# The Bw-Tree: A B-Tree On Steroids

Justin Levandoski
David Lomet
Sudipta Sengupta

# The Bw-Tree: What is it?

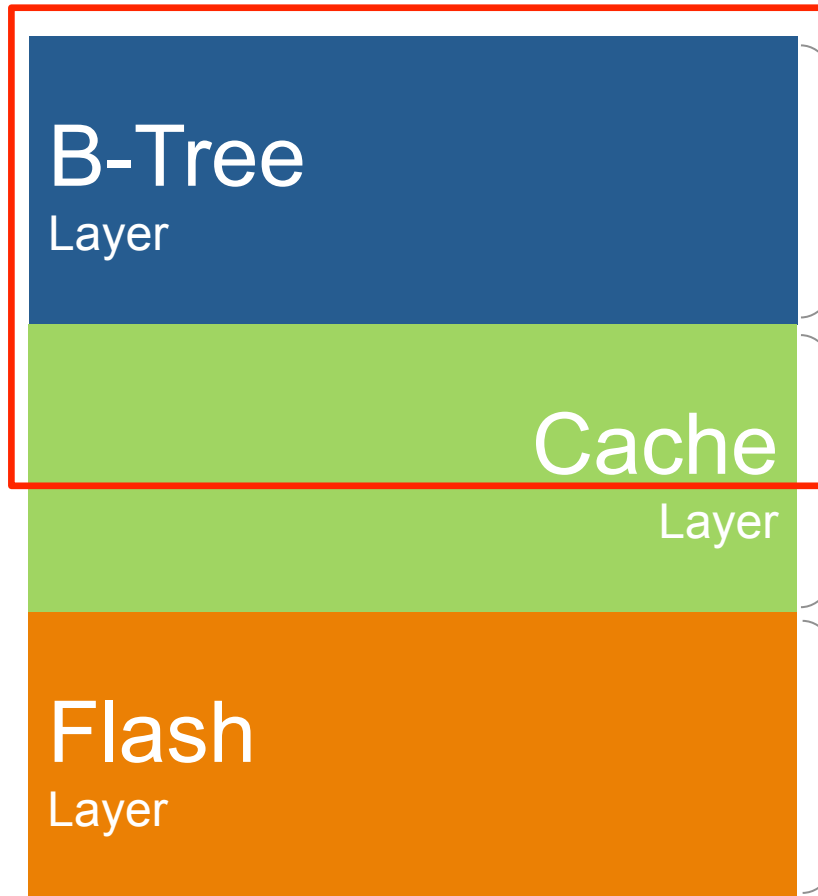"A Latch-free, Log-structured B-tree for Multi-core Machines with Large Main Memories and Flash Storage"

**Bw = "Buzz Word"**

# The Buzz Words: Attacking Two Trends

- Multi-core + large main memories
  - Latch (lock) free
    - Worker threads do not set latches for any reason
    - No latch contention
  - "Delta" updates
    - No updates in place
    - Reduces cache invalidation

