FOEDUS: OLTP Engine for a Thousand Cores and NVRAM

Hideaki Kimura HP Labs, Palo Alto, CA

Slides By : Hideaki Kimura Presented By : Aman Preet Singh

Next-Generation Server Hardware?



HP The Machine



UC Berkeley

Firebox



Differences

- HP Photonics? Infiniband? QPI?
- HP Memristor? PCM? Flash?

Commonalities

- 1,000s of CPU Cores
- Fast Interconnect
- Huge Low-Latency NVRAM
- Endurance will be an issue



GAP ?

- Traditional OLTP databases do not scale to large number of cores.
- Recent main-memory databases achieve better scalability but they do not efficiently handle NVRAM.

GAP?

- Software should place data wisely.
 - System on a chip.
 - Systems capable of handling cache in-coherrent architectures.
 - Non-uniform Memory access (NUMA) aware systems.
- Aim should be to try to make the data local.

Solution ?

- To provide these capabilities, FOEDUS was built.
 - Ground up database
 - Open source
 - Fully ACID
 - Designed to scale up to thousands of cores
 - Makes the best use of DRAM for OLTP,
 - And NVRAM for OLAP queries.

FOEDUS Key Principles

- Bring in-memory speed to NVRAM
 - SILO-like Lightweight Optimistic Concurrency Control
- Dual-Page: physically independent, logically equivalent
 - Mutable Volatile Pages in DRAM
 - Immutable Snapshot Pages in NVRAM
- Master-Tree: Simple and Scalable OCC for NVRAM
 - Masstree + Foster B-Tree + Foster-Twin

SILO: Lightweight/Decentralized OCC [Tu et al]





SILO: Decentralized Logger with Epoch [Tu et al]





SILO's OCC Benefits

Extremely Lightweight and Scalable

Block-free (lock-free with retries)

No Write for Reads (cf., read-lock)

In-page Lock Mechanism

No centralized lock manager

Cache Friendly

But, how can we go beyond DRAM?



In-Memory DBMS Beyond DRAM





Logically Equivalent, Physically Independent Dual •Volatile Page: Mutable, on DRAM Snapshot Page: Immutable, on NVRAM



- : Volatile-Page is just made. No Snapshot yet.
- : Snapshot-Page is latest. No modification.
- : X is the latest truth, but Y may be equivalent.





Constructing Stratified Snapshot from Logical Logs



Benefits of Stratified Snapshots

\succ Serializable transactions check only a single version of data \leftrightarrow LSM-Trees

>Drastically more scalable and efficient construction of NVRAM-resident data pages

- >Separate from transactions/logging
- >Everything Batched, Processed in a tight loop
- ≻Large Sequential Write only. No frequent flush.



Master-Tree in a Nutshell

- •Mass-tree: Cache Crafty B-tree and OCC with inpage Lock Mechanism.
- •Foster B-tree: Single incoming pointer via fosterchild to ease page in/out.
- Foster-Twins: Drastically Simplifies OCC and Reduces Aborts/Retries.



OCC Problem 1: Unnecessary Retries/Aborts



Further, SILO must take Page-Version set and abort if changed by pre-commit.



Figure 6. Find the border node containing a key.



© Copyright 2015 Hewlett-Packard Development Company, L.P. The information contained herein is subject to change without notice.

Foster Twins here to Save You!

- Strong Invariants to drastically simplify OCC and eliminate unnecessary aborts/retries
- Still everything in-page
 - No per-tuple GC or compaction/migration
 - No centralized data structure







Foster-Twins Commit Protocol



Dual Pages: Benefits

✓ Snapshot Pages are <u>Immutable</u>

Drastically simplifies Caching/Replication/Traversal/Commit/etc

✓ Physically <u>Independent</u>

Transactions never interfere w/ construction/retrieval/eviction of Snapshot Pages

✓ Logically <u>Equivalent</u>

- No Bloom Filter or global data needed for Serializability ↔ LSM-Trees
- Volatile Pages <u>guaranteed to</u> contain latest data
- Snapshot Pages <u>quaranteed to</u> be complete as of snapshot-epoch



Experiments

TPC-C, Serializable

vs. H-Store, SILO, Shore-MT

- (Traditional DBs (e.g., MySQL) are too slow to compare with.)
- HP Superdome X (DragonHawk): <u>16 sockets, 240 cores, 12TB DRAM</u>.
- Emulated NVRAM





© Copyright 2015 Hewlett-Packard Development Company, L.P. The information contained herein is subject to change without notice.

Experiment 2: DBMS on NVRAM (emulated) FOEDUS

<u>Environments</u>

- tmpfs-based NVRAM emulation (varied latency)
- Data/xlog in NVRAM
- H-Store w/ anti-cache

<u>Observations</u>

- 1. FOEDUS 100x~ faster
- 2. Performs best when NVRAM read **<10us**



Experiment 3: OLAP Workload

<u>Workload</u>

- 30x Larger Data
- Cursor-Heavy

- Read-Only

- Volatile Pages Only vs Snapshot Pages Only

Observations

- 100x~ faster than H-Store.
- FOEDUS runs faster when database is cold (!!)



Ideas

• The key idea of FEODUS is to use the master-tree model to reduce contention, making OCC lightweight. Can this be extended to distributed memory systems with shared nothing topology, especially where reading remote memory is not very costly, like RDMA?

Reserved Slides



© Copyright 2015 Hewlett-Packard Development Company, L.P. The information contained herein is subject to change without notice.

Ongoing Work

✓ Try on 1,000 Cores~

✓ Further Resilience to Contentions

Combining Pessimistic Approach When Beneficial

✓SIMD, Xeon-Phi

Log-Gleaner's Sorting, Page Construction, etc

✓Non-C++ Interface

SQL/JDBC/ODBC for OLAP. But, for OLTP???



Easy OCC with Foster Twins

Simple, Robust, and Efficient Search:

Find a pointer that <u>probably</u> contains the key; Follow the page pointer; If the page's <u>key-range</u> contain the key: go on; Else: <u>locally</u> retry; No hand-over-hand verification/split counters. No unnecessary local retry. Absolutely no global retry.

 No Page-Version-set nor unnecessary aborts: All Records/TIDs always track-able via Foster-Twin!

