# Hekaton: SQL Server's Memory-Optimized OLTP Engine

Cristian Diaconu, Craig Freedman, Erik Ismert, Per-Åke Larson,
Pravin Mittal, Ryan Stonecipher, Nitin Verma, Mike Zwilling
Microsoft
{cdiaconu, craigfr, eriki, palarson, pravinm, ryanston, nitinver, mikezw}@microsoft.com
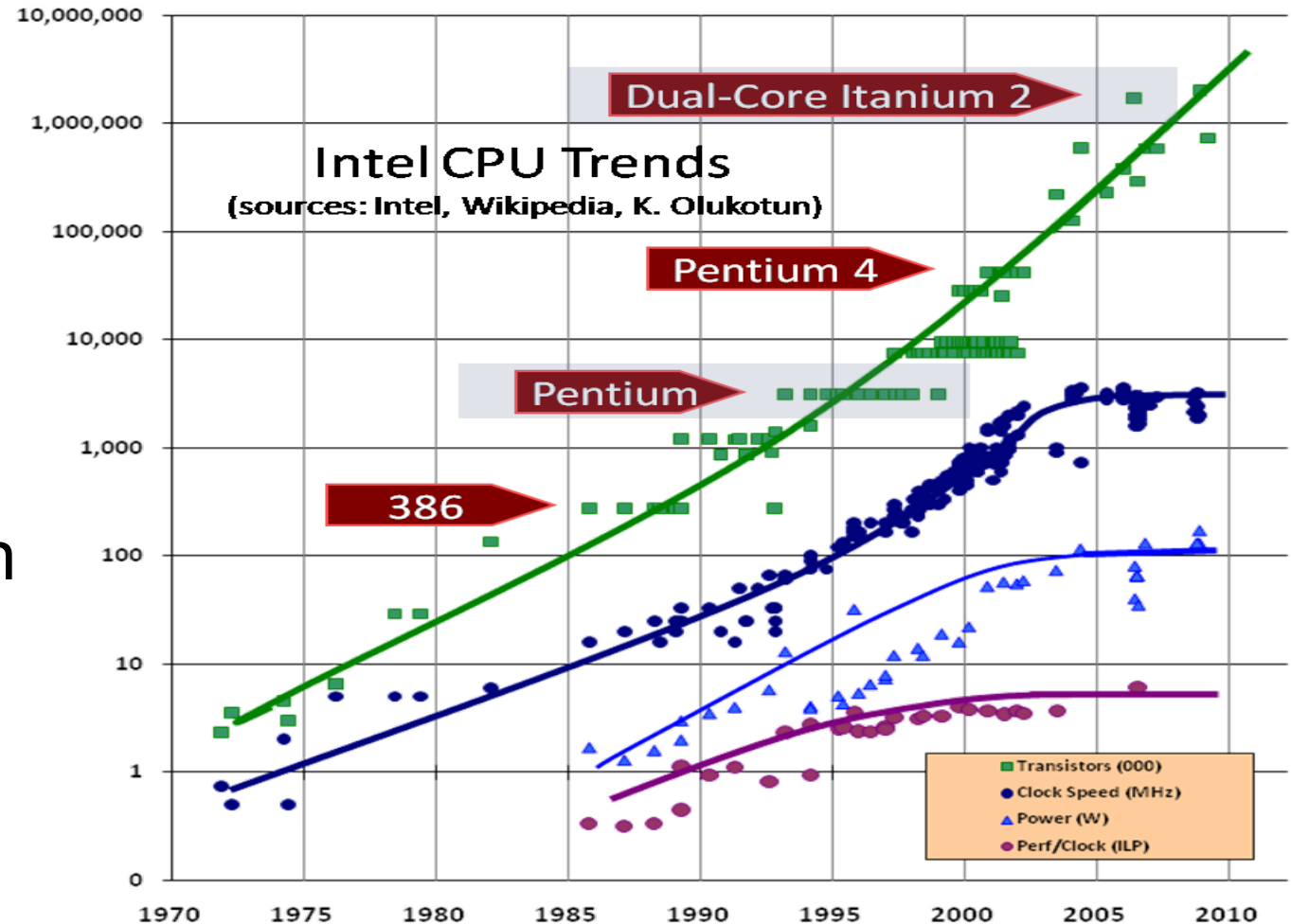
Presented by: Prateek Gulati

# Agenda

- Why do you need in-memory processing?
- Hekaton engine overview
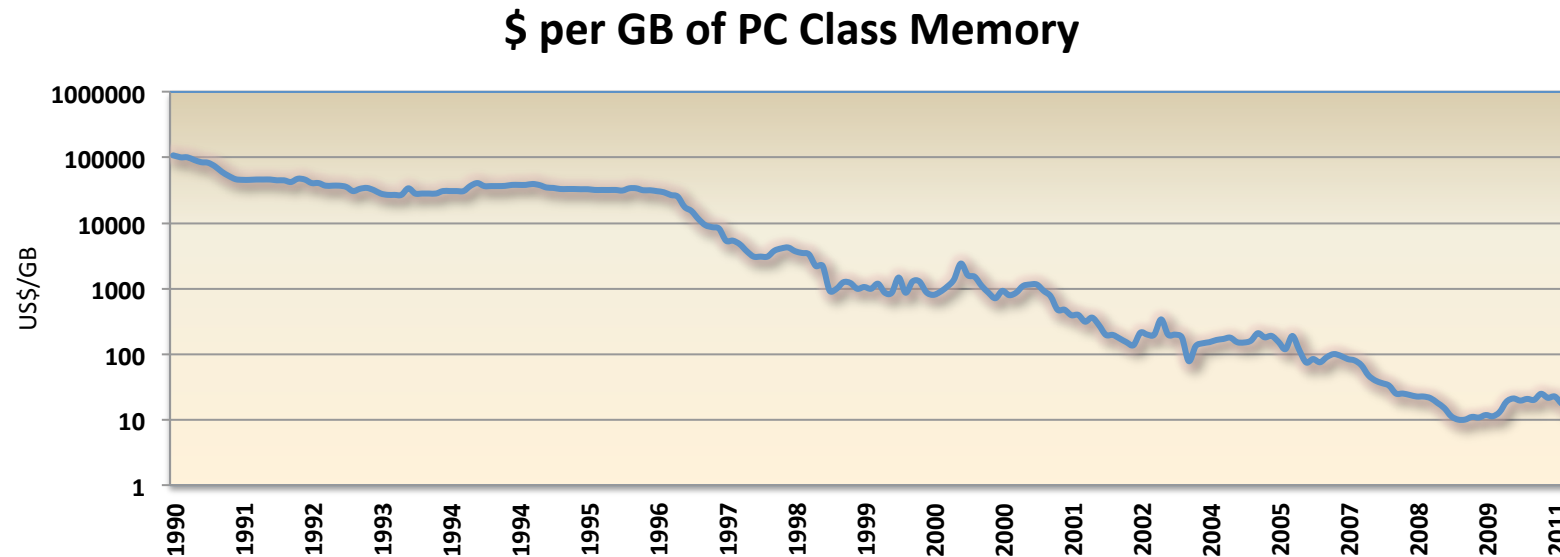- How it is done
- Benefits
- Limitations

