# OLTP Through the Looking Glass 16 Years Later
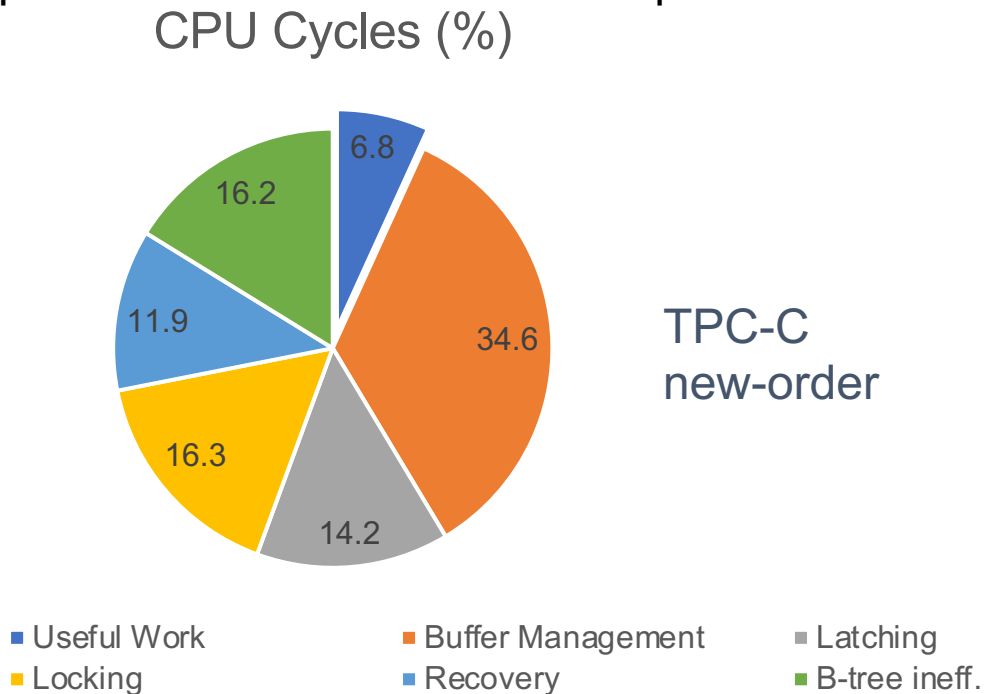
**Xinjing Zhou**, Viktor Leis, Xiangyao Yu, Michael Stonebraker

**MIT CSAIL**, TUM, UW-Madison, MIT CSAIL

# OLTP Looking Glass Back in 2008

- A performance study of a disk-based OLTP system - Shore
- Bottlenecks were spread across various core components when data fits in memory

CPU Cycles (%)

TPC-C
new-order

- Useful Work
- Buffer Management
- Latching
- Locking
- Recovery
- B-tree ineff.

6.8
34.6
14.2
16.3
11.9
16.2

# Many New OLTP Engines since then
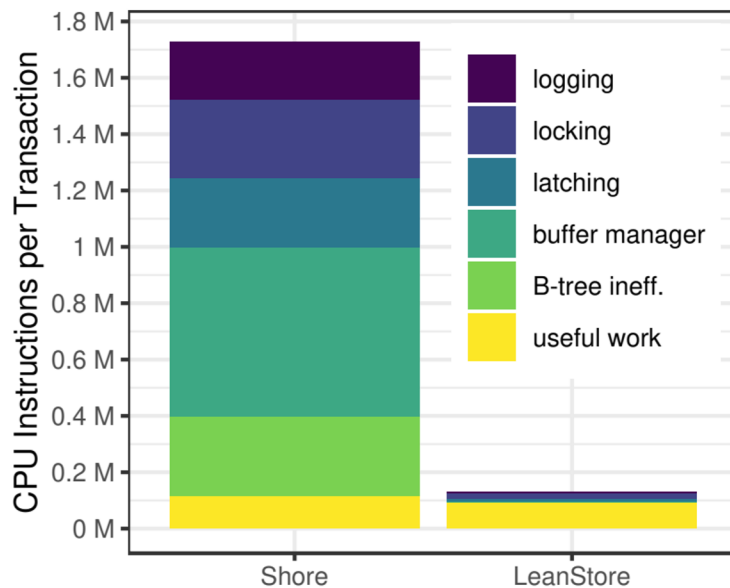
**H-Store** **VOLT**DB **leanstore**
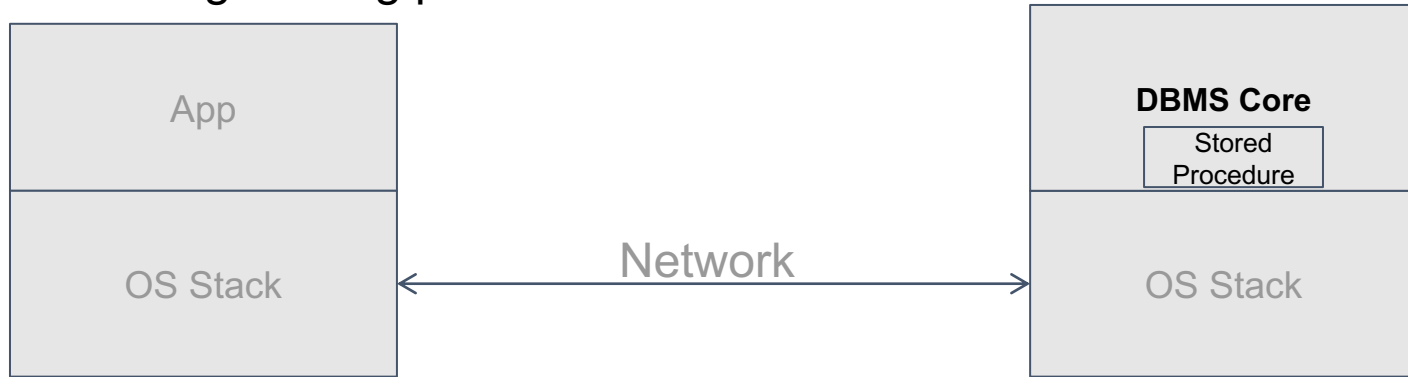
Silo **HyPer** **UMBRA**

