# Self-Organizing Data Containers
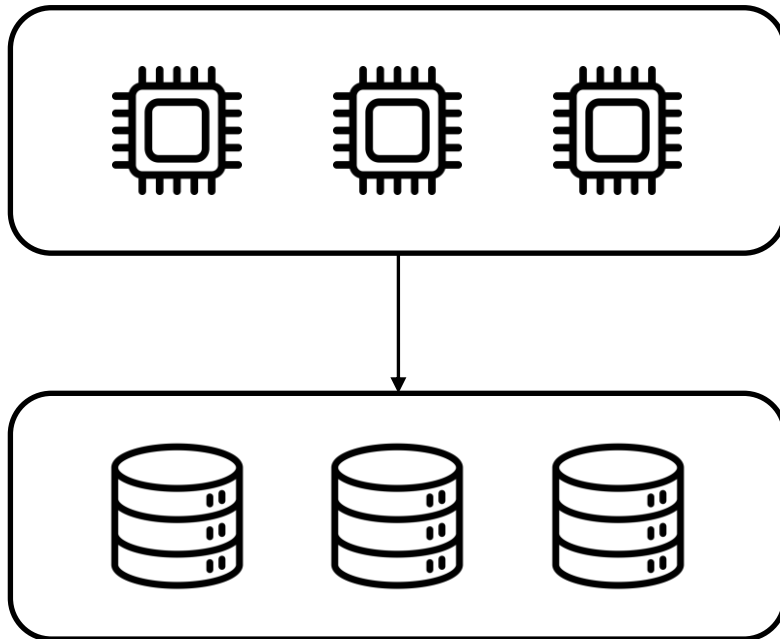
Thomas Glas

Technische Universität München

TUM School of Computation, Information and Technology

Lehrstuhl für Datenbanksystem
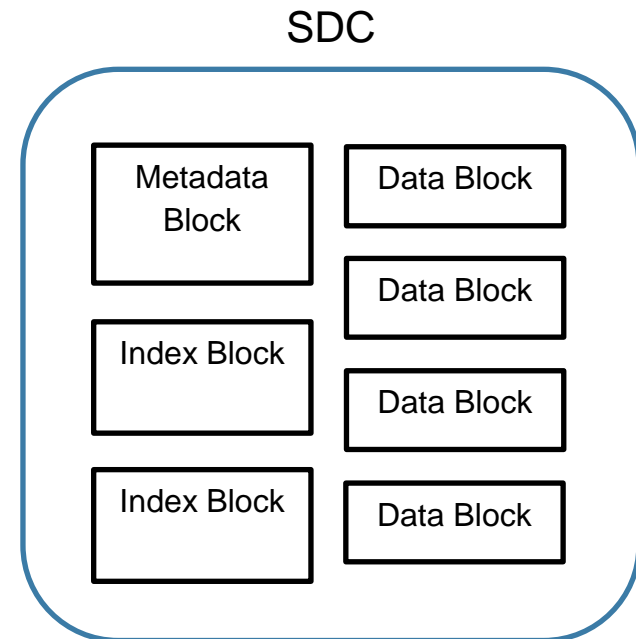
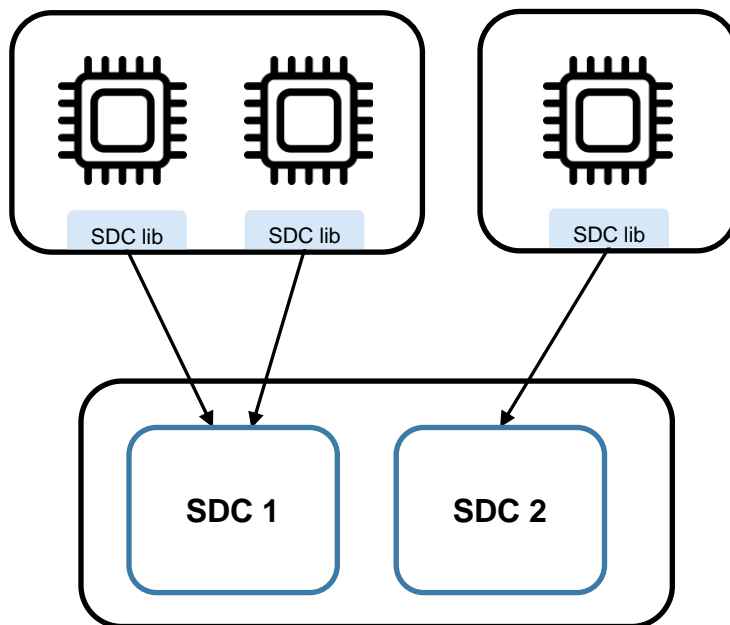München, 14. März 2023

