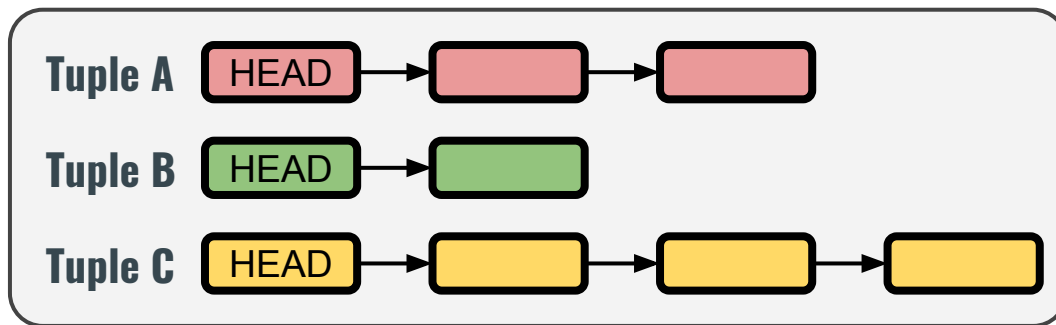# An Experimental Evaluation of In-Memory Multi-Version Concurrency Control

Yingjun Wu, Joy Arulraj, Jiexi Lin, Ran Xian, Andrew Pavlo

**NUS** National University of Singapore

**Carnegie Mellon University**

CARNEGIE MELLON DATABASE GROUP
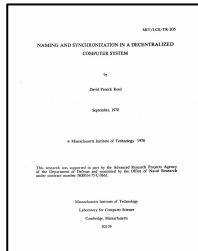
# MULTI-VERSION CONCURRENCY CONTROL

❏ **Avoid read-write conflicts**
❏ **Support time-travel queries**
❏ **Enable snapshot isolation**



**VERSION CHAINS**

# A BRIEF HISTORY OF MVCC

## 1979: FIRST MENTION

DAVID
**REED**

Naming and synchronization in a decentralized computer system
*Ph.D. Thesis, 1979*

# A BRIEF HISTORY OF MVCC

- **1979:** FIRST MENTION
- **1981:** FIRST IMPLEMENTATION

**InterBase/Firebird**

# A BRIEF HISTORY OF MVCC

- **1979: FIRST MENTION**
- **1981: FIRST IMPLEMENTATION**
- **1984:** Oracle
- **1985:** Postgres

# A BRIEF HISTORY OF MVCC

- **1979: FIRST MENTION**
- **1981: FIRST IMPLEMENTATION**
- **1984:** Oracle
- **1985:** Postgres
- **2001:** MySQL-InnoDB

# A BRIEF HISTORY OF MVCC

- **1979: FIRST MENTION**
- **1981: FIRST IMPLEMENTATION**
- **1984:** Oracle
- **1985:** Postgres
- **2001:** MySQL-InnoDB
- **2010-2017:** Hyrise, Hekaton, MemSQL, SAP HANA, NuoDB, HyPer...
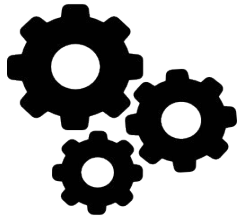
7

# A BRIEF HISTORY OF MVCC



Search for the best **MVCC** scheme
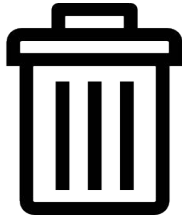for multi-core main-memory DBMSs

# DESIGN DECISIONS OF MVCC

**Concurrency Control Protocol**

**Version Storage**

**Garbage Collection**

**Index Management**

# CONCURRENCY CONTROL PROTOCOL

| SCHEME | DBMS |
|---|---|
| Timestamp Ordering (MVTO) | N/A |
| Optimistic Concurrency Control (MVOCC) | HYRISE   HEKATON Microsoft SQL Server   HyPer   MEMSQL |
| Two-phase Locking (MV2PL) | ORACLE   MySQL   SAP HANA   NUODB |
| Serialization Certifier | PostgreSQL |

10

# CONCURRENCY CONTROL PROTOCOL

❏ **Approach #1: Timestamp Ordering (MVTO)**
   ❏ The DBMS assigns transactions timestamps that determine serial order.

| | TXN-ID | READ-TS | BEGIN-TS | END-TS |
|---|---|---|---|---|
| $A_x$ | 0 | Tid | 10 | 20 |
| $B_x$ | Tid | 17 | 15 | 30 |
| $B_{x+1}$ | Tid | 0 | - | - |

READ(A)

WRITE(B)

# CONCURRENCY CONTROL PROTOCOL

## Approach #3: Optimistic Concurrency Control (MVOCC)

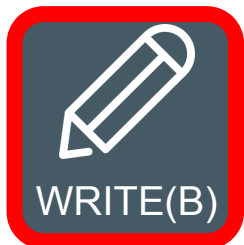Transactions optimistically access physical versions before validating the read-set consistency.

