# Dances with SAP HANA:
## Innovations Towards In-Memory Transactional and Analytical Big Data Processing

Juchang Lee (juc.lee@sap.com, SAP Labs Korea)
December 2016

# Agenda
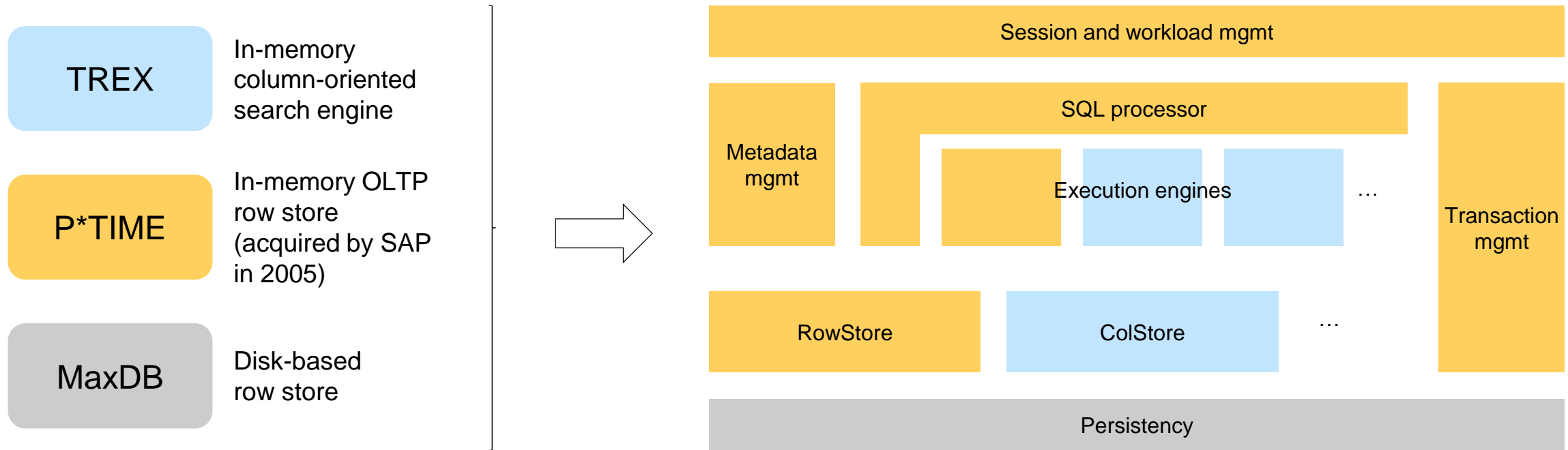
**Brief history of SAP HANA**

**Challenges and Innovations with SAP HANA**

- #1. *In-memory* database
- #2. In-memory *Big Data* processing
- #3. In-memory *Transactional and Analytical* Big Data processing
- #4. *HW+SW* co-innovation
- #5. Development and delivery *process* for high-performance DB engine

**Further opportunities and call for idea-to-market**

# Brief History of SAP HANA

**TREX** — In-memory column-oriented search engine

**P*TIME** — In-memory OLTP row store (acquired by SAP in 2005)

**MaxDB** — Disk-based row store

⇒

| Session and workload mgmt |
|---|

| Metadata mgmt | | SQL processor | | | | Transaction mgmt |
| | | Execution engines | | … | |

| RowStore | ColStore | … | |

| Persistency |

SAP HANA In-memory Database
for OLTP/OLAP workloads
(first product shipment in late Nov/2010)

# Acknowledgement

# #1. In-Memory Database Engine

**Paradigm shift of the DB engine optimization focus:**

### Cache-Conscious Concurrency Control of Main-Memory Indexes on Shared-Memory Multiprocessor Systems

Sang K. Cha    Sangyong Hwang    Kihong Kim    Keunjoo Kwon

School of Electrical Engineering and Computer Science
Seoul National University
{chask, syhwang, next, icdi}@kdb.snu.ac.kr

**Abstract**

Recent research addressed the importance of optimizing L2 cache utilization in the design of main memory indexes and proposed the so-called cache-conscious indexes such as the CSB+-tree. However, none of these indexes took account of concurrency control, which is crucial for running the real-world main memory database applications involving index updates and taking ... DBMS (DRDBMS) in many problem domains. MMDBMS can show orders-of-magnitude higher performance than DRDBMS not only for read transactions but also for update transactions. However, such a significant performance gain of MMDBMS over DRDBMS does not come automatically by just placing the database in memory but requires MMDBMS-specific optimization techniques. For example, the so-called differential logging scheme improves the update and recovery performance of MMDBMS significantly by

