



Implications of Storage Class Memories (SCM) on Software Architectures

C. Mohan, IBM Almaden Research Center, San Jose

mohan@almaden.ibm.com <http://www.almaden.ibm.com/u/mohan>

Suparna Bhattacharya, IBM Systems and Technology Lab, India

bsuparna@in.ibm.com

HPCA 2010 Workshop on the use of Emerging Storage and memory Technologies (WEST)
Bangalore, India, January 2010

Acknowledgments and References

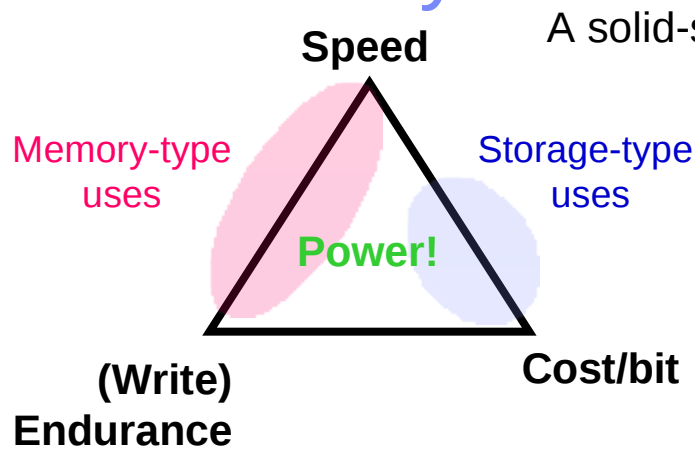
Thanks to

- Colleagues in various IBM research labs in general
- Colleagues in IBM Almaden in particular

References:

- "Storage Class Memory, Technology, and Uses", Richard Freitas, Winfried Wilcke, Bülent Kurdi, and Geoffrey Burr, Tutorial at the 7th USENIX Conference on File and Storage Technologies (FAST '09), San Francisco, February 2009, <http://www.usenix.org/events/fast/tutorials/T3.pdf>
- "Storage Class Memory - The Future of Solid State Storage", Richard Freitas, Flash Memory Summit, Santa Clara, August 2009, http://www.flashmemorysummit.com/English/Collaterals/Proceedings/2009/20090812_T1B_Freitas.pdf
- "Storage Class Memory: Technology, Systems and Applications", Richard Freitas, 35th SIGMOD international Conference on Management of Data, Providence, USA, June 2009.

