



Main-Memory Centric Data Management – Open Problems and Some Solutions Wolfgang Lehner

Technische Universität Dresden

New Realities

- TB disks < \$100
- Everything is data
- Rise of data-driven culture
 - CERN's LHC generates 15 PB a year

> The Reality: The Petabyte Age

Sloan Digital Sky Survey (200 GB/night)

The quest for knowledge used to begin with grand theories. Now it begins with massive amounts of data.

Welcome to the Petabyte Age.







The Web is a huge source of information: search engines (Google, Yahoo!) collect and store billions of documents and click streams

- 20 PB processed every day at Google (2008)
- 200 million photos are uploaded to Facebook every day → 2,314 photos/second (2010)
- 60 hours of video uploaded to YouTube every minute (2012)
- By 2015, 1 Zettabyte of data will flow over the internet per day (Cisco Visual Networking Index, June 2011)
- One zettabyte = stack of books from Earth to Pluto 20 times





chnolog

TECH | 2/16/2012 @ 11:02AM | 1,530,629 views



Tweet

2.4k

28207

Share

How Target Figured Out A Teen Girl Was Pregnant Before Her Father Did



in Share

9241

reddit 💮

329

Submit

⊈ +1 _

💵 💹 🔜 🏙 🚺 📕 🔍 258 comments, 149 called-out

Every time you go shopping, you share intimate details about your consumption patterns with retailers. And many of those retailers are studying those details to figure out what you like, what you need, and which coupons are most likely to make you happy. <u>Target</u>, for example, has figured out how to data-mine its way into your womb, to figure out whether you have a baby on the way long before you need to start buying diapers.

Charles Duhigg outlines in the <u>New York</u> <u>Times</u> how Target tries to hook parents-to-be at that crucial moment before they turn into rampant — and loyal — buyers of all things



+ Comment now

Target has got you in its aim

pastel, plastic, and miniature. He talked to Target statistician Andrew Pole before Target freaked out and cut off all communications — about the clues to a customer's impending bundle of joy. Target assigns every customer a Guest ID number, tied to their credit card, name, or email address that becomes a

_1anagement

4



Main Memory Centric Data Management





RAM locality is king



- Increased CPU calculation speed
- Increased memory bandwith
- Stagnating memory latency
- Avoid CPU idle time due to missing data
- RAM locality very important



> Memory Performance Comparison





Results for a quad-core i7 2.66GHz, DDR3 1666. 32GB data accessed total. © Tim Kaldewey

Is memory the new disk ???

- Some characteristics are very similar, e.g. random vs. sequential
- Memory architecture complicates things !

cessed total.		K
Aspect		
Rand vs. seq	1-2 orders of magnitude	3 orders of magnitude
Access granuarity	Byte addressable in theory, Caches get in the way	1 disk block, usually 4KB
Writes	Read-modify write (CL)	Read-modify-write (block)
Concurrency	Parallel memory access for peak performance	Multiple seq. streams ➔ random access

→ Motivation for CPU-cache aware data structures, e.g. Vectorwise/X11/Ingres, SAP HANA, Greenplum, Vertica Flexstore , Oracle 11g R2 (PAX)

HAEC – Collaborative Research Center



Energy-Adaptive High-Speed Computing Platform



Optical Interconnect

- adaptive analog/digital circuits for e/o transceiver
- embedded polymer waveguide
- packaging technologies (e.g. 3D stacking of Si/III-V hybrids)
- 90° coupling of laser

Radio Interconnect

- on-chip/on-package antenna arrays
- analog/digital beamsteering and interference minimization
- 100Gb/s
- 100-300GHz channel
- 3D routing & flow management



Project A01

Millimeter-Wave Integrated Circuits for Ultra High-Speed Wireless Boardto-Board Computer Communications

Antennas and Wave Propagation for Adaptive Wireless Backplane Communication





Network Coding for Wireless and Wired Onboard and Backplane Colution high efficient ion:

© Prof. Dr.-Ing. Wolfgang Lehner

Low-resolution design: phase modulations (A02) relax amplitude accuracy, require more bandwidth

Parallelization of low-res ADC for max bandwidth



[SWB+10b]	BiCMOS 130 nm	122 GHz	3 GHz (2.5%)
[POZ-10]	BiCMOS 130 nm	160 GHz	~7 GHz (4%)
[BGV+09]	CMOS 65 nm	60 GHz	~6 GHz (10%)
CCN (Tx above)	BiCMOS 250 nm	61 GHz	10 GHz (16%)
HAEC	BiCMOS >130 nm	>150 GHz	> 20 GHz (>13%)

