# Instant Recovery for Main-Memory Databases

Ismail Oukid*°, Wolfgang Lehner*, Thomas Kissinger*, Peter Bumbulis°, and Thomas Willhalm +

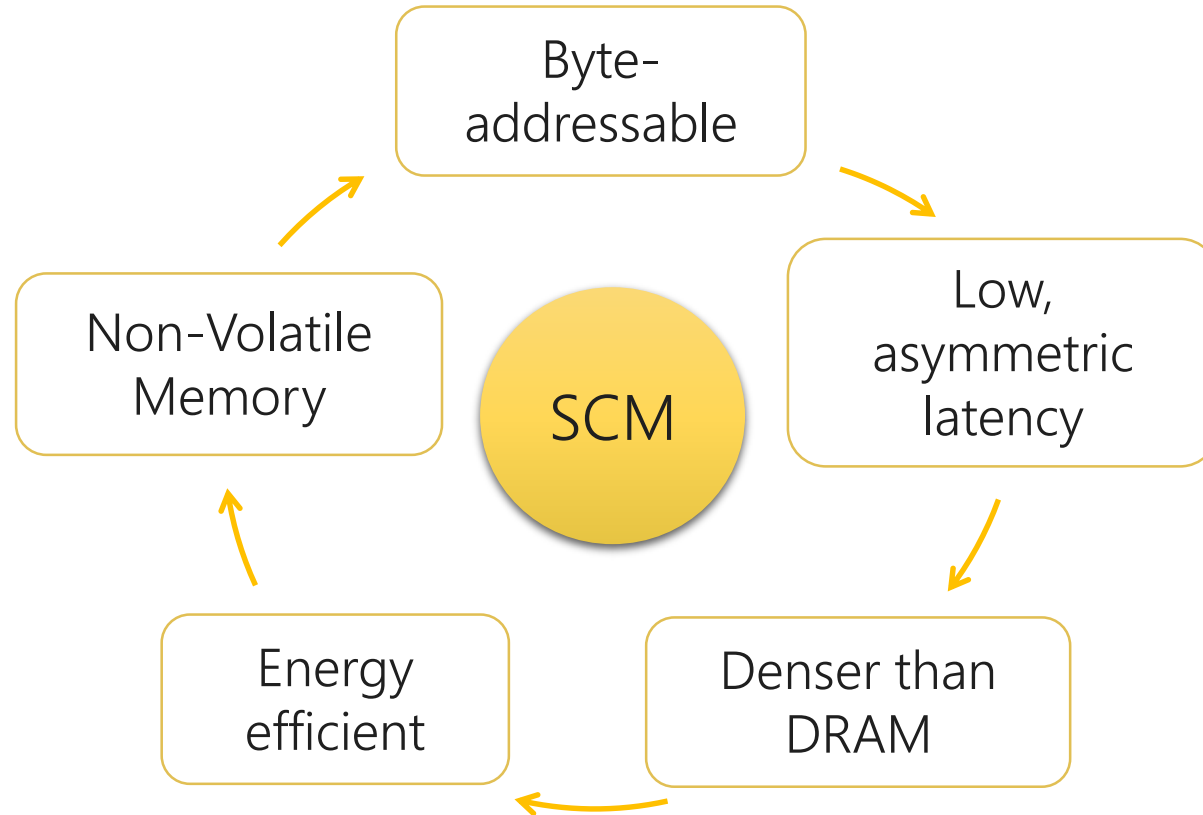*TU Dresden          °SAP SE          + Intel GmbH
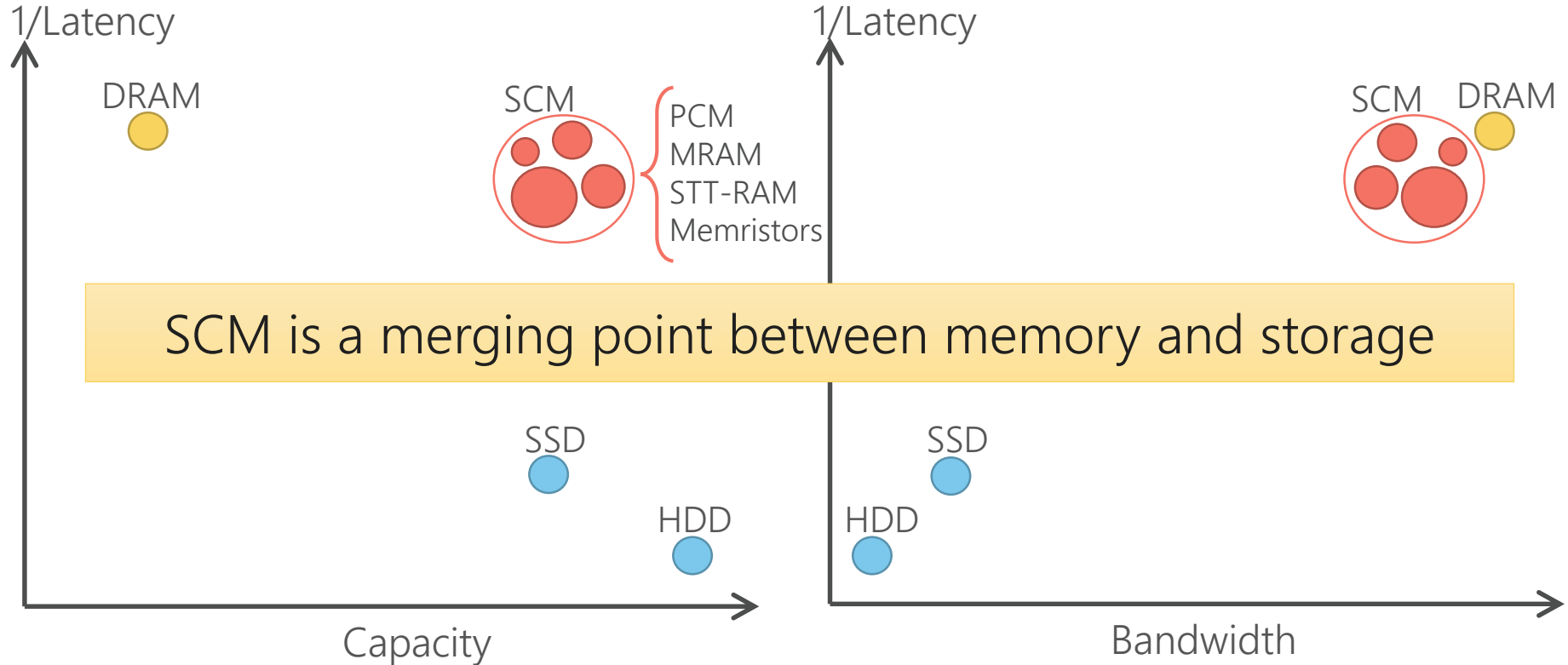
CIDR 2015, Asilomar, California, USA,
January 5, 2015