READ(A)

WRITE(B)

| | TXN-ID | BEGIN-TS | END-TS |
|---|---|---|---|
| $A_x$ | 0 | 10 | 20 |
| $B_x$ | Tid | 15 | 30 |
| $B_{x+1}$ | Tid | - | - |

12

# CONCURRENCY CONTROL PROTOCOL

## ❏ Approach #3: Two-Phase Locking (MV2PL)

❏ Transactions acquire appropriate lock on physical version before they can read/write a logical tuple.

|  | TXN-ID | READ-CNT | BEGIN-TS | END-TS |
|---|---|---|---|---|
| $A_x$ | 0 | 2 | 10 | 20 |
| $B_x$ | Tid | 0 | 15 | 30 |
| $B_{x+1}$ | Tid | 0 | - | - |

READ(A)

WRITE(B)

# CONCURRENCY CONTROL PROTOCOL

❏ **Approach #4: Serialization Certifier**

  ❏ The DBMS maintains a serialization graph for detecting and removing "dangerous structures" formed by concurrent transactions.

|  | TXN-ID | BEGIN-TS | END-TS |
|---|---|---|---|
| $A_X$ | 0 | 10 | 20 |
| $B_X$ | Tid | 15 | 30 |
| $B_{X+1}$ | Tid | - | - |

READ(A)

WRITE(B)

# VERSION STORAGE

| SCHEME | DBMS |
|--------|------|
| Append-Only | PostgreSQL   HYRISE   Microsoft HEKATON SQL Server   NUODB   MEMSQL |
| Time-Travel | SAP HANA |
| Delta | ORACLE   MySQL   HyPer |

# VERSION STORAGE

❏ **Approach #1: Append-Only Storage (Oldest-to-Newest)**
  ❏ New versions are appended to the same table space.
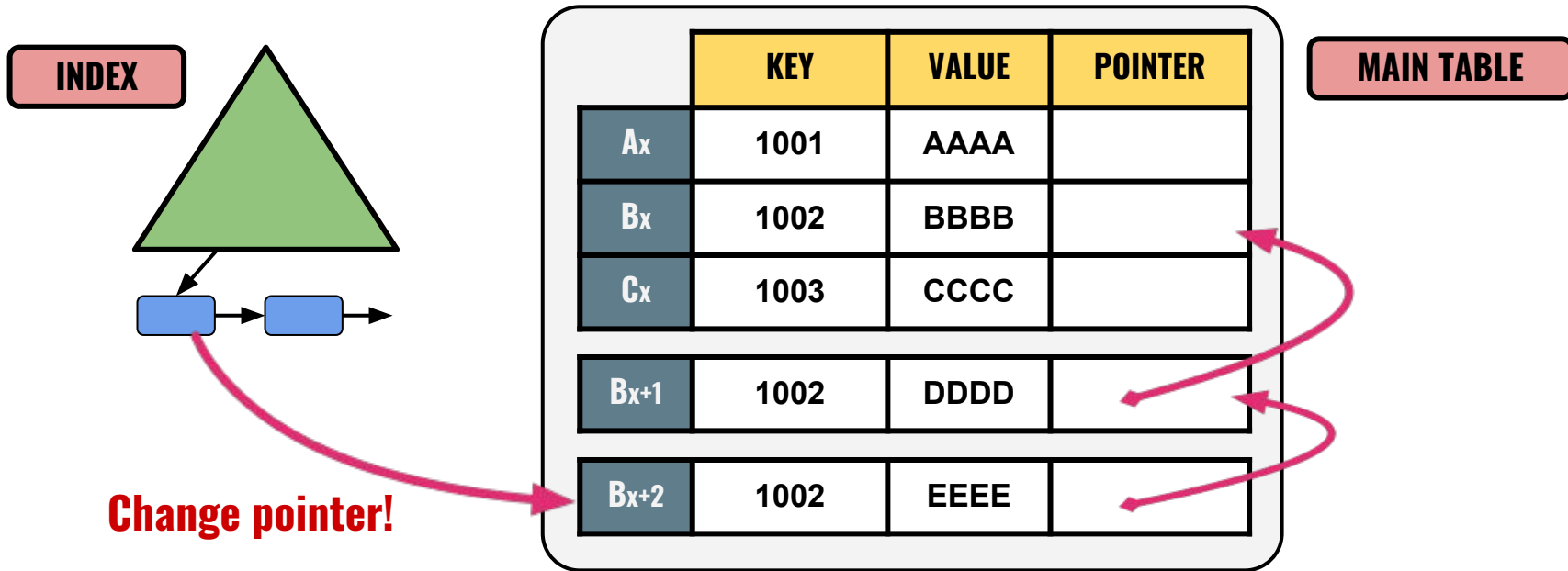


INDEX

MAIN TABLE

| | KEY | VALUE | POINTER |
|---|---|---|---|
| A$_x$ | 1001 | AAAA | |
| B$_x$ | 1002 | BBBB | |
| C$_x$ | 1003 | CCCC | |
| B$_{x+1}$ | 1002 | DDDD | |
| B$_{x+2}$ | 1002 | EEEE | |

Unchanged!