Storage Class Memory

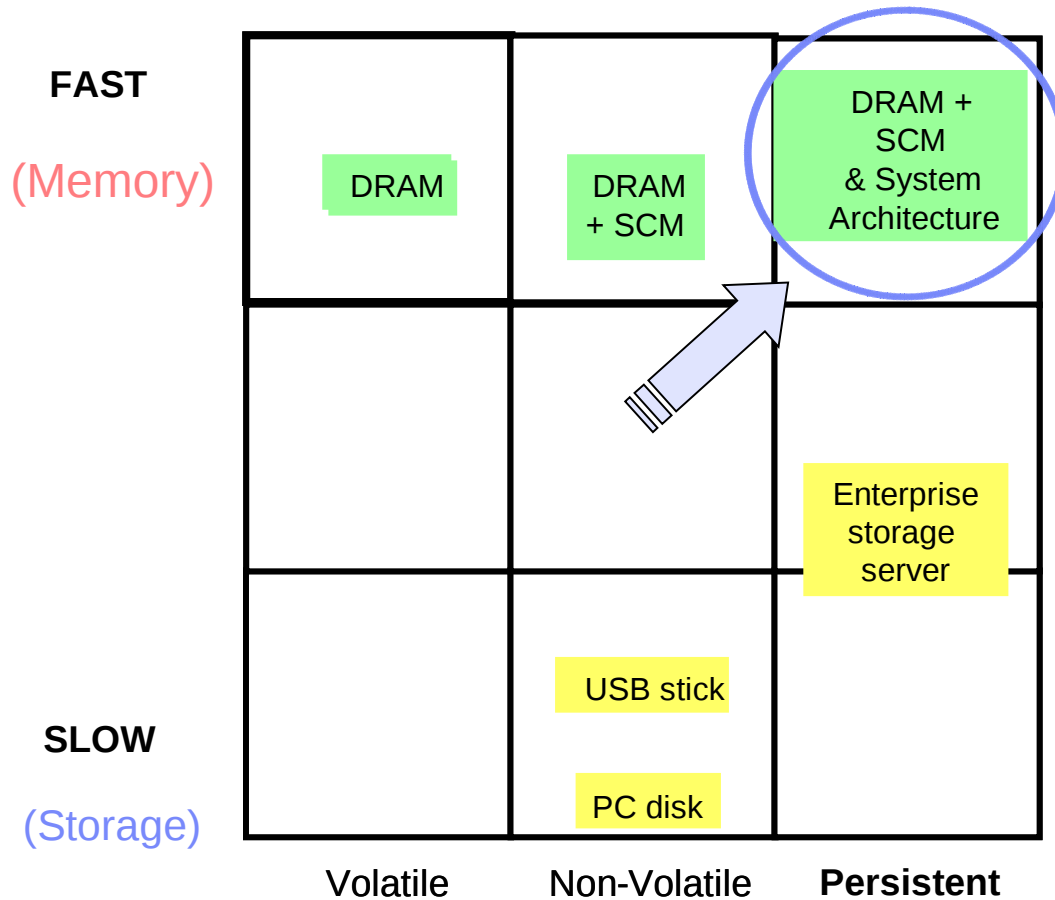


A solid-state memory that **blurs the boundaries** between storage and memory by being **low-cost, fast, and non-volatile.**

SCM system requirements for Memory (Storage) apps

- No more than 3-5x the **Cost** of enterprise HDD (< \$1 per GB in 2012)
- **<200nsec (<1μsec)** Read/Write/Erase time