Minimizing L2 cache misses
on concurrent index accesses
VLDB 2001

### Differential Logging: A Commutative and Associative Logging Scheme for Highly Parallel Main Memory Database

Juchang Lee    Kihong Kim    Sang K. Cha

Graduate School of Electrical Engineering and Computer Science
Seoul National University
{juch, next, chask}@kdb.snu.ac.kr

**Abstract**

With a gigabyte of memory priced at less than $2,000, the main-memory DBMS (MMDBMS) is emerging as an economically viable alternative to the disk-resident DBMS (DRDBMS) in many problem domains. The MMDBMS can show significantly higher performance than the DRDBMS by reducing disk accesses to the sequential form of log writing and the occasional checkpointing. Upon the system ... processing the disk-resident data indirectly through the main-memory buffer. With a gigabyte of memory priced at less than $2,000 these days, the MMDBMS emerges as an economically viable alternative to the DRDBMS with the data structures and algorithms optimized for the in-memory access.

The MMDBMS can show higher performance than the DRDBMS not only for read transactions but also for update transactions by orders of magnitude. This is because the disk access in the MMDBMS is reduced to the sequential ...

Resolving the bottleneck at log I/O during
transaction commit and database restart
ICDE 2001

# #2. Big Data Processing with In-Memory Database Engine

**Approach 1: compressed column store and additional acceleration on table scan**



SIGMOD 2009



VLDB 2009

# #2. Big Data Processing with In-Memory Database Engine

## Approach 2: scaling out



Practical optimizations under shared-nothing
distributed database systems
ICDE 2013



Optimizations for delayed snapshot isolation
VLDB Journal 2014

# #2. Big Data Processing with In-Memory Database Engine

## Approach 3: offloading cold data to disks

### Page As You Go:
### Piecewise Columnar Access In SAP HANA

Reza Sherkat, Colin Florendo, Mihnea Andrei, Anil K. Goel, Anisoara Nica, Peter Bumbulis
Ivan Schreter, Günter Radestock, Christian Bensberg, Daniel Booss, Heiko Gerwens
SAP SE
<*firstname.lastname*>@sap.com

## ABSTRACT

In-memory columnar databases such as SAP HANA achieve extreme performance by means of vector processing over logical units of main memory resident columns. The core in-memory algorithms can be challenged when the working set of an application does not fit into main memory. To deal with memory pressure, most in-memory columnar databases evict candidate columns (or tables) using a set of heuristics gleaned from recent workload. As an alternative approach, we propose to reduce the unit of load and eviction from column to a contiguous portion of the in-memory columnar representation, which we call a *page*. In this paper, we adapt the core algorithms to be able to operate with partially loaded columns while preserving the performance benefits of vector processing. Our approach has two key advan-

## 1. INTRODUCTION

Modern in-memory database systems achieve extreme performance for analytical workloads by taking advantage of recent trends in hardware technology, including affordable large dynamic random-access memory (DRAM) and the instruction level parallelism of register vector processing, e.g. Advanced Vector Extensions (AVX) processor capabilities [7, 23, 27]. For in-memory columnar databases, dictionary based compression schemes produce uniform representation and in-memory layout for each column. The memory layout can be optimized for extremely fast basic database operations, such as scan and search [8, 17, 26].

With the rapid growth in the variety of applications that want to combine business data with the emerging Internet of Things (IoT) and/or social media data, there is an important

SIGMOD 2016

# #3. OLTP and OLAP together in a single Database Engine

## Hybrid Garbage Collection for Multi-Version Concurrency Control in SAP HANA

Juchang Lee #
juc.lee@sap.com

Hyungyu Shin ‡
hgshin@dblab.postech.ac.kr

Chang Gyoo Park #
chang.gyoo.park@sap.com

Seongyun Ko ‡
syko@dblab.postech.ac.kr

Jaeyun Noh #
jaeyun.noh@sap.com

Yongjae Chuh #
yongjae.chuh@sap.com

Wolfgang Stephan †
wolfgang.stephan@sap.com

Wook-Shin Han ‡*
wshan@dblab.postech.ac.kr

#SAP Labs Korea, Korea
‡Pohang University of Science and Technology (POSTECH), Korea
†SAP SE, Germany