# Storage Class Memory

Byte-addressable

Non-Volatile Memory

SCM

Low, asymmetric latency

Energy efficient

Denser than DRAM

# SCM Compared with Today's Technologies



SCM is a merging point between memory and storage

Left chart: axes 1/Latency (vertical) and Capacity (horizontal). DRAM, SCM (PCM, MRAM, STT-RAM, Memristors), SSD, HDD.

Right chart: axes 1/Latency (vertical) and Bandwidth (horizontal). SCM, DRAM, SSD, HDD.

# SCM and Databases

## Improving the logging infrastructure, e.g.:

- Fang et al. High performance database logging using Storage Class Memory. ICDE'11
- Pelley et al. Storage management in the NVRAM era. VLDB'13
- Huang et al. NVRAM-aware Logging in Transaction Systems. VLDB'14

## Improving specific database algorithms, e.g.:

- Chen et al. Rethinking Database Algorithms for Phase Change Memory. CIDR'11
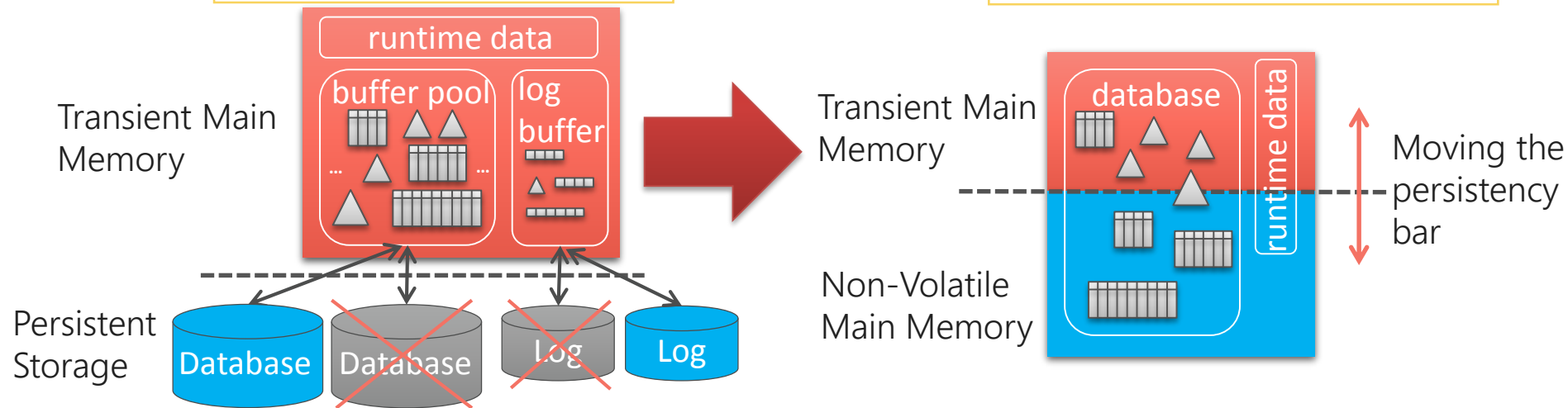- Stratis D. Viglas. Write-limited sorts and joins for persistent memory. VLDB'14

It takes a greenfield approach to measure the full potential of SCM
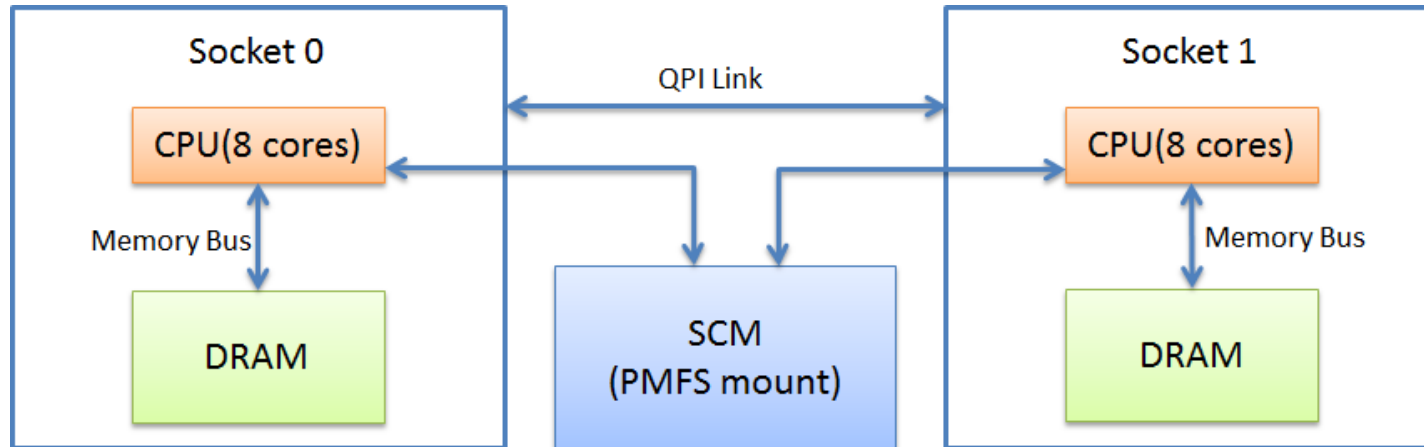
# SCM-enabled Architecture



