## The New Era of Solid State Drives: Technical Opportunities and Challenges

**October 20, 2008** 

Daniel Donggi Lee, Ph.D Flash Solution, Memory Division Samsung Electronics Co., Ltd.



## Contents



IV Samsung SSD



## **A New Hurdle in Computing Platform**

Need to bridge performance gap in computing platform

- CPU cache bridges CPU and system memory (DRAM) gap
- Need to bridge the gap between system memory and HDD



## **Existing Approaches - I**

#### NVC (Non-Volatile Cache)

- Flash memory is added as a storage cache
  - Pros
    - No HDD SKU proliferation for PC OEMs
  - Cons
    - Data separation problem: Inconsistency may happen if either HDD or flash memory is detached from the system
    - System synchronization required in every power on cycle



## **Existing Approaches - II**

#### **SSD (Solid State Drive)**

- Data storage device that uses sold-state memory
- Utilizing existing storage device interfaces
  - Pros
    - No additional system variation for computing platform
    - Fast TAT time for production launch
  - Cons
    - Cost, System optimization, variety of SKUs (over 83 vendors)

-5-



Texas Memory System – RamSan-440



4500MB/s sustained throughputs 600k random IOPS Fiber Channel I/F, DRAM SSD

NVRAMOS 2008 Source: www.superssd.com

Violin Memory Inc – Violin 1010



Scalable Memory Architecture (VXM) 84 VIMMs (Violin Intelligent Memory Modules) 1M random IOPS PCIe x4/x8 I/F, DRAM/Flash SSD

Source: www.violin-memory.com



## Why NAND Flash is Good Solutions ?

#### Flash complements HDD for Fast Seek and Low Power

- Consistent performance at sequential and random access
- Low power and instant standby and activation



\*Using Interleave, Flash can support the best read and write performance



## NAND Flash like Toy of Lego<sup>™</sup>

17MB/s with a single 50nm 8Gb SLC in one plane mode
Performance = F {multi-plane, multi-way/channel...}



## **No More NAND Price is Limit!**

□ Corp. Notebook: 30~40GB/Home User: 48GB (eWeek.com '06 Aug.)

- OS (XP: 3~4GB, Vista: 10GB) + App. (6~7GB) + User data (~20GB)

2010: 1.8" HDD 43Mpcs, 2B\$ forecasted





## Contents



## **II** SSD Opportunities

## **III SSD Technical Issues**

## IV Samsung SSD



## **SSD Overview**



#### □ Access Time Comparison



- > SSD is superior to HDD considering all aspects..!!
- > Especially, SSD has ultra fast access time thanks to no rotation platters..!!



## **SSD Lineup**



Major Target	Consumer/Laptop Storage	Enterprise/Server Storage	Military/ Industrial Storage
Optimized For	Lowest \$\$/GB	Performance, Power	Reliability, Durability, Rugged, Security
Applications	Mixed	OLTP, Office Productivity, Data Mining	Data Recording, Mission Media
Trend	Moving to MLC	Higher MTBF	Encryption



## **Fast Startup**



## **Fast User Experiences**



-13-

## **Fast Gaming Experiences**



## **Power Savings**

#### □ SSD has 9X to 14X power savings compared HDD

- > 15K RPM 3.5" HDD, operating 14W, idle 9W
- SSD, under 1W at both operating & idle



\* HDD DATA back ground

RPM	Drive
4200	Hitachi GST Travelstar 4k120
5400	Toshiba MK1032GSX
7200	Hitachi GST Travelstar 7k100
10K	Seagate Savvio 10K.1
15K	Seagate Cheetah 15K.4

SAMSU

## **More Lifetime**

#### MLC SSD Estimated Notebook Lifetime



## **Server Opportunity**

## SSD's Performance/Power overcomes HDD's



\*SSD: IOPS: random read @ 512B, Data (SEC Marketing)

\*HDD: Price and Data from storagereview.com; IOPS: average read speed

#### NVRAMOS 2008



	15K HDD 3.5"	15K HDD 2.5"	SSD 2.5"
Capacity	300GB	144GB	64GB
Read	93 MB/s	108 MB/s	100 MB/s
IOPS	320	380	800
Price Est.	~ \$350	~ \$420	~ \$600
Avg. Watts	16 W	7W	1W
GB / \$	0.86GB <b>21.2</b>	0.34GB	₹ <b>1.0GB</b>
IOPS / W	20 240	54 ξ <mark>15</mark> χ	<b>800</b>

SAMSUNG

**HARMANNAN** 

## **Better IOPS/\$, IOPS/W in Enterprise**



## Contents



**II** SSD Opportunities

## **III SSD Technical Issues**

## IV Samsung SSD





## SSD Related Technical Issues – Big Pictures



## **Storage System Architecture**

Violin-Memory Inc. (www.violin-memory.com)

#### System Architecture

- > VIMM + VXM architecture (SSD unit -> system level)
- > Multi-level Flash RAID Hierarchy
  - > Tightly integrated RAID/flash controller