Hekaton  .......



TPC-C
new-order

# Problems of Previous Research

- Benchmarks ignore OS stack and communication
- Most assume stored-procedure as the core technique to reduce network overhead.
- The reality [1-3]: many apps prefer interactive transactions due to better software engineering practices

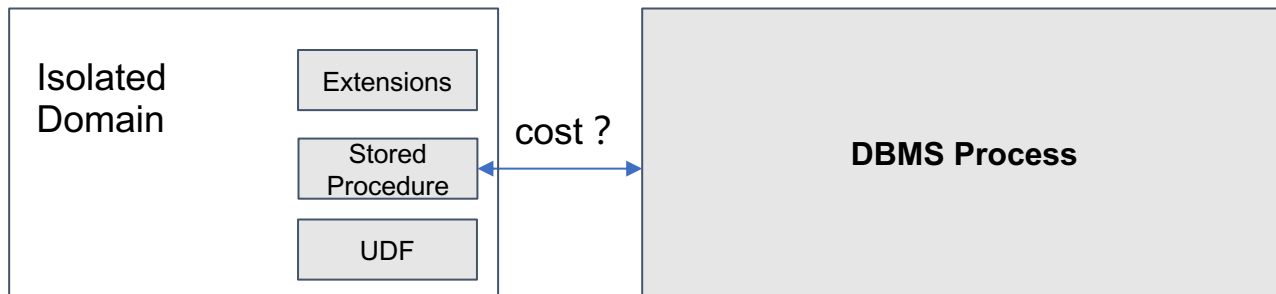| App | Network | **DBMS Core** |
|-----|---------|---------------|
| OS Stack | | Stored Procedure |
| | | OS Stack |

[1] Pavlo, Andrew. "What are we doing with our lives? Nobody cares about our concurrency control research." *SIGMOD 2017*.
[2] Gupta, Surabhi, and Karthik Ramachandra. "Procedural Extensions of SQL: Understanding their usage in the wild." *VLDB 2021*
[3] Hu, Gansen, et al. "WeBridge: Synthesizing Stored Procedures for Large-Scale Real-World Web Applications." *SIGMOD 2024.*

# Security of Stored-procedure

- Procedure run in the same address space of DBMS process for performance
  - written in various languages: PL/SQL, C/C++, Java, Python
- Malicious/errant procedures could read unauthorized data or crash DBMS
- DBMSs are becoming more multi-tenant as people move to the cloud
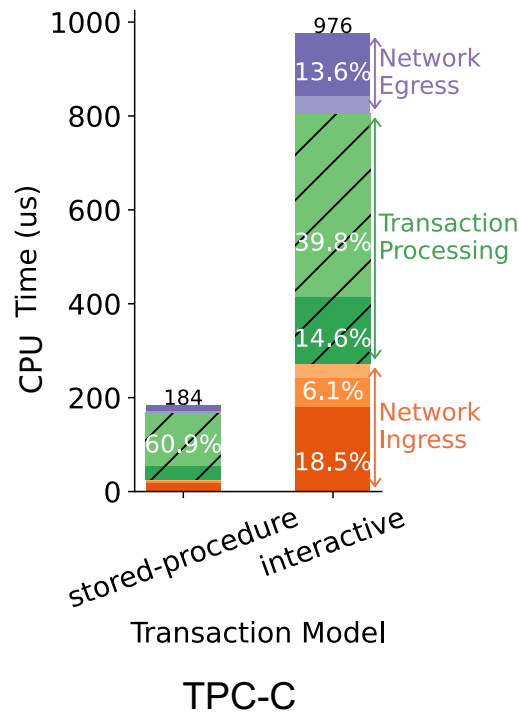- This applies to other extensibility mechanisms: UDF and extensions

