



WHAT'S REALLY NEW WITH NEWSQL

@ANDY_PAVLO



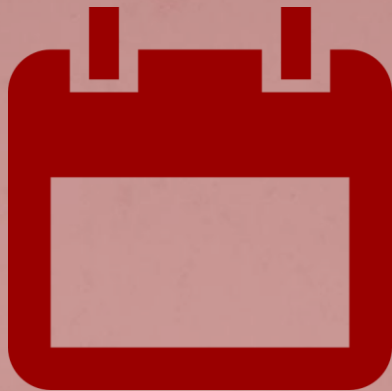
Fast



Repetitive



Small



The Last Decade of DATABASE SYSTEMS

Early 2000s — Sharding Middleware

- Custom middleware to combine multiple nodes together into a logical database.
- Route queries to correct node.

Middleware Problems

- Developers spend time writing middleware rather than working on core applications.
- Some features are implemented in application code.

Late 2000s — NoSQL

- Forgo transactional guarantees in order to achieve high-availability and high-scalability.
- Non-relational data models.

NoSQL Problems

- Developers write code to handle eventually consistent data, lack of transactions, and joins.
- Not all applications can give up strong transactional semantics.



The Rise of NEWSQL SYSTEMS

Aslett White Paper

[Systems that] deliver the scalability and flexibility promised by NoSQL while retaining the support for SQL queries and/or ACID, or to improve performance for appropriate workloads.

Matt Aslett – 451 Group (April 4th, 2011)

<https://www.451research.com/report-short?entityId=66963>

the 451 group

451 TECHDEALMAKER

4 April 2011 – Sector IQ

How will the database incumbents respond to NoSQL and NewSQL?

Analyst: Matt Aslett

The acquisition of **MySQL AB** by **Sun Microsystems** in January 2008 appeared to signal that open source databases were on the brink of opening up a new battleground against the proprietary database giants. In announcing the deal, Sun signaled its intention to provide the support and development resources required for MySQL to challenge the established vendors in supporting mission-critical, high-performance applications on Web-based architectures. Needless to say, reality was somewhat different as Sun faced wider problems of its own and eventually succumbed to takeover by **Oracle** (Nasdaq: ORCL) in April 2009, in doing so handing ownership of the leading commercial open source database to the database heavyweight.

We had previously argued that MySQL was very much the crown jewel of the open source database world thanks to its focus on Web applications, its lightweight architecture and its fast read capabilities, which made it potentially complementary technology for all of the established database players. Additionally, if Oracle's major rivals were seeking an obvious alternative to MySQL in 2009, they were out of luck.

Just two years later, however, the database market is awash with open source databases with lightweight architectures targeted at Web applications. Not only have the likes of **Monty Program** and **SkySQL** emerged to provide alternative support for MySQL and its forks, but there are also a large number of products available under the banner of NoSQL, which emerged in mid-2009 as an umbrella term for a loosely affiliated collection of non-relational database projects. We have also seen the emergence of what we have termed 'NewSQL' database offerings, with companies promising to deliver the scalability and flexibility promised by NoSQL while retaining the support for SQL queries and/or ACID (atomicity, consistency, isolation and durability), or to improve performance for appropriate workloads to the extent that the advanced scalability promised by some NoSQL databases becomes irrelevant.

From MySQL to NoSQL

Despite being a good match for many read-intensive applications, MySQL does not provide predictable performance at scale, particularly with a few writes thrown into the mix. The memcached distributed memory object-caching system can be used – and has been widely adopted – to improve performance but does not provide any persistence and lacks consistency. To some extent, the rise of NoSQL has been driven by the inadequacies of

451 TechDealmaker
Copyright 2011 The 451 Group

4 April 2011
Page 1 of 5

Stonebraker Article

*SQL as the primary interface.
ACID support for transactions
Non-locking concurrency control.
High per-node performance.
Shared-nothing architecture.*

Mike Stonebraker – Blog@CACM (June 16th, 2011)
<http://cacm.acm.org/blogs/blog-cacm/109710>

COMMUNICATIONS OF THE ACM

TRUSTED INSIGHTS FOR COMPUTING'S LEADING PROFESSIONALS

[Home](#) » [Blogs](#) » [BLOG@CACM](#) » [New SQL: An Alternative to NoSQL and Old SQL for New...](#) » [Full Text](#)

[BLOG@CACM](#)

New SQL: An Alternative to NoSQL and Old SQL for New OLTP Apps

Michael Stonebraker

June 16, 2011



Historically, Online Transaction Processing (OLTP) was performed by customers submitting traditional transactions (order something, withdraw money, cash a check, etc.) to a relational DBMS. Large enterprises might have dozens to hundreds of these systems. Invariably, enterprises wanted to consolidate the information in these OLTP systems for business analysis, cross selling, or some other purpose. Hence, Extract-Transform-and-Load (ETL) products were used to convert OLTP data to a common format and load it into a data warehouse. Data warehouse activity rarely shared machine resources with OLTP because of lock contention in the DBMS and because business intelligence (BI) queries were so resource-heavy that they got in the way of timely responses to transactions.

This combination of a collection of OLTP systems, connected to ETL, and connected to one or more data warehouses is the gold standard in enterprise computing. I will term it "Old OLTP." By and large, this activity was supported by the traditional RDBMS vendors. In the past I have affectionately called them "the elephants"; in this posting I refer to them as "Old SQL."

