

ExpoDB: Towards a Unified OLTP and OLAP over a Secure Platform

Mohammad Sadoghi

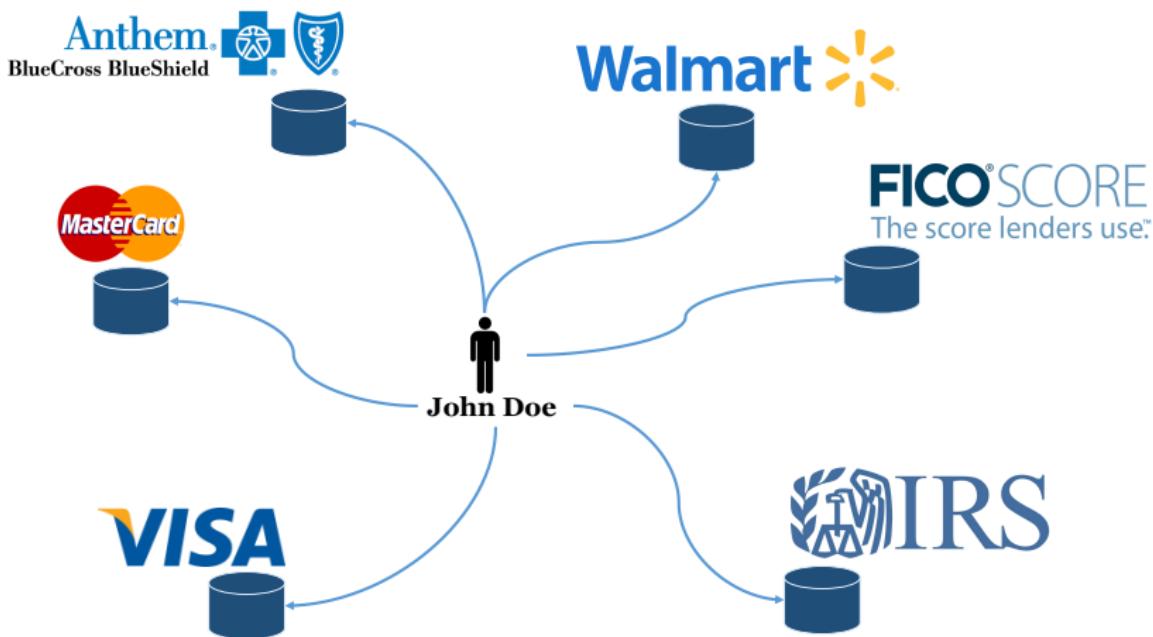
Exploratory Systems Lab
University of California, Davis

University of Waterloo

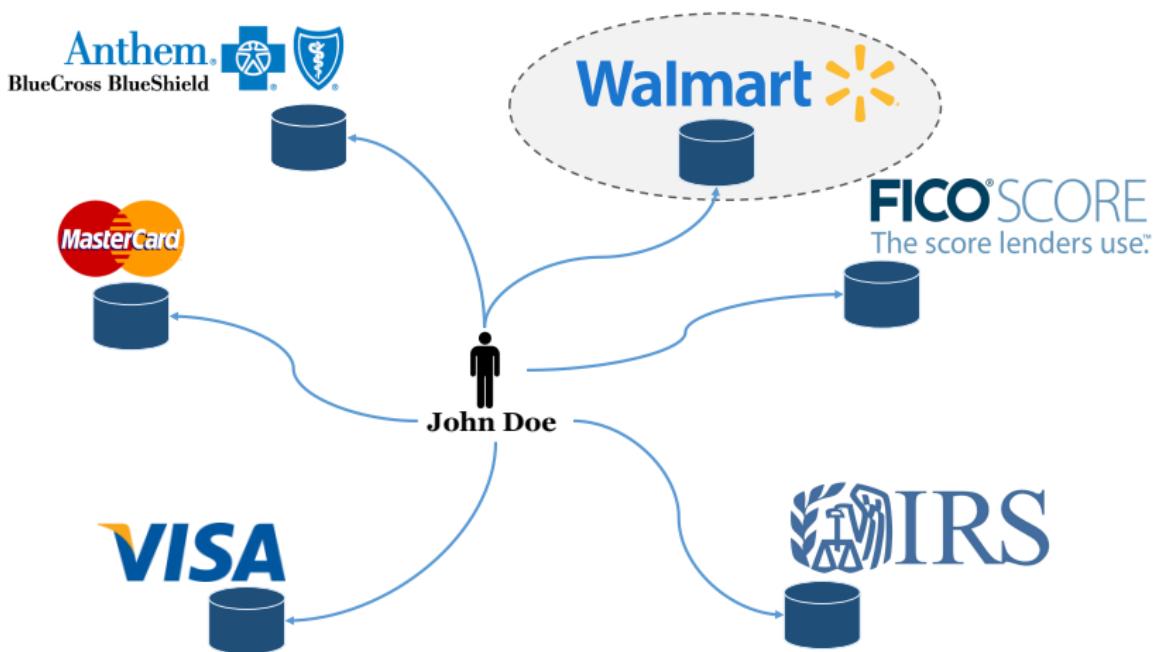
October 15, 2018



Data Management at Macroscale



Data Management at Macroscale

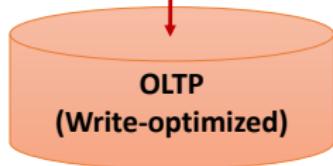


Data Management at Microscale: Volume & Velocity

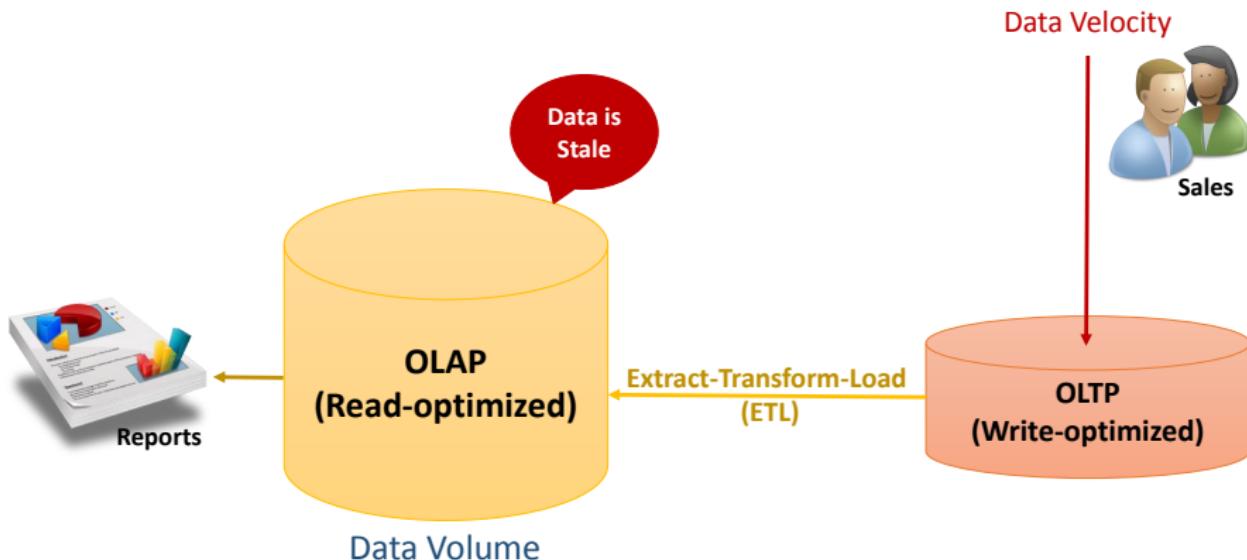
Data Velocity



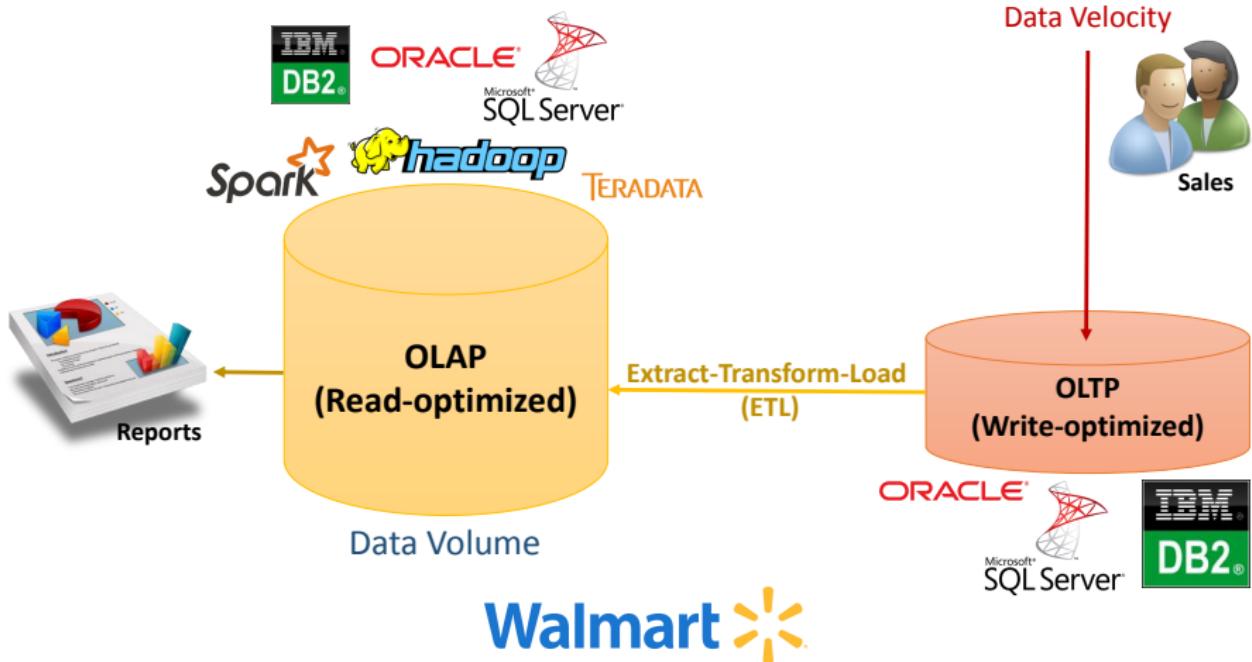
Sales



Data Management at Microscale: Volume & Velocity

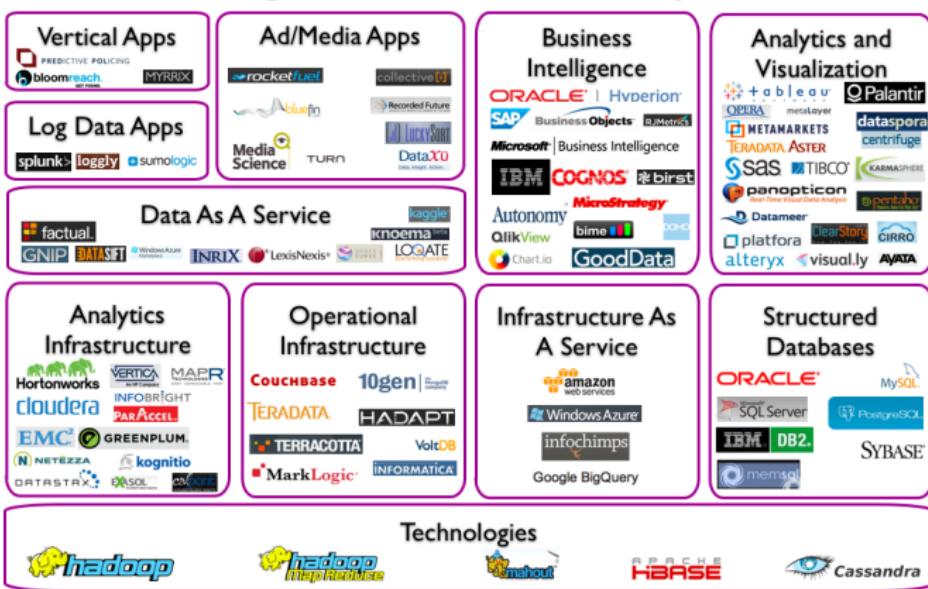


Data Management at Microscale: Volume & Velocity



One Size Does not Fit All As of 2012

Big Data Landscape

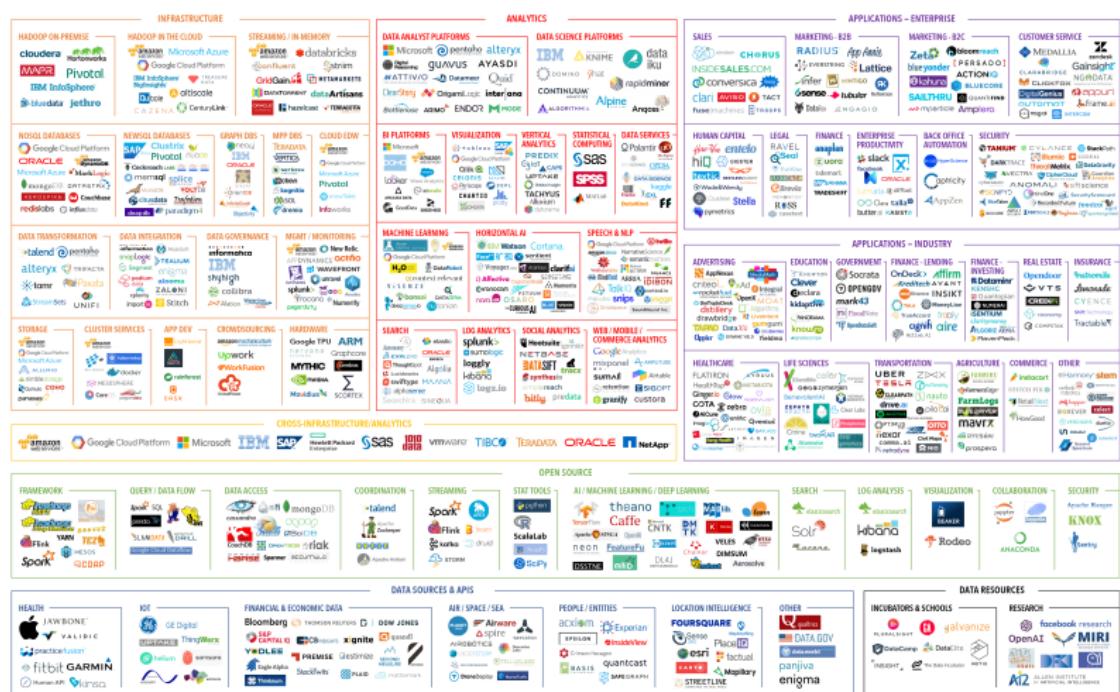


Copyright © 2012 Dave Feinleib

dave@vcddave.com

blogs.forbes.com/davefeinleib

One Size Does not Fit All As of 2017



142 Last run date: 15/12/2017

© Matt Turek (@mattturek), Jim Hao (@jimhao), & FirstMark (@firstmarkcap).