## ABSTRACT

While multi-version concurrency control (MVCC) supports fast and robust performance in in-memory, relational databases, it has the potential problem of a growing number of versions over time due to obsolete versions. Although a few TB of main memory is available for enterprise machines, the memory resource should be used carefully for economic and practical reasons. Thus, in order to maintain the necessary number of versions in MVCC, versions which will no longer be used need to be deleted. This process is called garbage collection. MVCC uses the concept of *visibility* to define garbage. A set of versions for each record is first identified as candidate if their version timestamps are lower than the minimum value of snapshot timestamps of active snapshots in the system. All such candidates, except the one which group garbage collection, table garbage collection, and interval garbage collection. Through experiments using mixed OLTP and OLAP workloads, we show that HYBRIDGC effectively and efficiently collects garbage versions with negligible overhead.

## Keywords

Garbage collection; multi-version concurrency control; SAP HANA

## 1. INTRODUCTION

Commercial database management systems (DBMSs) including SAP HANA employ multi-version concurrency control (MVCC) due to robust and better performance for var

SIGMOD 2016

# Problem: Garbage Collection under OLTP/OLAP-mixed Workloads

OLTP
workload

OLAP
workload

**OLTP** generates versions
with high speed from highly
concurrent write
transactions

**OLAP** generates long-
running queries holding old
read timestamps which
blocks the garbage collector
(e.g. >1hr of cursor life time
in a commercial system)

a single
database
system

**Consequence**
- Version space grows – more critical especially in in-memory DBMS
- Version lookup cost increases, i.e. performance drops

Figure 9: Hybrid garbage collection.



Figure 10: The number of record versions with a long-duration cursor.



Figure 12: TPC-C throughput with a long-duration cursor.

# #4. HW+SW Co-Innovation

## FPTree: A Hybrid SCM-DRAM Persistent and Concurrent B-Tree for Storage Class Memory

Ismail Oukid[†‡]    Johan Lasperas[‡]    Anisoara Nica[‡]

Thomas Willhalm[*]    Wolfgang Lehner[†]

[†]TU Dresden    [‡]SAP SE    [*]Intel Deutschland GmbH
first.last@tu-dresden.de    first.last@sap.com    first.last@intel.com