As noted by most pundits, "the Web changes everything," and I have noticed a very different collection of OLTP requirements that are emerging for Web properties, which I will term "New OLTP." These sites seem to be driven by two customer requirements:

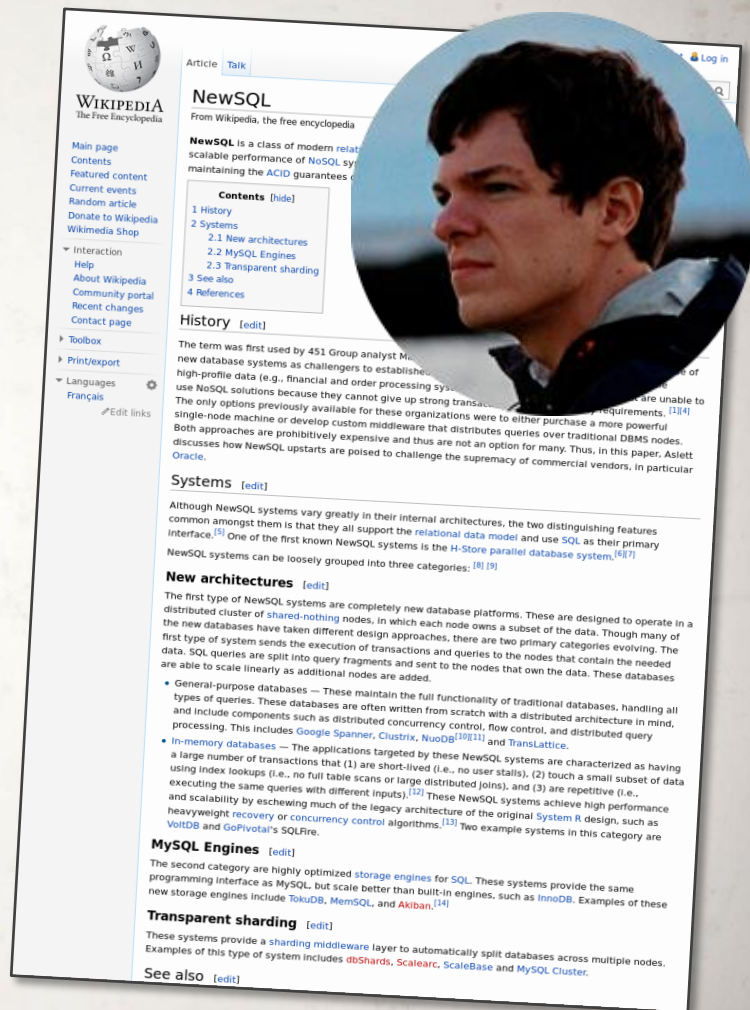
The need for far more OLTP throughput. Consider new Web-based applications such as multi-player games, social networking sites, and online gambling networks. The aggregate number of interactions per second is skyrocketing for the successful Web properties in this category. In addition, the explosive growth of smartphones has created a market for applications that use the phone as a geographic sensor and provide location-based services. Again, successful

Wikipedia Article

A class of modern relational database systems that provide the same scalable performance of NoSQL systems for OLTP workloads while still maintaining the ACID guarantees of a traditional database system.

Wikipedia (October 19th, 2013)

<http://en.wikipedia.org/wiki/NewSQL>



SCALABILITY

HIGH
(Many Nodes)

NOSQL

NEWSQL

LOW
(One Node)

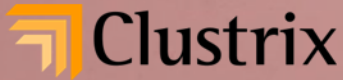
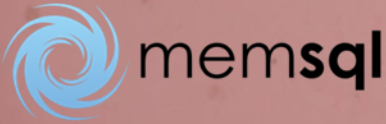
TRADITIONAL

WEAK
(None/Limited)

GUARANTEES

STRONG
(ACID)

New Design



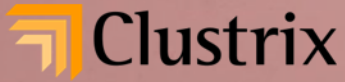
MySQL Engines



Middleware



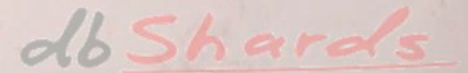
New Design



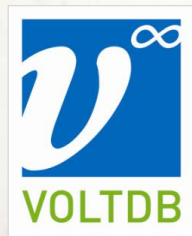
MySQL Engines



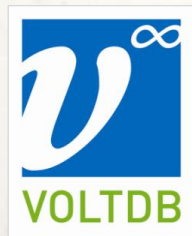
Middleware



Distributed Concurrency Control



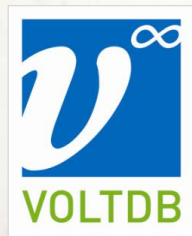
Distributed Concurrency Control



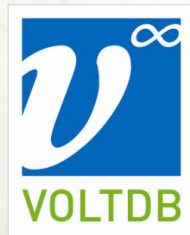
Distributed Concurrency Control



Distributed Concurrency Control



Distributed Concurrency



Phil Bernstein, et al. – Computing Surveys (June 1981)

<http://dl.acm.org/citation.cfm?id=356846>

Concurrency Control in Distributed Database Systems

PHILIP A. BERNSTEIN AND NATHAN GOODMAN

Computer Corporation of America, Cambridge, Massachusetts 02139

In this paper we survey, consolidate, and present the state of the art in distributed database concurrency control. The heart of our analysis is a decomposition of the concurrency control problem into two major subproblems: read-write and write-write synchronization. We describe a series of synchronization techniques for solving each subproblem and show how to combine these techniques into algorithms for solving the entire concurrency control problem. Such algorithms are called "concurrency control algorithms." We describe 48 principal methods, including all practical algorithms that have appeared in the literature plus several new ones. We concentrate on the structure and correctness of concurrency control algorithms. Issues of performance are given only secondary treatment.

Keywords and Phrases: concurrency control, deadlock, distributed database management systems, locking, serializability, synchronization, timestamp ordering, timestamps, two-phase commit, two-phase locking

CR Categories: 4.33, 4.35

INTRODUCTION

The Concurrency Control Problem

Concurrency control is the activity of coordinating concurrent accesses to a database in a multuser database management system (DBMS). Concurrency control permits users to access a database in a multiprogrammed fashion while preserving the illusion that each user is executing alone on a dedicated system. The main technical difficulty in attaining this goal is to prevent database updates performed by one user from interfering with database retrievals and updates performed by another. The concurrency control problem is exacerbated in a distributed DBMS (DDBMS) because (1) users may access data stored in many different computers in a distributed system, and (2) a concurrency control mechanism at one computer cannot instantaneously know about interactions at other computers.

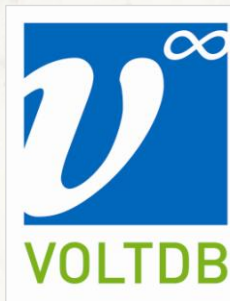
Concurrency control has been actively investigated for the past several years, and the problem for nondistributed DBMSs is well understood. A broad mathematical theory has been developed to analyze the problem, and one approach, called *two-phase locking*, has been accepted as a standard solution. Current research on nondistributed concurrency control is focused on evolutionary improvements to two-phase locking, detailed performance analysis and optimization, and extensions to the mathematical theory.