Isolated Domain

Extensions

Stored Procedure

UDF

cost ?

**DBMS Process**

# OLTP Looking Glass 2.0

- Consider OS stacks and communication
- Consider procedure isolation
- Assume previous bottlenecks were solved after more than a decade of research - We use VoltDB as the testbed.
- Assume single-partition transactions
and single-node setup

Transaction Processing

Network Ingress

Network Egress

Partitioned OLTP Engine 1 — 6

Partitioned OLTP Engine N

5

...

Stored Procedure Executor 1 — 4

Stored Procedure Executor N

3  7

**Userspace**

Network Thread 1

...

Network Thread M

2

TCP/IP Stack

1  NiC Driver & Interrupt  8

Ingress

Egress  **Kernel**

# No-Isolation – Server-side CPU-time Breakdown, Communication is the bottleneck



**Legend:** Network Receive, Socket Read, Request Queuing, Isolation Overhead, Procedure Execution, Query Execution, Response Queuing, Network Send

YCSB-C

Voter

TPC-C

# Isolating Procedure

| | Stored Procedure Process | Stored Procedure Process in Container | Stored Procedure Process | App with interactive transactions |
|---|---|---|---|---|
| | | | Guest OS Kernel | OS Kernel |
| | ↕ IPC | ↕ IPC | ↕ IPC | ↕ Network |
| DBMS Process | DBMS Process | DBMS Process | DBMS Process | DBMS Process |
| Stored Procedure | | | | |
| OS Kernel | OS Kernel | OS Kernel | OS Kernel | OS Kernel |

No-Isolation    Process Isolation    Container Isolation    VM Isolation    Client-server Isolation

# Isolated Stored Procedure Execution, Communication for Isolation is the bottleneck



**Legend:**
- Network Receive
- Socket Read
- Request Queuing
- Isolation Overhead
- Procedure Execution
- Query Execution
- Response Queuing
- Network Send

**No Isolation**

**Process Isolation**

**Container Isolation**

**VM Isolation**

CPU Time (us)

| Isolation Mechanisms | none | tcp | shm | shm domain socket | docker shm | docker tcp | vm shm | vm tcp |
|---|---|---|---|---|---|---|---|---|

- none: 178, 60.7%
- tcp: 616, 23.6%, 63.2%
- shm: 296, 48.5%, 26.1%
- shm domain socket: 563, 27.2%, 57.1%
- docker shm: 297, 49.4%, 25.3%
- docker tcp: 641, 22.4%, 64.6%
- vm shm: 340, 47.4%, 30.8%
- vm tcp: 3250, 90.9%

TPC-C

# Wish #1: Usable Kernel Bypass

- DPDK + User space TPC/IP stack (F-Stack)
  - Reduces kernel network stack overhead of VoltDB by 85%
- Only two DBMS vendors support kernel-bypass: Yellowbrick and ScyllaDB
- Three Problems
  - **Interface-Mismatch**: DPDK is a layer-2 stack – no transport/routing layer support
  - **Design Limitation**: A DPDK app requires complete control of a NIC
    - Linux tooling are not available on DPDK-managed NIC, making debugging and deployment hard.
  - **Engineering and Maintenance:** User-space TCP/IP stacks often require DBMS to rewrite their network layer code due to API differences.

# Wish #2: More Exploration in the Trade-off Space



Stored Procedure

Interactive Transaction

Security

Is there a better approach that attend to all three?

Debuggability
Testing
Language Flexibility
DBMS-agnostic
Version-Control

Ease-of-Use

Performance

# Conclusion

- Communication and OS stack should be more focused by the community for future DB systems research.
- We need more usable and efficient kernel bypass abstractions to make larger impact on DBMS.
- We should revisit the debate about stored-procedure and interactive transaction, factoring in security and usability. The trade-off space is under-explored.

CIDR 2025 Preprint