- Flash storage
  - Good at random reads and sequential reads/writes
  - Bad at random writes
  - Use flash as append log
  - Implement log-structured storage layer over flash
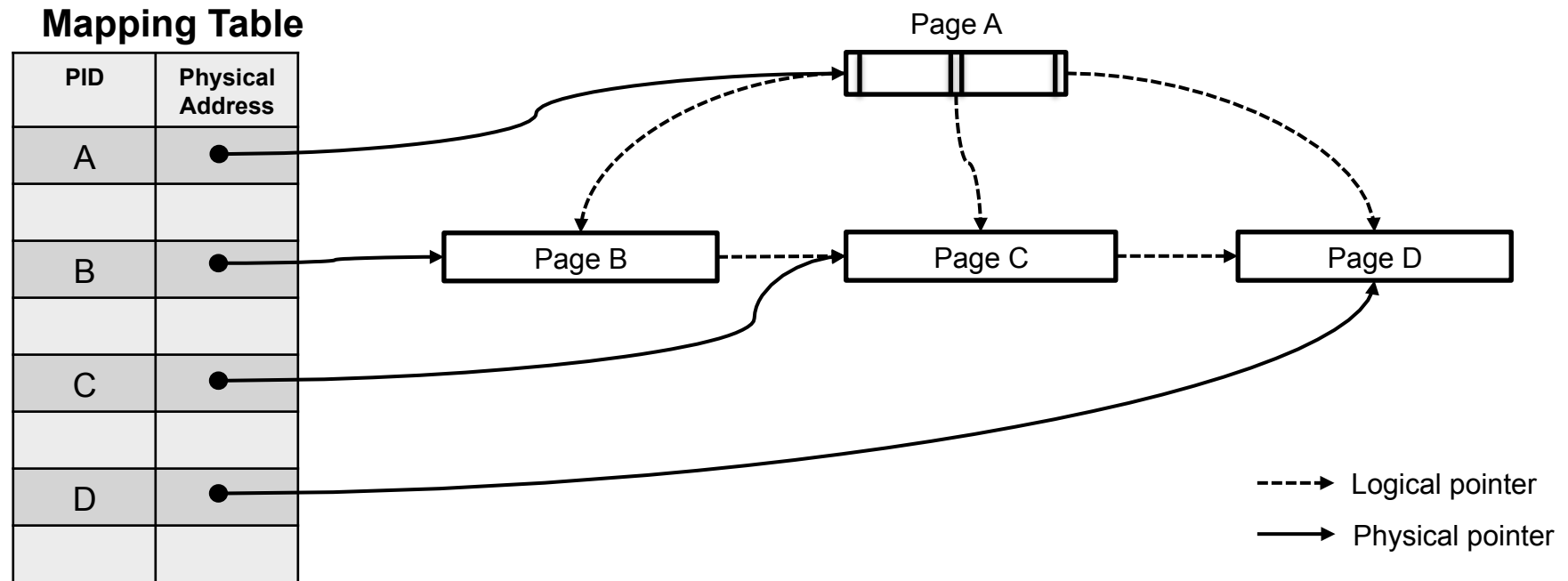  - Must run efficiently on both expensive AND cheap devices

**Let's talk about memory first…**

| | |
|---|---|
| **B-Tree** Layer | • "CRUD" API<br>• B-tree search/update logic<br>• In-memory pages only |
| **Cache** Layer | • Logical page abstraction for B-tree layer<br>• Brings pages from flash to RAM as necessary |
| **Flash** Layer | • Sequential writes to log-structured storage<br>• Flash garbage collection |

# Logical Pages and Mapping Table

**Mapping Table**

| PID | Physical Address |
|-----|------------------|
| A | ● |
| | |
| B | ● |
| | |
| C | ● |
| | |
| D | ● |
| | |