Distributed concurrency control, by contrast, is in a state of extreme turbulence. More than 20 concurrency control algorithms have been proposed for DDBMSs, and several have been, or are being, implemented. These algorithms are usually complex, hard to understand, and difficult to prove correct (indeed, many are incorrect). Because they are described in different terminologies and make different assumptions

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its copy otherwise, or to republish, requires a fee and/or specific permission.

© 1981 ACM 0010-4892/81/0600-0185 \$00.75

Main Memory Systems



memsql

IMPLEMENTATION TECHNIQUES FOR MAIN MEMORY DATABASE SYSTEMS

David J. DeWitt¹, Randy H. Kats², Frank Olken³,
 Leonard D. Shapiro⁴, Michael R. Hodoraker⁴, David Wood²

¹Computer Sciences Department, University of Wisconsin
²ICS Department, University of California at Berkeley
³CSAM Department, Lawrence Berkeley Laboratory
⁴Department of Computer Science, North Dakota State University

ABSTRACT With the availability of very large, relatively inexpensive main memories, it is becoming possible to keep large databases resident in main memory. In this paper we consider the changes necessary to permit a relational database system to take advantage of large amounts of main memory. We evaluate AVL, B+-tree, and other access methods for main memory databases, hash-based query processing, and other techniques. We also discuss the need for new indexing techniques and other techniques for main memory. As expected, B+-trees are the preferred storage mechanism unless more than 80-90% of the database fits in main memory. A somewhat surprising result is that hash based query processing strategies are advantageous for large memory situations.

Key Words and Phrases: Main Memory Databases, Access Methods, Join Algorithms, Access Planning, Recovery Mechanisms

1. Introduction

Throughout the past decade main memory prices have plummeted and are expected to continue to do so. At the present time, memory for supercomputers such as the VAX 11/780 costs chips will be commonplace and should further reduce prices by another order of magnitude. Thus, in 1980 a 64Kbyte of memory should cost less than \$200,000. If a megabyte of memory is available, the price might be as low as \$50,000.

With the availability of large amounts of main memory, it becomes possible to contemplate the storage of databases in main memory. In fact, DAS Fast Path [DAS82] has supported such databases for some time. In this paper we consider the changes that might be needed to a relational database system if most (or all) of a database fits in main memory.

In Section 2, the performance of alternative access methods for main memory database systems are considered. Algorithms for relational database operations in the environment are presented and evaluated in Section 3. In Section 4, we describe how access planning will be affected by the availability of large amounts of main memory for query processing. Section 5 discusses recovery in main memory resident databases. Our conclusions and suggestions for future research are presented in Section 6.

The research was partially supported by the National Science Foundation under grant MCS-81-0106, MCS-82-0476, by the Department of Energy under contract DE-AC02-81OR21000, DE-AC02-82OR00001, and by the Air Force Office of Scientific Research under Grant DMR-81-0106.

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

© 1984 ACM 0-89791-128-8/84/006-0001 \$00.75

future research are contained in Section 6.

2. Access Methods for Memory Resident Databases

The standard access method for data on disk is the B+-tree [COM67], providing both random and sequential key access. A B+-tree is especially designed to provide fast access to disk-resident data and makes fundamental use of the page size of the device. On the other hand, if a key is known to be in memory, then a B+-tree is less efficient. In this section we analyze the performance of both structures for a relation R with the following characteristics:

- [N] number of tuples in relation R
- [K] width of the key for R in bytes
- [L] width of a tuple in bytes
- [P] page size in bytes
- [4] size of a pointer in bytes

We have analyzed two cases of interest. The first is the cost of retrieving a single tuple using a random key value. An example of this type of query is:

retrieve (emp salary) where emp name = "Jones"

The second case analyzed is the cost of making N records sequentially. Consider the query:

retrieve (emp salary, emp name) where emp name = "J"

which requests data on all employees whose names begin with J. To execute this query, the database system would locate the first employee with a name beginning with J and then read sequentially. This second case analyzes the sequential access portion of such a command.

For both cases (random and sequential access), there are two costs that are specific to the access method:

[page read] the number of pages read to execute the query

[comparison] the number of record comparisons required to isolate the particular data of interest

The number of comparisons is indicative of the CPU time required to process the command while the number of page reads approximates the I/O cost.

To compare the performance AVL and B+-trees, we propose the following cost function:

cost = 2 * [page-reads] + [comparisons]

Since a page read consumes perhaps 200 instructions of operating system overhead and 30 milliseconds of elapsed time while a compare can easily be done in 200, we expect relative values of 2 to be in the range of 10 to 30. Later in the section we will see several values in the range.

Moreover, it is possible (although not very likely) that an AVL-tree comparison will be cheaper than a B+-tree comparison. The reasoning is that the B+-tree record must be located within a page while an AVL tree does not contain any page structure and a record can be directly located. Consequently, we assume that an AVL-tree comparison costs Y times a B+-tree comparison for some

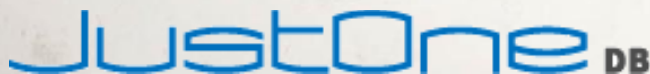
Hybrid Architectures



deap[®]db

JustOne[®] DB

Hybrid Architectures



Anastassia Ailamaki, et al. – VLDB (2001)
<http://dl.acm.org/citation.cfm?id=672367>