### ABSTRACT

The advent of Storage Class Memory (SCM) is driving a rethink of storage systems towards a single-level architecture where memory and storage are merged. In this context, several works have investigated how to design persistent trees in SCM as a fundamental building block for these novel systems. However, these trees are significantly slower than DRAM-based counterparts since trees are latency-sensitive and SCM exhibits higher latencies than DRAM. In this paper we propose a novel hybrid SCM-DRAM persistent and concurrent $B^+$-Tree, named *Fingerprinting Persistent Tree* (FPTree) that achieves similar performance to DRAM-based counterparts. In this novel design, leaf nodes are persisted in SCM while inner nodes are placed in DRAM and rebuilt upon recovery. The FPTree uses *Fingerprinting*, a technique that limits the expected number of in-leaf probed keys to one. In addition, we propose a hybrid concurrency scheme for the FPTree that is partially based on Hardware Transactional Memory. We conduct a thorough performance evaluation and show that the FPTree outperforms state-of-the-art persistent trees with different SCM latencies by up to a factor of 8.2. Moreover, we show that the FPTree scales very well on a machine with 88 logical cores. Finally, we integrate the evaluated trees in *memcached* and ... We show that the FPTree ...
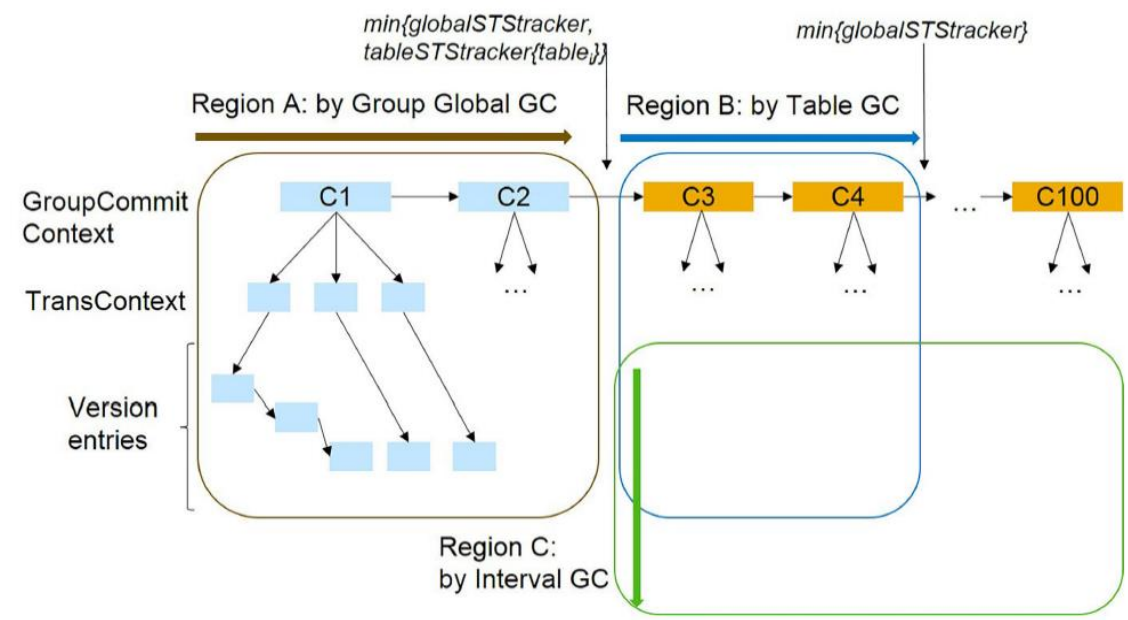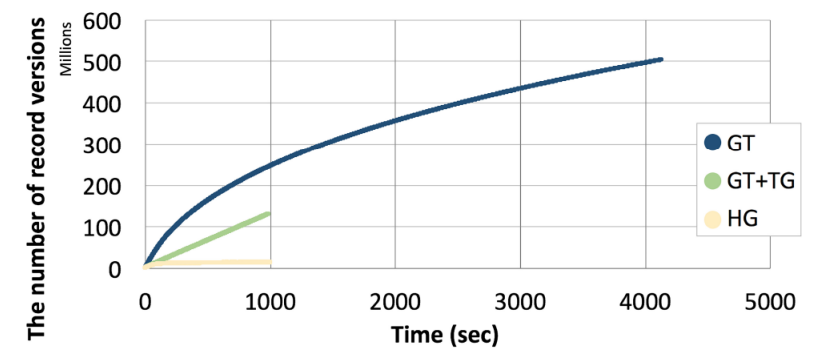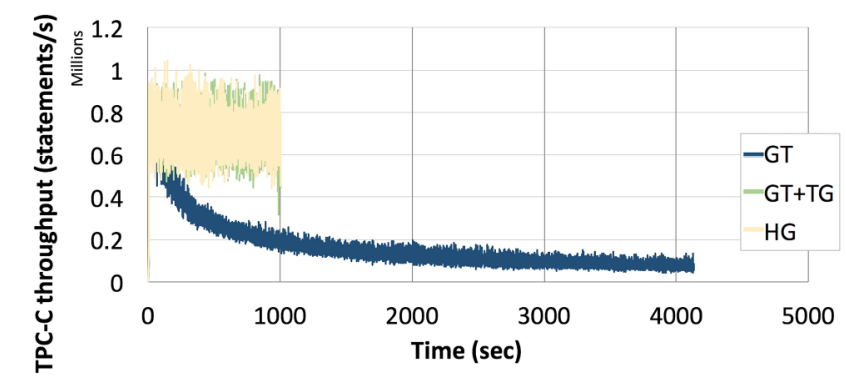
performance. This potential has been acknowledged by recent works and new systems that adopt this novel paradigm are emerging [5, 21, 18]. In this paper we investigate the design of index structures as one of the core database structures, motivated by the observation that traditional main memory B-Tree implementations do not fulfill the consistency requirements needed for such a use case. Furthermore, while expected in the range of DRAM, SCM latencies are slower and asymmetric with writes noticeably slower than reads. We argue that these performance differences between SCM and DRAM imply that the assumptions made for previous well-established main memory B-Tree implementations might not hold anymore. We therefore see the need to design a novel, persistent B-Tree that leverages the capabilities of SCM while exhibiting performance similar to that of traditional transient B-Tree.

Designing persistent data structures presents unprecedented challenges for data consistency since SCM is accessed via the volatile CPU cache over which software has only little control. Several works proposed SCM-optimized B-Trees such as the CDDS B-Tree [24], the wBTree [8], and the NV-Tree [28], but they fail to match the speed of an optimized DRAM-based B-Tree. Additionally, they do not address all SCM programming challenges we identify in Section 2, especially those of persistent memory leaks and data recovery.