[GAB+10]	CMOS 65 nm	40 GS/s	18 GHz	6 bit	1.5 W	Time interl. 16 x SAR ADC
[LC10]	BiCMOS 180 nm	50 GS/s	20 GHz	5 bit	5.4 W	Time interl. 2 x flash ADC
HAEC	BiCMOS or CMOS	> 50 GS/s	> 20 GHz	1-2 bit (?)	Min.	Time interl. (?)

Sub-THz radio receiver:

Technology: BiCMOS for higher efficiency at target performance

Design approach: positive feedback for $>f_t/2$ operation

Secure Network Coding for Board-to-Board Communication

25 GHz bandwidth ADC:



Grou

10

Low-resolution for both energy efficiency and fastest



...a plea for specialized DB systems



The End of an Architectural Era (It's Time for a Complete Rewrite)

Michael Stonebraker Samuel Madden Daniel J. Abadi Stavros Harizopoulos Nabil Hachem AvantGarde Consulting, LLC nhachem@agdba.com Pat Helland Microsoft Corporation phelland@microsoft.com

Implications for the Elephants

- They are selling "one size fits all"
- Which is 30 year old legacy technology that good at nothing



Implications for the Community











Some Techniques

- Compression
- Indexing
- Resilience
- (Hardware) Transactional Memory



Many different compression schemes



Challenges:

- What procedure to apply?
- When to apply compression at all?



- For each unique value create dictionary entry
- Dictionary can be per-block or percolumn
- Column-stores have the advantage that dictionary entries may encode multiple values at once



Database Technology

- Encodes values as b bit offset from chosen frame of reference
- Special escape code (e.g. all bits set to 1) indicates a difference larger than can be stored in b bits
- After escape code, original (uncompressed) value is written







QUERY:
SELECT custID,SUM(price)
WHERE (prodID = 4) AND
(storeID = 1) AND GROUP BY custID

Create rows first

Disadvantages:

- Need to construct all tuples
- Need to decompress data
- Poor memory bandwith utilization
- Loose opportunity for vectorized operation

> Late Materialization Example



QUERY:
SELECT custID,SUM(price)
FROM table
WHERE (prodID = 4) AND
(storeID = 1) AND
GROUP BY custID

	4	2	2	7
	4	1	3	13
	4	3	3	42
1	4	1	3	80
p	rodID	storeID	custID	price



- Database Technology
- Each operator adjusts its output to the requirements of the successive operator





Some Techniques

- Compression
- Indexing
- Resilience
- (Hardware) Transactional Memory

> The need for Indexing: Scan vs. Index

- About 40 GB/s scan performance with in-memory Databases
- Real-Time Analytics requires low response time



 \rightarrow Indexing is still necessary for ordering, grouping, point- and range-queries

Database Technology

Group



> KISS-Tree Overview

Database Technology

Properties

- Specialized version for 32bit keys
- Latch-free updates
- Order-preserving
- 2-3 memory accesses per key
- → Comparable fast to reported order-preserving in-memory indexes for read access

 \rightarrow BUT:

High update performance

- Heterogeneous in-memory index structure
- Combination of direct and indirect addressing
- Takes advantage of virtual memory management
- Enables different compression mechanisms



KISS-Tree: Smart Latch-Free In-Memory Indexing OAEMON 2012 on Modern Architectures

> Thomas Kissinger, Benjamin Schlegel, Dirk Habich, Wolfgang Lehner Database Technology Group Technische Universität Dresden 01082 Dresden. Germany (firsthame lastname)@tu-dresden.de

ABSTRACT

Growing main memory capacities and an increasing number of hardness threads in modern severe systems lad to fundamental changes in database architectures. Most importantly, query processing is nowadays performed on data that is often completely stored in main memory. Despite of a high main memory scan performance, index structures arcstill important components, but they have to be designed from scratch to cope with the specific characteristics of main memory and to exploit the high degree of parallelism. Current research mainly focused on adapting block-optimized B+-Trees, but these data structures were designed for secondary memory and involve comprehensive structural maintenance for updates.