# Industry Trends: CPU

- Computing power holds Moore Law due to parallelism
- CPU clock frequency stalled
- Parallel processing has its limits due to lock contention



Intel CPU Trends
(sources: Intel, Wikipedia, K. Olukotun)

Dual-Core Itanium 2

Pentium 4

Pentium

386

- Transistors (000)
- Clock Speed (MHz)
- Power (W)
- Perf/Clock (ILP)

# Industry Trends: RAM

- RAM prices continue to fall
- Servers have HUGE memory
- DDR4 expected to hit mainstream in 2014-2015
- Traditional page based architecture has limitations, even when all pages are in memory

**$ per GB of PC Class Memory**

# Hekaton-In-memory OLTP engine Architecture

**Architectural Pillars**

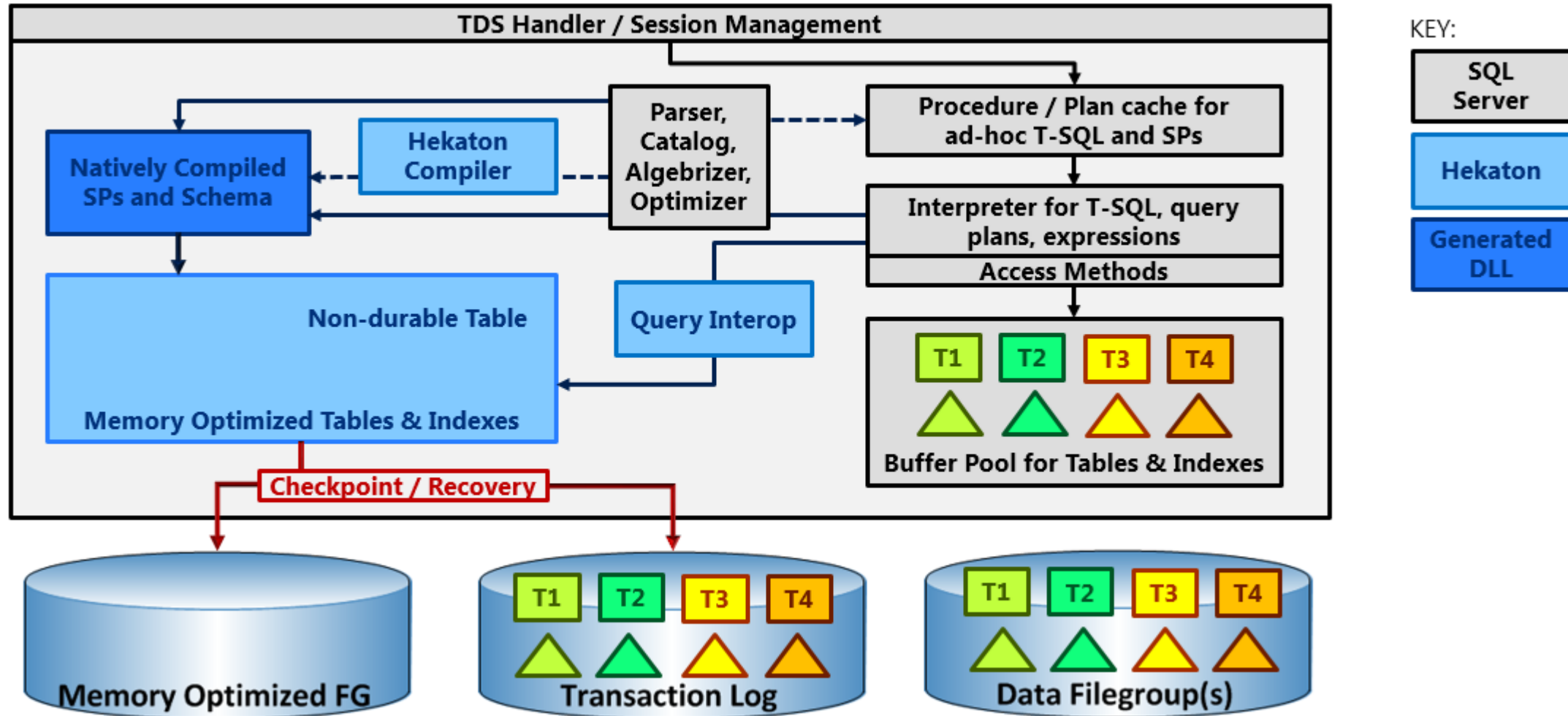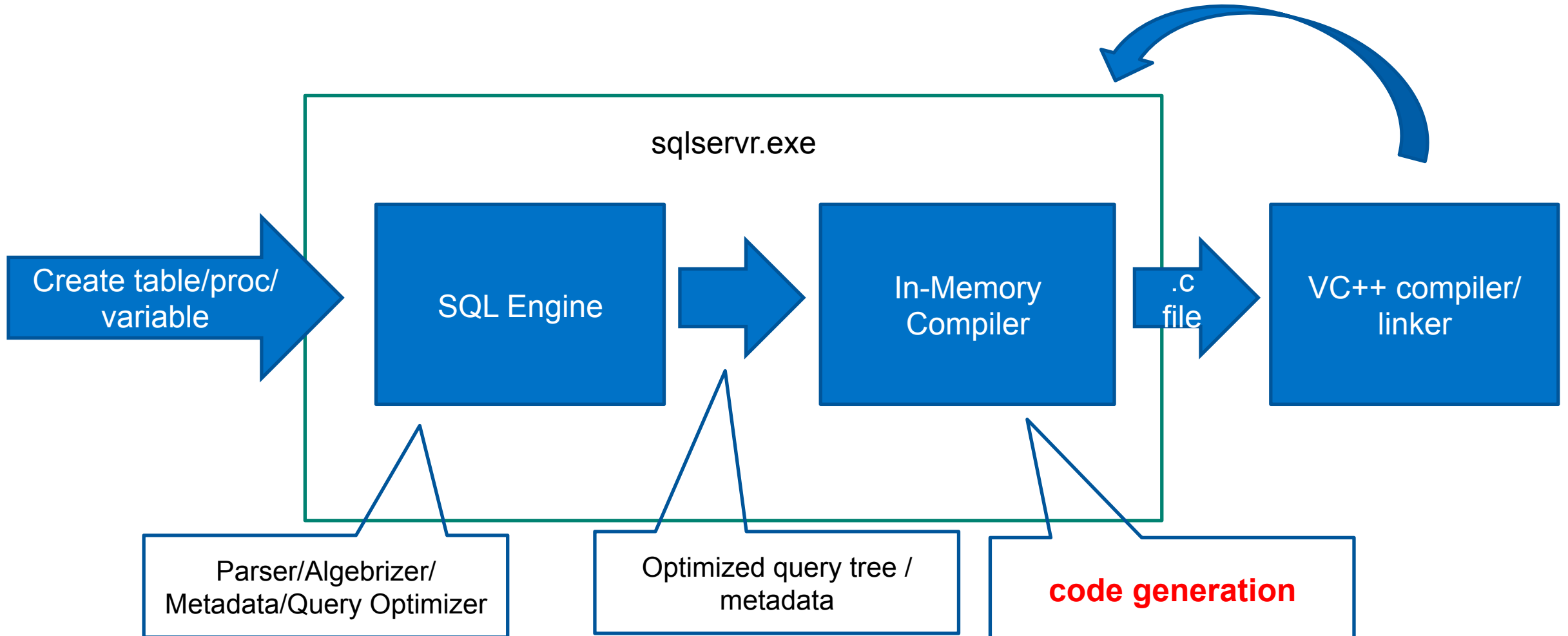| Main-Memory Optimized | T-SQL Compiled to Machine Code | High Concurrency | SQL Server Integration |
|---|---|---|---|
| • Optimized for in-memory data<br>• Memory optimized Indexes (hash and range) exist only in memory<br>• No buffer pool, B-trees<br>• Stream-based storage<br>• Transaction log optimization (block writes, no undo) | • T-SQL compiled to machine code via C code generator and VC<br><br>• Invoking a procedure is just a DLL entry-point<br><br>• Aggressive optimizations at compile-time | • Multi-version optimistic concurrency control (MVCC) with full ACID support<br>• Core engine uses non blocking lock-free algorithms<br>• No lock manager, latches or spinlocks<br>• No TempDB | • Same manageability, administration & development experience<br><br>• Integrated queries & transactions<br><br>• Integrated backup/restore<br><br>• If SQL Server crashes data is fully recoverable. |