16

# VERSION STORAGE

❏ **Approach #1: Append-Only Storage (Newest-to-Oldest)**

    ❏ New versions are appended to the same table space.



INDEX

MAIN TABLE

| | KEY | VALUE | POINTER |
|---|---|---|---|
| $A_x$ | 1001 | AAAA | |
| $B_x$ | 1002 | BBBB | |
| $C_x$ | 1003 | CCCC | |
| $B_{x+1}$ | 1002 | DDDD | |
| $B_{x+2}$ | 1002 | EEEE | |

**Change pointer!**

# VERSION STORAGE

## ❏ Approach #2: Time-Travel Storage

❏ Old versions are copied to separate table space.



INDEX

MAIN TABLE

| | KEY | VALUE | POINTER |
|---|---|---|---|
| A$_x$ | 1001 | AAAA | |
| B$_{x+2}$ | 1002 | EEEE | |
| C$_x$ | 1003 | CCCC | |

| | | | |
|---|---|---|---|
| B$_x$ | 1002 | BBBB | |
| B$_{x+1}$ | 1002 | DDDD | |

TIME-TRAVEL TABLE

18

# VERSION STORAGE

❑ **Approach #3: Delta Storage**
  ❑ The original values of the modified attributes are copied into a separate delta space.
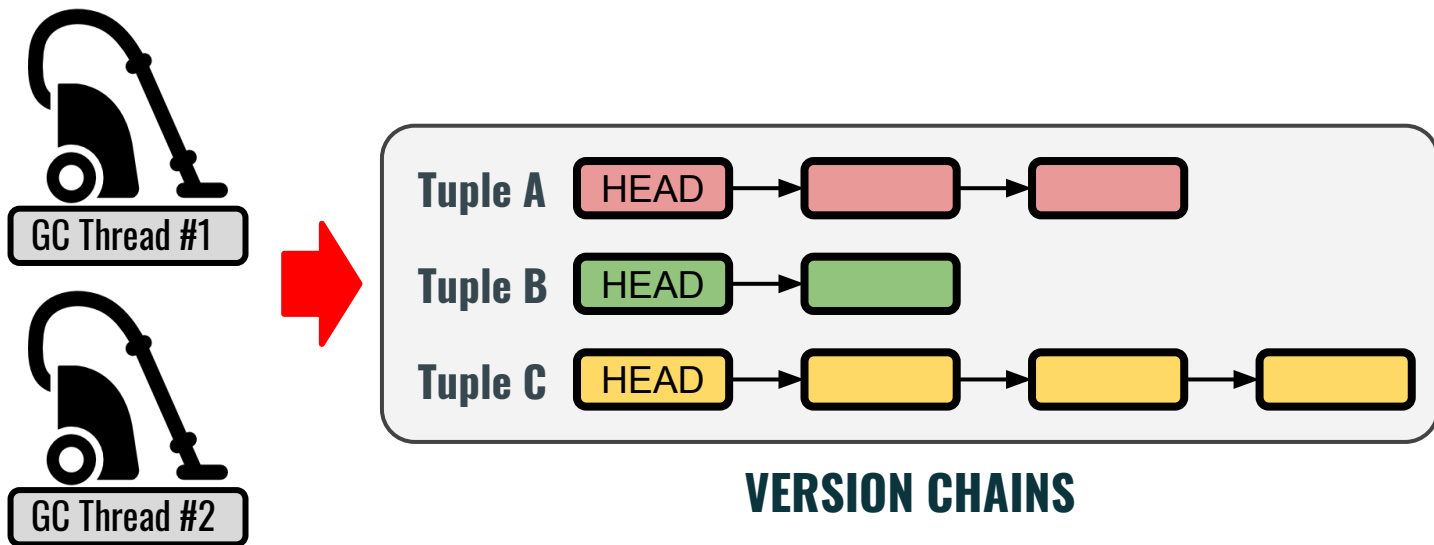


INDEX

MAIN TABLE

| | KEY | VALUE | POINTER |
|---|---|---|---|
| $A_x$ | 1001 | AAAA | |
| $B_{x+2}$ | 1002 | EEEE | |
| $C_x$ | 1003 | CCCC | |

| $B_x$ | | | $B_{x+1}$ | |
|---|---|---|---|---|
| BBBB | | | DDDD | |

DELTA STORAGE

19

# GARBAGE COLLECTION

| SCHEME | DBMS |
|--------|------|
| Tuple-Level | ORACLE MySQL PostgreSQL HEKATON Microsoft SQL Server SAP HANA MEMSQL HYRISE NUODB |
| Transaction-Level | HyPer |

# GARBAGE COLLECTION

❑ **Approach #1: Tuple-level**
   ❑ Find old versions by examining tuples directly.