- >100,000 **Read I/O operations** per second
- **>1GB/sec (>100MB/sec)**
- **Lifetime** of 10^9 – 10^{12} write/erase cycles
- 10x lower **power** than enterprise HDD (100mW ON power, 1mW standby)

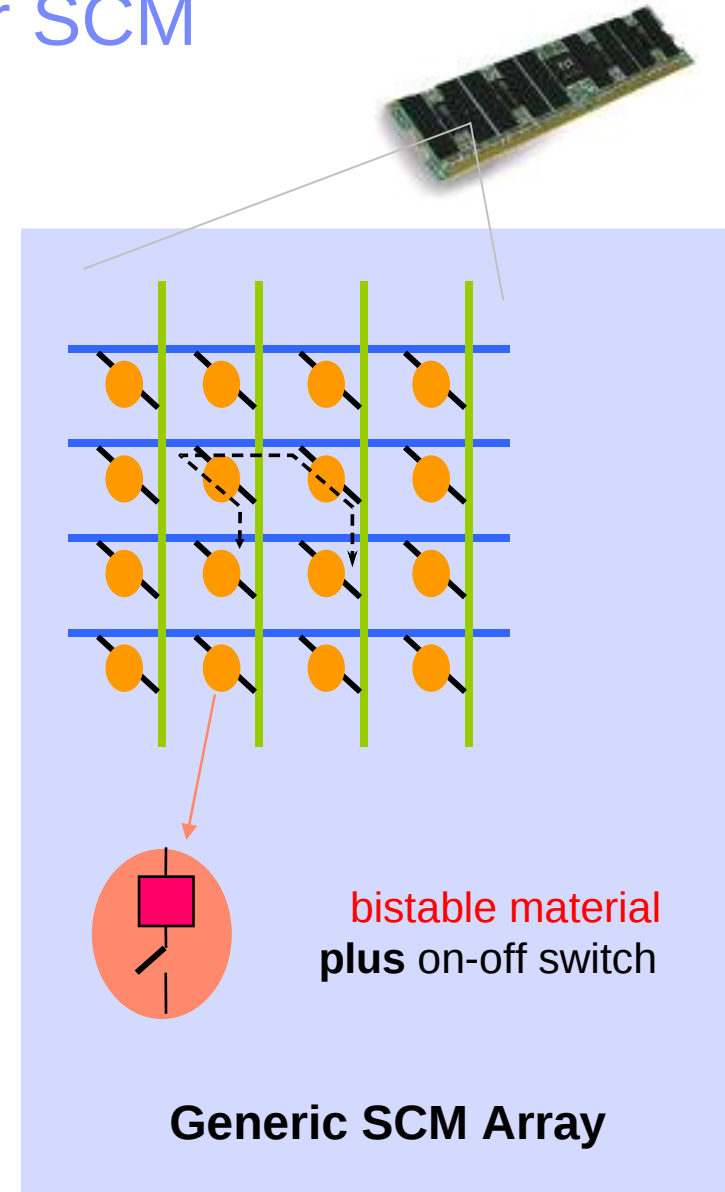
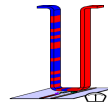
Speed/Volatility/Persistence Matrix