<http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3737201/>

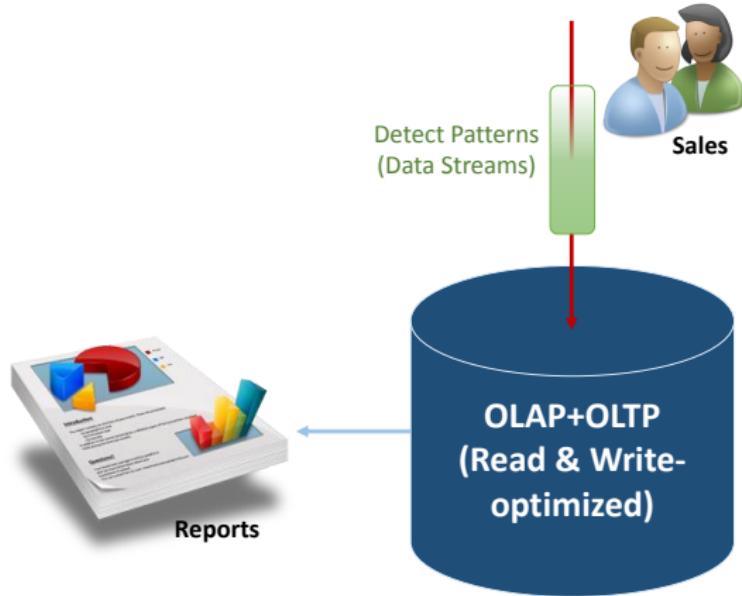
Data Management at Microscale: Volume & Velocity



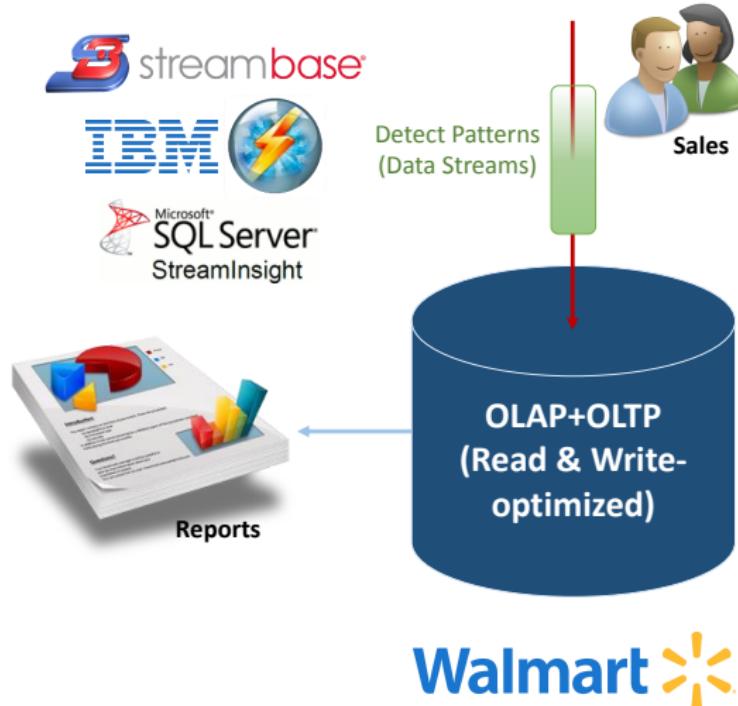
Data Management at Microscale: Volume & Velocity



Data Management at Microscale: Volume & Velocity



Data Management at Microscale: Volume & Velocity



Walmart

1 Data Management at Microscale

2 Data Management at Microscale

3 Data Velocity: Index Maintenance

4 Data Volume: MVCC Concurrency

5 Data Volume: Coordination-free Concurrency

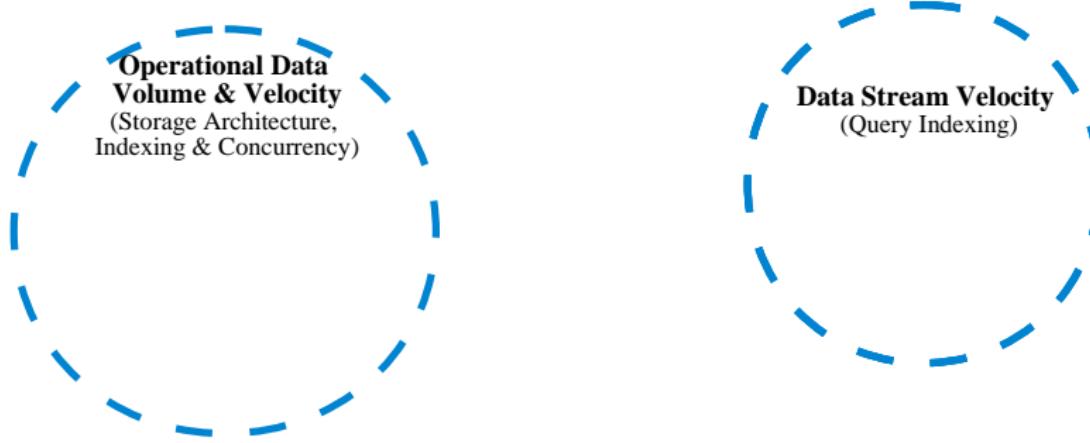
6 Combining Volume & Velocity: Lineage-based Storage Architecture