# Hekaton Integration with SQL Server

# Native Compilation Process

Compile T-SQL statements *and* table data access logic into machine code



sqlservr.exe

Create table/proc/ variable → SQL Engine → In-Memory Compiler → .c file → VC++ compiler/ linker

Parser/Algebrizer/ Metadata/Query Optimizer

Optimized query tree / metadata

**code generation**

# Native Compiled Stored Procedures

**Interpreted T-SQL Access**

- Access both memory- and disk-based tables
- Less performant
- Virtually all T-SQL functions supported
- **When to use**
  - Ad hoc queries
  - Reporting-style queries
  - Speeding up app migration

**Natively Compiled Procs**

- Access only memory optimized tables
- Maximum performance
- Limited T-SQL functions supported
- **When to use**
  - OLTP-style operations
  - Optimize performance critical business logic
  - More the logic embedded, better the performance improvement
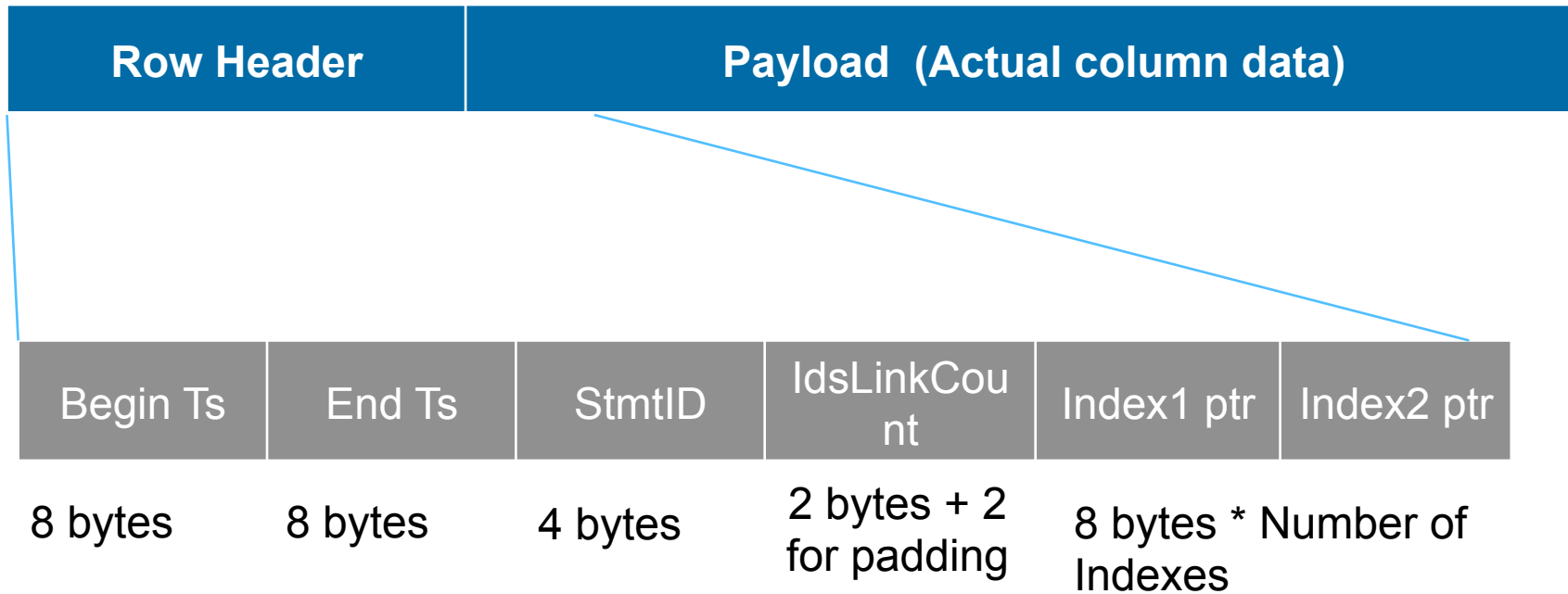
# In-Memory OLTP Structures summary