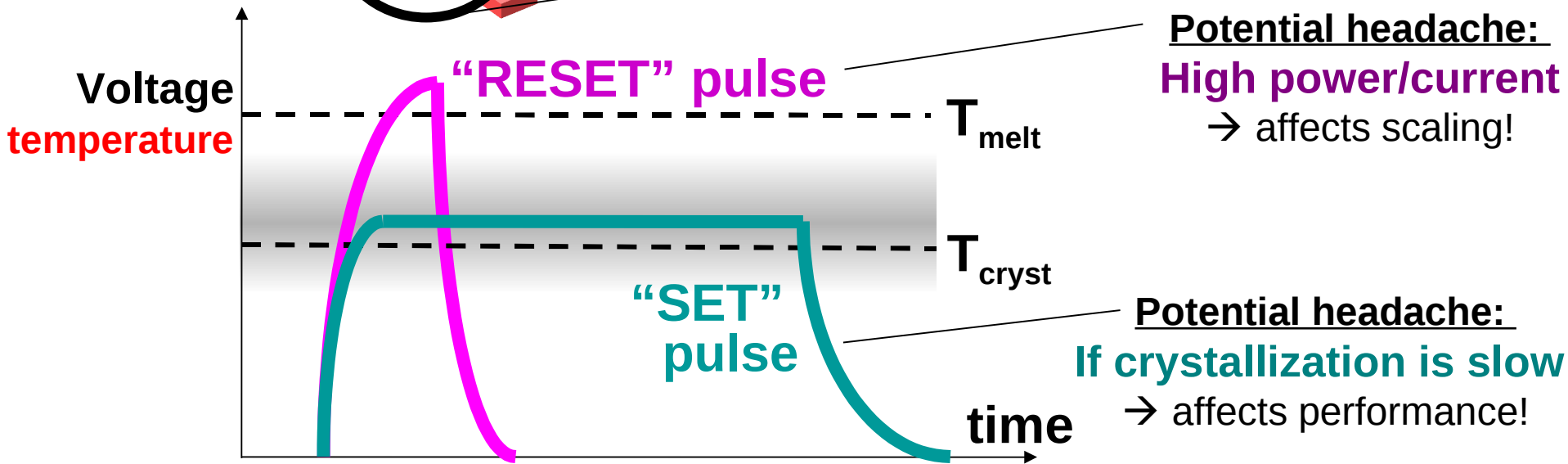
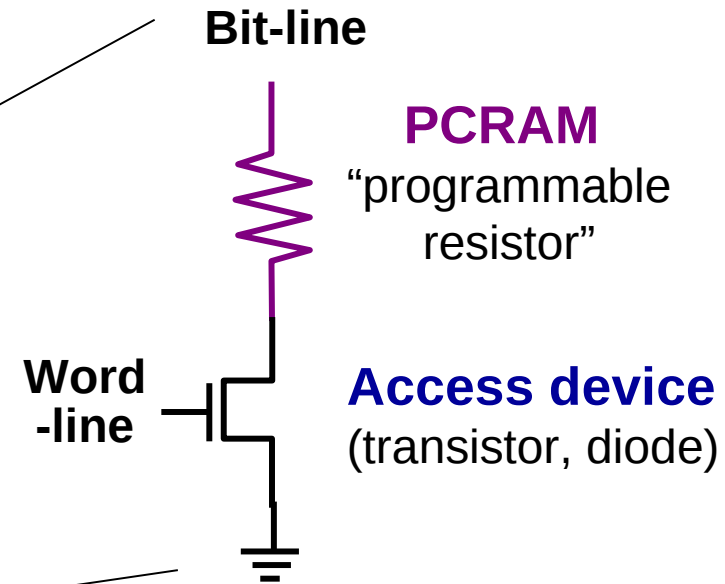
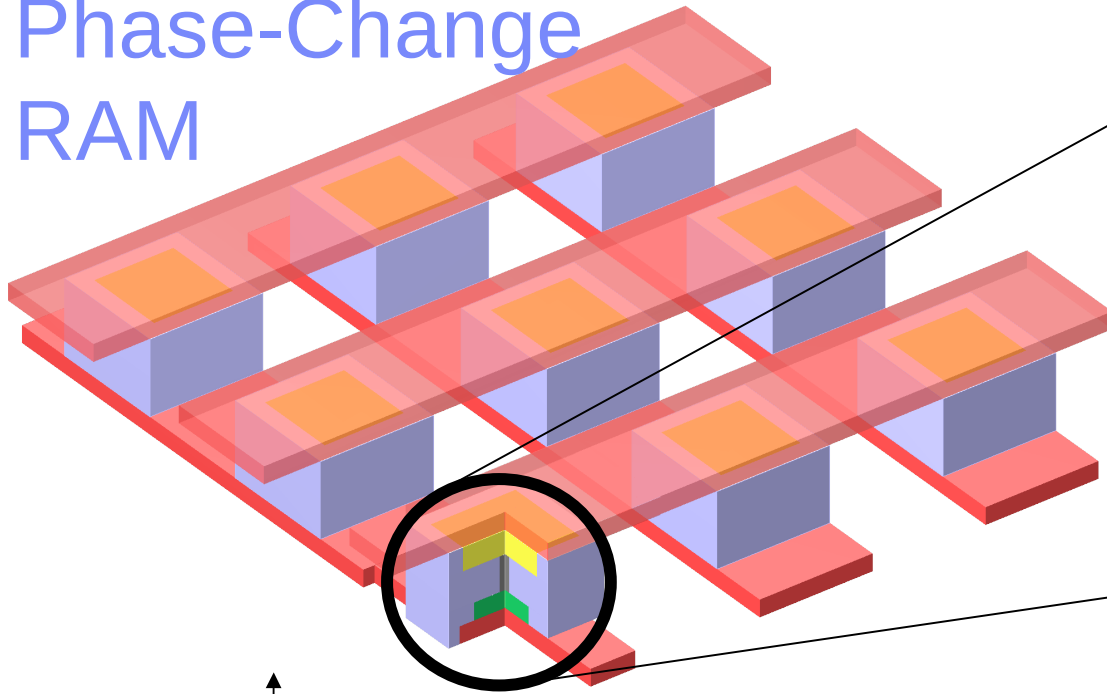
Persistent storage
will not lose data