Industry Standard Interconnects

> Asynchronous IO over multi-VIMMs



## **Operating System**

#### □ OS has been optimized for HDD over decades

- > Especially, file system and memory management
- > Ex. delayed write, sequentially-optimized allocation, read-ahead, swap operation

#### □ HDD optimization = SSD optimization?

#### Common

- Favors large I/O
- Favors delayed write
- Segmented FS layout

#### Different

- Page/block boundary
- Large read-ahead is a waste
- No seek time, but write perf. is dependent on position & fragmentation

#### **CPU utilization vs Storage Utilization**

- > OS kernel's efficiency related to resource balancing
- > Resource scheduling optimization issue

#### □ Not only OS, but also Applications

- > Application should utilize the internal I/O parallelism of SSD (like CPU parallelism)
- > App-level I/O parallelism: multi-tasking I/O, asynchronous I/O



## **File System**

#### Linux FS Benchmark

Running postmark on various Linux FS

Log-structured and COW FS performs better than traditional FS on SSD

Workload	File size	# of file (work-set)
SS	9-15K	10,000
SL	9-15K	100,000
LS	0.1-3M	1,000
LL	0.1-3M	4,250



## **Internal SSD Architecture**

#### Storage Buffer Management

- > SSD has larger buffer than HDD (similar to RAID buffer)
- > Efficient buffer coordination between host and storage device will be helpful

#### □ How to coordinate with existing RAID system?

- > Multi-layered hierarchy of RAID system with (duplicated) buffers
- > Translated host I/O pattern into SSD-friendly pattern (ex. Pre-fetch, deferred write)

#### □ How to utilize internal parallelism?

- > Interface for multiple I/O command T10 TCQ (common practice : 32 ~ 256)
- > Cf. number of flash chips in SSD 64 (Gen3)







## **Host Interface**

#### NAND Flash is no more performance bottleneck

- Host interface would be performance limit
- However, complexity in protocol and SI/PI increase



## **NAND Internals**

**Endurance and Write speed is reverse-proportional** 



## **SSD Resource Management**

#### □ How I/O is processed in SSD?

- > Internal processing is a sequence of pipelined stages
- > I/O stage is parallelized into multiple hardware resources (ex. Chip, channel)

#### □ Resource scheduling of SSD

- > How to speed up the entire I/O processing throughput?
- Similar to super-scalar pipeline problem (without inter-request dependency)



Resources

## **Low Power Management**

#### Power consumption directly proportional to Storage Parallelism



# of NAND / Rate of Parallelism (channels/ways..)



## **Related Works (Research-side)**

#### □ Flash Translation Layer (FTL)

- A space-efficient flash translation layer for CompactFlash Systems (IEEE Transactions on Consumer Electronics, 2002)
- > A superblock-based flash translation layer for NAND flash memory (EMSOFT, 2006)
- A Log Buffer Based Flash Translation Layer Using Fully-Associative Sector Translation (ACM Transaction on Embedded Computing System, 2007)
- A Re-Configurable FTL(Flash Translation Layer) Architecture for NAND Flash based Applications (RSP, 2007)

#### Flash-optimized Buffer Management

- > CFLRU: a replacement algorithm for flash memory (CASES, 2006)
- FAB: flash-aware buffer management policy for portable media players (IEEE Transactions on Consumer Electronics, 2006)
- BPLRU: A Buffer Management Scheme for Improving Random Writes in Flash Storage (USENIX Conference on File and Storage Technologies, 2008)

#### Flash-optimized File System

- > A New Type of NAND Flash-based File System: Design and Implementation (WiCOM, 2006)
- An Efficient NAND Flash File System for Flash Memory Storage (IEEE Transactions on Computers, 2006)
- Embedded NAND Flash file System for Mobile Multimedia Device (IEEE Transactions on Computer, 2008)



## **Related Works (Industry-side)**

Btrfs	5	
Developer	Oracle Corporation	
Full name	Btrfs	
Introduced	June 12, 2007 (Linux)	
Structu		
Directory contents	btree	
File allocation	extents	
Limit		
Max file size 16 EiB		
Max number of files	284	
Max filename length	255 bytes	
Max volume size	16 EiB	
Allowed characters in filenames	All bytes except NUL and '/'	
Feature	es	
Attributes	POSIX	
File system permissions	POSIX, ACL	
Transparent encryption	No	
Supported operating systems	Linux	

#### NVRAMOS 2008 -30-

	ZFS	
Developer	Sun Microsystems	
Full name	ZFS	
Introduced	November 2005 (OpenSolaris)	
St	ructures	
Directory contents	Extensible hash table	
	Limits	
Max file size	16 EiB	
Max number of files	2 <sup>48</sup>	
Max filename length	255 bytes	
Max volume size	16 EiB	
Features		
Forks	Yes (called Extended Attributes)	
Attributes	POSIX	
File system permissions	POSIX, ACL	
Transparent compression	Yes	
Transparent encryption	Yes (currently beta) <sup>[1]</sup>	
Supported	Sun Solaris, Apple Mac	
operating	OS X Server 10.5, FreeBSD Linux via EUSE	