7 Data at Macroscale: Decentralized & Democratic Data Platform

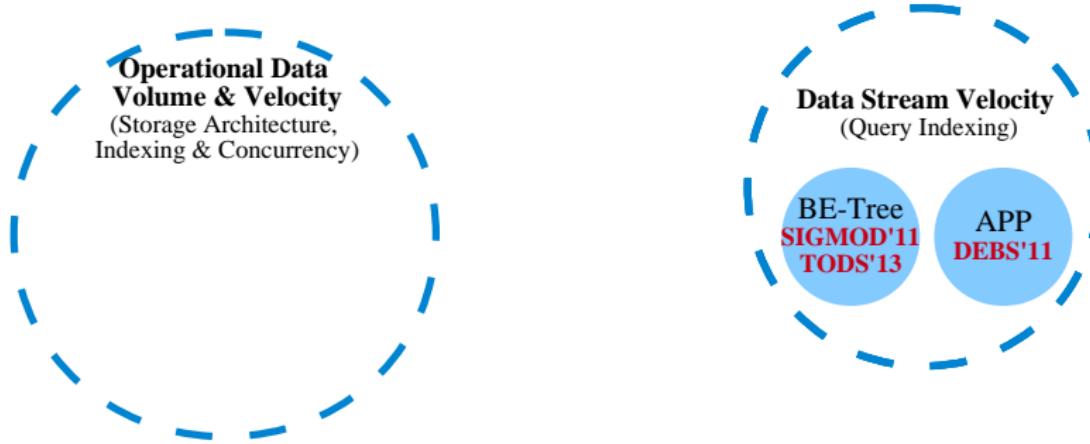
8 Conclusions

9 References

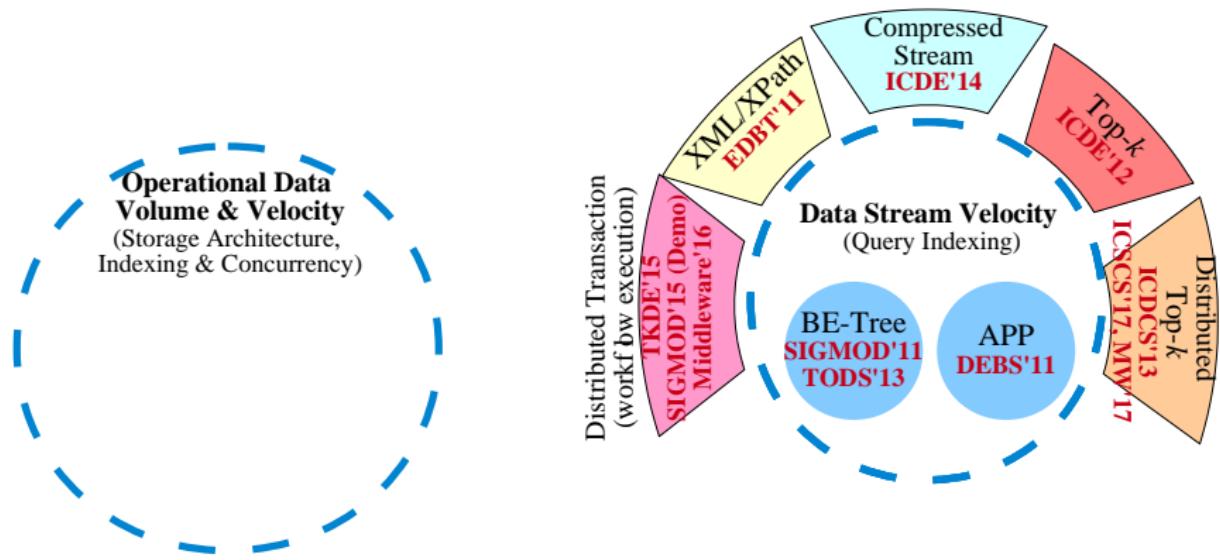
Big Picture



Big Picture



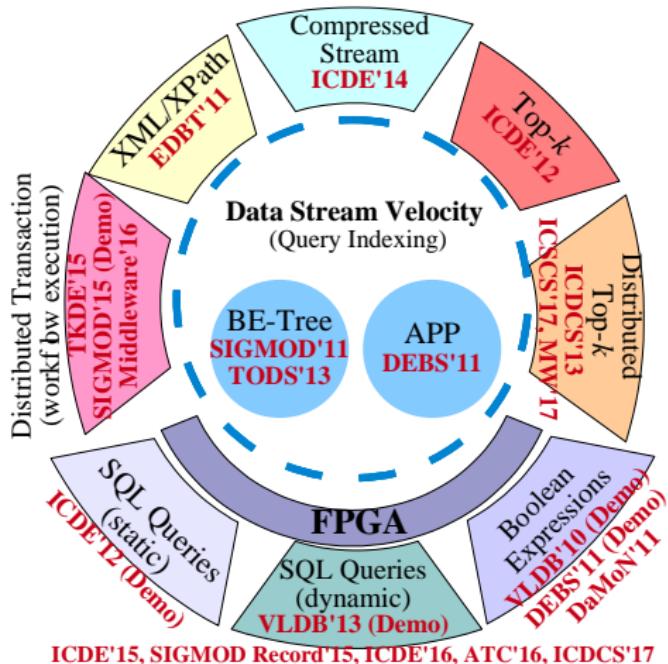
Big Picture



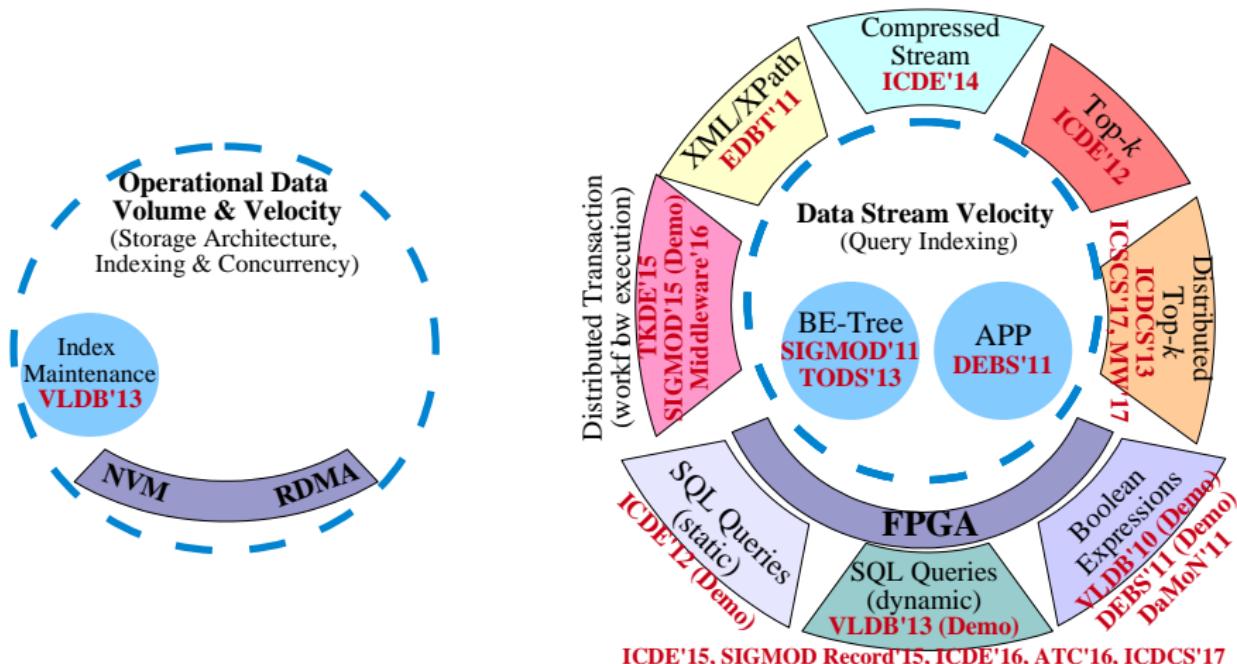
Big Picture



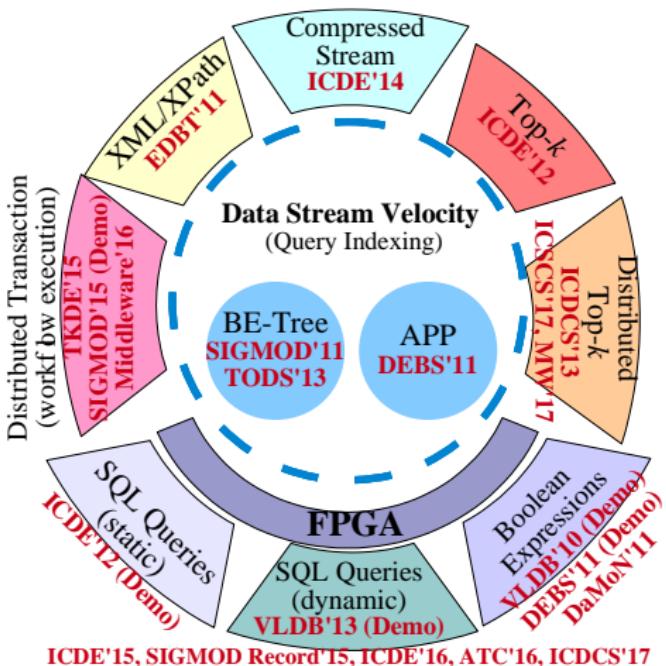
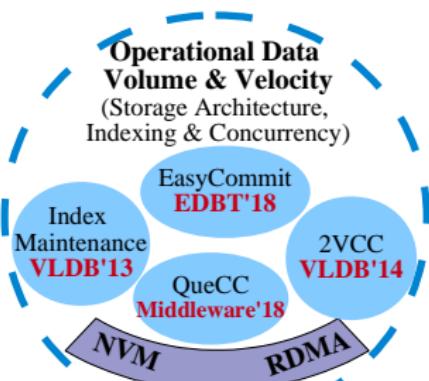
Operational Data Volume & Velocity
(Storage Architecture,
Indexing & Concurrency)



Big Picture

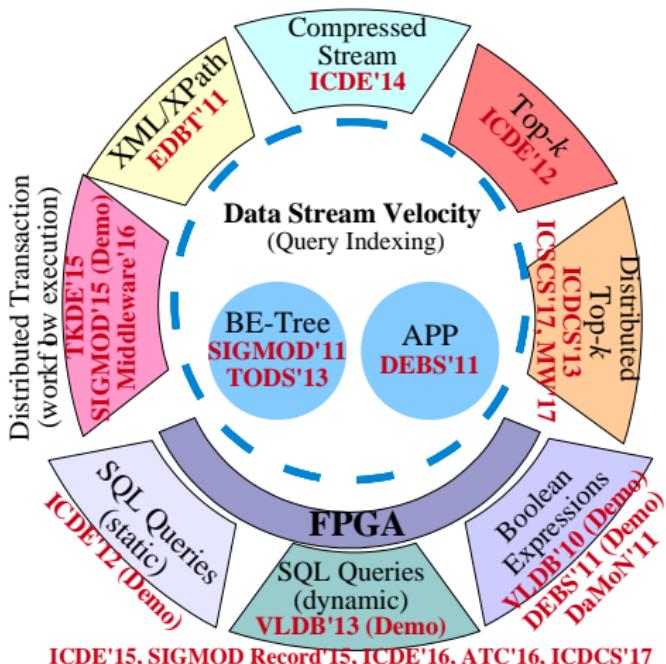
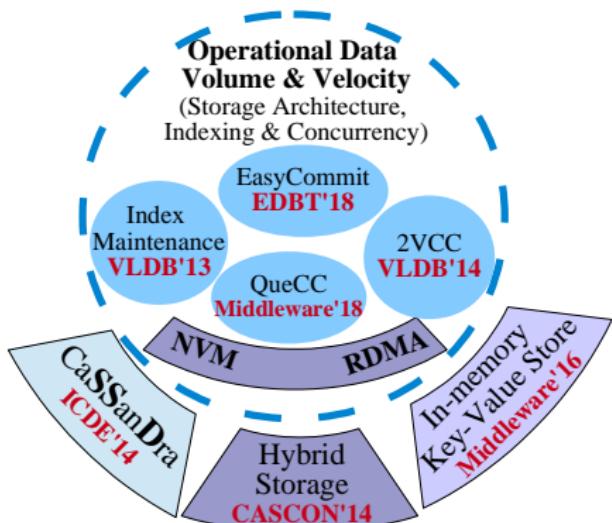


Big Picture

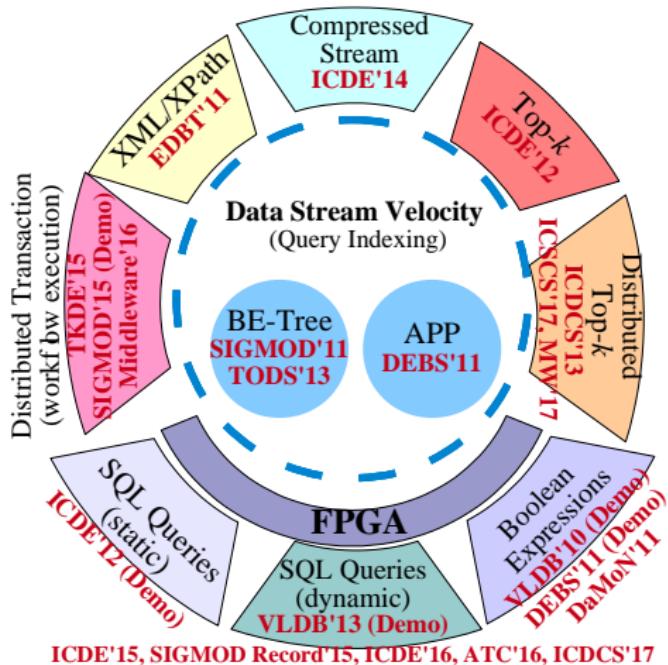
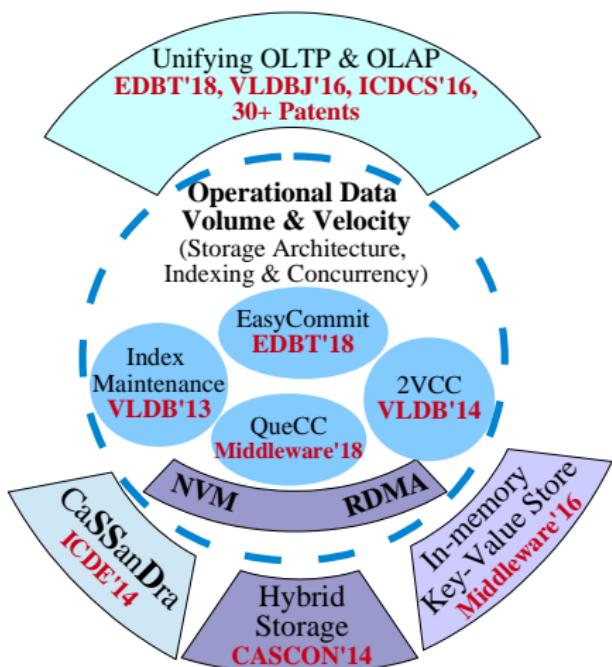


ICDE'15, SIGMOD Record'15, ICDE'16, ATC'16, ICDCS'17

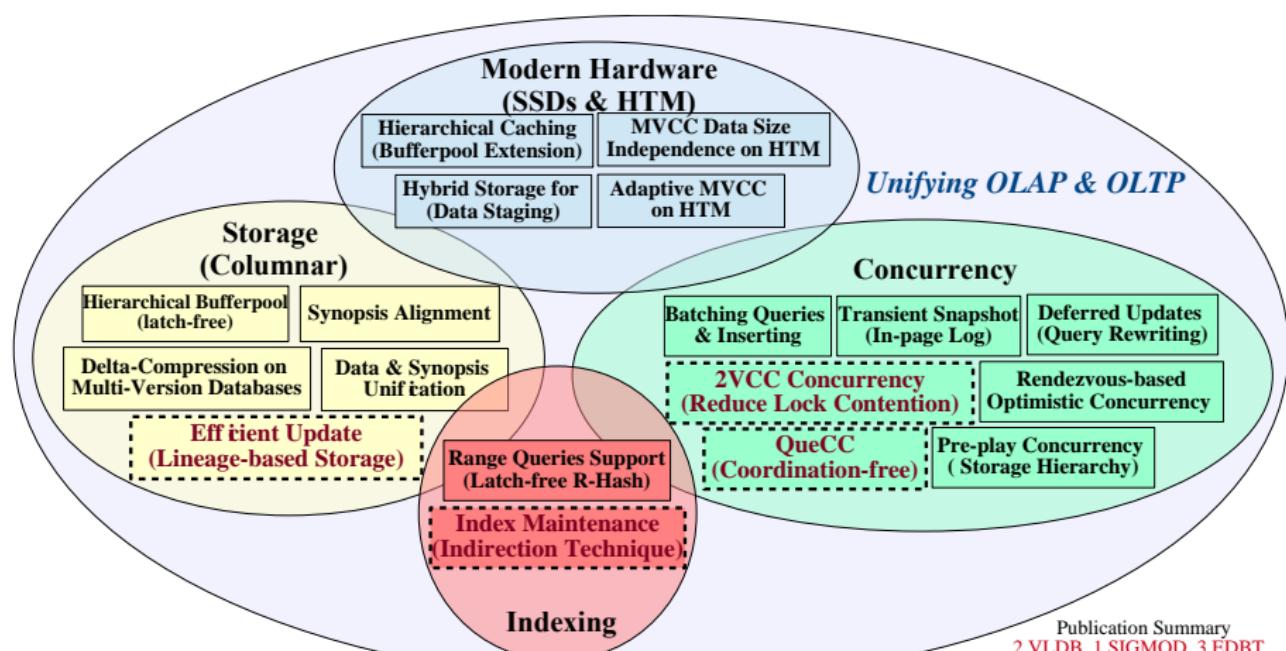
Big Picture



Big Picture



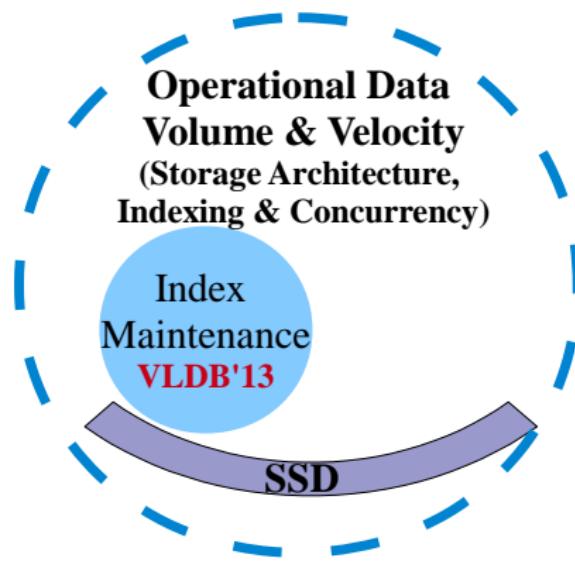
Deep Dive: Unifying OLAP & OLTP



Publication Summary
2 VLDB, 1 SIGMOD, 3 EDBT,
1 VLDBJ, 1 ICDE, 1 ICDCS,
2 Middleware, 30+ Patents

- 1 Data Management at Microscale
- 2 Data Management at Microscale
- 3 Data Velocity: Index Maintenance
- 4 Data Volume: MVCC Concurrency
- 5 Data Volume: Coordination-free Concurrency
- 6 Combining Volume & Velocity: Lineage-based Storage Architecture
- 7 Data at Macroscale: Decentralized & Democratic Data Platform
- 8 Conclusions
- 9 References

Extending Storage Hierarchy with Indirection Layer



Reducing Index maintenance: Velocity Dimension

Observed Trends

In the absence of in-place updates in operational multi-version databases, the cost of index maintenance becomes a major obstacle to cope with data velocity.

Reducing Index maintenance: Velocity Dimension

Observed Trends

In the absence of in-place updates in operational multi-version databases, the cost of index maintenance becomes a major obstacle to cope with data velocity.

Extending storage hierarchy (using fast non-volatile memory) with *an extra level of indirection* in order to

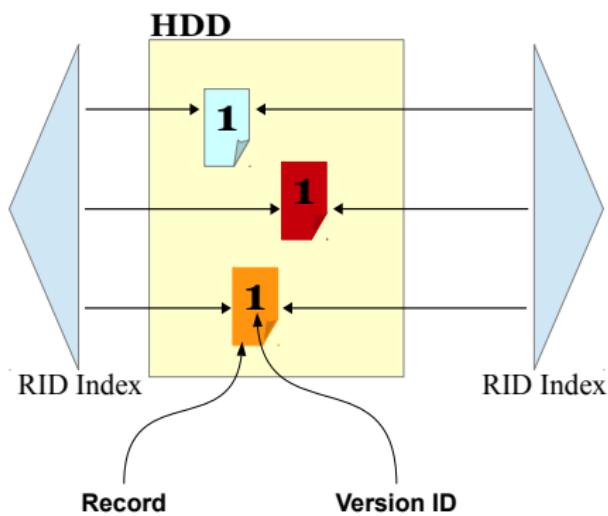
Reducing Index maintenance: Velocity Dimension

Observed Trends

In the absence of in-place updates in operational multi-version databases, the cost of index maintenance becomes a major obstacle to cope with data velocity.