Many Competing Technologies for SCM

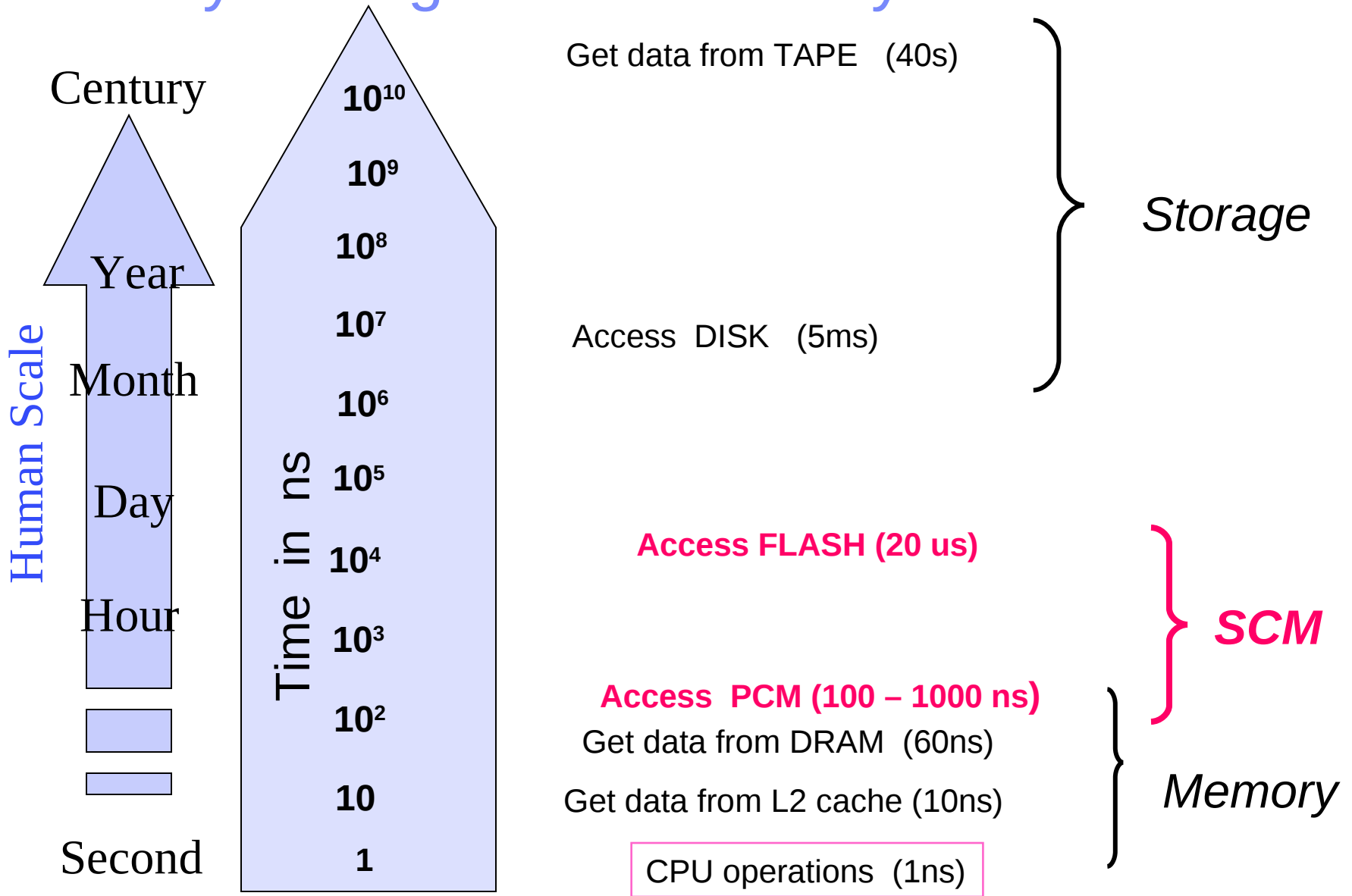
- **Phase Change RAM**
 - most promising now (scaling)
- **Magnetic RAM**
 - used today, but poor scaling and a space hog
- **Magnetic Racetrack**
 - basic research, but very promising long term
- **Ferroelectric RAM**
 - used today, but poor scalability
- **Solid Electrolyte and resistive RAM (Memristor)**
 - early development, maybe promising
- **Organic, nano particle and polymeric RAM**
 - many different devices in this class, unlikely
- **Improved FLASH**
 - still slow and poor write endurance