Rows
- Row structure is optimized for memory access
- There are no Pages
- Rows are versioned and there are no in-place updates
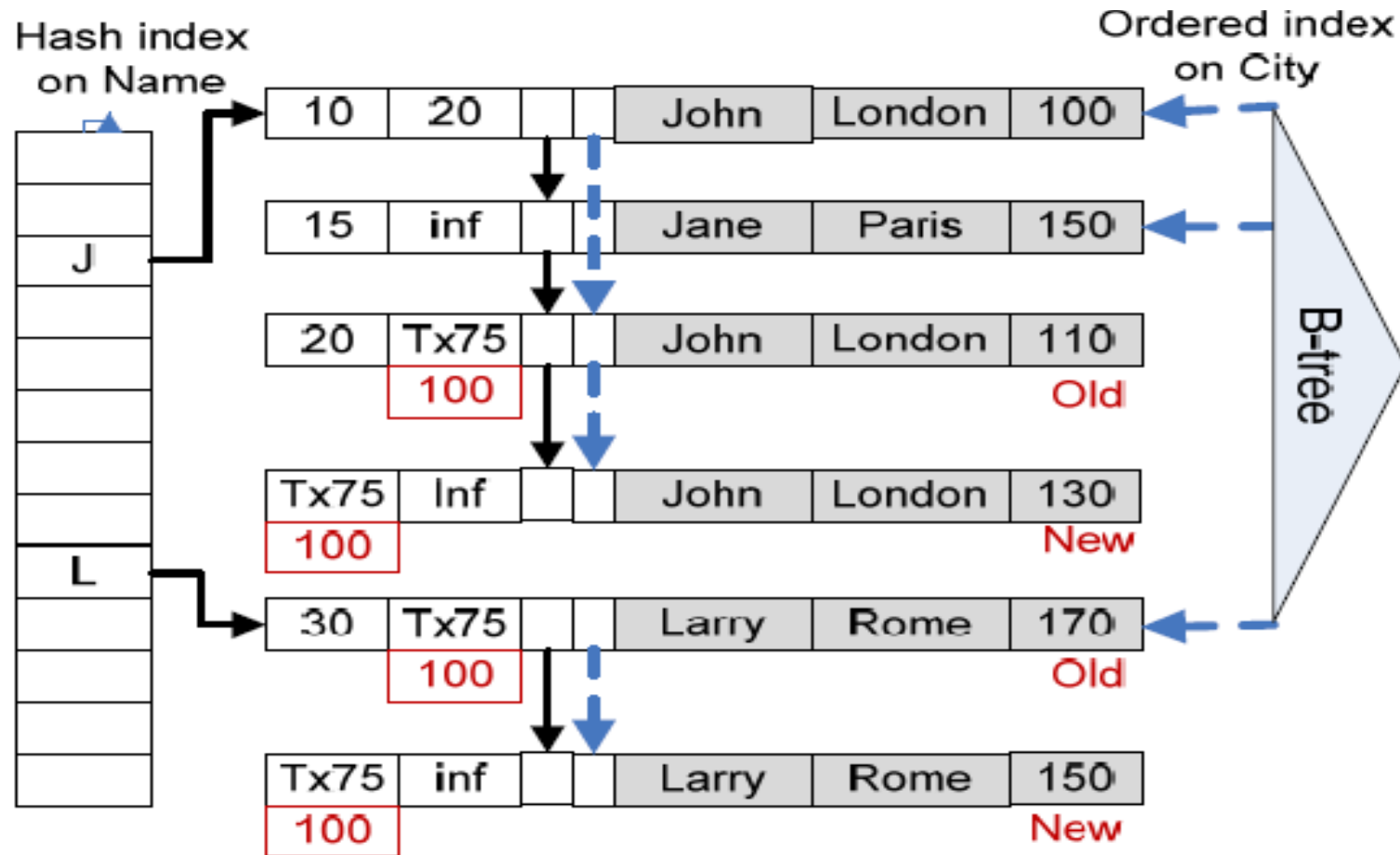- Fully durable by default (but they don't have to be)

Indexes
- There is no clustered index, only non-clustered indexes
- Indexes point to rows, access to rows is via an index
- Indexes do not exist on disk, only in memory, recreated during recovery
- Hash indexes for point lookups
- Range indexes for ordered scans and Range Scans

# In-Memory Row Format

| Row Header | Payload  (Actual column data) |
|---|---|

| Begin Ts | End Ts | StmtID | IdsLinkCount | Index1 ptr | Index2 ptr |
|---|---|---|---|---|---|
| 8 bytes | 8 bytes | 4 bytes | 2 bytes + 2 for padding | 8 bytes * Number of Indexes | |