Weaving Relations for Cache Performance

Anastassia Ailamaki[†]
Carnegie Mellon University
ailamaki@cs.cmu.edu

David J. DeWitt
Univ. of Wisconsin-Madison
dewitt@cs.wisc.edu

Mark D. Hill
Univ. of Wisconsin-Madison
markhill@cs.wisc.edu

Marios Skounakis
Univ. of Wisconsin-Madison
marios@cs.wisc.edu

Abstract

Relational database systems have traditionally optimized for I/O performance and organized records sequentially on disk pages. Recent research, however, indicates that cache utilization and performance is becoming increasingly important on modern platforms. In this paper, we first demonstrate that in-page data placement is the key to high cache performance and that NSM exhibits low cache utilization on modern platforms. Next, we propose a new data organization model called PAX (Partitioned Attributes Across), that significantly improves cache performance by grouping together all values of each attribute within each page. Because PAX only affects layout inside the pages, it incurs no storage penalty and does not affect I/O behavior. According to our experimental results, when compared to NSM (a) PAX exhibits superior cache and memory bandwidth utilization, saving at least 75% of NSM's stall time due to data cache accesses, (b) range selection queries and updates on memory-involving I/O execute 11-48% faster, and (c) TPC-H queries

1 Introduction

The communication between the CPU and the secondary storage (I/O) has been traditionally recognized as the major database performance bottleneck. To optimize data transfer to and from mass storage, relational DBMSs have long organized records in slotted disk pages using the N-ary Storage Model (NSM). NSM stores records contiguously starting from the beginning of each disk page, and uses an offset (slot) table at the end of the page to locate the beginning of each record [27].

Unfortunately, most queries use only a fraction of each record. To minimize unnecessary I/O, the Decomposition Storage Model (DSM) was proposed in 1985 [10]. DSM partitions an n -attribute relation vertically into n sub-relations, each of which is accessed only when the corresponding attribute is needed. Queries that involve multiple attributes from a relation, however, must spend

tremendous additional time to join the participating sub-relations together. Except for Sybase IQ [33], today's relational DBMSs use NSM for general-purpose data placement [20][29][32].

Recent research has demonstrated that modern database workloads, such as decision support systems and spatial applications, are often bound by delays related to the processor and the memory subsystem rather than I/O [20][5][26]. When running commercial database systems on a modern processor, data requests that miss in the cache hierarchy (i.e., requests for data that are not found in any of the caches and are transferred from main memory) are a key memory bottleneck [1]. In addition, only a fraction of the data transferred to the cache is useful to the query: the item that the query processing algorithm requests and the transfer unit between the memory and the processor are typically not the same size. Loading the cache with useless data (a) wastes bandwidth, (b) pollutes the cache, and (c) possibly forces replacement of information that may be needed in the future, incurring even more delays. The challenge is to repair NSM's cache behavior without compromising its advantages over DSM.

This paper introduces and evaluates Partitioned Attributes Across (PAX), a new layout for data records that combines the best of the two worlds and exhibits performance superior to both placement schemes by eliminating unnecessary accesses to main memory. For a given relation, PAX stores the same data on each page as NSM. Within each page, however, PAX groups all the values of a particular attribute together on a minipage. During a sequential scan (e.g., to apply a predicate on a fraction of the record), PAX fully utilizes the cache resources, because on each miss a number of a single attribute's values are loaded into the cache together. At the same time, all parts of the record are on the same page. To reconstruct a record one needs to perform a mini-join among minipages, which incurs minimal cost because it does not have to look beyond the page.

We evaluated PAX against NSM and DSM using (a) predicate selection queries on numeric data and (b) a variety of queries on TPC-H datasets on top of the Shore storage manager [7]. We vary query parameters including selectivity, projectivity, number of predicates, distance between the projected attribute and the attribute in the predicate, and degree of the relation. The experimental results show that, when compared to NSM, PAX (a) incurs 50-75% fewer second-level cache misses due to data

[†] Work done while author was at the University of Wisconsin-Madison.
Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the VLDB copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Very Large Data Base Endowment. To copy otherwise, or to republish, requires a fee and/or special permission from the Endowment.
Proceedings of the 27th VLDB Conference,
Roma, Italy, 2001

Query Code Compilation



Google
Spanner

Query Code Compilation



Google
Spanner

M. M. Astrahan, et al. – ACM TODS (1976)

<http://dl.acm.org/citation.cfm?id=320457>

System R: Relational Approach to Database Management

M. M. ASTRAHAN, M. W. BLASGEN, D. D. CHAMBERLIN,
K. P. ESWARAN, J. N. GRAY, P. P. GRIFFITHS,
W. F. KING, R. A. LORIE, P. R. MCJONES, J. W. MEHL,
G. R. PUTZOLU, I. L. TRAIGER, B. W. WADE, AND V. WATSON
IBM Research Laboratory

System R is a database management system which provides a high level relational data interface. The system provides a high level of data independence by isolating the end user as much as possible from underlying storage structures. The system permits definition of a variety of relational views on common underlying data. Data control features are provided, including authorization, integrity assertions, triggered transactions, a logging and recovery subsystem, and facilities for maintaining data consistency in a shared-update environment.

This paper contains a description of the overall architecture and design of the system. At the present time the system is being implemented and the design evaluated. We emphasize that System R is a vehicle for research in database architecture, and is not planned as a product.

Key Words and Phrases: database, relational model, nonprocedural language, authorization, locking, recovery, data structures, index structures
CR categories: 3.74, 4.22, 4.33, 4.35

1. INTRODUCTION

The relational model of data was introduced by Codd [7] in 1970 as an approach toward providing solutions to various problems in database management. In particular, Codd addressed the problems of providing a data model or view which is divorced from various implementation considerations (the data independence problem) and also the problem of providing the database user with a very high level, nonprocedural data sublanguage for accessing data.

To a large extent, the acceptance and value of the relational approach hinges on the demonstration that a system can be built which can be used in a real environment to solve real problems and has performance at least comparable to today's existing systems. The purpose of this paper is to describe the overall architecture and design aspects of an experimental prototype database management system called System R, which is currently being implemented and evaluated at the IBM San Jose Research Laboratory. At the time of this writing, the design has been

Copyright © 1976, Association for Computing Machinery, Inc. General permission to republish, but not for profit, all or part of this material is granted provided that ACM's copyright notice is given and that reference is made to the publication, to its date of issue, and to the fact that reprinting privileges were granted by permission of the Association for Computing Machinery. Authors' address: IBM Research Laboratory, San Jose, CA 95103.

ACM Transactions on Database Systems, Vol. 1, No. 2, June 1976, Pages 97-137.