SOFORT is a **single-level** column-store, i.e., the working copy **is** the durable copy

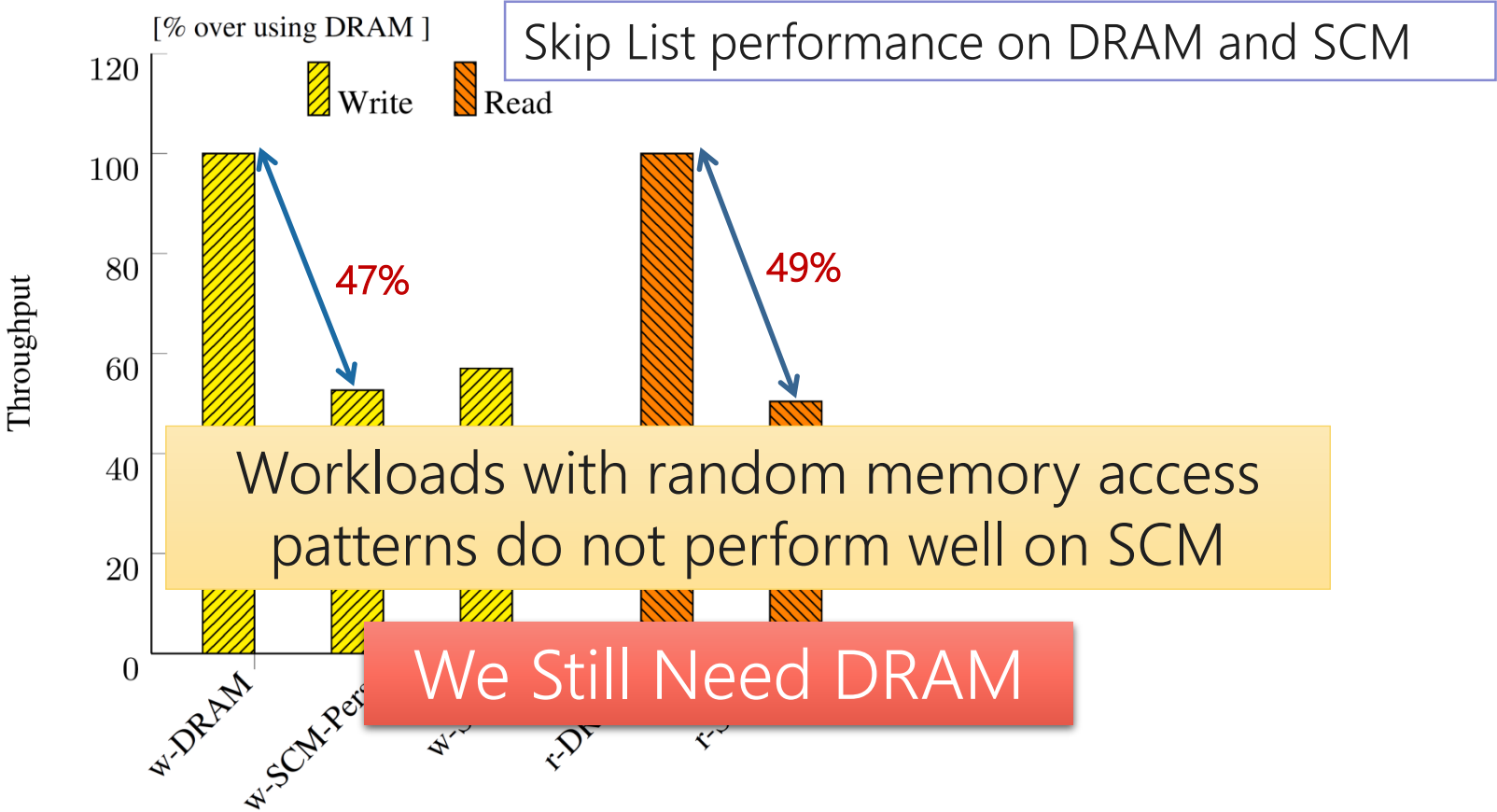# Understanding SCM through Microbenchmarks