SIGMOD 2016



Samsung-SAP Research Center
(since 29/Sep/2016)

# #5. Development and Delivery Process
# for High-Performance DB Engine

## A Performance Anomaly Detection and Analysis Framework for DBMS Development

Donghun Lee, Sang K. Cha, *Member, IEEE*, and Arthur H. Lee

**Abstract**—Detecting performance anomalies and finding their root causes are tedious tasks requiring much manual work. Functionality enhancements in DBMS development as in most software development often introduce performance problems in addition to bugs. To detect the problems as soon as they are introduced, which often happens during the early phases of a development cycle, we adopt performance regression testing early in the process. In this paper, we describe a framework that we developed to manage performance anomalies after establishing a set of conditions for a problem to be considered an anomaly. The framework uses Statistical Process Control (SPC) charts to detect performance anomalies and differential profiling to identify their root causes. By automating the tasks within the framework we were able to remove most of the manual overhead in detecting anomalies and reduce the analysis time for identifying the root causes by about 90 percent in most cases. The tools developed and deployed based on the framework allow us continuous, automated daily monitoring of performance in addition to the usual functionality monitoring in our DBMS development.

**Index Terms**—CUSUM chart, differential profiling, statistical process control (SPC), performance anomaly, DBMS.

✦

## 1 INTRODUCTION

PERFORMANCE concerns in a software product are usually serious enough to cause delays in deployment and sometimes reported as one of the most critical problems of deployed systems as documented in [10]. As we have been developing the DBMS engine P*TIME [28], one of the most critical challenges has been fighting occasional performance degradations. With the advent in hardware technology and the anomalies and finding their root causes efficiently with as little of disruption to the development process as possible. What is worse, the development continues thus possibly affecting some of the performance metrics continuously often in some unexpected places and directions. Furthermore, performance values show inevitable variations affected by various factors such as workload by

**IEEE TKDE 2012**

## Performance monitoring in SAP HANA's continuous integration process

Kim-Thomas Rehmann
SAP SE
Walldorf, Germany
kim-thomas.rehmann@sap.com

Changyun Seo
SAP SE
Seoul, South Korea
changyun.seo@sap.com

Dongwon Hwang
SAP SE
Seoul, South Korea
dong.won.hwang@sap.com

Binh Than Truong
SAP SE
Ho Chi Minh City, Viet Nam
binh.thanh.truong@sap.com

Alexander Boehm
SAP SE
Walldorf, Germany
alexander.boehm@sap.com

Dong Hun Lee
SAP SE
Seoul, South Korea
dong.hun.lee@sap.com

### ABSTRACT

Development principles such as continuous integration [2] and continuous delivery [5] become increasingly popular in the software industry. They allow for the quick and auto-mated build, test, and delivery of software, thereby signif-icantly improving the overall quality assurance and release processes.

In this paper, we show how to apply the ideas of continu-ous delivery to complex system software, as exemplified by the SAP HANA database platform. We discuss the inte-gration of performance testing early in the delivery process and the construction of micro-services to detect and report performance anomalies in a continuous integration process.

the pre-commit checks.

In contrast to functional tests, performance regression benchmarks are much more complex to handle in a CI pipe-line, because they are susceptible to influences of the test en-vironment and produce large amounts of test results, which have to be analyzed in an efficient manner.

Over the last years, we devised a method to monitor data-base performance early and at various stages of the CI pro-cess. In the following sections, we discuss how

• to integrate performance tests into the pre-commit tests of a CI pipeline, introducing the concept of pre-commit *performance barriers*

• to deal with long-running benchmarks and user sce-

**SIGMETRICS 2016**

# Further Opportunities and Call for Idea-to-Market

**Industry-University co-research**

**Open job positions**

- DB kernel developers (transaction, replication, session, metadata, in-memory store):
  - https://jobs.sap.com/job/Seoul-C%2B%2B-%28Senior%29-Developer-for-SAP-HANA-Database-Kernel-Team-Job-11/321176501/
- SQL engine developers (SQL/SQLScript executor, optimizer):
  - https://jobs.sap.com/job/Bangalore-C%2B%2B-%28Senior%29-Developer-for-SAP-HANA-SQL-engine-team-Job-11/321330101/
- Quality specialist, development infrastructure, development support
  - https://jobs.sap.com/job/Seoul-Quality-Specialist-Job-11/320983201/
  - https://jobs.sap.com/job/Seoul-SAP-HANA-Platform-Dev-Support-Job-11/321336901/