1	WAFL	
Developer	NetApp	
Full name	Write Anywhere File Layout	and the
	Limits	
Max file size	16TB (limited by containing aggregate size)	
Max volume size	16TB (limited by containing aggregate size)	
Allowed characters in filenames	selectable (UTF-8 default)	NetA
F	eatures	
Dates recorded	atime, ctime, mtime	
File system permissions	UNIX permissions and ACLs	
Transparent compression	No	
Transparent encryption	No (possible with 3rd party appliances like Decru DataFort)	
Single Instance Storage	Yes ( <i>FAS Dedup</i> : periodic offline scans, block based; <i>VTL Dedup</i> : online byte-range based)	-

Source: Wikipedia



## **Research Hints**

#### Operating System

- > File system alignment, allocation policy, design (ex. COW)
- > Block layer SSD-optimized I/O scheduler
- > Volume manager: alignment, allocation
- > Virtual memory: read-ahead

#### Storage System

- > On-storage buffer management
- > RAID-like optimization algorithm
- > Exchange hint info with host

#### □ Lessons from flash FS

- > Sequential writing at multiple logging points
- > Wandering tree
  - > Trace-off between sequentiality vs. amount of write
  - > Cf. space map (Sun ZFS)
- > Need to optimize garbage collection overhead
  - Either FS itself or FTL in SSD



## **Research Issue Summary**

	Technology	Revisiting existing technology	New technology for flash storage
S <sup>r</sup> fc	ystem Architecture or Flash Storage	Host RAID System Storage management	Violin, TMS
(	Operating System	N/A	T13: Trim command
	File System	Log-structured FS COW FS	N/A
В	uffer Management	N/A	CFLRU, FAB, BPLRU
	RAID	RAID Buffer management	N/A
	Low Power Management	Dynamic power management	N/A
/RAMC	<b>S 2008</b>	-32-	SA

## Contents





## Samsung SSD Roadmap



**NVRAMOS 2008** 

SAMSUNG

## The 1<sup>st</sup> SSD Only NotePC - Lenovo X300<sup>™</sup>



# First "Mainstream" SSD Only Notebook 64GB Solid State Drive

No HDD Option or "Opt Out"

□ Represents an Irreversible Industry Trend...



## **Samsung SSD Press Release - Client**

#### SAMSUNG Develops World's Fastest and Largest Capacity 2.5-inch, MLC-based (256GB) SSD with SATA II Interface

Taipei, Taiwan on May 26, 2008



Taipei, Taiwan - May 26, 2008 : Samsung Electronics, the world leader in advanced memory technology, announced today that it has developed the world's fastest, 2.5-inch, 256 Gigabyte (GB) multi-level cell (MLC) based solid state drive (SSD) using a SATA II interface at the fifth annual Samsung Mobile Solution Forum held at the Westin Taipei Hotel. Samsung's new 256GB SSD is also the thinnest drive with the largest capacity to be offered with a SATA II interface.

#### Samsung Introduces High-performance, Low-density SATA II SSDs for Low-priced PC Market

Half Slim

Client

MLC

Seoul, Korea on Aug 27, 2008



Seoul, Korea - Aug. 27, 2008 - Samsung Electronics Co., Ltd., the world leader in advanced semiconductor technology, announced today that it has begun sampling low-density, higher-performance solid state drives (SSDs) that are only 30 percent of the size of 2.5inch SSDs and highly cost-efficient to manufacture. With the introduction of these smaller, low-capacity SSDs, Samsung now offers an attractive replacement for existing hard drives used in lowcost PCs. Available in densities of 8GB, 16GB and 32GB, the new multi-level-cell SSDs will be mass produced beginning next month.



## Samsung SSD Press Release - Server

#### Samsung Solid State Drive Selected for New HP ProLiant Virtualization Blade Server

Blade Launch

By: HP News Desk Oct. 10, 2008 01:45 PM



Samsung Electronics Co., Ltd., the world leader in advanced <u>semiconductor</u> technology, announced today that its 32 gigabyte (GB) and 64 GB solid state drives (SSDs) have been selected, after extensive testing, for use in the HP ProLiant BL495c virtualization blade server. This represents the first time that an SSD has been qualified for use in a server optimized for virtualization. The server is the

world's first server blade designed specifically to host virtual machines, and is designed for use in virtualized environments that require significant memory, data storage and network connections to optimize server performance.

#### High Endurance SLC

#### SAMSUNG Collaborates with Sun Microsystems to Develop New Ultra-Endurance Flash Memory for SSD Products in Server Applications

Seoul, Korea on Jul 17, 2008



Seoul, Korea - July 17, 2008 : Samsung Electronics Co., Ltd., the world leader in advanced semiconductor technology, announced today that it has collaborated with Sun Microsystems to develop a single-level-cell NAND flash memory device for use in solid state drives that offers much higher endurance levels than any other flash memory device on the market today.



## **World Recognized Samsung SSD**

#### Why Samsung ?

The Samsung SSD has been voted one of the ten most brilliant and bold ideas that are set to change the worldrecognized for pushing the envelope of technology and for design and engineering innovation.





# THANK YOU