In this paper, we present the *KISS*-Tree, a latch-free inmemory index that is optimized for a minimum number of memory accesses and a high number of concurrent updates. More specifically, we ainfor the same performance as modern hash-based algorithms but keeping the order-preserving nature of trees. We achieve this by using a perfix tree that incorporates virtual memory management functionality and compression schemes. In our experiments, we evaluate the *KISS*-Tree on different workloads and hardware platforms and compare the results to existing in-memory indexes. The *KISS*-Tree offset he highest performance on current architectures, a balanced real/write performance, and has a low memory footprint.



Figure 1: KISS-Tree Overview.

indexes) in meranzy and use secondary memory (e.g., disk) only for persistence. While disk block accesses constituted the botthenck for disk-based systems, modern in-memory databases shift the memory hierarchy closer to the CPU and face the "memory wall" [13] as new bottleneck. Thus, they have to care for CPU disk caches, TLBs, main memory accesses and access patterns. This sessurilal change also facts indexes and forces us to design new index structures that are new optimized for these new design goals. Morover, the movement from disks to main memory duramatically increased data access bandwidth and reduced latency. In combination with the increasing row, because the overhead true imposes a new challenges for us, because the overhead





Requirement for Level 2: All nodes have to be stored sequentially in memory











© Prof. Dr.-Ing. Wolfgang Lehner |

Main-Memory Centric Data Management | 30





Thread-Local Memory Management Subsystem

> Duplicate Handling

- Database Technology Group
- Efficient duplicate handling necessary for query processing
 - Scanning a linked list results in random memory accesses



- Page boundaries are a barrier for hardware prefetchers
 - store values sequentially in 4KB blocks
 - blocks grow exponentially until reaching 4KB
 - \rightarrow trade-off between scan performance and memory consumption

> Read Performance

Database Technology



© Prof. Dr.-Ing. Wolfgang Lehner |

Main-Memory Centric Data Management

> Update Performance





© Prof. Dr.-Ing. Wolfgang Lehner |



Some Techniques

- Compression
- Indexing
- Resilience
- (Hardware) Transactional Memory



Database Technolog

Grour

Increasing Component Error Rates

- Decreasing feature sizes (new tech generations)
- Reduced voltage supply
- Static (hard) vs. dynamic (soft) errors
- 8% increase error rate per tech generation [Borkar05]
- 25,000 70,000 FIT / Mbit [Schroeder09]

Increasing System Error Rates

- Increasing scale
 - # of components (core, transistor)
 - Memory capacities
- Example:
 - Fixed error rate / component



Errors and error-prone behavior will become the normal case

P(**()**)=**0.039**


Implicit (silent) vs. Explicit (detected/corrected) Errors

State-of-the-art: error detection and correction at HW/OS level

State-of-the-Art: Resilient Memory

ECC / parity bits / memory scrubbing / full data redundancy

ECC Extended Hamming(7+1,4)

d1 d2 d3 d4 🗭 p1 p2 d1 p3 d2 d3 d4 P

State-of-the-Art: Resilient Computing

Computation redundancy

Double Modular Redundancy (DMR):





Such resiliency mechanisms cause "resiliency costs"

37

© Prof. Dr.-Ing. Wolfgang Lehner |

Resiliency Costs Categories

> Motivation: Resiliency Costs (2)

- Performance overhead (throughput, latency)
- Memory overhead
- Energy consumption
- Monetary HW costs

Resiliency Costs @ OS-Level

- **Memory overhead** (capacity, bandwidth)
- **Computation overhead**
- Energy consumption (increased time)

Resiliency Costs @ HW-Level

- **Monetary HW costs** (Chipset, ECC RAM)
- **Energy consumption** (time, chip space)

Increasing error rates ~ increasing resiliency costs!