Recap

- Distributed Concurrency Control
- Main Memory Stores
- Hybrid Architectures
- Query Code Compilation



The Future of OLTP DBMS RESEARCH

Nearly Solved Problems



Fine-grain Elasticity

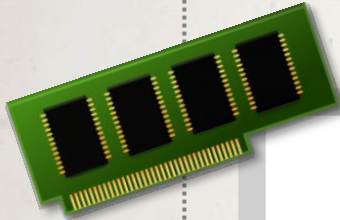


Hybrid Architectures



Larger-than Memory DBs

Application

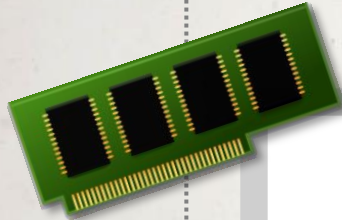


Primary Storage



Anti-Cache

Application



Primary Storage



Anti-Cache

Application



Primary Storage



Four blue arrows point downwards from the Primary Storage box to the Anti-Cache box.



Anti-Cache

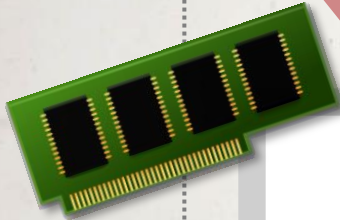
Application



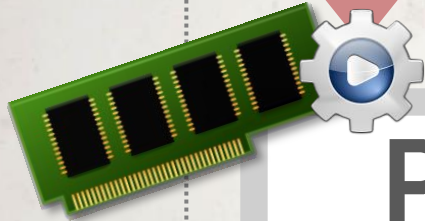
Primary Storage



Anti-Cache



Application

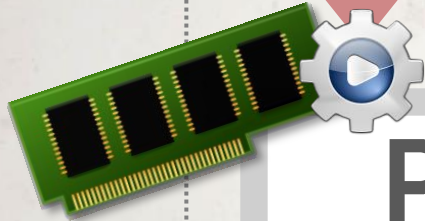


Primary Storage



Anti-Cache

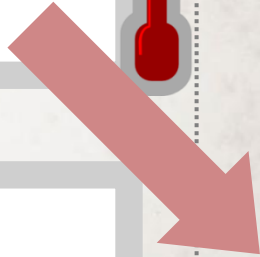
Application



Primary Storage



Anti-Cache

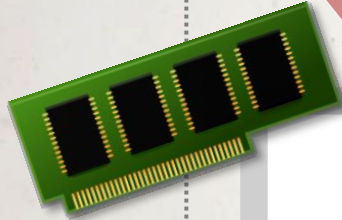


```
<script language=javascript>  
document.write("Hello World!")  
</script>  
<script language=javascript>  
window.prompt("Enter a number")  
</script>  
<script language=javascript>  
window.prompt("Enter a number")  
</script>
```

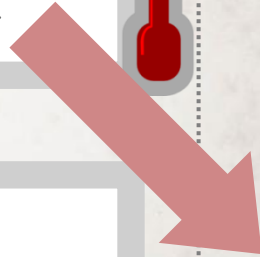
Application



Primary Storage



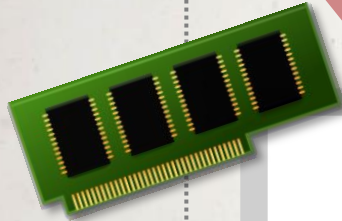
Anti-Cache



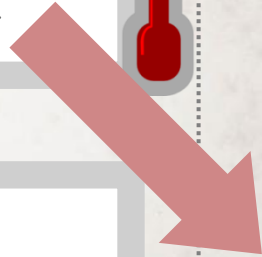
Application



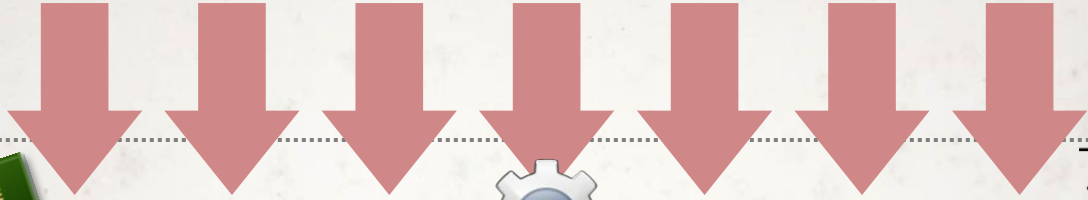
Primary Storage



Anti-Cache



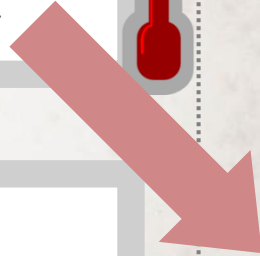
Application



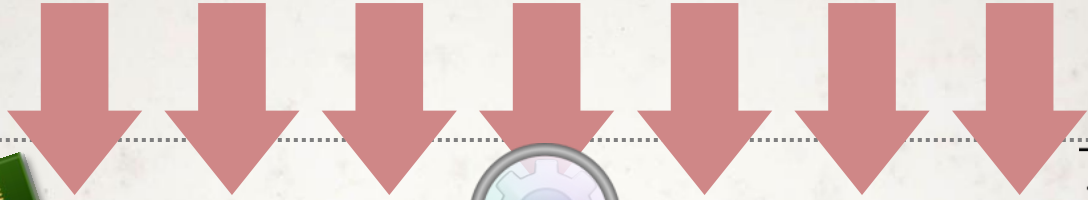
Primary Storage



Anti-Cache



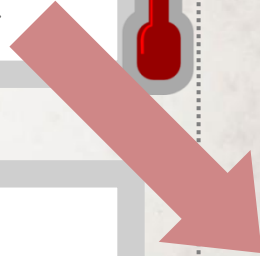
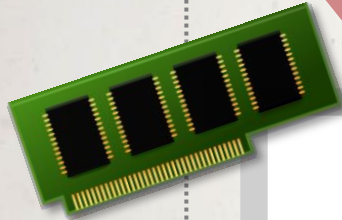
Application