Page A

Page B     Page C     Page D

- - - - → Logical pointer

⎯⎯→ Physical pointer

- Logical pages identified by mapping table index
- Isolates update to a single page
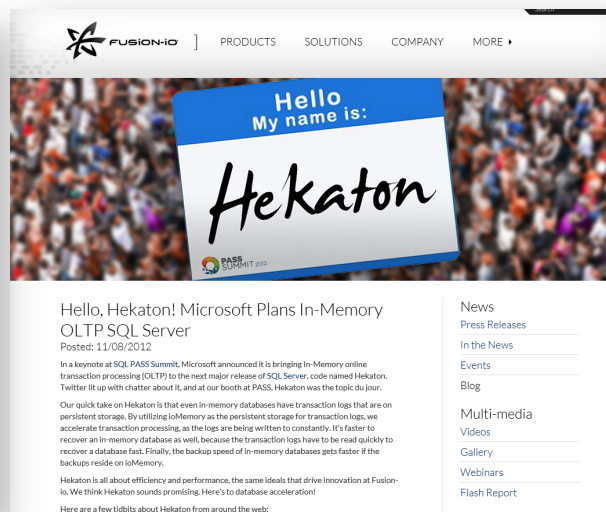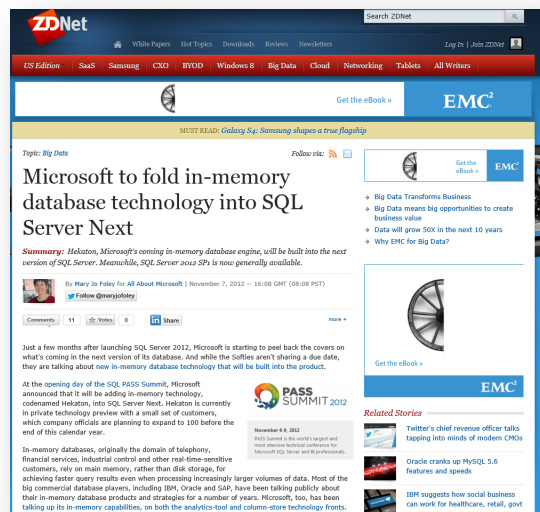- Important for latch-free behavior and log-structuring

# Delta Updates



- Each page update produces a new address (the delta).
- Install new page address in map using compare-and-swap.
- Only one winner on concurrent update to the same address.
- Eventually install new consolidate page with deltas applied.
- Single-page updates are easy, solved node splits and deletes.

# Microsoft SQL Server Hekaton

- Main-memory optimized OLTP engine
  - Engine is completely latch-free
  - Multi-versioned, optimistic concurrency control (VLDB 2012)
- Bw-tree is the ordered index in Hekaton



http://research.microsoft.com/main-memory_dbs/

# Architecture



- API
- B-tree search/update logic
- In-memory pages only

- Logical page abstraction for B-tree layer
- Brings pages from flash to RAM as necessary

- Sequential writes to log-structured storage
- Flash garbage collection

**B-Tree** Layer

**Cache** Layer

**Flash** Layer

# Handling pages located on flash

**Memory**

**Mapping Table**

| PID | Physical Address |
|-----|------------------|
| A | ● |
| | |
| B | ● |
| | |
| C | ● |
| | |
| D | ● |
| | |

Page A

Page B

Page D

| flash/mem flag | address |
|----------------|---------|
| 1 bit | 63 bits |

------→ Logical pointer

——→ Physical pointer

Page X

Page C

Page Y

**Log Structured Store (LSS) on Flash**

# Flushing pages

Flush Δ

Δ: Delete record 48

Δ: Insert record 50

Flush Δ

Page P

- Swapout drops page from memory.
- Install LSS offset in mapping table.
- Can also perform partial swapout.
- May require random read to retrieve page.

**Mapping Table**

| PID | Physical Address |
|-----|------------------|
|     |                  |
|     |                  |
| P   |                  |
|     |                  |

**Latch-free Write Buffer**

**Log Structured Store (LSS) on Flash**

| Page X | Page F | Page P | Page T | Page G | $\Delta_P$ 50 | $\Delta_P$ 48 | Page E | $\Delta_T$ 5 |

Write ordering in log

- ## LSS Garbage Collection

  - Cleans orphaned data unreachable from mapping table.

  - Relocates entire pages in sequential blocks (to reduce random reads from LSS).

- ## Access Method Recovery

  - Occasionally checkpoint mapping table.

  - Recover by:

    - Restoring mapping table.

    - Scan LSS forward from position recorded in checkpoint to the end of the log.

    - End result is latest LSS offset for pages in mapping table.

# The Big Picture

App Needing
Transactional
Key-Value Store

App Needing Atomic
Key-Value Store

App Needing High
Performance Log
Structured "Page" Store

**Transactional Component**

**Data Component**

Bw-Tree Latch Free
Ordered Index

Latch-Free
Linear Hashing

LLAMA Storage Engine
(Latch-Free, Log-Structured, Access-Method Aware)