**Hardware-based SCM simulation:**
- Special BIOS, tunable latency with means of a microcode patch
- Limitation: symmetric instead of asymmetric read/write latency
- Avoiding NUMA effects: benchmark run on a single socket
- DRAM Latency: 90ns    SCM latency: 200ns

SIMD-Scan performance on DRAM and SCM

[% over Scan-DRAM ]

Legend: DRAM, DRAM-NoPrefetch, SCM, SCM-NoPrefetch

Throughput (y-axis), Bit case (x-axis)

8% average slowdown

41% average slowdown

Workloads with sequential memory access patterns perform well on SCM

Skip List performance on DRAM and SCM

Workloads with random memory access patterns do not perform well on SCM

We Still Need DRAM

# SOFORT Column Structure



Persisted in SCM
Volatile in DRAM

On DRAM for better performance

SOFORT

Tx array

Tables

Column

(0, Asilomar)

(1, Dresden)   (2, Heidelberg)

Dict. Index

| 0 | Asilomar |
| 1 | Dresden |
| 2 | Heidelberg |

Dict. Array

| 2 |
| 0 |
| 1 |
| 2 |

Value IDs
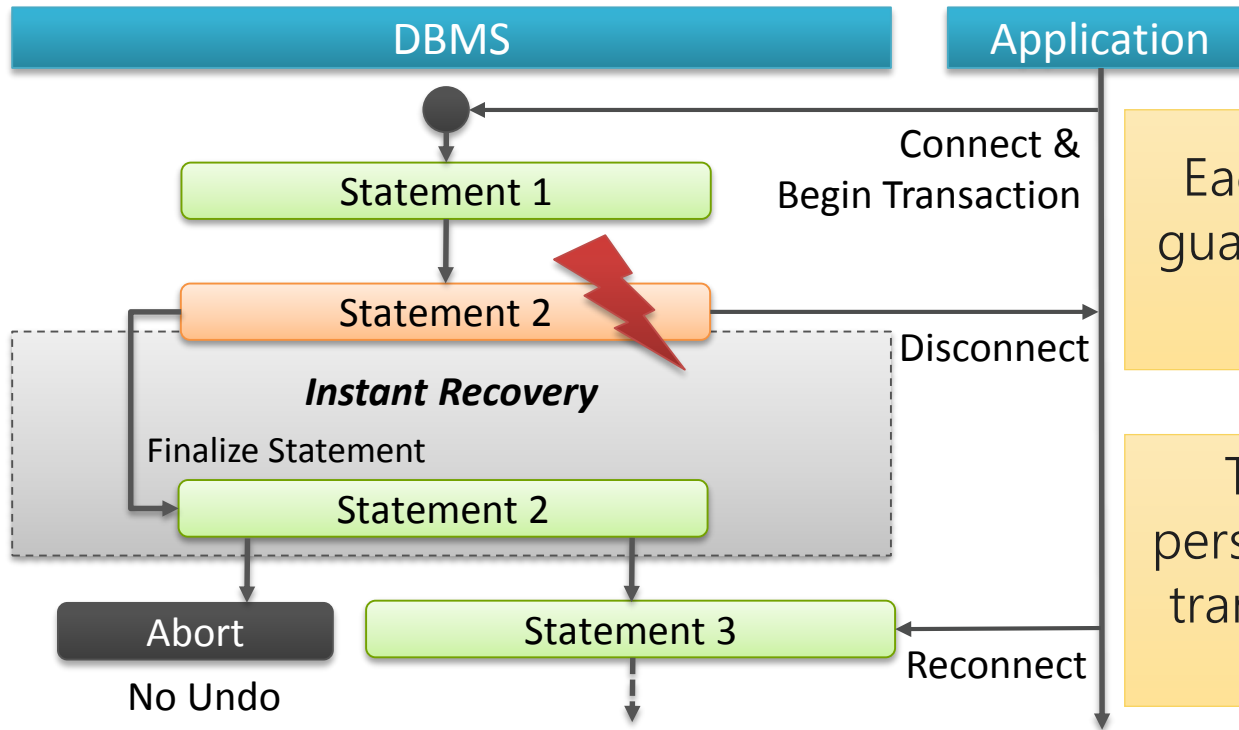
Persistent to enable continuing
unfinished transactions