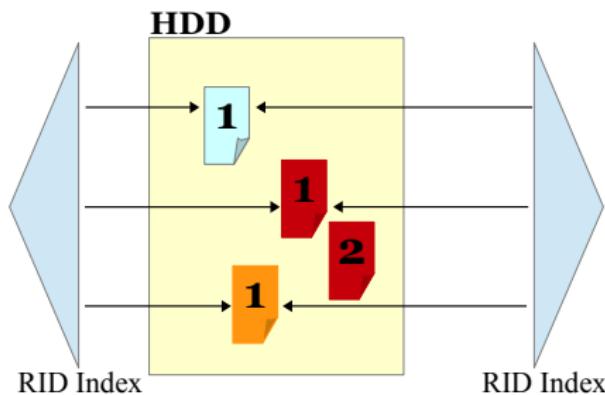
Extending storage hierarchy (using fast non-volatile memory) with *an extra level of indirection* in order to
Decouple Logical and Physical Locations of Records to Reduce Index Maintenance

Traditional Multi-version Indexing: Updating Records



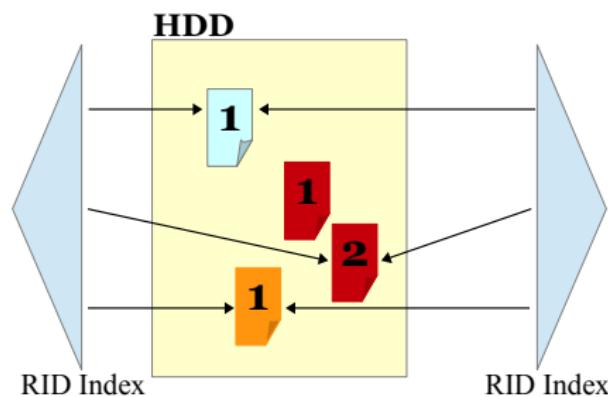
Updating random leaf pages

Traditional Multi-version Indexing: Updating Records



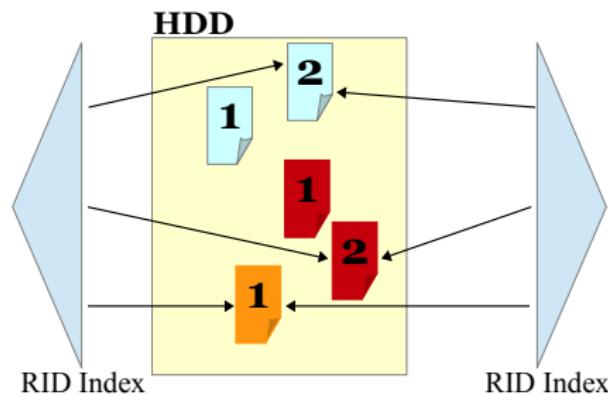
Updating random leaf pages

Traditional Multi-version Indexing: Updating Records



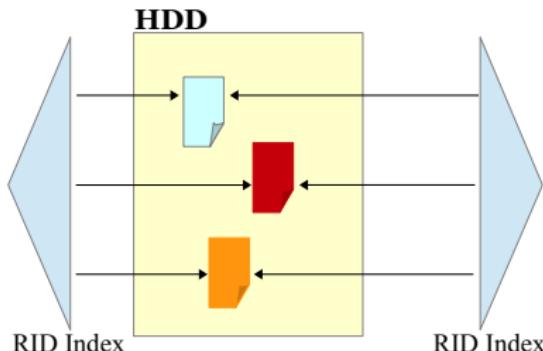
Updating random leaf pages

Traditional Multi-version Indexing: Updating Records

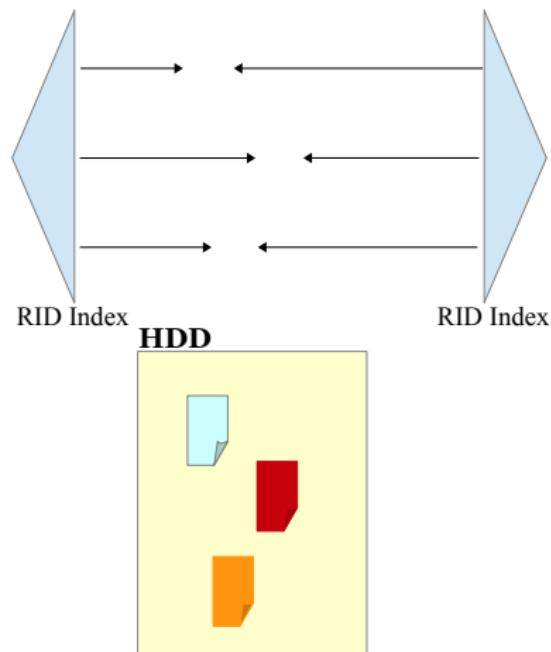


Updating random leaf pages

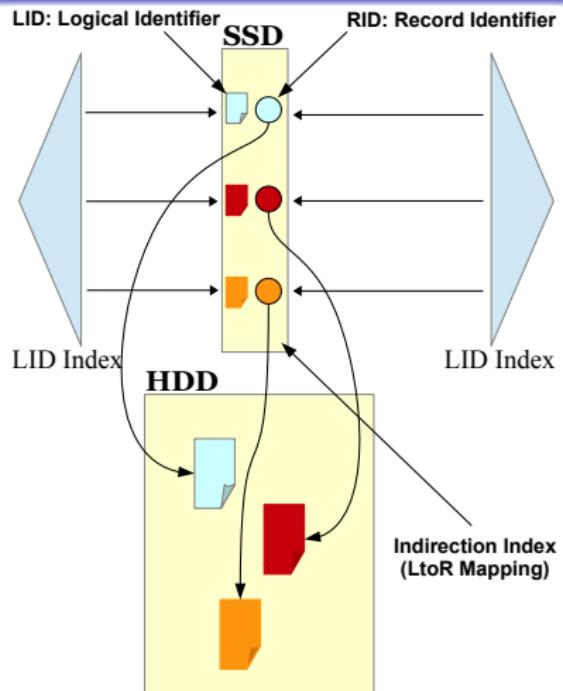
Indirection Indexing: Updating Records



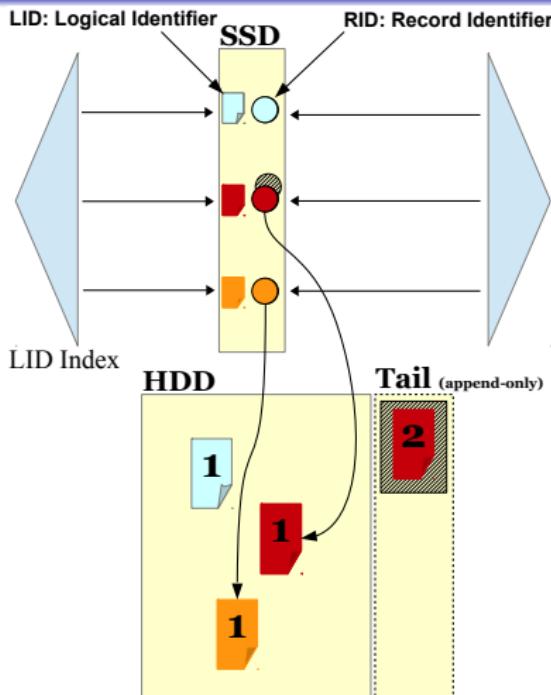
Indirection Indexing: Updating Records



Indirection Indexing: Updating Records

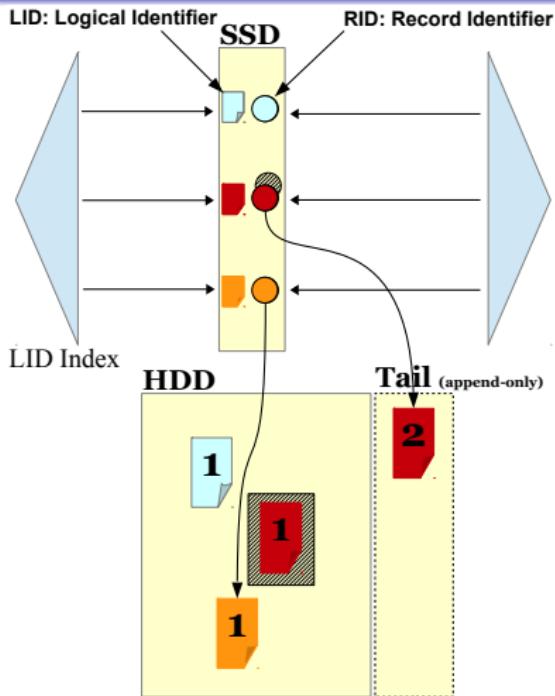


Indirection Indexing: Updating Records



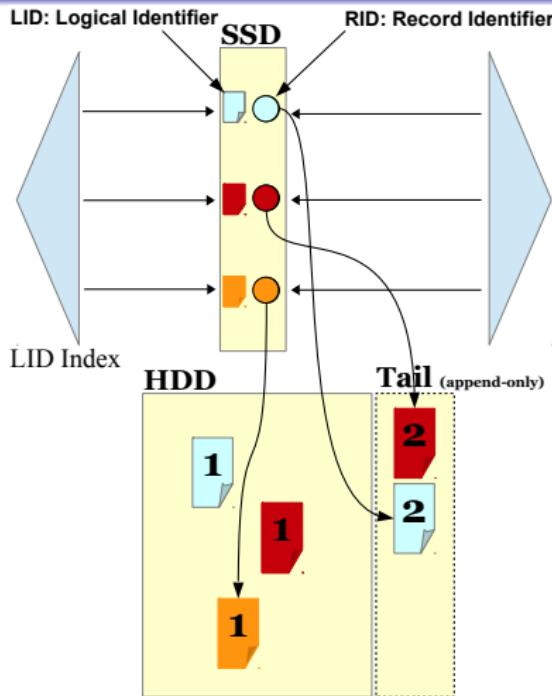
Eliminating random leaf-page updates

Indirection Indexing: Updating Records



Eliminating random leaf-page updates

Indirection Indexing: Updating Records



Eliminating random leaf-page updates

Analytical & Experimental Evaluations

Indirection Time Complexity Analysis

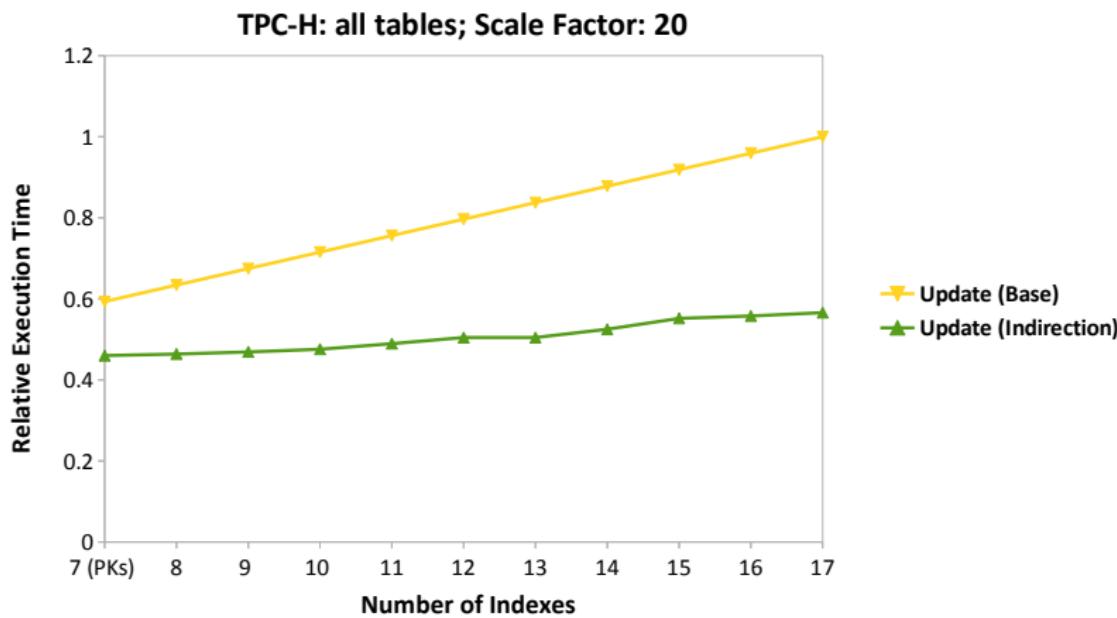
Legend	
K	Number of indexes
LB	LIDBlock size
M	Number of matching records

Method	Type	Imm. SSD	Def. SSD	Imm. HDD	Def. HDD
Base	Deletion	0	0	$2 + K$	$\leq 1 + K$
	Single-attr. update	0	0	$3 + K$	$\leq 2 + K$
	Insertion	0	0	$1 + K$	$\leq 1 + K$
	Search Uniq.	0	0	2	0
	Search Mult.	0	0	$1 + M$	0
Indirection	Deletion	2	0	2	≤ 3
	Single-attr. update	2	0	4	≤ 3
	Insertion	$2 + 2K$	$2K/LB$	1	$\leq 1 + 2K/LB$
	Search Uniq.	2	0	2	0
	Search Mult.	$1 + M$	0	$1 + M$	0

Experimental Setting

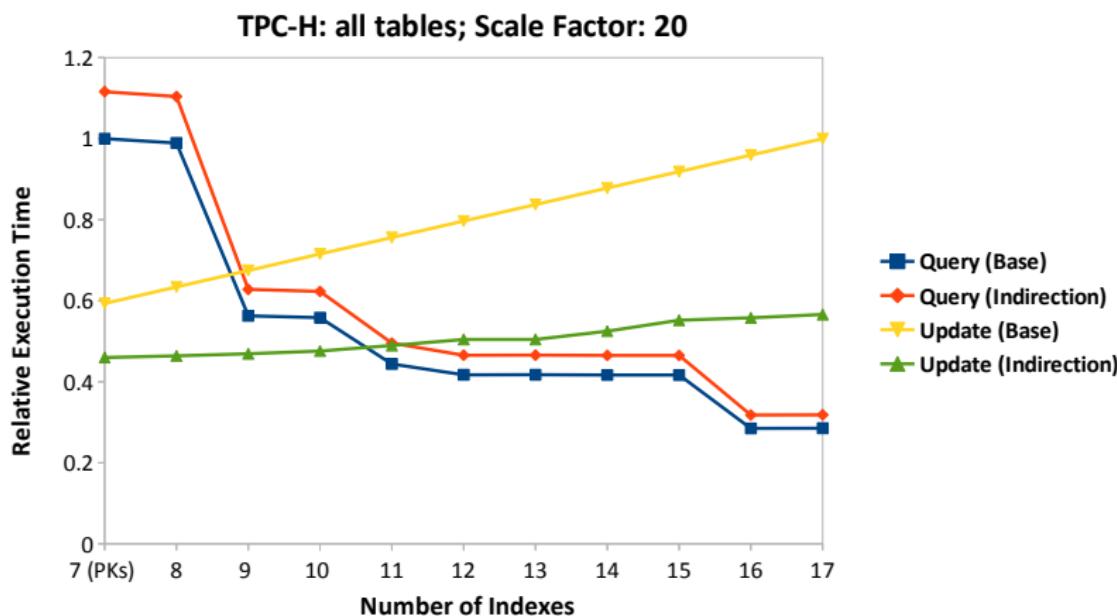
- Hardware:
 - (2 × 8-core) Intel(R) Xeon(R) CPU E7-4820 @ 2.00GHz, 32GB, 2 × HDD, SSD Fusion-io
- Software:
 - Database: IBM DB2 9.7
 - Prototyped in a commercial proprietary database
 - Prototyped in Apache Spark by UC Berkeley
 - LIBGist v.1.0: Generalized Search Tree C++ Library by UC Berkeley (**5K LOC**)
(Predecessor of Generalized Search Tree (GiST) access method for PostgreSQL)
 - **LIBGist^{mv} Prototype:** Multi-version Generalized Search Tree C++ Library over LIBGist supporting Indirection/LIDBlock/DeltaBlock (**3K LOC**)
- Data:
 - TPC-H benchmark
 - Microsoft Hekaton micro benchmark

Indirection: Effect of Indexes in Operational Data Stores



Substantially improving the update time ...

