Beyond Main Memory Database Systems

Justin DeBrabant Brown University



Talk Overview

1. anti-caching architecture

- larger than memory datasets in main memory DBMS
- 2. anti-caching + persistent memory
 ditching the disk (finally!)



A bit of history...



1974 – System R query optimization recovery Itransaction serialization Allows concurrent execution of transactions Iots of locks





Change is Good





Great, that's what the buffer pool is for...right?



More Memory \rightarrow Higher Throughput?

- all data resides in memory (i.e. in buffer pool)
 - ► No disk stalls
- ▶ must still...
 - maintain buffer pool
 - ► lock/latch data
 - maintain ARIES-style recovery logs
- question: What is the overhead of all these things?





OLTP Through the Looking Glass, and What We Found There SIGMOD '08



Ok, no buffer pool, how about a distributed cache?



Distributed Caches ▶ e.g. memcached just in-memory key-value pair no inherent persistence > application programmer must maintain consistency ▶ or not!



So, we need a hybrid of these two architectures.



Fiorestand

Parallel Main Memory Transaction Processing System



H-Store: A High-Performance, Distributed Main Memory Transaction Processing System VLDB 2008.



H-Store Architecture

partitioned, shared-nothing
 data is sharded across nodes
 single-threaded main memory execution

no need for locks and latches



H-Store Architecture (cont'd) stored procedures > no ad hoc queries in OLTP Command logging **recovery** Snapshots + command log







H-Store Assumptions

- 1. OLTP workload
 - short-lived transactions that touch only a few records at a time
- 2. mostly single-site transactions
 - distributed transactions need multi-node coordination
- 3. data fits in memory
 - virtual memory is bad!



YCSB, update-heavy, data < memory





Assumptions: Revisited

1. OLTP workload

- "One Size Fits All": An Idea Whose Time Has Come and Gone
 - ▶ ICDE '05
- The End of An Architectural Era: (Its Time for A Complete Rewrite)
 - ► VLDB '07



Assumptions: Revisited

2. mostly single-site transactions

- Skew-Aware Automatic Database
 Partitioning In Shared-Nothing, Parallel
 OLTP Systems
 - *SIGMOD '12*
- On Predictive Modeling For Optimizing Transaction Execution in Parallel OLTP Systems
 - **VLDB '11**



Assumptions: Revisited





Workload Skew Exists!

hot data in memory
cold data to disk
goals
maintain transactional consistency
avoid blocking



Anti-Caching



Anti-Caching Phases evict pre-pass fetch

>merge



Evict

- 1. data > anti-cache threshold
- 2. dynamically construct anti-cache blocks of coldest tuples
- 3. asynchronously write to disk



Pre-Pass

- 1. a transaction enters pre-pass when evicted data is accessed
- 2. continues execution, creating list of evicted blocks
- 3. abort, queue blocks to be fetched



Fetch

data is fetched asynchronously from disk avoids blocking copied into merge buffer



Merge

- 1. previously aborted transaction is restarted
- 2. moves data from merge buffer to normal table
- 3. transaction executes normally



Multiple Restarts

- in-memory data for restarted transaction is relatively cold
 - mark tuples in pre-pass phase as hot
- In the data dependencies with evicted tuples
 - mark recently merged tuples as hot
- Iarger-than-memory queries still an issue
 not in OLTP



Other Design Points LRU chain \blacktriangleright embedded in tuple headers \rightarrow 0(1) updates EvictedTable stores <tuple id, block id> pairs for evicted tuples > anti-cache BerkeleyDB BROWN 吚 Database Group





Sounds like swapping...



Anti-Caching vs. Swapping

fine-grained eviction
 blocks constructed dynamically
 non-blocking fetches
 remove disk from critical path



Sounds like caching...



Anti-Caching vs. Caching data exists in exactly one location Caching architectures have multiple copies, must maintain consistency I data is moved, not copied ▶ goal is data size, not throughput



Benchmarking **YCSB** Zipfian skew >data > memory >read/write mix MySQL, MySQL + memcached BROWN 吚 Database Group

YCSB, read-only, data 8X memory



Conclusions ▶8-17X improvement for skewed workloads at 8X memory > avoids blocking for disk > fine-grained eviction **disk becomes the bottleneck**



Anti-Caching: A New Approach to Database Management System Architecture

| Justin DeBrabant | Andrew Pavlo | Stephen Tu |
|-----------------------------------|-------------------|----------------------------|
| Brown University | Brown University | MIT CSAIL |
| debrabant@cs.brown.edu | pavlo@cs.brown.ec | lu stephentu@csail.mit.edu |
| Michael Stonebraker Stan | | Stan Zdonik |
| MIT CSAIL Brown I | | Brown University |
| stonebraker@csail.mit.edu sbz@cs. | | 27@cs.brown.edu |

ABSTRACT

The traditional windom for building disk-based relational database management systems (DBMS) is to organize data in heavily-encoded blocks stored on disk, with a main memory block cache. In order to improve performance given high disk latency, these systems use a multi-threaded architecture with dynamic record-level locking that allows multiple transactions to access the database at the same time. Previous research has shown that this results in substantial overbaced for on-high transactions (JCIP) somelicitum. (JSI)

allows multiple transactions to access the database at the same time. Previous research has shown that this results in substantial overhead for on-line transaction processing (OLTP) applications [15]. The next generation DBMS seek to verecome these limitations with architecture based on main memory resident data. To overcome the restriction that id data fit in main memory, we reprose to disk in a transactionally-side rammer as the database grows in the Research of the second second second second second second to disk in a transactionally-side rammer, and at-caching architecture reverses the traditional storage hierarchy of disk-based systems. Main memory is now the primary storage device.