Implementation details in "SOFORT: A Hybrid SCM-DRAM Storage Engine for Fast Data Recovery", DaMoN'14

# Continuing Unfinished Transactions



| DBMS | Application |
|---|---|

Connect & Begin Transaction

Statement 1

Statement 2

Disconnect

**Instant Recovery**

Finalize Statement

Statement 2

Abort
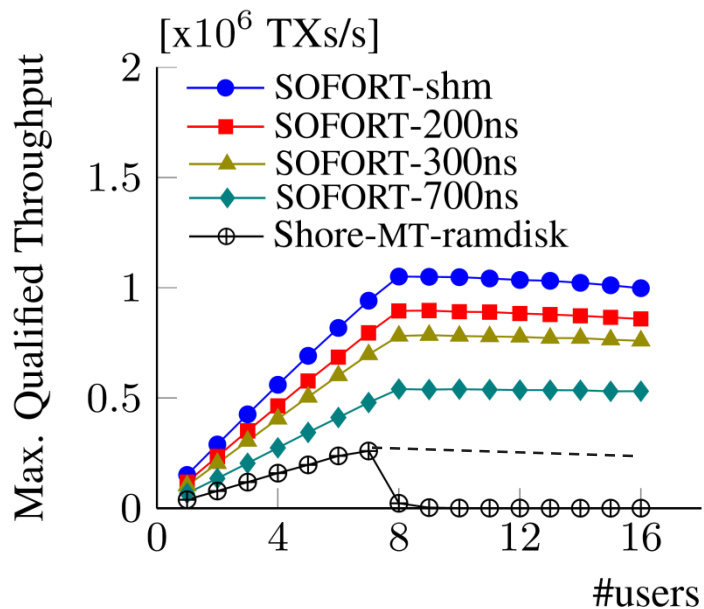
No Undo

Statement 3

Reconnect

Each executed statement is guaranteed to have persisted its changes in SCM.

The Transaction array is persistent allowing unfinished transactions at crash time to continue.
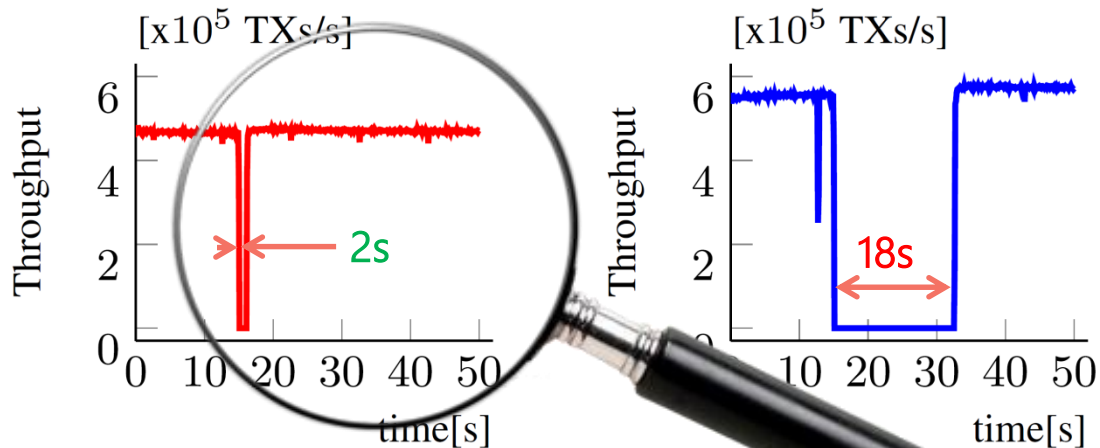
# Performance Overview