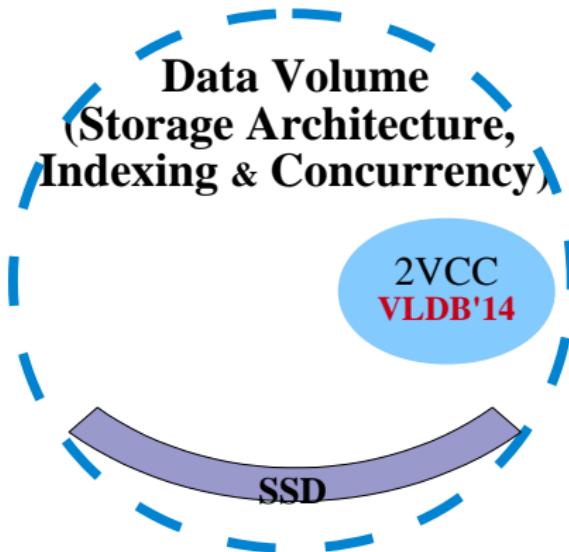
Indirection: Effect of Indexes in Operational Data Stores



... Consequently affording more indexes and significantly reducing the query time

- 1 Data Management at Microscale
- 2 Data Management at Microscale
- 3 Data Velocity: Index Maintenance
- 4 Data Volume: MVCC Concurrency
- 5 Data Volume: Coordination-free Concurrency
- 6 Combining Volume & Velocity: Lineage-based Storage Architecture
- 7 Data at Macroscale: Decentralized & Democratic Data Platform
- 8 Conclusions
- 9 References

Introducing Multi-version Concurrency Control



Generalized Concurrency Control: Volume Dimension

Observed Trends

In operational multi-version databases, there is a tremendous opportunity to avoid clashes between readers (scanning a large volume of data) and writers (frequent updates).

Generalized Concurrency Control: Volume Dimension

Observed Trends

In operational multi-version databases, there is a tremendous opportunity to avoid clashes between readers (scanning a large volume of data) and writers (frequent updates).

Introducing a (latch-free) *two-version concurrency control (2VCC)* by extending indirection mapping (i.e., central coordination mechanism) and exploiting existing two-phase locking (2PL) in order to

Generalized Concurrency Control: Volume Dimension

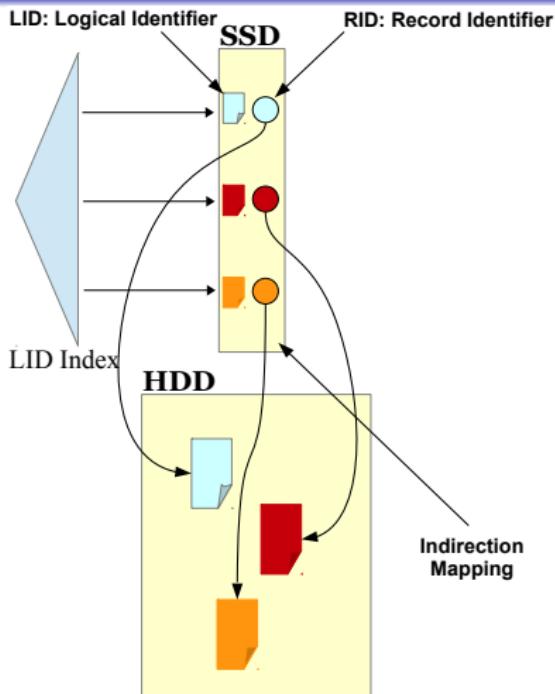
Observed Trends

In operational multi-version databases, there is a tremendous opportunity to avoid clashes between readers (scanning a large volume of data) and writers (frequent updates).

Introducing a (latch-free) *two-version concurrency control (2VCC)* by extending indirection mapping (i.e., central coordination mechanism) and exploiting existing two-phase locking (2PL) in order to

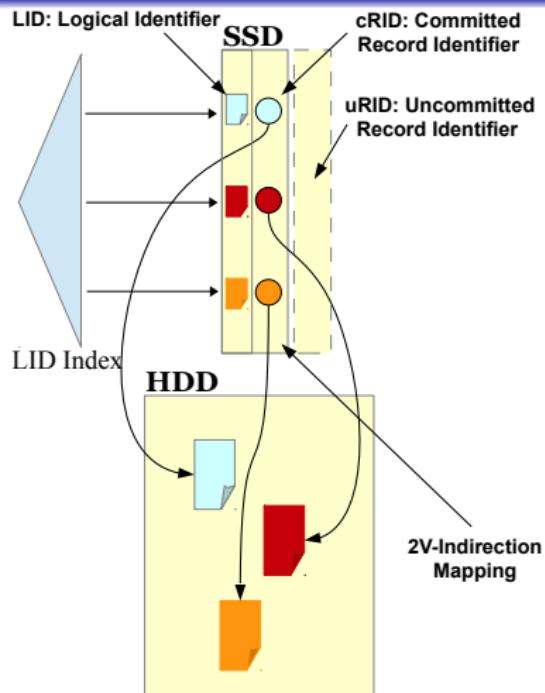
Decouple Readers/Writers to Reduce Contention
(Pessimistic and Optimistic Concurrency Control Coexistence)

2V-Indirection Indexing: Updating Records



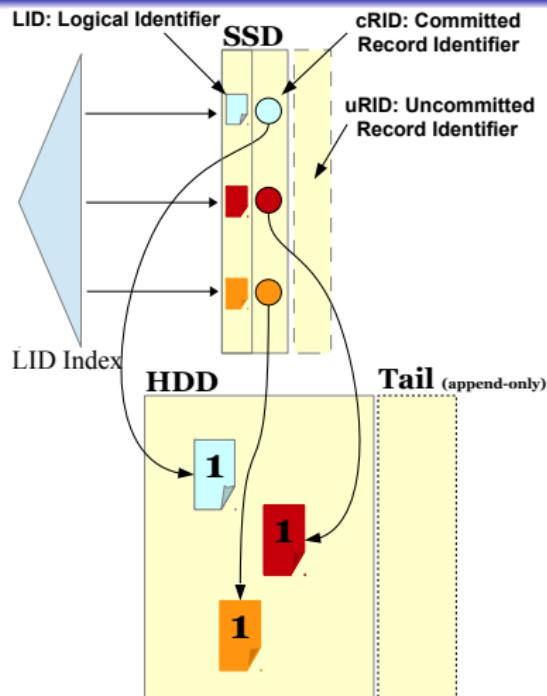
Recap: Indirection technique for reducing index maintenance

2V-Indirection Indexing: Updating Records



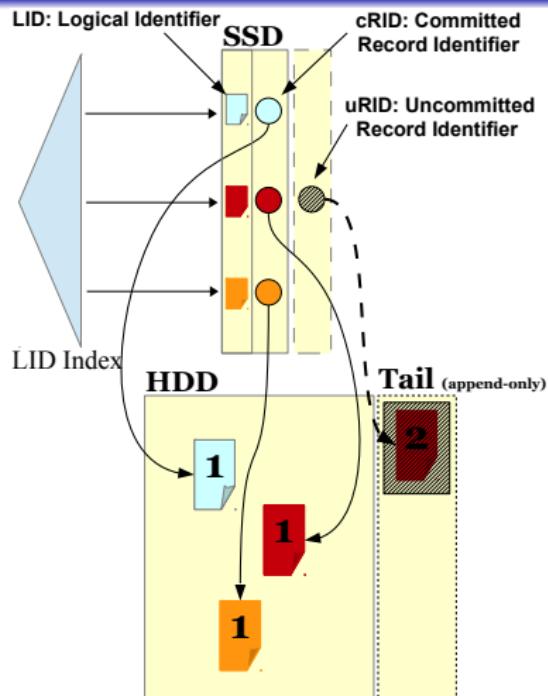
Extending the indirection to committed/uncommitted records

2V-Indirection Indexing: Updating Records



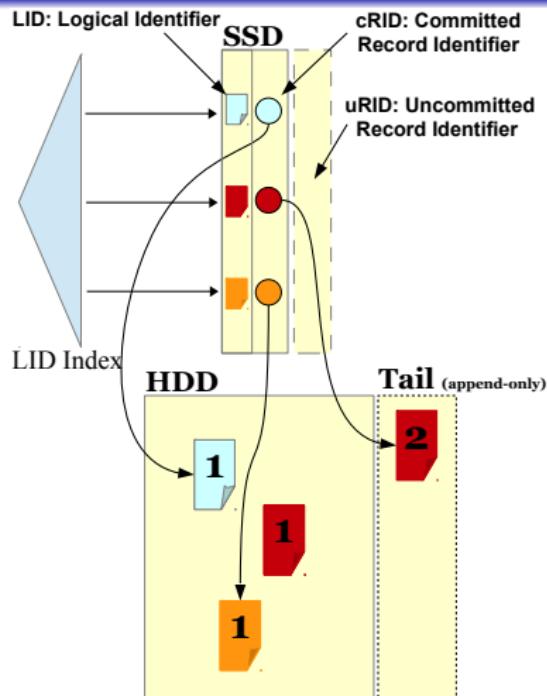
Extending the indirection to committed/uncommitted records

2V-Indirection Indexing: Updating Records



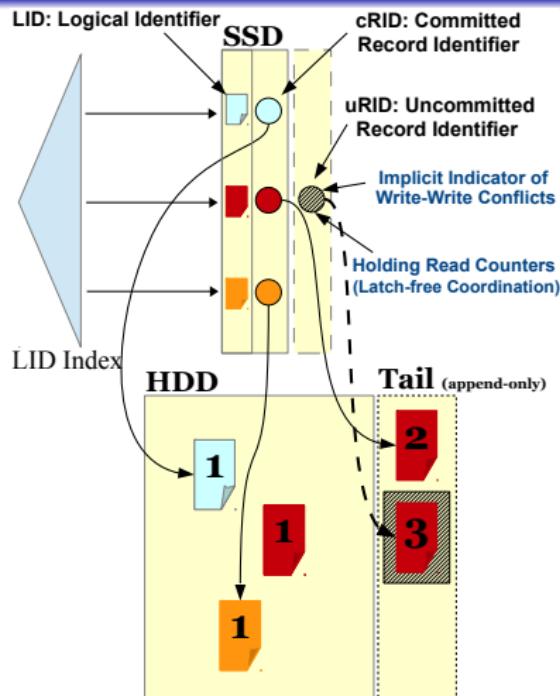
Decoupling readers/writers to eliminate contention

2V-Indirection Indexing: Updating Records



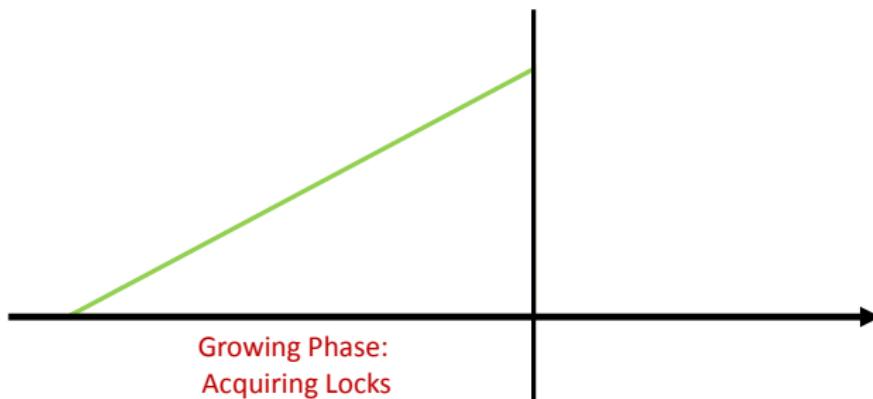
Decoupling readers/writers to eliminate contention