- Begin/End timestamp determines row's version validity and visibility
- No concept of data pages, only rows exist
- Row size limited to 8060 bytes (@table create time) to allow data to be moved to disk-based tables
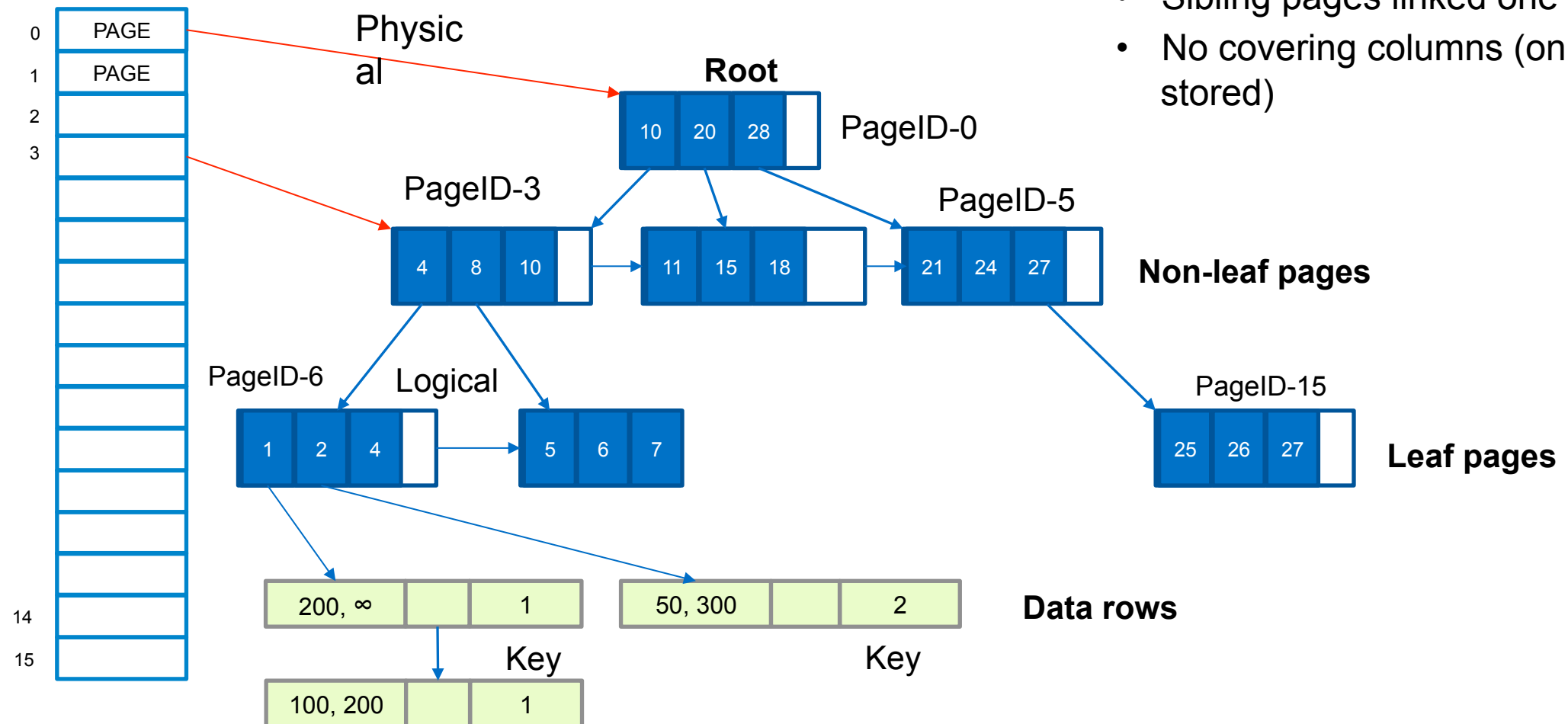- Not every SQL table schema is supported (Ex: LOB and SqlVariant)

# Hash Indexes

# Non Clustered (Range) Index

- No latch for page updates
- No in-place updates on index pages
- Page size- up to 8K. Sized to the row
- Sibling pages linked one direction
- No covering columns (only the key is stored)

Page Mapping Table

| | |
|---|---|
| 0 | PAGE |
| 1 | PAGE |
| 2 | |
| 3 | |
| | |
| | |
| | |
| | |
| | |
| | |
| | |
| | |
| | |
| | |
| 14 | |
| 15 | |

Physical

**Root**

| 10 | 20 | 28 | |

PageID-0

PageID-3

| 4 | 8 | 10 | |

| 11 | 15 | 18 | |

PageID-5

| 21 | 24 | 27 | |

**Non-leaf pages**

PageID-6    Logical

| 1 | 2 | 4 | |

| 5 | 6 | 7 |

PageID-15

| 25 | 26 | 27 | |

**Leaf pages**

| 200, ∞ | | 1 |

| 50, 300 | | 2 |

**Data rows**

Key

| 100, 200 | | 1 |

Key

# Memory Optimized Table Insert



| Timestamps | Chain ptrs | Name | City |
|---|---|---|---|

Hash index on Name

Hash index on City

*f*(John)

*f*(Prague)

| 50, ∞ | | Jane | Prague |
|---|---|---|---|

| 100, ∞ | | John | Prague |
|---|---|---|---|

| 90, ∞ | | Susan | Bogota |
|---|---|---|---|

T100: INSERT (John, Prague)

13

# Memory Optimized Table Update



| Timestamps | Chain ptrs | Name | City |
|---|---|---|---|

Hash index on Name

Hash index on City

f(John)

f(Beijing)

| 50, ∞ | | Jane | Prague |

| 100, 200 | | John | Prague |

| 200, ∞ | | John | Beijing |

| 90, 150 | | Susan | Bogota |

T200: UPDATE (John, Prague) to (John, Beijing)

14

# Memory Optimized Table Delete

| Timestamps | Chain ptrs | Name | City |
|---|---|---|---|

**Hash index on Name**

**Hash index on City**

| 50, ∞ | | Jane | Prague |
|---|---|---|---|

| 100, ∞ | | John | Prague |
|---|---|---|---|

| 90, 150 | | Susan | Bogota |
|---|---|---|---|

T150: DELETE (Susan, Bogota)