Computation overhead



ECC RAM Memory

ECC RAM





Challenge

- Problem: data loss/corruption (explicit/implicit)
- Goal: data stability by data redundancy and error correction



Optimization

- Exploit the multiple replicas → (complementary) layouts
- E.g., different sorting orders, partitioning schemes, compression schemes, etc



Plan Scheduling

Challenge

- Problem: missing/invalid tuples (explicit/implicit)
- Goal: reliable query results by error correction / error-tolerant algorithms



© Prof. Dr.-Ing. Wolfgang Lehner | UTECHNISCHE

- Assumption of valid raw data (single point of truth)
- Corrupt raw data are propagated through the model

Alternative

- Based on a single query (\rightarrow three different QEPs)
- Majority Gate on a single tuple basis
- Different OEPs based on
 - different operators and
 - in particular on redundant data sources
- Properties
 - no single point of truth
 - data replication is important

Redundant Query Processing (Triple Modular Redundancy)

- Based on a single query (\rightarrow three times execution of an optimal QEP)
- Majority Gate on a single tuple basis

> Resilient Query Processing (2)





Gate





Some Techniques

- Compression
- Indexing
- Resilience
- (Hardware) Transactional Memory



Concurrent Execution

- No serialization, no communication if no data conflicts
- Data conflicts are atomatically termined by the underlying TM infrastructure



Example: Hash Table







Main Memory Centric Data Management









let us map the situation of data analytics to ...

Phillip G. Armour

The Five Orders of Ignorance

Viewing software development as knowledge acquisition and ignorance reduction.



Oth Order Ignorance (OOI)—Lack of Ignorance

- I know something and can demonstrate my lack of ignorance (00I is knowledge)
- When I have 00I, I have the answer to the problem.



- Data Consumption Side
 - no need to look at the data
- Data Cube / Data Space Selection
 - no need to identify, tap, and combine potentially interesting data cubes
- Data Provisioning
 - don't bother about data analytics





1st Order Ignorance (10I)— Lack of Knowledge

- I don't know something and can readily identify that fact (10I is basic ignorance)
- When I have 10I, I have the question.
- Usually, having a good question makes it fairly easy to find the answer

Mapping to data analytics

- Data Consumption Side
 - Reporting!
 you know the dimension values

you are looking specifically for a measure of the specific dimension values

- Data Cube / Data Space Selection
 - no need to actively search for possibly interesting data cubes place where to search your information is well known
- Data Provisioning
 - data sources are well-known and understood; integration/transformation steps exist





2nd Order Ignorance (201)— Lack of Awareness

- I don't know that I don't know something (not only am I ignorant of something, I am unaware of this fact. I don't know enough to know that I don't know enough.)
- Not only do I not have the answer I need, I don't even have the question.

Mapping to data analytics

- Data Consumption Side
 - Guided Navigation / Explorer-Style
- Data Cube Selection
 - ??? how can I find the ,right' data
- Data Provisioning
 - I don't know my sources, but I suspect that there is helpful data out there...





3rd Order Ignorance (301) — Lack of Process

- I have 3OI when I don't know a suitably efficient way to find out I don't know that I don't know something → lack of process
- If I have 3OI, I don't know of a way to find out there are things I don't know that I don't know

Mapping to data analytics

- Data Consumption Side
 - Data mining ("Tell me something interesting!")
 - Data visualization
- Data Cube Selection
 - I don't even know how to design my BI entities
- Data Provisioning
 - I don't know anything about my data sources



Main Memory Centric Data Management

What are the main benefits?

Help business people with an *xth Order Ignorance*

> Situation in a Typical Organization

Database Technology



Data is Everywhere !!!

Corporate EDW(s)

- IT governed infrastructure
- well defined with standardized processes, taxonomies, and KPIs
- Top-down control for the enterprise-wide business data

Dozens of data marts, 100s of local databases, 1000s of spreadsheets

- locally defined analytics
- Bottom-up generation and ad-hoc access pattern





Data Marts and Shadow' Databases ~90% of data

EDW ~10% of data



Self-Service/ Agile BI

Magnetic

- "Attract data and practitioners"
- Usage of all data source independet of their data quality

Agile

- "Rapid iteration: ingest, analyze, productionalize"
- Continous evolution of the logical and physical structures
- ELT (Extraction, Loading, Transformation)

Deep

- "Sophisticated analytics in Big Data"
- Extended algorithmic run-time
- Ad-hoc advanced analytics and statistics

92 up.

To be able to do/perform amazing/unexpected things

I gots me mad skills, yo.

To be said after performing an extraordinairy feat.