# Motivation



How can we minimize data transfer?
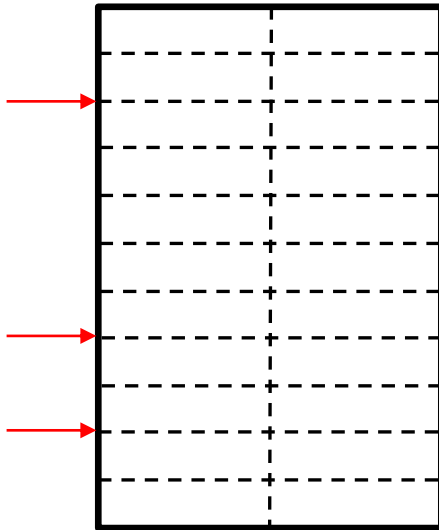➡️ only retrieve data needed for query

How can we create a storage layer with rich metadata to support this?

# Self-Organizing Data Containers

# Partitioning Strategies

Column Range Partitioning

Qd-Tree

# Self-Organizing Data Containers

## Client

1. Run workload on primary replica
2. Update workload statistics in SDC metadata
3. Compute new optimized layout for SDC
4. Create index and data blocks for new layout
5. Run workload on optimized layout

SDC lib

## SDC

### Metadata Block
- Table metadata
- Workload statistics
- Existing indexes

### Index Block
- Data blocks
- Data ranges contained in data block

### Data Block
- Parquet file

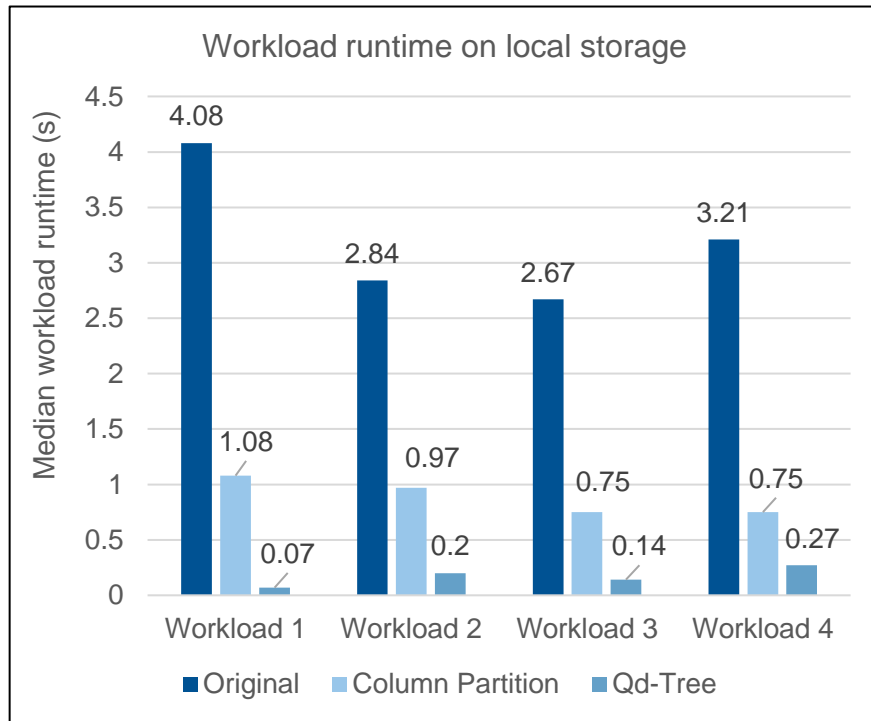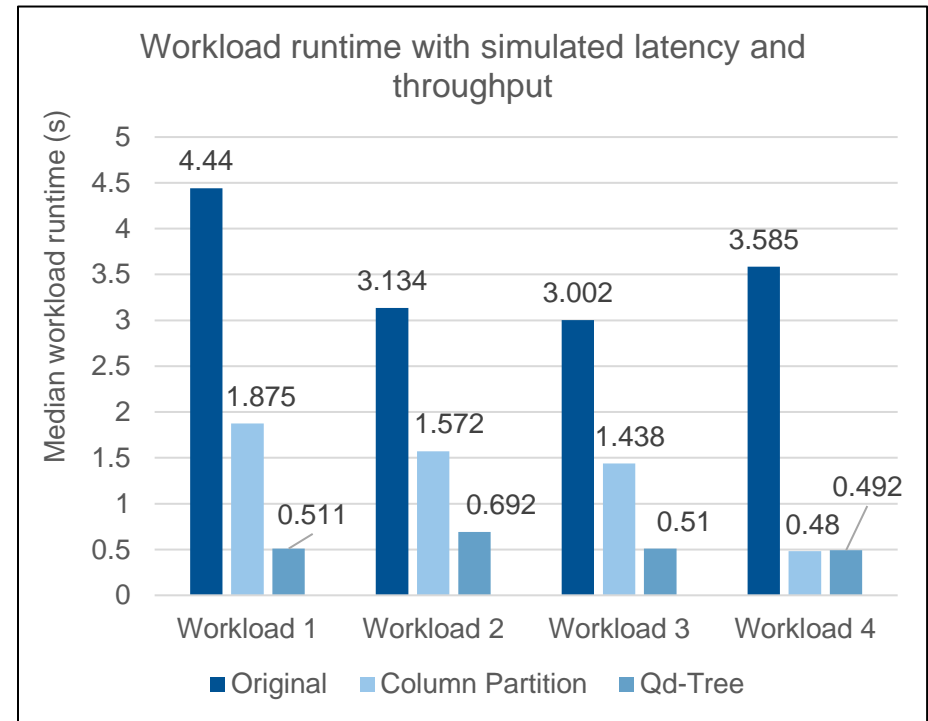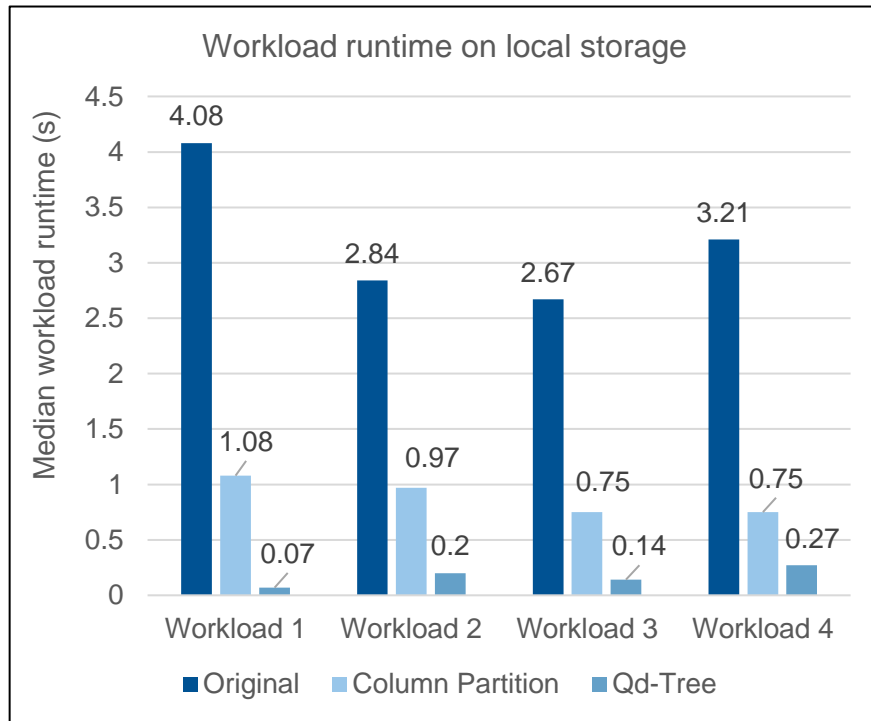# Demo of my SDC prototype