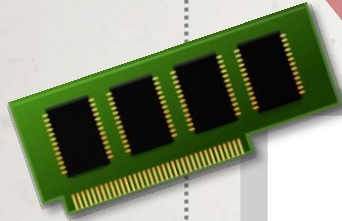
Primary Storage



Anti-Cache



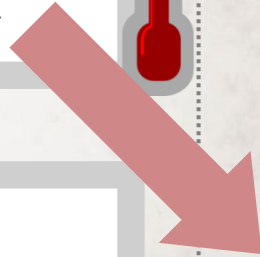
Application



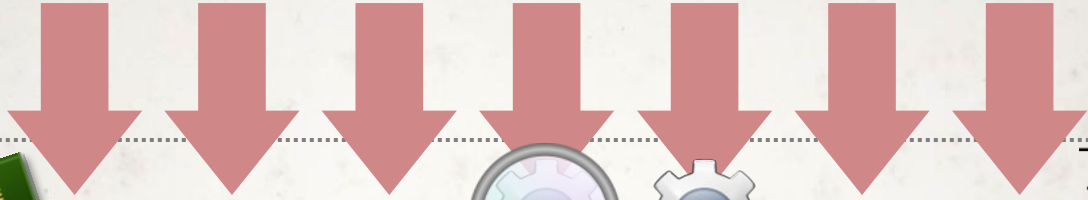
Primary Storage



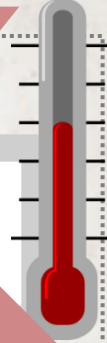
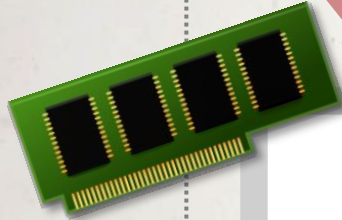
Anti-Cache