2V-Indirection Indexing: Updating Records



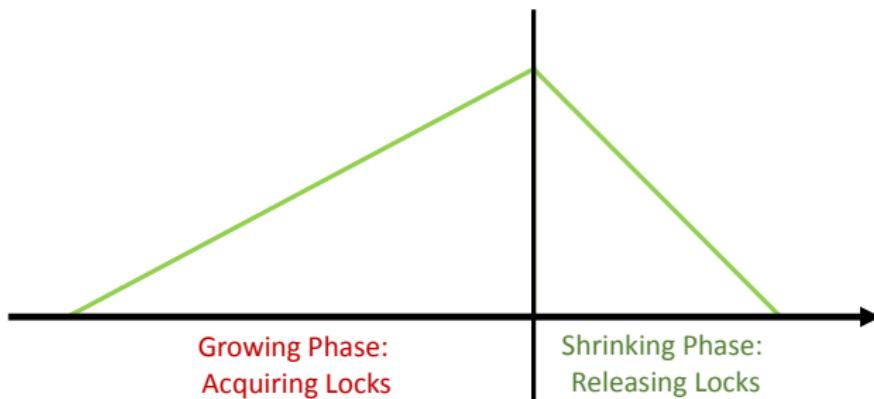
Decoupling readers/writers to eliminate contention

Overview of Two-version Concurrency Control Protocol



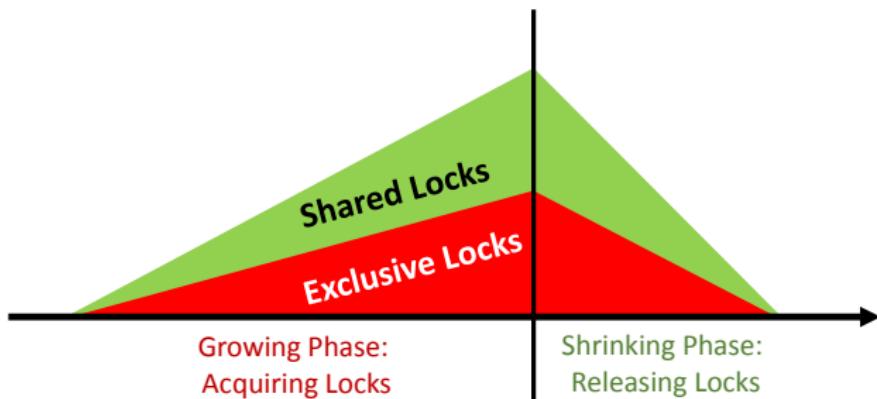
Two-phase locking (2PL) consisting of growing and shrinking phases

Overview of Two-version Concurrency Control Protocol



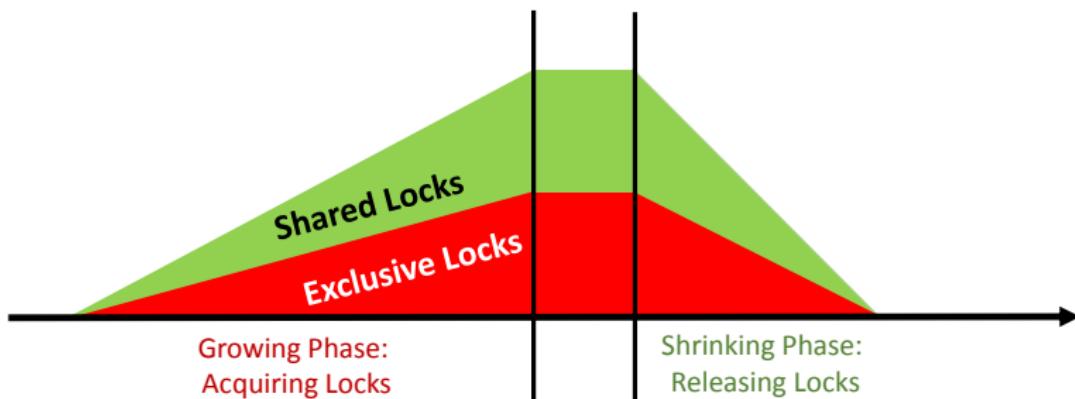
Two-phase locking (2PL) consisting of growing and shrinking phases

Overview of Two-version Concurrency Control Protocol



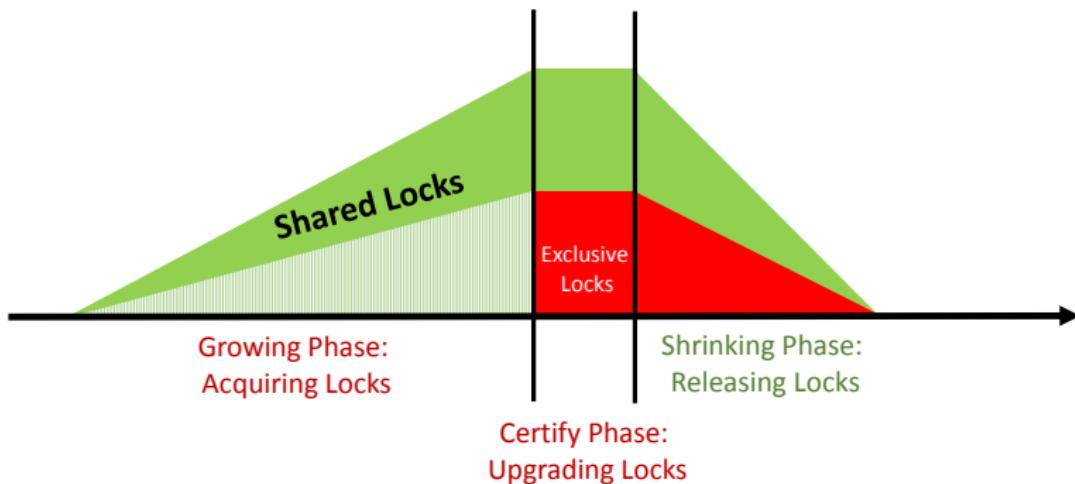
Two-phase locking (2PL) consisting of growing and shrinking phases

Overview of Two-version Concurrency Control Protocol



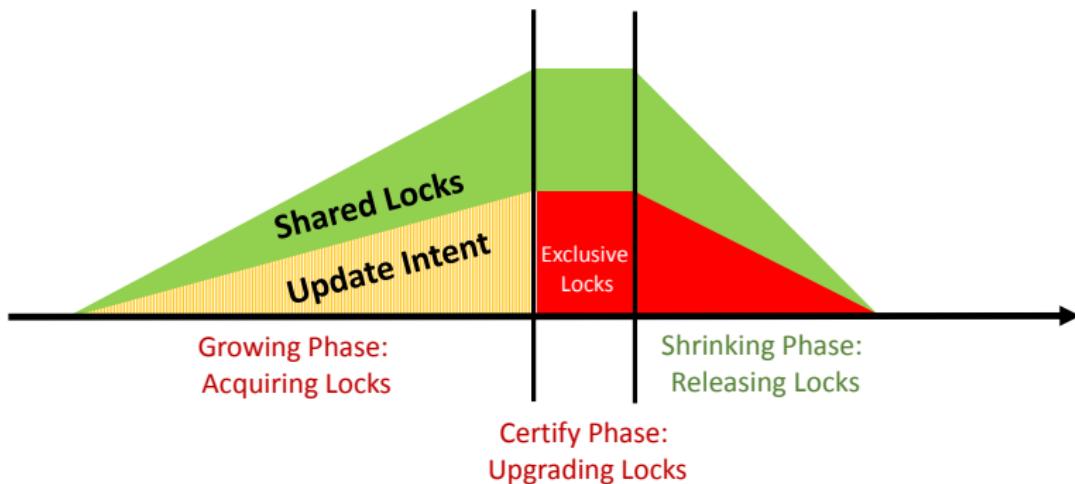
Extending 2PL with certify phase

Overview of Two-version Concurrency Control Protocol



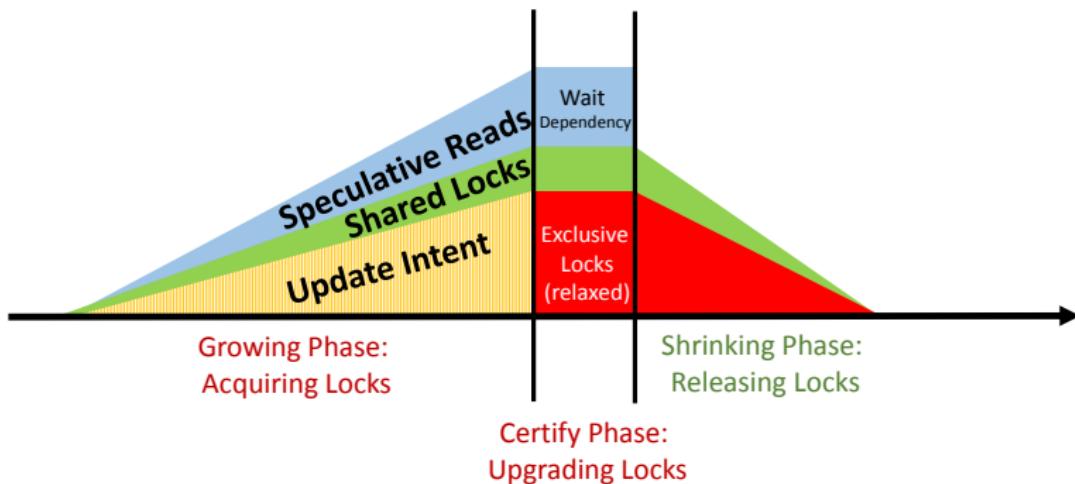
Exclusive locks held for shorter period (inherently optimistic)

Overview of Two-version Concurrency Control Protocol



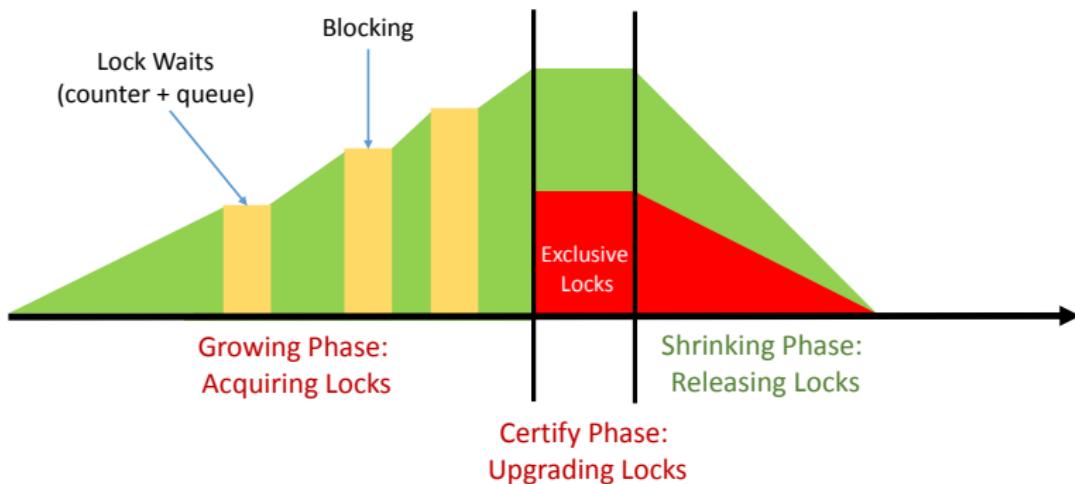
Exclusive locks held for shorter period (inherently optimistic)

Overview of Two-version Concurrency Control Protocol



Relaxed exclusive locks to allow speculative reads (increased optimism)

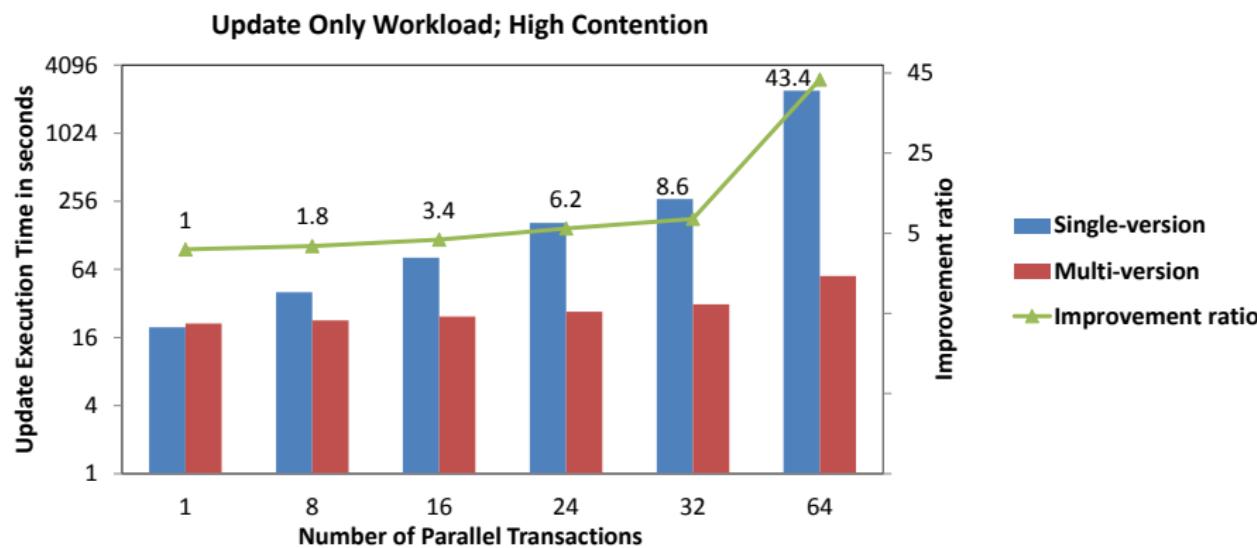
Overview of Two-version Concurrency Control Protocol



Trade-offs between blocking (i.e., locks) vs. non-blocking (i.e., read counters)