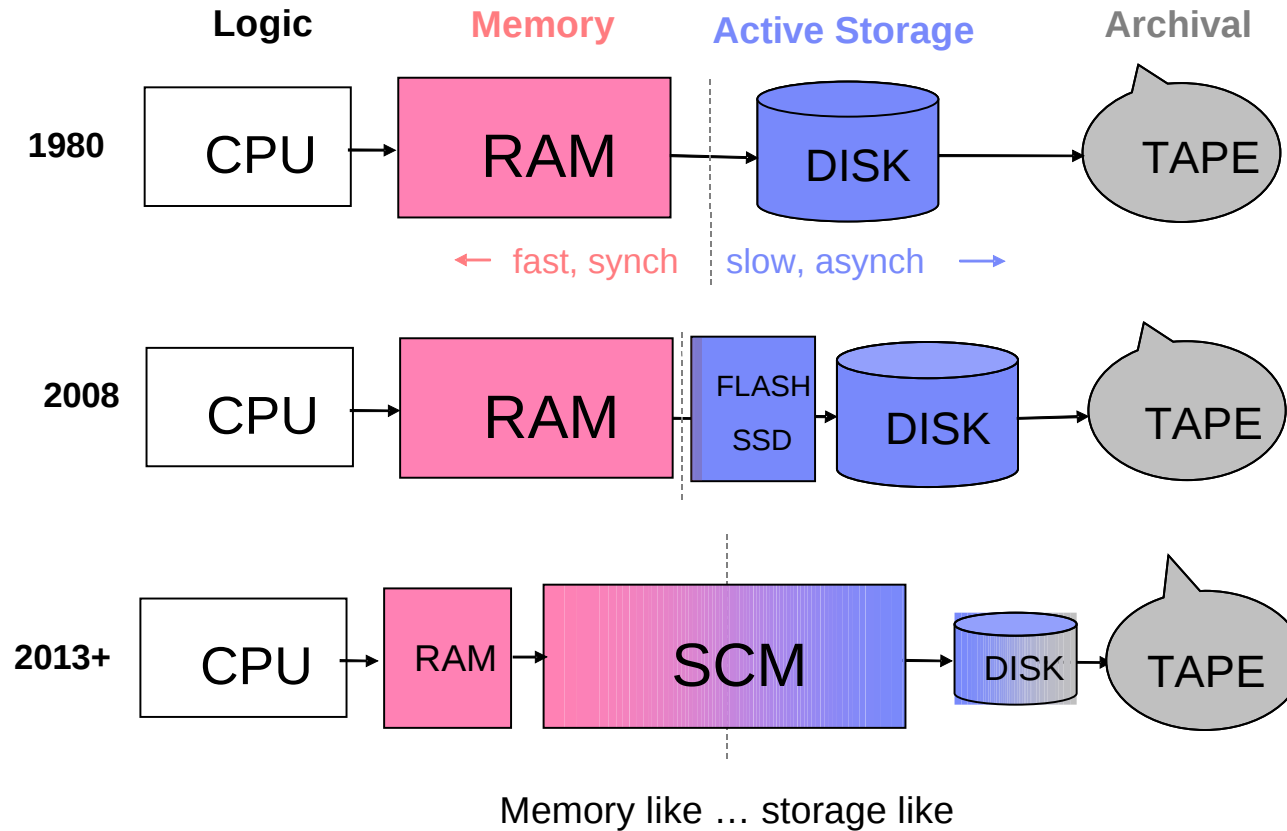
Phase-Change RAM



Memory/Storage Stack Latency Problem



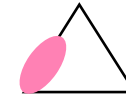
SCM as Part of Memory/Storage Solution Stack