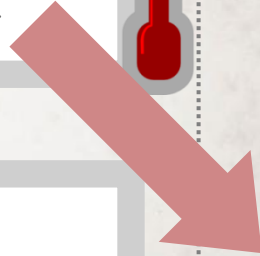
Application



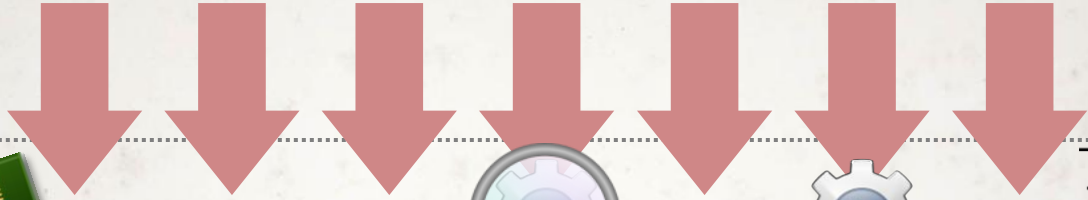
Primary Storage



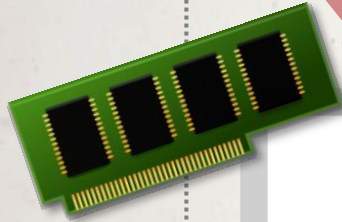
Anti-Cache



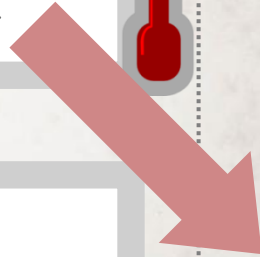
Application



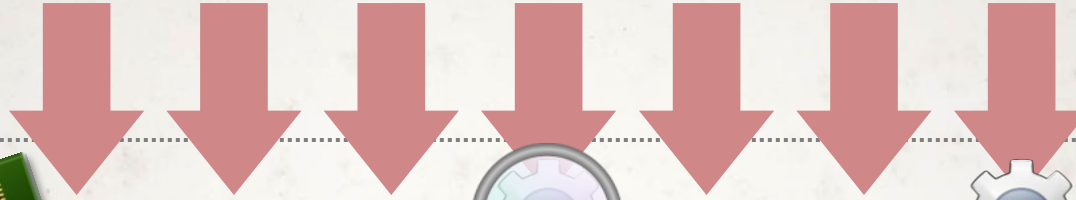
Primary Storage



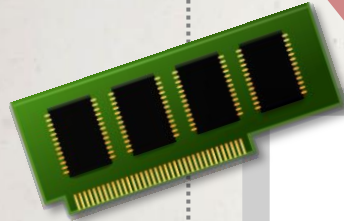
Anti-Cache



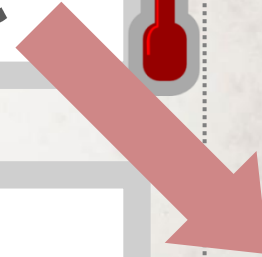
Application



Primary Storage



Anti-Cache



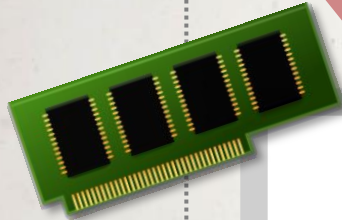
Application



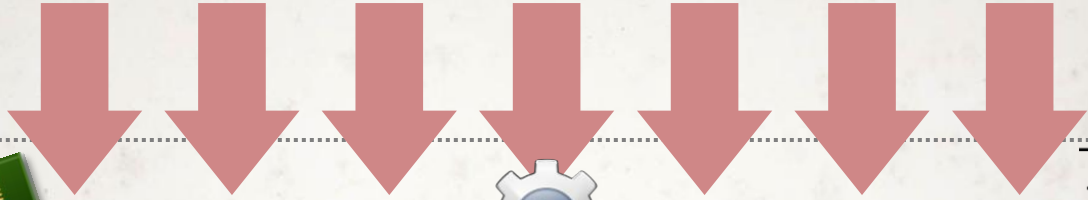
Primary Storage



Anti-Cache



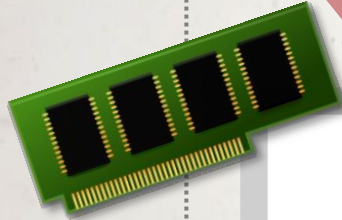
Application



Primary Storage

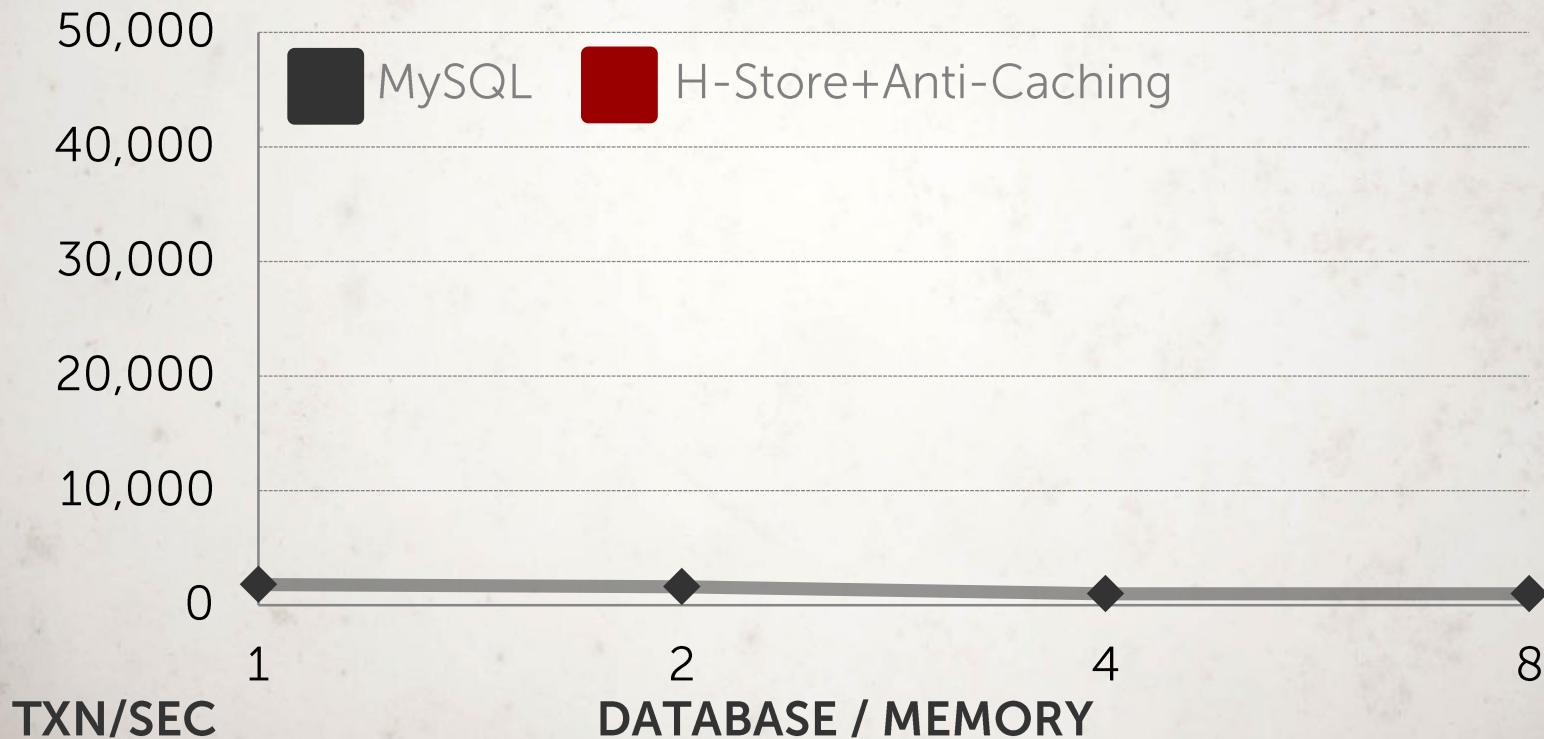


Anti-Cache



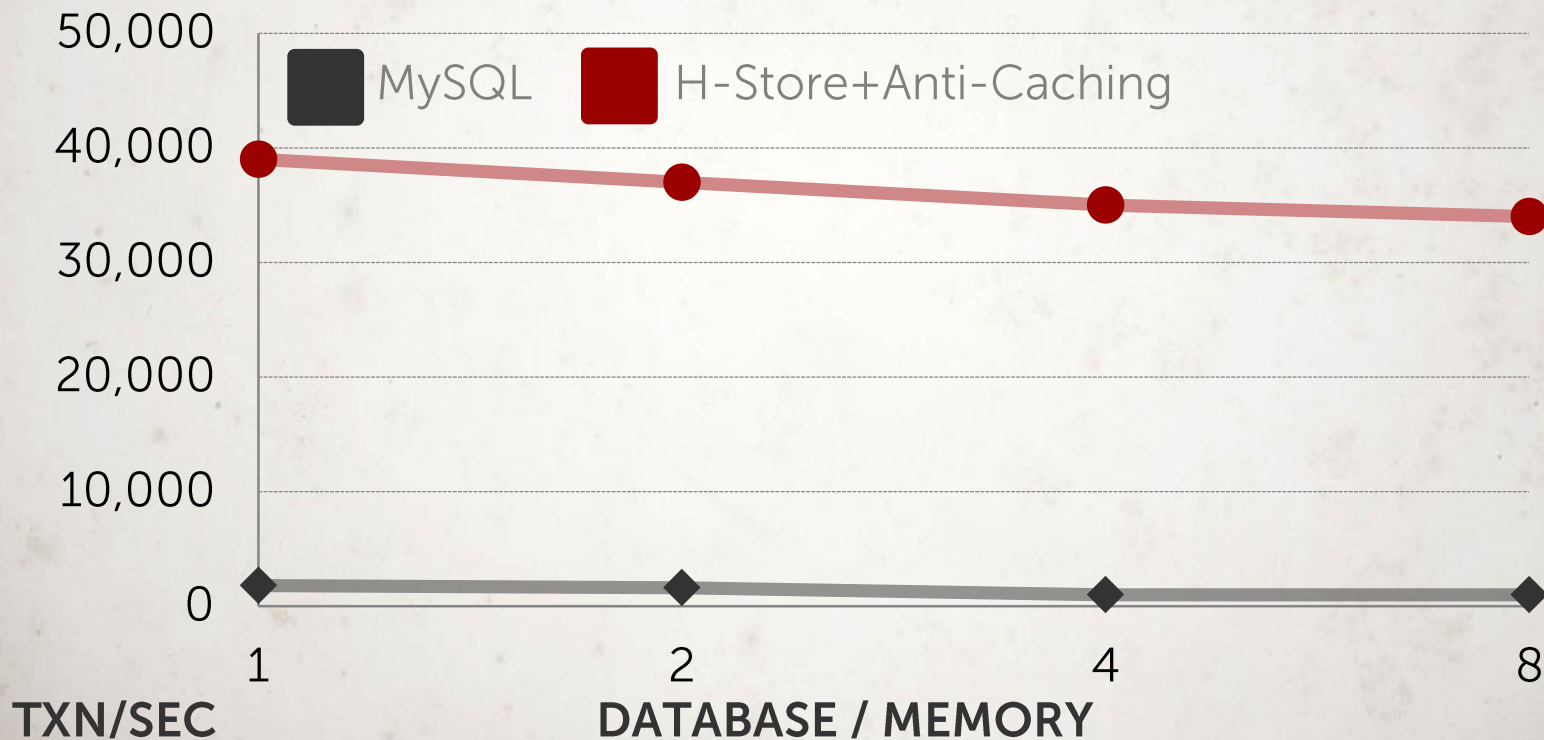
TPC-C Benchmark

100% Single-Partition Transactions – 100 Warehouses



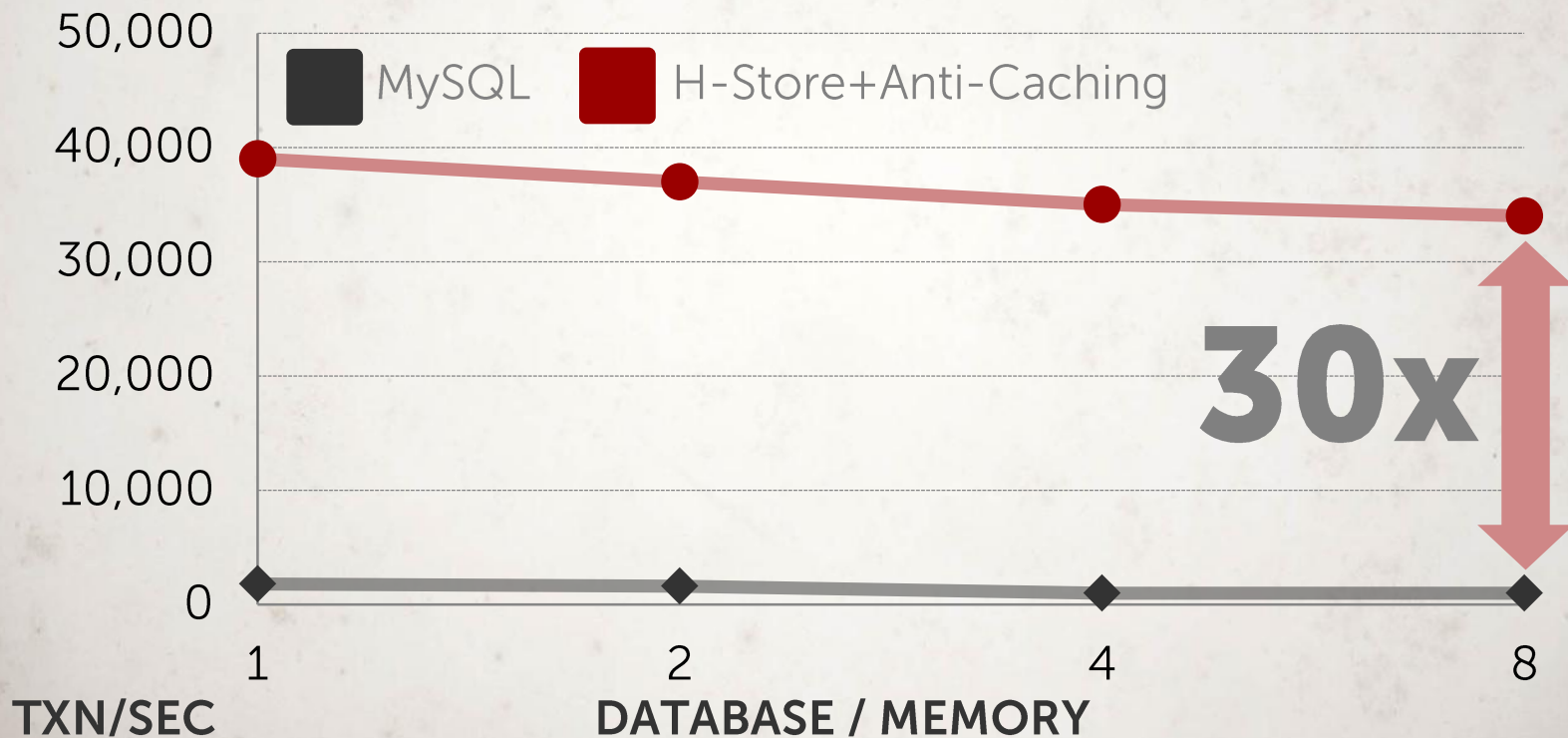
TPC-C Benchmark

100% Single-Partition Transactions – 100 Warehouses

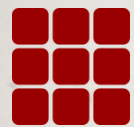


TPC-C Benchmark

100% Single-Partition Transactions – 100 Warehouses



Future Research @ CMU



Many-Core Concurrency



Non-Volatile Memory



Geo-replicated DBMS

Spänner

Spä
X
ner

F₂



BillHoweDB

BillH~~**X**~~**veDB**

OpauSQLTM

Oracle, Please Acquire Us

OpauSQL™ — Design Principles

- Don't treat the DB as a black box.
- Machine learning to understand intra- & inter-txn dependencies.
- Introspection of integrity constraints.

OLTP Application Library

- Examine open source software to create a catalog of application properties.
- Automatically infer optimizations by examining DB access patterns.

Conclusion

- Most NewSQL systems are using known ideas to achieve high-performance for OLTP workloads.
- Database systems research is back at CMU.

END

@ANDY_PAVLO