# Transaction Durability

- Transaction durability is ensured to allows system to recover memory-optimized table after a failure.
- Log streams contain the effects of committed transactions logged as insertion and deletion of row versions
- Checkpoint streams come in two forms:
  - a) data streams which contain all inserted versions during a timestamp interval,
  - b) delta streams, each of which is associated with a particular data stream and contains a dense list of integers identifying deleted versions for its corresponding data stream

- Hekaton table can be durable or non-durable

- Stored in a single memory-optimized FILEGROUP based on FILESTREAM implementation

- Sequential IO pattern (no random IO)
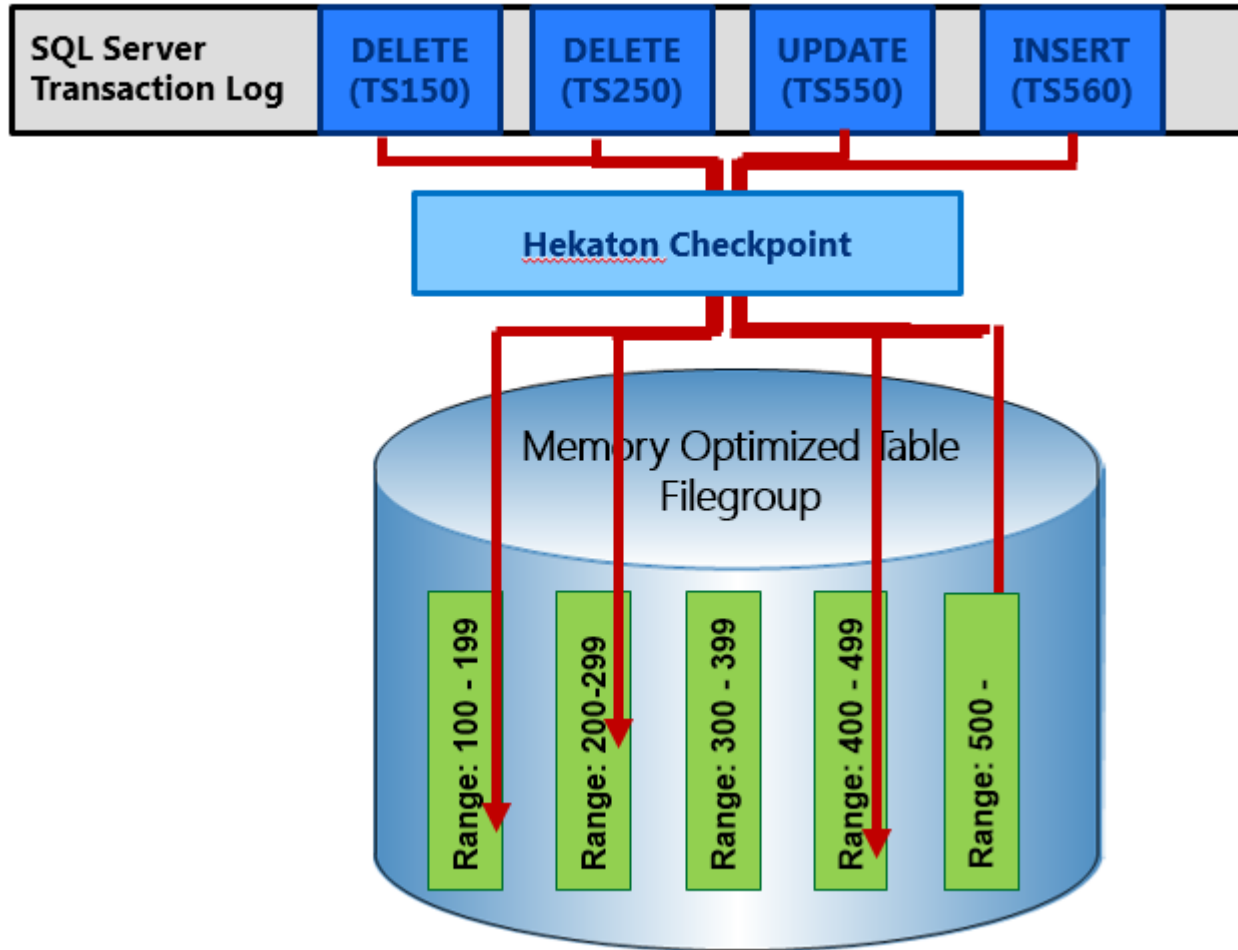
# Transaction Logging

- Uses database's transaction log to store content

- Each Hekaton log record contains a transaction log record header, followed by Hekaton-specific log content

- All logging in Hekaton is logical
  - No physical log records for physical structure modifications
  - No index-specific / index-maintenance log records
  - Redo-only log records in transaction log

# Checkpoints

## Hekaton Checkpoint

- Not tied to recovery interval or SQL checkpoint. Has its own log truncation
- Gets triggered when generated log exceeds a threshold (1GB) or internal min time-threshold  has crossed since last checkpoint or manual checkpoint
- Checkpoint  is a  "set of {Data, Delta} files and checkpoint file inventory to apply transaction log from"

# Populating Data / Delta files



Data files:

- Pre-allocated size (128 MB)
- Hekaton Engine switches to new data file when it estimates that current set of log records will fill the file
- Stores only the inserted rows
- Indexes exist only in memory, not on disk
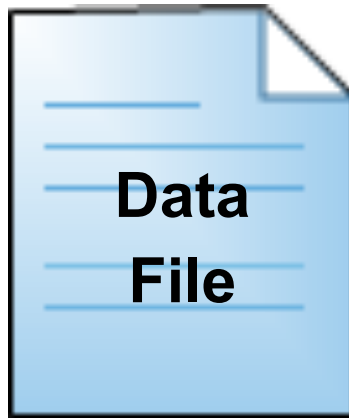- Once a data file is closed, it becomes read-only

Delta files:

- File size is not constant, write 4KB pages over time
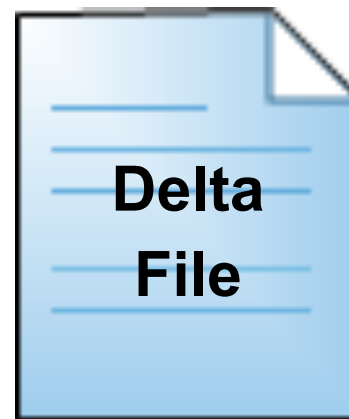- Stores IDs of deleted rows

# Data / Delta Files

**Transaction Timestamp Range**

0 → 100

**Checkpoint File Pair**



| Timestamp (INSERT) | TableID | RowID | Payload |
|---|---|---|---|
| Timestamp (INSERT) | TableID | RowID | Payload |
| Timestamp (INSERT) | TableID | RowID | Payload |

Data file contains rows inserted within a given transaction range

| Timestamp (INSERT) | Timestamp (DELETE) | RowID |
|---|---|---|
| Timestamp (INSERT) | Timestamp (DELETE) | RowID |
| Timestamp (INSERT) | Timestamp (DELETE) | RowID |

Delta file contains deleted rows within a given transaction range

# Merge Operation

- Merges 2+ adjacent data / delta files pairs into 1 pair
- Need for merge - Deleting rows causes data files to have stale rows
- Manual checkpoints closes file before it is "full"
- Reduces storage required to "store" active data rows
- Improves the recovery time
- Stored Procedure provided to invoke merge manually

# Merge Operation



Timestamp: 500

Memory Optimized Table Filegroup

Range: 100 - 199
Range: 200-299
Range: 300 - 399
Range: 400 - 499

Merge 200-399

Timestamp: 600

Memory Optimized Table Filegroup

Range: 100 - 199
Range: 200 - 299
Range: 200 - 399
Range: 300- 399
Range: 400 - 499

# Garbage Collection

| Timestamps | Chain ptrs | Name | City |
|---|---|---|---|

Hash index on Name

*f*(Jane)

*f*(John)

| 50, ∞ | | Jane | Prague |
|---|---|---|---|

| 100, 200 | | John | Prague |
|---|---|---|---|

| 200, ∞ | | John | Beijing |
|---|---|---|---|

| 90, 150 | | Susan | Bogota |
|---|---|---|---|

T250: *Garbage collection*

# Cooperative Garbage Collection

- Scanners can remove expired rows when they come across them

- Offloads work from GC thread

- Ensures that frequently visited areas of the index are clean of expired rows

TX4:  Begin = 210
Oldest Active Hint = 200

| | | |
|---|---|---|

| 100 | 200 | 1 | | John | Smith | Kirkland |
|---|---|---|---|---|---|---|

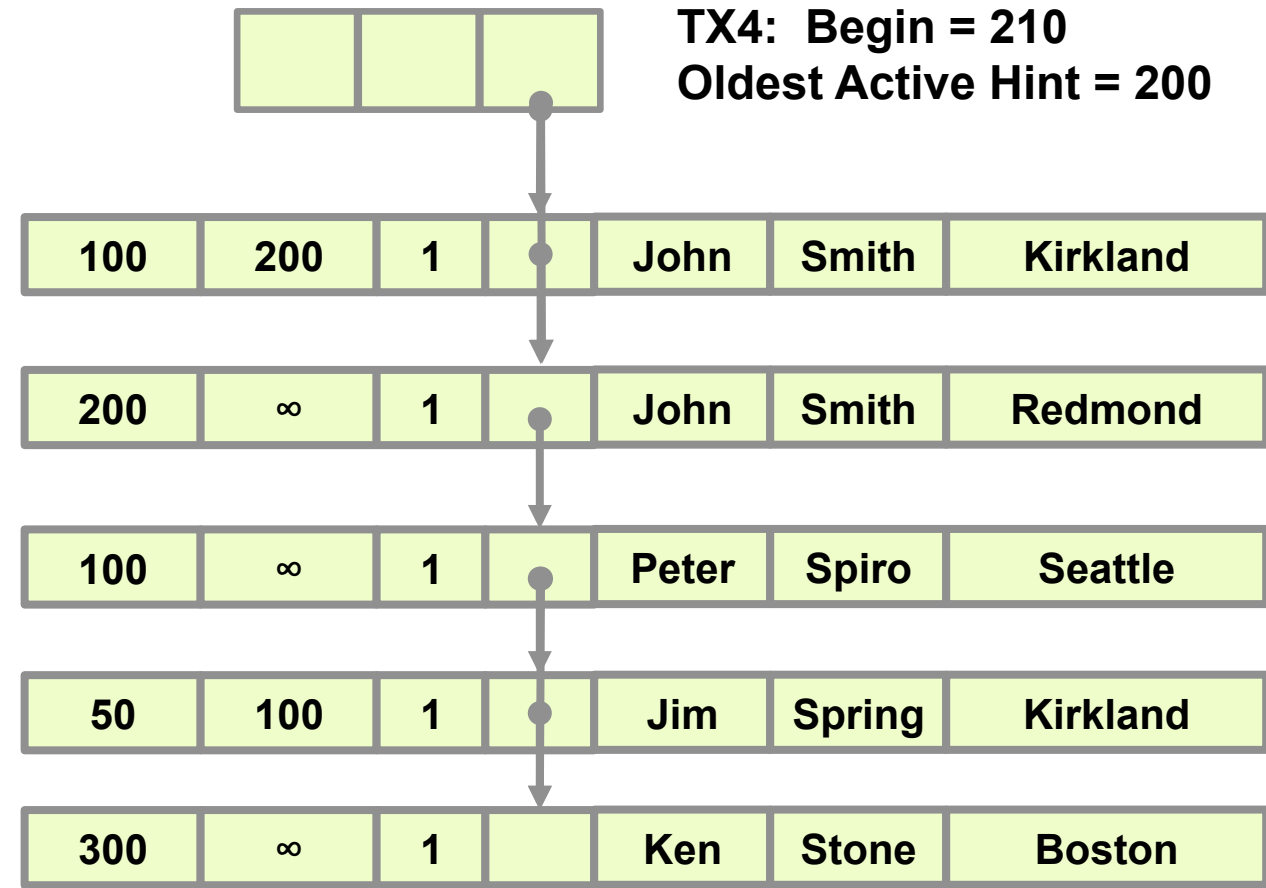| 200 | ∞ | 1 | | John | Smith | Redmond |
|---|---|---|---|---|---|---|