THROUGHPUT

[x10^6 TXs/s]

- SOFORT-shm
- SOFORT-200ns
- SOFORT-300ns
- SOFORT-700ns
- Shore-MT-ramdisk

Max. Qualified Throughput

#users

TATP Mix

RESTART TIME

[x10^5 TXs/s]

Throughput

2s

time[s]

(a) SOFORT-PMFS-200ns

[x10^5 TXs/s]

Throughput

18s

time[s]

(b) SO...isk

Competitive performance even in high latency environment

Fast restart time. No need to fetch data stored in SCM

Still not instant

# Improving Recovery Performance

## Synchronous Recovery

- Step 1: Recovery memory management
- Step 2: Recover primary data
- Step 3: Continue unfinished statements
- Step 4: Rebuild secondary data structures on DRAM
- Step 5: Start accepting user queries

Primary data already "loaded"

Restart time depends on the size of secondary data structures to be rebuilt

## Instant Recovery

- Idea 1:
  - Use primary data to answer queries
  - Rebuild secondary data structures asynchronously

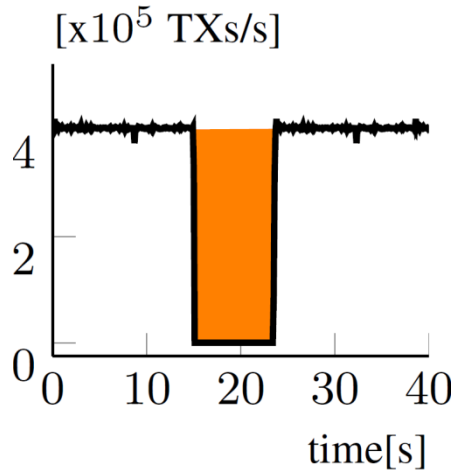Instant responsiveness

- Idea 2:
  - Persist part of or all secondary data structures in SCM

Instant recovery at peak performance
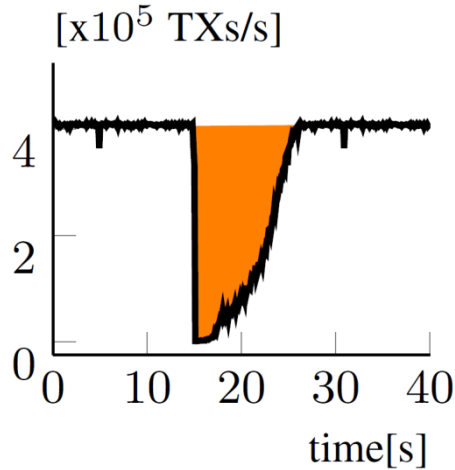
Perf. Penalty on throughput

# Evaluation: Recovery Time



| Synchronous Recovery | Instant Recovery | | |
|---|---|---|---|
| | 0% indexes in SCM | 40% indexes in SCM | 100% indexes in SCM |

First query accepted after ~8s, i.e., Recovery delta = 8s

Throughput: -0%
Recovery area: -16%
Recovery delta: ~8s

Throughput: -14%
Recovery area: -82%
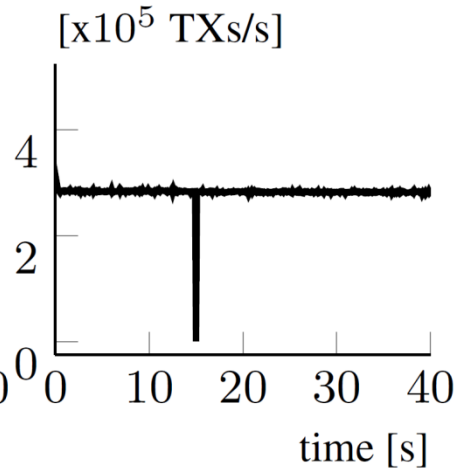Recovery delta: <2s

