# SIGMOD 2020

TECHNISCHE
UNIVERSITÄT
DARMSTADT

THE BEST RUN SAP

## ROBUST PERFORMANCE OF MAIN MEMORY DATA STRUCTURES BY CONFIGURATION

**Norman May**

SAP SE

**Ismail Oukid**

(SAP SE)

**Ilia Petrov**

Reutlingen
University

**Carsten Binnig**

TU Darmstadt

**Tiemo Bang**

TU Darmstadt
& SAP SE

# PLETHORA OF MODERN OLTP ARCHITECTURES



Many OLTP architectures proposed following evolution of modern hardware

# CANDIDATE 1: H-STORE

## Kallman et al.: Fine-grained shared-nothing architecture[1]

Gen-Z Sandbox: 2000 cores, 200 TB memory[3]



**Partitioning per core**

**Scalable OLTP architecture:** Applies to broad hardware

**Depends on partitioning:** Sensitive to skew!

[1] Kallman et al. 2008. H-Store: a High-Performance, Distributed Main Memory Transaction Processing System. Proc. VLDB Endow.
[2] Pavlo. 2011. Magical Parallel OLTP Databases. https://hstore.cs.brown.edu/slides/hstore-houdini-nov2011.pdf
[3] Burts. 2018. HPE BOOTS UP SANDBOX OF THE MACHINE FOR EARLY USERS. https://www.nextplatform.com/2018/06/21/hpe-boots-up-sandbox-of-the-machine-for-early-users/

# CANDIDATE 2: HEKATON

**Diaconu et al.: Efficient shared-everything architecture[4]**

Joint operation by all resources on all data:

+ Non-partitionable workloads

+ Fluctuating workloads

- Physical contention

- NUMA effects

**E.g.: TPC-C on many-core hardware**

Few warehouses

[4] Diaconu et al. .2013. Hekaton: SQL Server's Memory-Optimized OLTP Engine. SIGMOD.

TECHNISCHE
UNIVERSITÄT
DARMSTADT

THE BEST RUN SAP

# DESIGN SPACE OF OLTP ARCHITECTURES



> ➢ **How to achieve robust performance for entire design space?**

[5] Drake et al. 2015. Microsoft's Open Cloud Server. http://download.microsoft.com/download/B/1/7/B179029E-7AE8-447A-B8C9-B823B3DFC727/Microsofts_Open_CloudServer_Strategy_Brief.pdf

# IDEA:
# CONFIGURE OLTP ARCHITECTURE

**One size does not fit entire design space!**

**Simply resize to fit point in design space!**

**Flexible resource partitions**

Partitionable on many-core

**Optimal instantiation**

Non-partitionable
on multi-socket

**= Optimal architecture
for any workload
on any hardware**

...

TECHNISCHE
UNIVERSITÄT
DARMSTADT

THE BEST RUN SAP

# EXAMPLE:
# CONFIGURATION FOR WORKLOADS



> ➤ **Simply optimise architecture for given workload!**

# CONFIGURATION OPPORTUNITIES

Configuration = flexible resource partitions + optimal instantiation

**Existing architectures**

**Mimic:**

e.g. NUMA-aware for read-only workload

**Fitting partitioning**

**Optimally sized:**

e.g. reduce contention for read-update workload

**For discrete components**

**Individually sized:**

e.g. H-Store for hot + NUMA-aware for cold indexes in single architecture

➢ **Configure optimal DBMS architecture without redesign!**

TECHNISCHE UNIVERSITÄT DARMSTADT  THE BEST RUN SAP

# CONFIGURATION APPROACH

## Configurable Virtual Domains



**Virtual** resource partition +
Contention and locality **Domain**

## Configuration procedure



Linear program
instantiates virtual domains
and assigns data structures

## Deploy on runtime

Establish efficient operation
between partitions…

## Initial Calibration



Capture general behaviour
of data structures

## Actual workload and hardware

E.g. hot and cold indexes
on 4-socket server

[6] Chiarandini. Linear and Integer Programming Lecture Notes. https://imada.sdu.dk/~marco/Teaching/AY2014-2015/DM554/Notes/dm554-main.pdf

# CONFIGURATION RUNTIME

Allow efficient execution across individually optimised partitions

Access partitions via Data-aware Task:
- **without interfering**
- **with maximal utilisation**
- **with minimal overhead**

**Delegation** + **async. execution**
of **Data-aware Tasks**



**In-memory Messaging**

➢ Robust performance by efficient execution across optimal partitions!

TECHNISCHE UNIVERSITÄT DARMSTADT — THE BEST RUN — SAP

# EVALUATION:
# SIMPLE WORKLOAD (YCSB)

Configuration of 2 Key-Value Stores for differing YCSB workloads

# EVALUATION:
# SIMPLE WORKLOAD (YCSB)

Throughput under increasing system size for differing workload



> **Robust across workloads, system sizes & data structures!**

# EVALUATION:
# COMPLEX WORKLOAD (TPC-C)

Configuration for non-partitionable TPC-C workload, 8 warehouses
(*details in the paper*)



➢ **Robust performance for complex workload by configuration!**

# ALL DETAILS…

## Robust Performance of Main Memory Data Structures by Configuration

Tiemo Bang*
TU Darmstadt & SAP SE

Ismail Oukid†
Snowflake Inc.