**VERSION CHAINS**

# GARBAGE COLLECTION

❏ **Approach #2: Transaction-level**
   ❏ Transactions keep track of their old versions so the DBMS does not have to scan tuples to determine visibility.
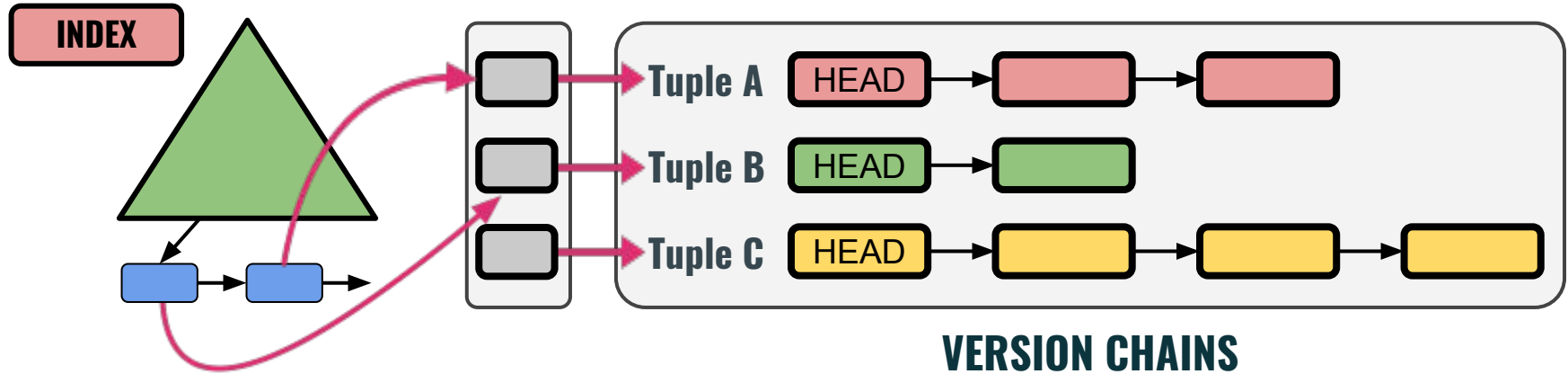


GC Thread #1

GC Thread #2

TXN1 WRITE-SET

TXN2 WRITE-SET

TXN3 WRITE-SET

⋮

**TRANSACTIONS**

Tuple A  HEAD → → 

Tuple B  HEAD → 

Tuple C  HEAD → → → 

**VERSION CHAINS**

# INDEX MANAGEMENT

| SCHEME | DBMS | | | | |
|---|---|---|---|---|---|
| Logical Pointers | ORACLE | MySQL | SAP HANA | HyPer | NUODB |
| Physical Pointers | PostgreSQL | HEKATON Microsoft SQL Server | HYRISE | MEMSQL | |

# INDEX MANAGEMENT

❏ **Approach #1: Logical Pointers**
   ❏ Use a fixed identifier per tuple that does not change.



**VERSION CHAINS**

# INDEX MANAGEMENT

❏ **Approach #2: Physical Pointers**
  ❏ Use the physical address to the version chain head.



**VERSION CHAINS**

# EVALUATION

- ❏ **Benchmarks**
  - ❏ **YCSB**
  - ❏ **TPC-C**
- ❏ **Configuration**
  - ❏ **4X Intel Xeon E7-4820 (40 cores)**
  - ❏ **128 GB DRAM**

PELOTON

# EVALUATION

❏ **Concurrency Control Protocol**

🔴━🔴 **MVTO**   🟦━🟦 **MVOCC**   🔻━🔻 **MV2PL**   🔺━🔺 **SI+SSN**

# EVALUATION

❑ **Version Storage**

# EVALUATION

❏ **Garbage Collection**

# EVALUATION

❏ **Index Management**



30

# EVALUATION

❏ **MVCC Configurations**

# CONCLUSION

❏ **Choosing the best MVCC scheme is challenging**

    ❏ **Four design aspects**

    ❏ **Multiple design decision combinations**

    ❏ **Optimize for different objectives**

# END

@yingjunwu