Experimental Analysis

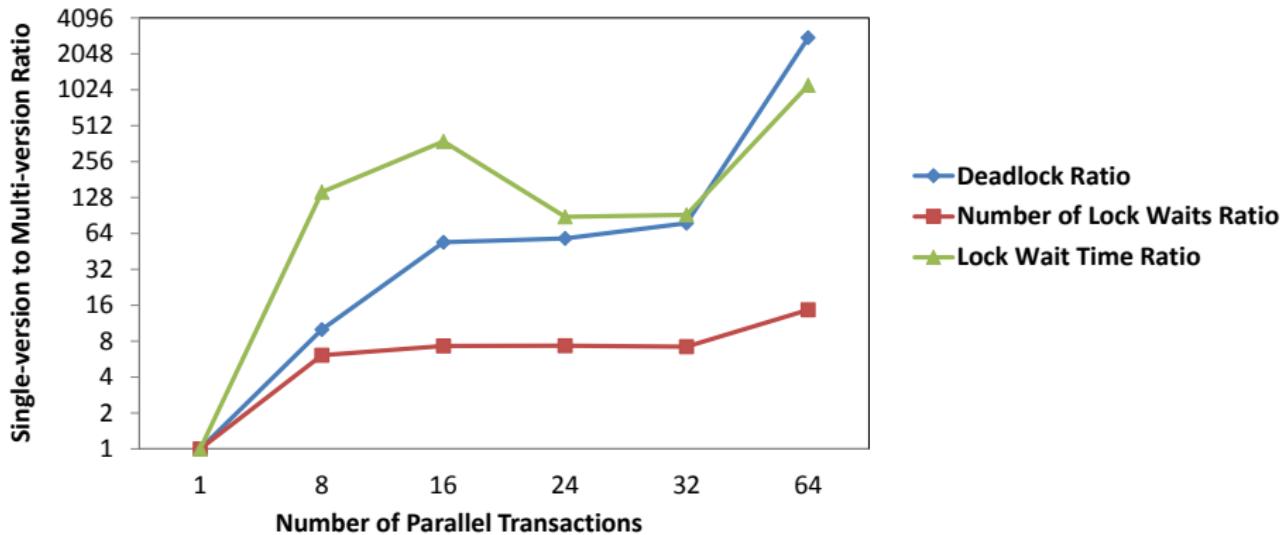
2VCC: Effect of Parallel Update Transactions



Substantial gain by reducing the read/write contention & using non-blocking operations

2VCC: Effect of Parallel Update Transactions

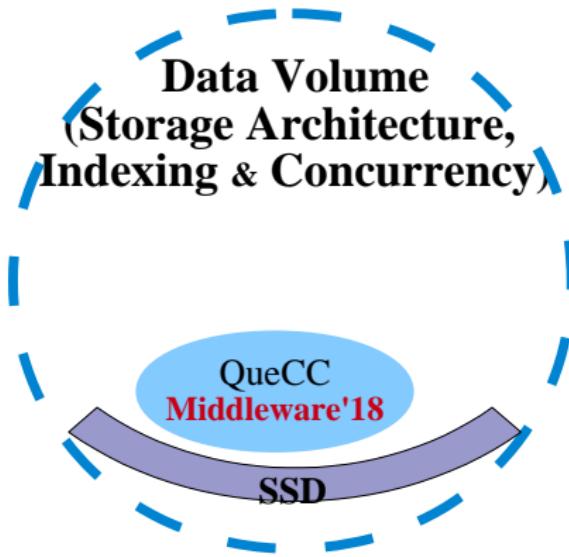
Lock Statistics Comparison; High Contention



Substantial gain by reducing the read/write contention & using non-blocking operations

- 1 Data Management at Microscale
- 2 Data Management at Microscale
- 3 Data Velocity: Index Maintenance
- 4 Data Volume: MVCC Concurrency
- 5 Data Volume: Coordination-free Concurrency
- 6 Combining Volume & Velocity: Lineage-based Storage Architecture
- 7 Data at Macroscale: Decentralized & Democratic Data Platform
- 8 Conclusions
- 9 References

Introducing Coordination-free Concurrency Control



Confrontation-free Concurrency Control

Observed Trends

In operational databases, the use of pre-compiled stored procedures is predominant. There is a tremendous opportunity to exploit transaction prior knowledge to eliminate the need for coordination.

Confrontation-free Concurrency Control

Observed Trends

In operational databases, the use of pre-compiled stored procedures is predominant. There is a tremendous opportunity to exploit transaction prior knowledge to eliminate the need for coordination.

Is it possible to have concurrent execution over shared data (not limited to partitionable workloads) without having any concurrency controls?

Confrontation-free Concurrency Control

Observed Trends

In operational databases, the use of pre-compiled stored procedures is predominant. There is a tremendous opportunity to exploit transaction prior knowledge to eliminate the need for coordination.

Is it possible to have concurrent execution over shared data (not limited to partitionable workloads) without having any concurrency controls?

Introducing a *queue-oriented, control-free concurrency* (QueCC) based on two parallel & independent phases of priority-driven planning & execution.

Confrontation-free Concurrency Control

Observed Trends

In operational databases, the use of pre-compiled stored procedures is predominant. There is a tremendous opportunity to exploit transaction prior knowledge to eliminate the need for coordination.

Is it possible to have concurrent execution over shared data (not limited to partitionable workloads) without having any concurrency controls?

Introducing a *queue-oriented, control-free concurrency* (QueCC) based on two parallel & independent phases of priority-driven planning & execution.

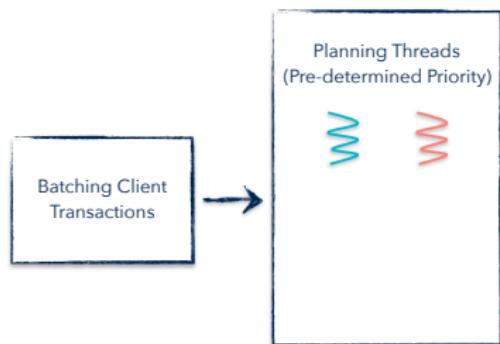
Execution and Synchronization Decoupling

Queue-oriented, Control-free Concurrency (QueCC)

Batching Client Transactions

Execution & Synchronization Decoupling: Deterministic priority-based planning followed by queue-oriented execution

Queue-oriented, Control-free Concurrency (QueCC)



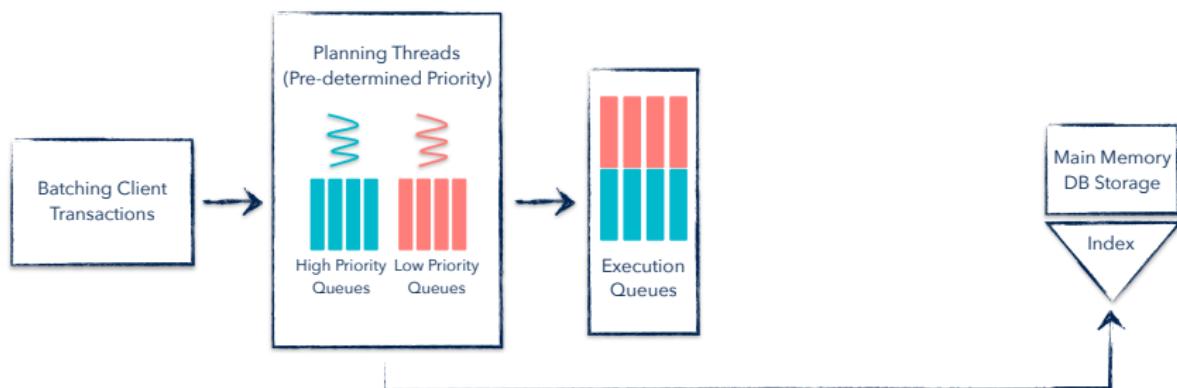
Execution & Synchronization Decoupling: Deterministic priority-based planning followed by queue-oriented execution

Queue-oriented, Control-free Concurrency (QueCC)



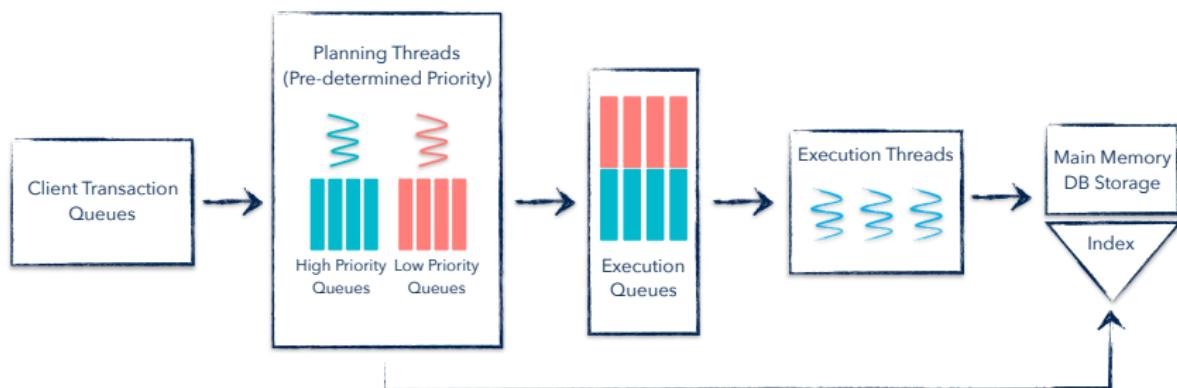
Execution & Synchronization Decoupling: Deterministic priority-based planning followed by queue-oriented execution

Queue-oriented, Control-free Concurrency (QueCC)



Execution & Synchronization Decoupling: Deterministic priority-based planning followed by queue-oriented execution

Queue-oriented, Control-free Concurrency (QueCC)

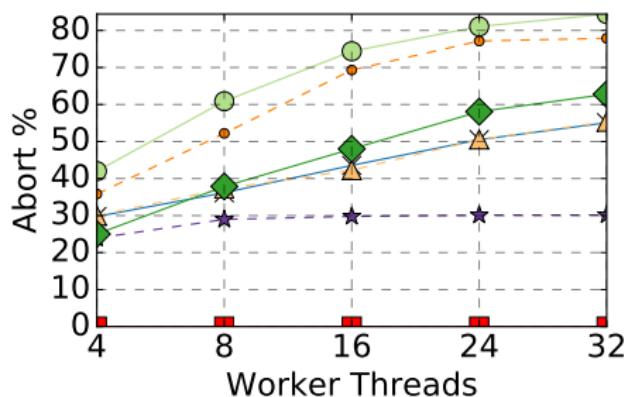
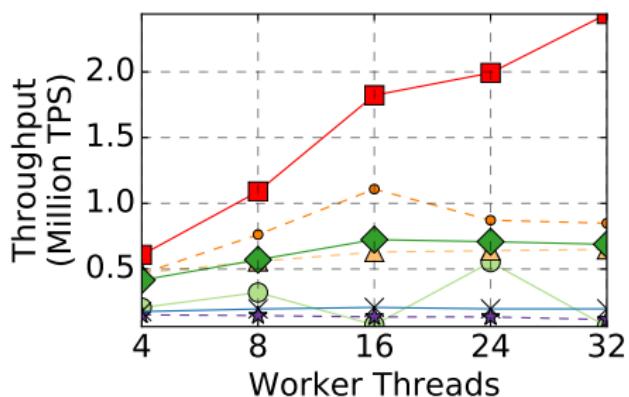


Execution & Synchronization Decoupling: Deterministic priority-based planning followed by queue-oriented execution

Experimental Analysis

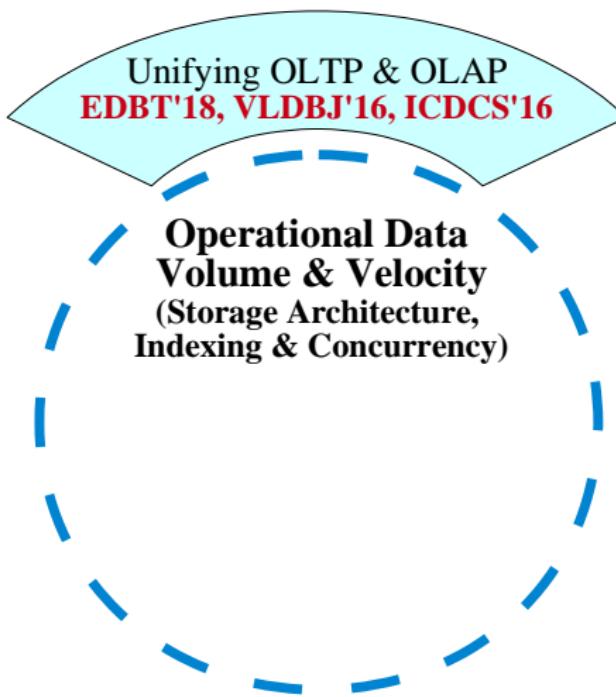
QueCC: Effect of Parallel Update Transactions

ERMIA-SI_SSNI FOEDUS-MOCC CICADA NO_WAIT QUECC SILO TICTOC



Avoiding thread coordination & eliminating all execution-induced aborts

Unifying OLTP and OLAP



Unifying OLTP and OLAP: Velocity & Volume Dimensions

Observed Trends

In operational databases, there is a pressing need to close the gap between the write-optimized layout for OLTP (i.e., row-wise) and the read-optimized layout for OLAP (i.e., column-wise).

Unifying OLTP and OLAP: Velocity & Volume Dimensions

Observed Trends

In operational databases, there is a pressing need to close the gap between the write-optimized layout for OLTP (i.e., row-wise) and the read-optimized layout for OLAP (i.e., column-wise).