- SDC Library API: Projections and filters on tables

- Indexes: Primary, column range partitioning,

  Qd-Tree

- Dataset: NYC TLC Trip Record Data

- Workloads: Single-table range queries

- Storage layer: Local disk vs cloud storage (simulated)

- Data blocks: Apache Parquet files

- Metadata blocks: JSON files

- Clients: Single client

# Benchmarks



Workload runtime on local storage

Median workload runtime (s)

- Original
- Column Partition
- Qd-Tree

| | Workload 1 | Workload 2 | Workload 3 | Workload 4 |
|---|---|---|---|---|
| Original | 4.08 | 2.84 | 2.67 | 3.21 |
| Column Partition | 1.08 | 0.97 | 0.75 | 0.75 |
| Qd-Tree | 0.07 | 0.2 | 0.14 | 0.27 |



(a) SDC: data on disk

Workload runtime (s)

- Default
- Range partitioned
- Optimized

| Dataset | Default | Range partitioned | Optimized |
|---|---|---|---|
| contributions | 197 | 17 | 6.41 |
| flights | 19 | 12 | 4.22 |
| taxi | 238 | 117 | 8.23 |
| tweets | 20.2 | 0.98 | 0.98 |

# Benchmarks



Workload runtime on local storage

Median workload runtime (s)

| | Workload 1 | Workload 2 | Workload 3 | Workload 4 |
|---|---|---|---|---|
| Original | 4.08 | 2.84 | 2.67 | 3.21 |
| Column Partition | 1.08 | 0.97 | 0.75 | 0.75 |
| Qd-Tree | 0.07 | 0.2 | 0.14 | 0.27 |