systems. Main memory is now the primary storage device. We implemented a prototype of our mit-aciding proposal in a high-performance, main memory OLTP DBMS and performed a series of experiments across as range of database stare, workload skews, and read/write mixes. We compared its performance with an open-source, data-base DBMS opticality fronted by a distributed main memory cache. One sheat DBMS opticality fronted by a distributed main memory cache. One sheat DBMS opticality is a sheat the sheat of workloads the anti-acading architecture has a performance advantage over either of the other architectures tested of up to 9× for a data size 6× larger thum memory.

1. INTRODUCTION

In TRKDUCE INTERM architecture of DBMSs has been predicated on the storage and management of data in heavily-encoded disk blocks. In most systems, there is a header at the beginning of each disk block to facilitate certain operations in the system. For eventy, this header structure is a line table? The the front of the block to support indirection to tuples. This allows the DBMS to recrepting to block, is read into main memory, it must then be translated into main memory format.

a disk block is read into män memory, it must then be translated into mäin memory format.
Permission su make ägind er band copies of all or part of this work for Permission su make ägind er band copies of all or parts of this work for entities of the strange of the strange of the strange of the optical band of the strange of the strange of the strange of the optical band of the strange of the strange of the strange of the strange band this optical band of the strange of the strange of the strange band the strange of the strange of the strange of the strange their results at The 39th International Conference on Wey Large Dan Basses, Against 26th - 312 or 320, River del Graft, Frence, Ialy, Proceedinge of the VLDE Endowment, Vol. 8, No. 14 DBMSk invariably maintain a buffer pool of Necks in main memory for faster access. When an exceeting query attempts to read a disk block, the DBMS first checks to see whether the block already exists in this buffer pool. If not, a block is exicted to make room for the needed one. There is substantial overhead to managing the system must maintain an ecicion order policy (e.g., least recently used). As noted block have to bo primoid in main memory, the covid mantening a DMS pool is nearly one-hind of all the CPU The excerce of manzing disk-resident data has for seen an

Cycle: listed by mir DBMS: or control of the second second second second second second second of new DBMS that pay the entire database in main memory and thus have no buffer pool [11]. Time:Ten was an early proposent of this approach [31] and more recent examples include H-Store [2, 18]. MemSQL [3], and KAMCloud [25]. H-Store (and its commercial version WithDI [14] performs arguincarily better than diskmain memory orientation, as well as from revising the overhead of concurrency control and heavy-weight data longging [22].

based DBMSs on standard OLTP benchmarks [29] because of this min memory orientiane, as well as from avoiding the overhead of concurrency control and heavy-weight data logging [22]. Is that this improve performance is only achievable when wents, is smaller than the amount of physical memory available in the system. If the database does not fit in memory, then the operating system will start to page virtual memory, available in the system. If the database does not fit in memory, then the operating system will start to page virtual memory, and main memory acsignificant poblem in a DBMS, like 14-Force, that execution of transactions is stalled while the page is fetched from disk. This is a significant poblem in a DBMS, like 14-Force, that executions in the size is the nucleon of heavy-weight locking and latchingceed the amount of real memory [2]. If memory is exceeded for if it might be at some point in the future), then a user must either (1) provision new hardware and mignare their database to a larger cluster, or (2) fall back to a traditional disk-based system, with its . The size divent prevents

One videly adopted performance enhancer is to use a main memory distributed cache, used na Mennechel [14], in front of a disbbased DBMS. Under this two ieir architecture, the application first look in the cache for the tuple of interest. If this tuple is not in the cache, then the application executes a query in the DBMS to fetch the desired data. Once the application receives this data from the DBMS, tupdues the cache for fast access in the future. Whenever a tuple is modified in the database, the application must invalidate is tache entry so that the next time it is accessed the application will retrieve the uncert vession from the DBMS. Many anothe web

Anti-Caching: A New Approach to Database Management System Architecture VLDB'13.









Memory Latencies (cycles) **SRAM: 1-30 DRAM: 100-300** Istash: 25,000-2,000,000 >disk: 5,000,000+



NVM ▶non-volatile random-access high write endurance except flash byte-addressable except flash



The Arms Race ▶ FeRAM high write endurance ► MRAM DRAM-like latency ▶ PCM (PRAM) DRAM-like capacity





NVM Emulation

goal: provide product-independent analysis

test wide range of latency profiles
automatically add specified latency
built by collaborators at Intel



H-Store + NVM Architectures

1. Anti-Cache to H-Store

- + fully utilizes memory hierarchy
- added memory overhead
- 2. H-Store on NVM
 - + no anti-cache overhead
 - wear-leveling necessary



Architecture 1



Architecture 2



Architecture 1, YCSB, read-only, data 8X memory



Work in Progress > implementation of Architecture 2 ▶ also in H-Store Further benchmarking YCSB, TPC-C, possibly others Architecture 1 vs. Architecture 2 > paper in preparation



Conclusions hardware has changed anti-caching to disk has 18-17X speedup over disk-based architecture better utilization of available memory next-generation persistent memories extend the benefits of anti-caching beyond skewed workloads



Collaborators





Questions?

debrabant@cs.brown.edu