Norman May
SAP SE

Ilia Petrov
Reutlingen University

Carsten Binnig
TU Darmstadt

**Abstract**

In this paper, we present a new approach for achieving robust performance of data structures making it easier to reuse the same design for d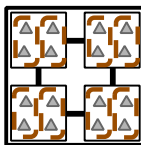ifferent hardware generations but also for different workloads. To achieve robust performance, the main idea is to strictly separate the data structure design from the actual strategies to execute access operations and adjust the actual execution strategies by means of so-called configurations instead of hard-wiring the execution strategy into the data structure. In our evaluation we demonstrate the benefits of this configuration approach for individual data structures as well as complex OLTP workloads.

## 1  Introduction

*Motivation:* Within the last decade, we have seen different hardware trends that significantly affected the design of single-node database systems: (1) Increases in main-memor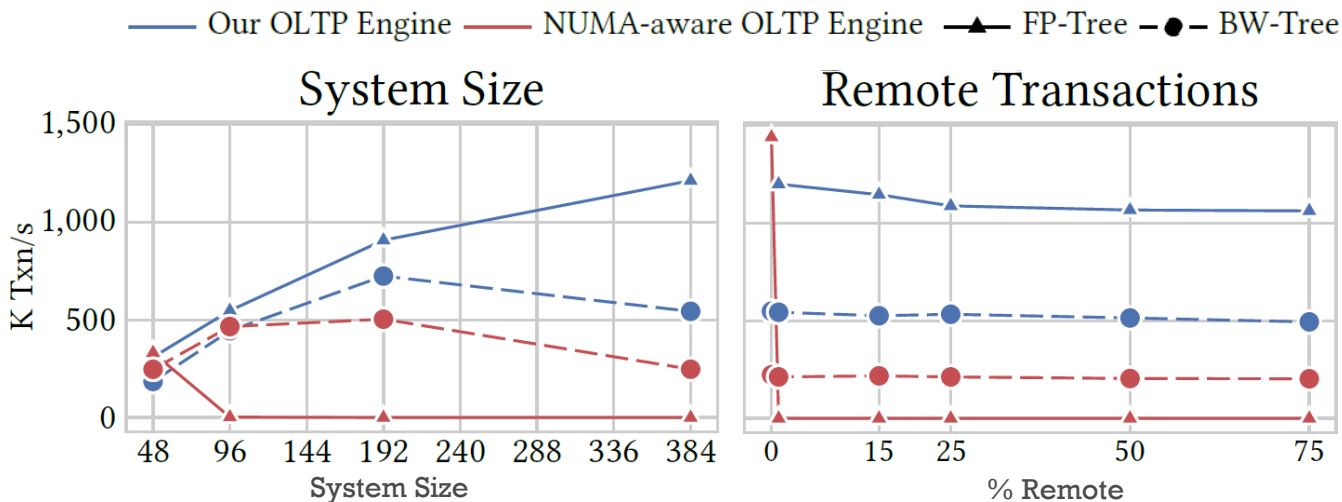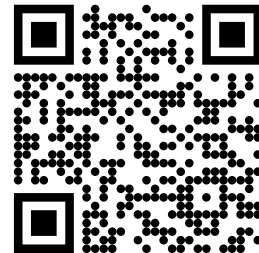y capacities made it possible to hold even larger data sets secondary storage (e.g., hard drives). (2) Moore's Law and Dennard Scaling required processor designers to move from single-socket and single-core designs to multi-socket and multi-core designs. As a result of these trends, we have seen a rapid evolution of hardware designs differing in essential characteristics not only memory capacities but also the underlying topology of how cores and memory are connected as well as cache sizes and coherence protocols.

A considerable body of existing work in DBMS research has thus focused on optimising the design of core DBMS data structures such as indexes for specific hardware configurations and workloads. For example, there have been various design alternatives proposed for classical B-trees to adapt them to modern memory hierarchies and make them more cache-conscious for read-heavy workloads [33, 34] or to optimise their behaviour for high-contention scenarios [25] under write-heavy workloads. A significant issue with this manual tailoring of core DBMS data structures is that not only their redesign involves high effort and reintegration into the DBMS but also that a design optimal for one hardware generation and one workload might induce severe performance degradation on another hardware generation when underlying assumptions change.
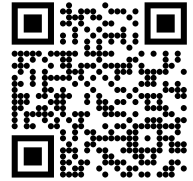
An alternative to this approach is designing data structures that can provide robust performance [18]. At its core, robust performance means the ability of a data structure to pro-

https://doi.org/10.1145/3318464.3389725

# CONCLUSION

> **Optimise DBMS architecture without redesign by configuration!**

*Stay healthy and see you next year in person!*

TECHNISCHE UNIVERSITÄT DARMSTADT

THE BEST RUN SAP