Throughput: -30%
Recovery area: -99,8%
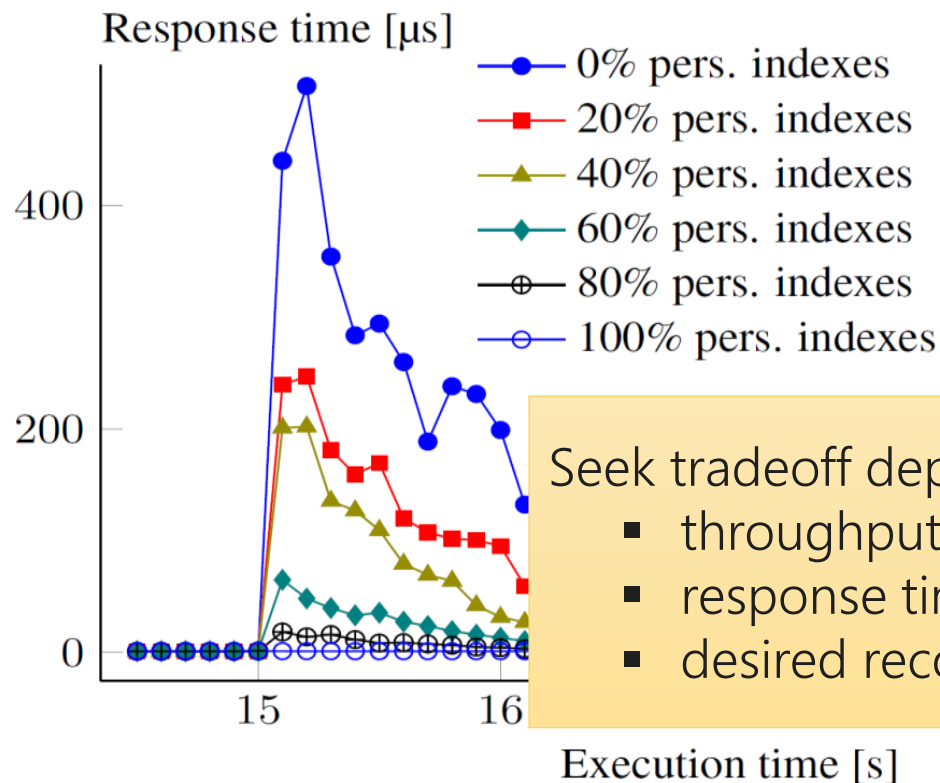Recovery delta: <5ms

# Evaluation: Throughput Vs. Recovery

Slowdown[%]

Area[$10^6$ x TXs]

Throughput drop limited to 30%

Curves are not linear: secondary data structures are not equally important for TATP

- Throughput drop
- Recovery area

Percentage of persistent indexes in SCM

Taking advantage of a workload's characteristics leads to an optimal tradeoff

# Evaluation: Average Response Time



Max. avg. (over 100ms) Response time:
- 0% pers. indexes: **506μs**
- 100% pers. indexes: **2μs**

Seek tradeoff depending on:
- throughput requirements
- response time requirements
- desired recovery performance

# Conclusion and Future Work

**WE SHOWED THAT SCM CAN HELP:**

- Achieve instant recovery for main-memory databases
- Continue unfinished transaction at crash time
- Simplify durability management
- Remove the need for a traditional transactional log

**CURRENT AND FUTURE WORK INCLUDE:**

- Improve recovery performance without compromising query performance
- Design new SCM-friendly persistent indexing structures
- Persistent, DRAM like memory management for SCM
- Testing tools for single-level store architectures

# Will SCM trigger a new rewrite of databases?

*Thank You! Questions? Comments?*

Ismail Oukid
ismail.oukid@sap.com
https://wwwdb.inf.tu-dresden.de/team/external-members/ismail-oukid/

## Instant Recovery for Main-Memory Databases

Ismail Oukid*°, Wolfgang Lehner*, Thomas Kissinger*, Peter Bumbulis°, and Thomas Willhalm[+]

*TU Dresden        °SAP SE        [+]Intel GmbH