Requirement 1: Balance Performance and Data Volume

- High query performance for interactive analysis / data exploration
- Challenge: Huge number of parallel users and ad-hoc query/data schemes
- Batch processing for offline tasks (hypothesis generation / consistency checks following complex business semantics)
- Challenge: Need to massage massive data volumes with complex business logic

Requirement 2: Support Pull and Push Processing

- Full spectrum of load granularity (batch trickle feed stream processing)
- **Challenge:** Need to integrate data streaming systems into DWH infrastructure

Requirement 3: Corporate Memory

- "Expect the unexpected"
- Challenge: Manage the organization's archives and individuals' memories
- Exist in hybrid environments
- Challenge: Provide cloud-enabled DWH environments

> Situation in a typical Organization





Data is Everywhere !!!

Corporate EDW(s)

- IT governed infrastructure
- well defined with standardized processes and measures
- Top-down control for the enterprise-wide business data

Dozens of data marts, 100s of local databases, 1000s of spreadsheets

- Locally defined analytics
- Bottom-up generation and ad-hoc access pattern





Data Marts and Shadow' Databases ~90% of data

EDW ~10% of data

Extend DWH Infrastructure



Approach: Enable business analysts to

- For business expert users, provide infrastructure to create agile data spaces
- For experienced users, easily run **ad-hoc queries** on complex data warehouses

Challenges

- Find the right input language for business people (simplified natural language)
- Work on very large and complex database schemas (hundreds of tables, inheritance patterns)
- Incrementally improve meta-model by experience gained from previous ad-hoc queries
- Balance well-structured and IT-governed data flow with ad-hoc-analytics requirements
- Reduce time to consumption no time left for tuning mechanisms



"Drill Beyond": Extending OLAP using Open Data

■ Open World" approach: Include all external open data sets → open schema LINEITEM



> The Big Wedding ahead!!!







data crunching meets number crunching



> The NetFlix Competition





> The NetFlix Competition (2)

444 4 CA CA

Netflix' star rating system helps determine personalized movie

those recommendations.

recommendations. Now the company is looking to outside developers to improve

1

61

BUSINESS

The \$1 Million Netflix Challenge

FRIDAY, OCTOBER 6, 2006 BY KATE GREENE

VP Jim Bennett discusses how recommendation systems suggest your next movie and the challenges of building a better one.

🖾 E-mail 🜿 Audio » 🖹 Print

Earlier this week, Netflix, the online movie rental service, announced it will award \$1 million to anyone who can come up with an algorithm that improves the accuracy of its movie recommendation service.

In doing so, the company is putting out a call to researchers who specialize in mac learning--the type of artificial intelligence used to build systems that recommend m books, and movies. The entrant who can increase the accuracy of the Netflix recommendation system, which is called Cinematch, by 10 percent by 2011 will w prize.

Recommendation systems such as those used by Netflix, Amazon, and other Well retailers are based on the principle that if two people enjoy the same product, they likely to have other favorites in common too.

But behind this simple premise is a complex algorithm that incorporates millions of ratings, tens of thousands of items, and ever-changing relationships between user preferences.



Main-Memory Centric Data Management





> The NetFlix Competition (3)













Database Technology

Group

> The NetFlix Competition (6)



> A simple experiment ...







color code := user rating



© Prof. Dr.-Ing. Wolfgang Lehner | UNIVERSITAT



Phase 1: drop 75% of all pixels





Phase 2: Random permutation of rows and columns



> The Experiment ...

Phase 3: Determine the latent factors



© Prof. Dr.-Ing. Wolfgang Lehner |



Database Technology Group



Phase 4: Reconstruction



© Prof. Dr.-Ing. Wolfgang Lehner | UTECHNISCHE



Phase 5: Final Result Generation



> The Experiment ...




> Summary and Conclusion

Database Technology

The quest for knowledge used to begin with grand theories. Now it begins with massive amounts of data.

Welcome to the Petabyte Age.

Main Memory Centric Data Management

- requires a holistic picture from system level to end-user experience
- data management layer is still the name of the game
- novel applications will combine the complexity of number-crunching and data-crunching

Huge impact on

- design of analytics infrastructures
- design of database systems
- implementation of analytical applications

to finally conclude: Recap

- 5 Orders of Ignorance
 - 00I Lack of Ignorance
 - 10I Lack of Knowledge
 - 20I Lack of Awareness
 - 30I Lack of Process
 - **4**01?

4th Order Ignorance (40I) — Meta Ignorance

I have 40I when I don't know about the Five Orders of Ignorance.

Main-Memory Centric Data Management

... we hopefully passed that stage!







Main-Memory Centric Data Management – Open Problems and Some Solutions Wolfgang Lehner

Technische Universität Dresden