2013 Possible Device Specs

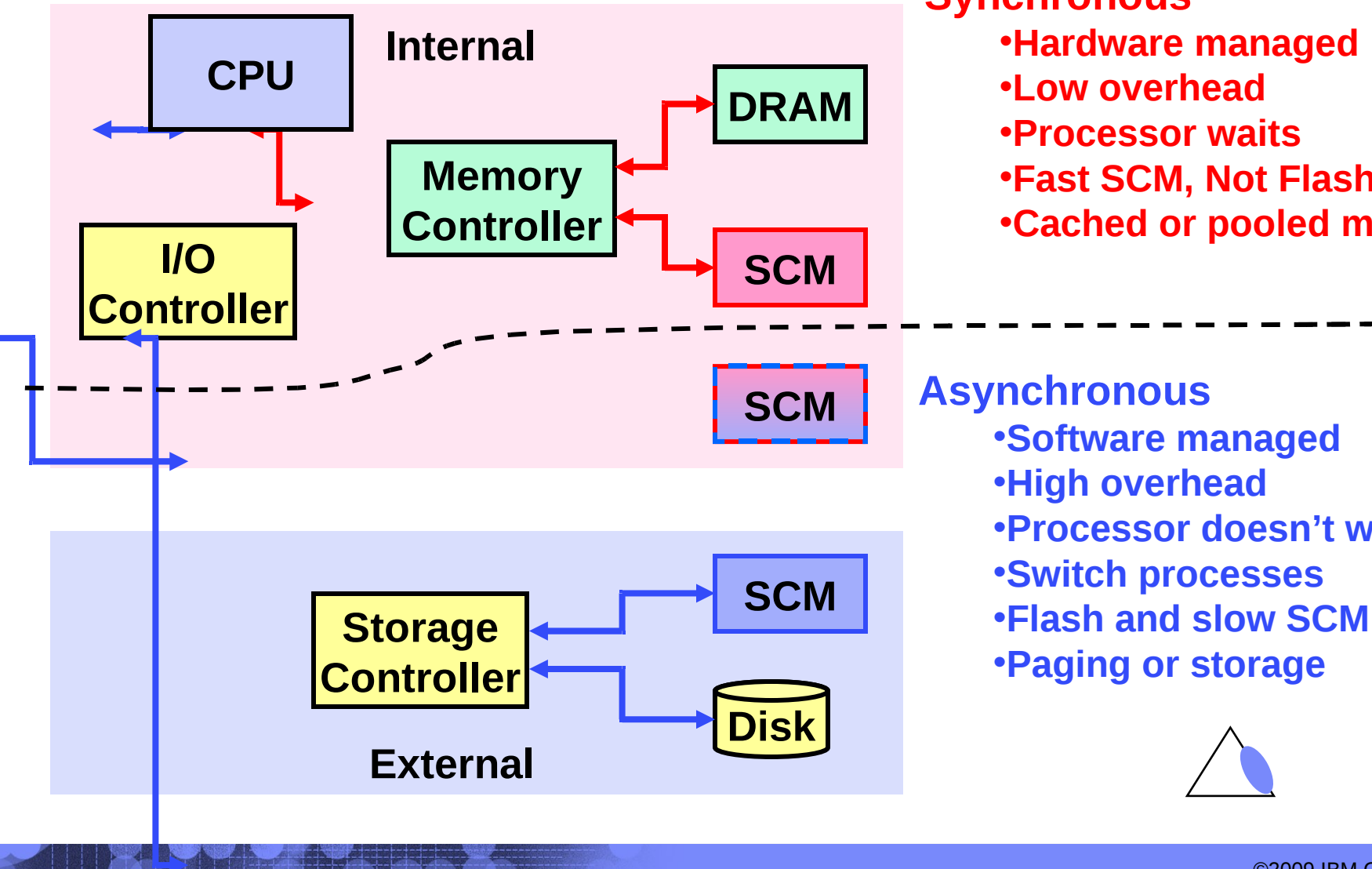
Parameter		PCM-S	PCM-M
Capacity	DRAM	128 Gbits	16 Gbits
Feature Size F	32nm	32nm	32nm
Effective cell size	$6 F^2$	$0.5 F^2$	$2 F^2$
Read latency	60ns	800ns	300ns
Write latency	60ns	1400ns	1400ns
Retention time	ms	2-10 years	Strongly temp. dependent

Architecture



Synchronous

- Hardware managed
- Low overhead
- Processor waits
- Fast SCM, Not Flash
- Cached or pooled memory



Asynchronous

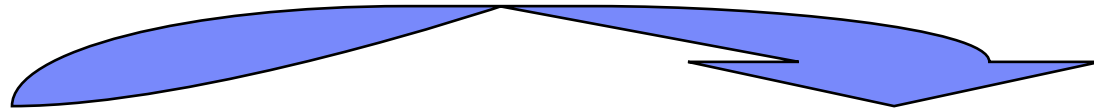
- Software managed
- High overhead
- Processor doesn't wait
- Switch processes
- Flash and slow SCM
- Paging or storage



Challenges with SCM

- Asymmetric performance
 - Flash: writes much slower than reads
 - Not as pronounced in other technologies
- Bad blocks
 - Devices are shipped with bad blocks
 - Blocks wear out, etc.
- The “fly in the ointment” is write endurance
 - In many SCM technologies writes are cumulatively destructive
 - For Flash it is the program/erase cycle
 - Current commercial flash varieties
 - Single level cell (SLC) → 10^5 writes/cell
 - Multi level cell (MLC) → 10^4 writes/cell
 - Coping strategy → Wear leveling, etc.