Workload runtime with simulated latency and throughput

Median workload runtime (s)

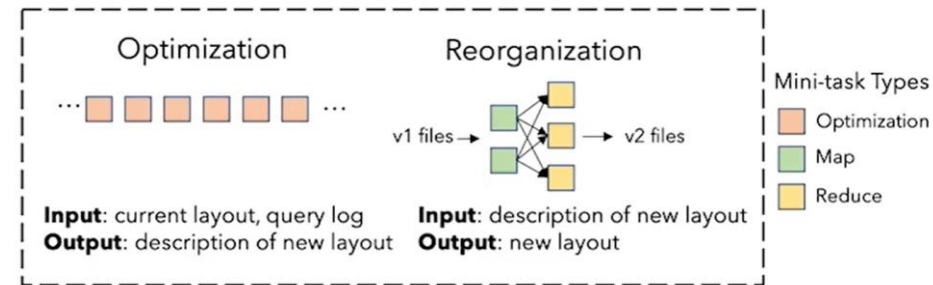| | Workload 1 | Workload 2 | Workload 3 | Workload 4 |
|---|---|---|---|---|
| Original | 4.44 | 3.134 | 3.002 | 3.585 |
| Column Partition | 1.875 | 1.572 | 1.438 | 0.48 |
| Qd-Tree | 0.511 | 0.692 | 0.51 | 0.492 |

# Conclusion & further research

- Simple yet effective self-learned storage optimization

- Easy to add indexes

- Easy to integrate into any type of applications (not just DBMS)

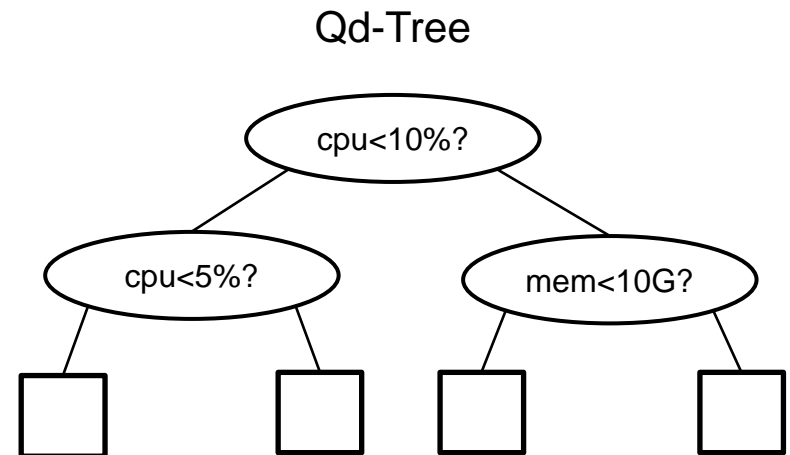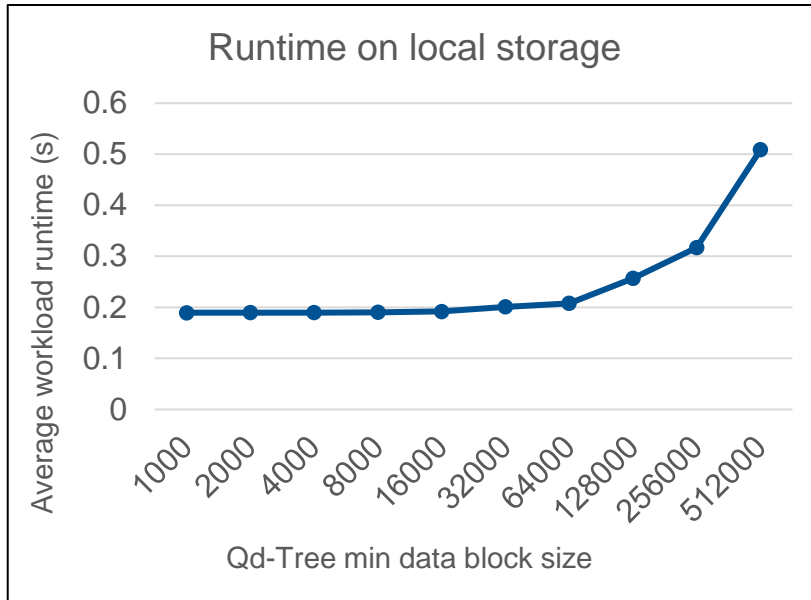- Use of data replication: trading off storage cost for query performance

Further research

- Distributing optimization work among clients

- Find query clusters in workload for effective

  indexes



**Code on Github:**

**https://github.com/thomasglas/SDCs**

# Benchmarking Qd-Tree min block sizes



Runtime on local storage

Qd-Tree

# Benchmarking Qd-Tree min block sizes