| 100 | ∞ | 1 | | Peter | Spiro | Seattle |
|---|---|---|---|---|---|---|

| 50 | 100 | 1 | | Jim | Spring | Kirkland |
|---|---|---|---|---|---|---|

| 300 | ∞ | 1 | | Ken | Stone | Boston |
|---|---|---|---|---|---|---|

**Client App**

TDS Handler and Session Management

No improvements in communication stack, parameter passing, result set generation

Hekaton Compiler

Parser, Catalog, Algebrizer, Optimizer

Proc/Plan cache for ad-hoc T-SQL and SPs

Natively Compiled SPs and Schema

10-30x more efficient

Interpreter for TSQL, query plans, expressions

Access Methods

Query Interop

In-Memory OLTP Engine for Memory_optimized Tables & Indexes

Reduced log bandwidth & contention. Log latency remains

Buffer Pool for Tables & Indexes

SQL Server.exe

Checkpoints are background sequential IO

Memory-optimized Table Filegroup

Transaction Log

Data Filegroup

**Key**

Existing SQL Component

Hekaton Component

Generated .dll

# Hekaton Engine's Scalability



## System throughput

| | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| Number of cores | 2 | 4 | 6 | 8 | 10 | 12 |
| SQL with contention | 984 | 1,363 | 1,645 | 1,876 | 2,118 | 2,312 |
| SQL without contention | 1,153 | 2,157 | 3,161 | 4,211 | 5,093 | 5,834 |
| Interop | 1,518 | 2,936 | 4,273 | 5,459 | 6,701 | 7,709 |
| Native | 7,078 | 13,892 | 20,919 | 26,721 | 32,507 | 36,375 |

Transactions Per Second

# Memory Optimized Table Limitations

## Optimized for high-throughput OLTP

- No XML and no CLR data types

## Optimized for in-memory

- Rows are at most 8060 bytes
- No Large Object (LOB) types like varchar(max)
- Durable memory-optimized tables are limited to 512 GB. (Non-durable tables have no size limit.)

## Scoping limitations

- No FOREIGN KEY and no CHECK constraints
- No schema changes (ALTER TABLE) – need to drop/recreate table
- No add/remove index – need to drop/recreate table
- No Computed Columns
- No Cross-Database Queries

# Natively Compiled Procedures Restrictions

- Not all operators/TSQLs are supported
- Only Nested Loop join, no TSQL MERGE or EXISTS, cursors, nested queries
- No CASE statement, CTEs, user-defined functions, UNION statement, DISTINCT statement
- Transaction isolation level
- SNAPSHOT, REPEATABLEREAD, and SERIALIZABLE
- READ COMMITTED and READ UNCOMMITED is not supported
- Cannot access disk-based tables
- No TEMPDB! Use In-Memory Table variables
- No automatic recompile on statistics changes
- Need to stop & start SQL or drop & create procedure

# Create Filegroup

FileGroup Container

```
CREATE DATABASE [Hekaton]
ON PRIMARY
(NAME = N'Hekaton_data', FILENAME = N'C:\Data\Data\Hekaton_data.mdf'),
FILEGROUP [Hekaton_InMemory] CONTAINS MEMORY_OPTIMIZED_DATA
(NAME = N'Hekaton_mem', FILENAME = N'C:\Data\Mem\Hekaton_Lun1.mdf')
LOG ON
(NAME = N'Hekaton_log', FILENAME = N'C:\Data\Log\Hekaton_log.ldf')


ALTER DATABASE [Hekaton]
ADD FILE (NAME = N'Hekaton_mem', FILENAME = N'C:\Data\Mem
\Hekaton_Lun2.mdf')
TO FILEGROUP [Hekaton_InMemory]
```

# Create Memory Optimized Table

```
CREATE TABLE [Customer] (
[CustomerID] INT NOT NULL
PRIMARY KEY NONCLUSTERED HASH WITH (BUCKET_COUNT = 1000000),
[AddressID] INT NOT NULL INDEX [IxName] HASH WITH (BUCKET_COUNT = 1000000),
[LName] NVARCHAR(250) COLLATE Latin1_General_100_BIN2 NOT NULL
INDEX [IXLName] NONCLUSTERED (LName)
)
WITH (MEMORY_OPTIMIZED = ON, DURABILITY = SCHEMA_AND_DATA);
```

Hash Index

Collation BIN2

Range Index

This table is memory optimized

This table is durable

# References

http://www.enterprisetech.com/2014/03/18/microsoft-turbocharges-transactions-hekaton-memory/

http://blogs.msdn.com/b/sqlcat/archive/2013/06/25/sql-server-in-memory-oltp-internals-overview-for-ctp1.aspx

http://blogs.msdn.com/b/arvindsh/archive/2013/07/03/sql-2014-in-memory-oltp-hekaton-training-videos-and-white-papers.aspx

https://www.simple-talk.com/sql/database-administration/exploring-in-memory-oltp-engine-(hekaton)-in-sql-server-2014-ctp1/

http://www.sqlskills.com/blogs/bobb/category/hekaton/

http://thomaslarock.com/2013/08/sql-server-2014-in-memory-oltp-hekaton-useful-links

http://blogs.msdn.com/b/carlnol/archive/2013/09/16/implementing-lob-storage-in-memory-optimized-tables.aspx

http://www.sqlpassion.at/archive/2013/08/12/extreme-transaction-processing-xtp-hekaton-the-solution-to-everything/

http://mattsql.wordpress.com/2013/07/08/in-memory-oltp-with-sql-server-2014/

http://research.microsoft.com/en-us/news/features/hekaton-122012.aspx

www.slideshare.net/RaviOkade/

# Thank you for your time!

# Q&A