Shift in Systems and Applications



■ *DRAM – Disk – Tape*

- Cost & power constrained
- Paging not used
- Only one type of memory:
volatile

Main Memory:

■ ***DRAM – SCM – Disk – Tape***

- Much larger memory space for same power and cost
- Paging viable
- Memory pools: different speeds, some persistent
- Fast boot and hibernate

- Active data on disk
- Inactive data on tape
- SANs in heavy use

Storage:

- Active data on SCM
- Inactive data on disk/tape
- Direct Attached Storage?

- Compute centric
- Focus on hiding disk latency

Applications:

- Data centric comes to fore
- Focus on efficient memory use and exploiting persistence
- Fast, persistent metadata

Potential early middle-ware exploiters of SCM

- Databases
 - Traditional RDBMS, specialized DBMS, in-memory DBMS
- Distributed object caching (in-memory grids)
 - Memcached, Websphere eXtreme Scale
- Alternate persistent stores
 - Document stores, key-value stores, DHTs
- Persistent message queues
- Data-intensive applications
 - Analytics, Genomics
- Long tail data access patterns (amenable to tiered store)

PCM Use Cases

- PCM as disk
- PCM as paging device
- PCM as memory
- PCM as extended memory

Let Us Explore DBMS as Middleware Exploiter of PCM

PCM as Logging Store – Permits > Log Forces/sec?

- Obvious one but options exist even for this one!
- Should log records be written directly to PCM or first to DRAM log buffers and then be forced to PCM (rather than disk)
- In the latter case, is it really that beneficial if ultimately you still want to have log on disk since PCM capacity won't be as much as disk – also since disk is more reliable and is a better long term storage medium
- In former case, all writes will be way slowed down!

PCM replaces DRAM? - Buffer pool in PCM?

- This PCM BP access will be slower than DRAM BP access!
- Writes will suffer even more than reads!!
- Should we instead have DRAM BPs backed by PCM BPs?

This is similar to DB2 z in parallel sysplex environment with BPs in coupling facility (CF)

But the DB2 situation has well defined rules on when pages move from DRAM BP to CF BP

- Variation was used in SafeRAM work at MCC in 1989

Assume whole DB fits in PCM?

- Apply old main memory DB design concepts directly?
- Shouldn't we leverage persistence specially?
- Every bit change persisting isn't always a good thing!
- Today's failure semantics lets fair amount of flexibility on tracking changes to DB pages – only some changes logged and inconsistent page states not made persistent!
- Memory overwrites will cause more damage!
- If every write assumed to be persistent as soon as write completes, then L1 & L2 caching can't be leveraged – need to do write through, further degrading perf

Assume whole DB fits in PCM? ...

- Even if whole DB fits in PCM and even though PCM is persistent, still need to externalize DB regularly since PCM won't have good endurance!
- If DB spans both DRAM and PCM, then
 - need to have logic to decide what goes where – hot and cold data distinction?
 - persistence isn't uniform and so need to book-keep carefully

What about Logging?

- If PCM is persistent and whole DB in PCM, do we need logging?
- Of course it is needed to provide at least partial rollback even if data is being versioned (at least need to track what versions to invalidate or eliminate); also for auditing, disaster recovery, ...

High Availability and PCM

- If PCM is used as memory and its persistence is taken advantage of, then such a memory should be dual ported (like for disks) so that its contents are accessible even if the host fails for backup to access
- Should locks also be maintained in PCM to speed up new transaction processing when host recovers

Start from Scratch?

- Maybe it is time for a fundamental rethink
- Design a DBMS from scratch keeping in mind the characteristics of PCM
- Reexamine data model, access methods, query optimizer, locking, logging, recovery, ...
- What are the killer apps for PCM? For flash, they are consumer oriented - digital cameras, personal music devices, ...