Introducing a *lineage-based storage architecture*, a contention-free update mechanism over a native columnar storage in order to

Unifying OLTP and OLAP: Velocity & Volume Dimensions

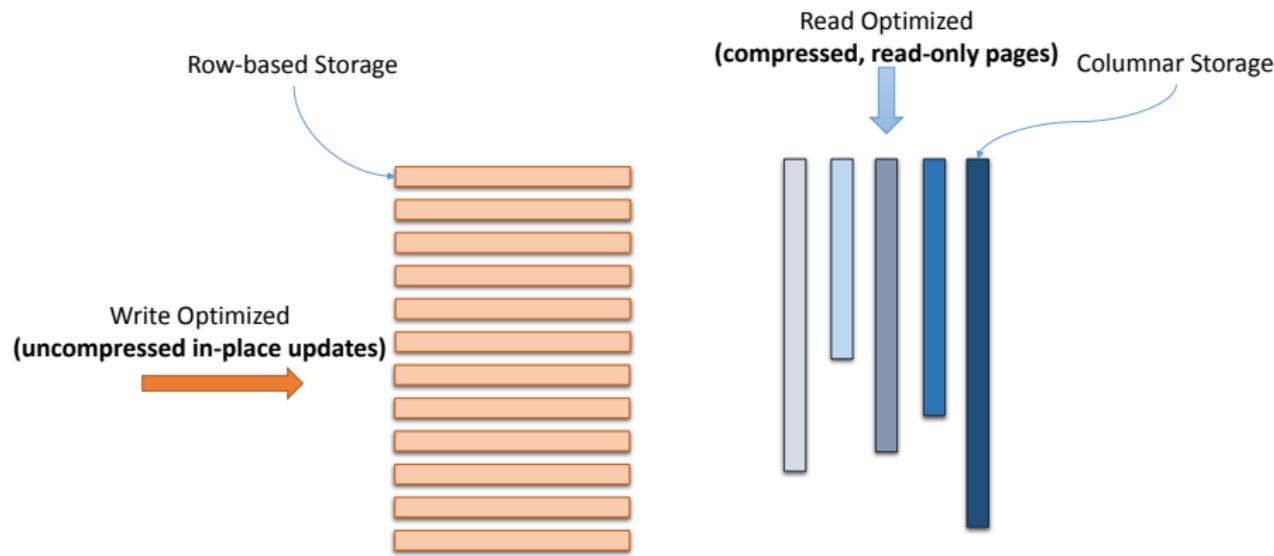
Observed Trends

In operational databases, there is a pressing need to close the gap between the write-optimized layout for OLTP (i.e., row-wise) and the read-optimized layout for OLAP (i.e., column-wise).

Introducing a *lineage-based storage architecture*, a contention-free update mechanism over a native columnar storage in order to

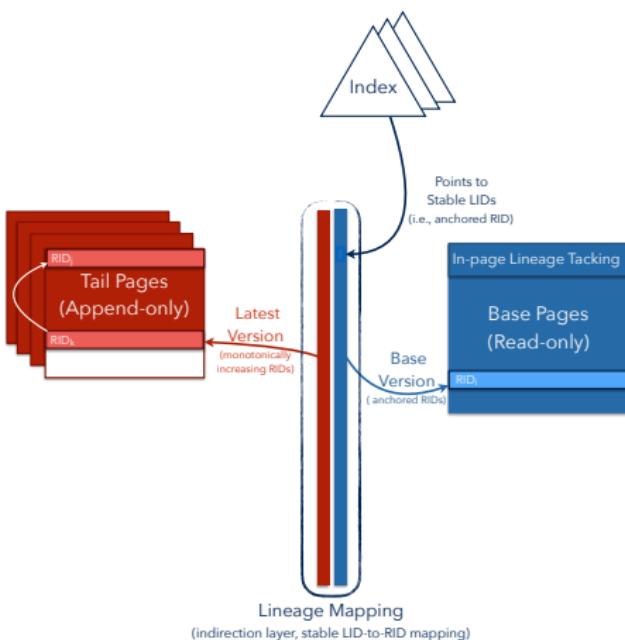
lazily and independently stage stable data from a write-optimized layout (i.e., OLTP) into a read-optimized layout (i.e., OLAP)

Storage Layout Conflict



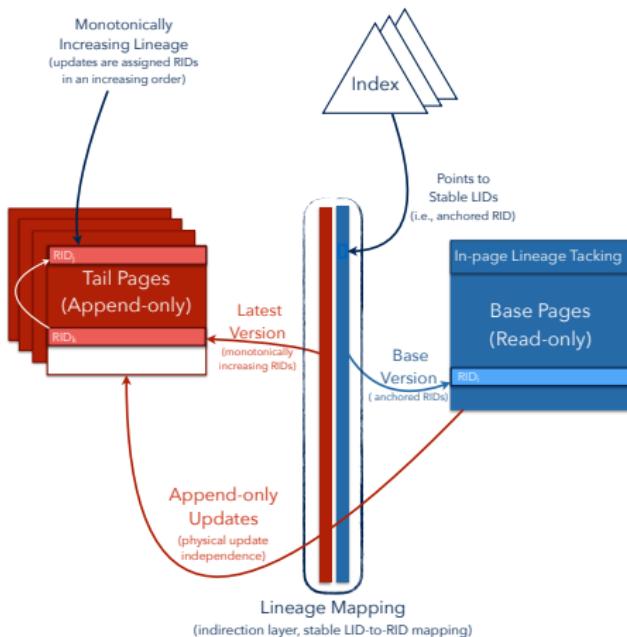
Write-optimized (i.e., uncompressed & row-based) vs. read-optimized (i.e., compressed & column-based) layouts

Lineage-based Storage Architecture (LSA): Intuition



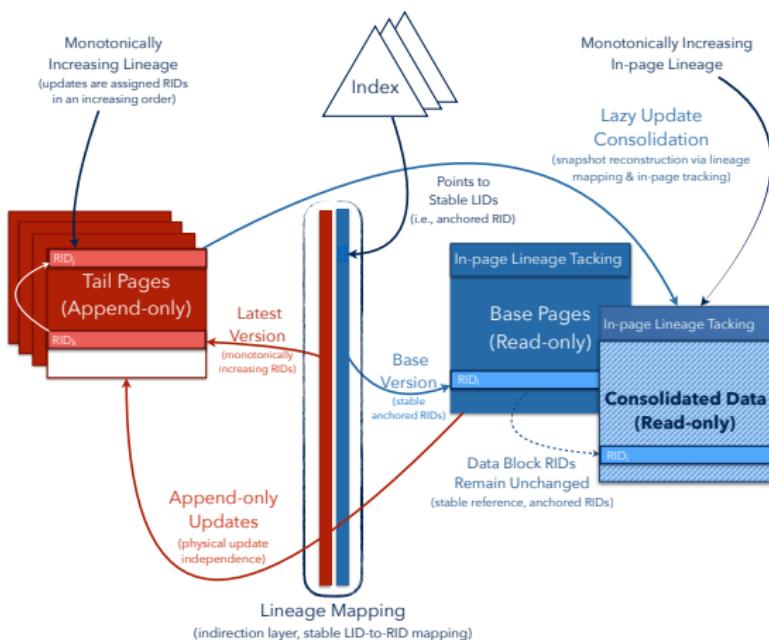
Physical Update Independence: De-coupling data & its updates
(reconstruction via in-page lineage tracking and lineage mapping)

Lineage-based Storage Architecture (LSA): Intuition



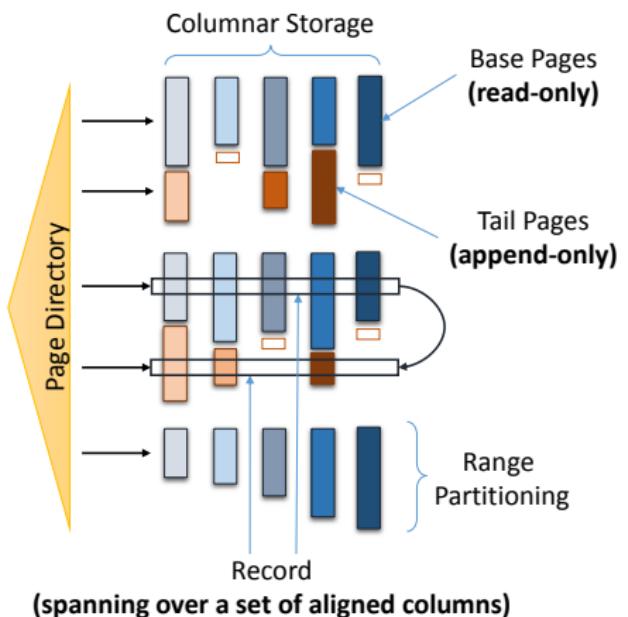
Physical Update Independence: De-coupling data & its updates
(reconstruction via in-page lineage tracking and lineage mapping)

Lineage-based Storage Architecture (LSA): Intuition



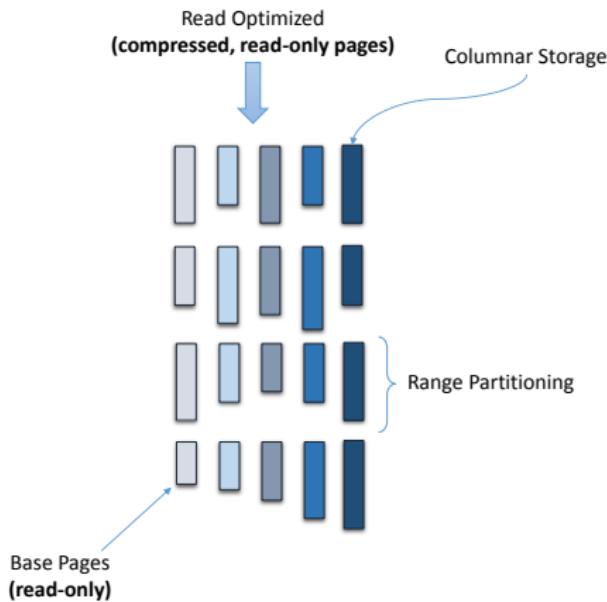
Physical Update Independence: De-coupling data & its updates
(reconstruction via in-page lineage tracking and lineage mapping)

Lineage-based Storage Architecture (LSA): Overview



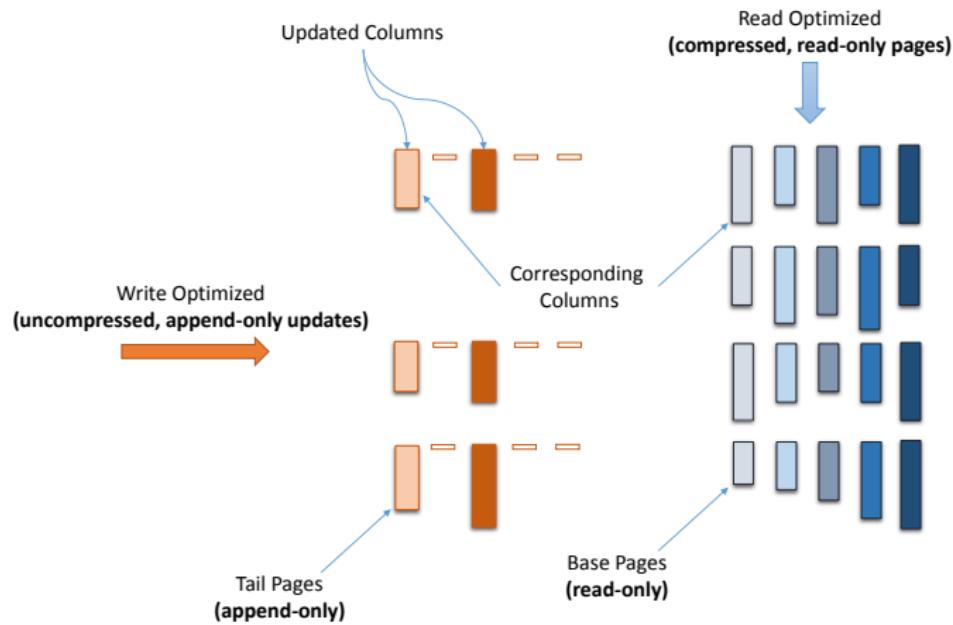
Overview of the lineage-based storage architecture
(**base pages** and **tail pages** are handled identically at the storage layer)

L-Store: Detailed Design



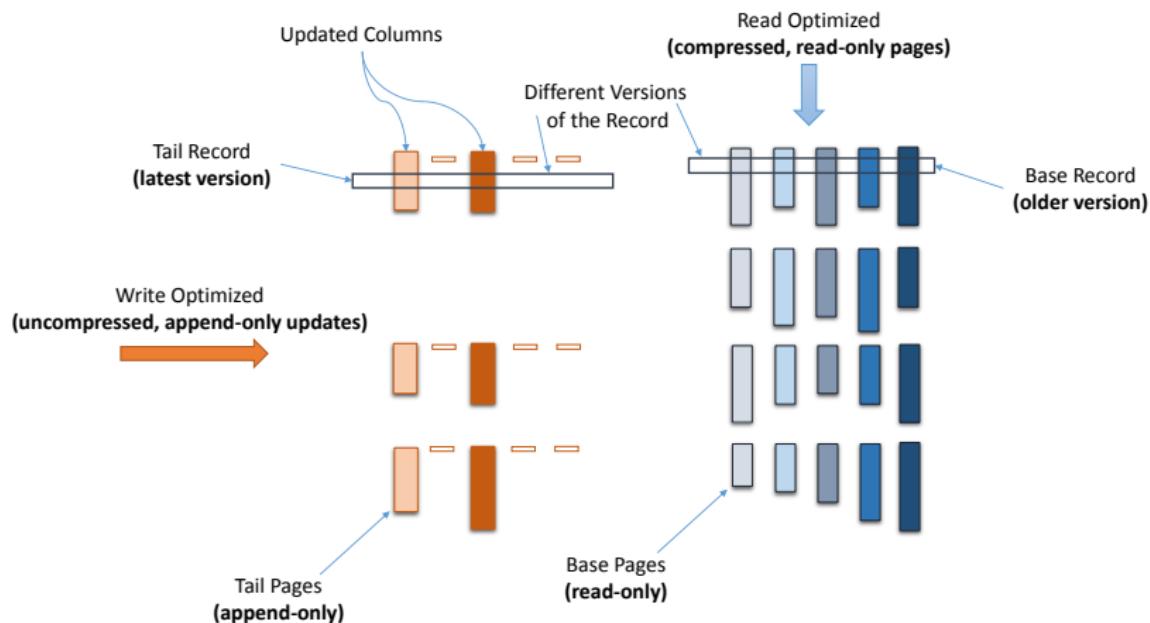
Records are range-partitioned and compressed into a set of ready-only **base pages** (accelerating analytical queries)

L-Store: Detailed Design



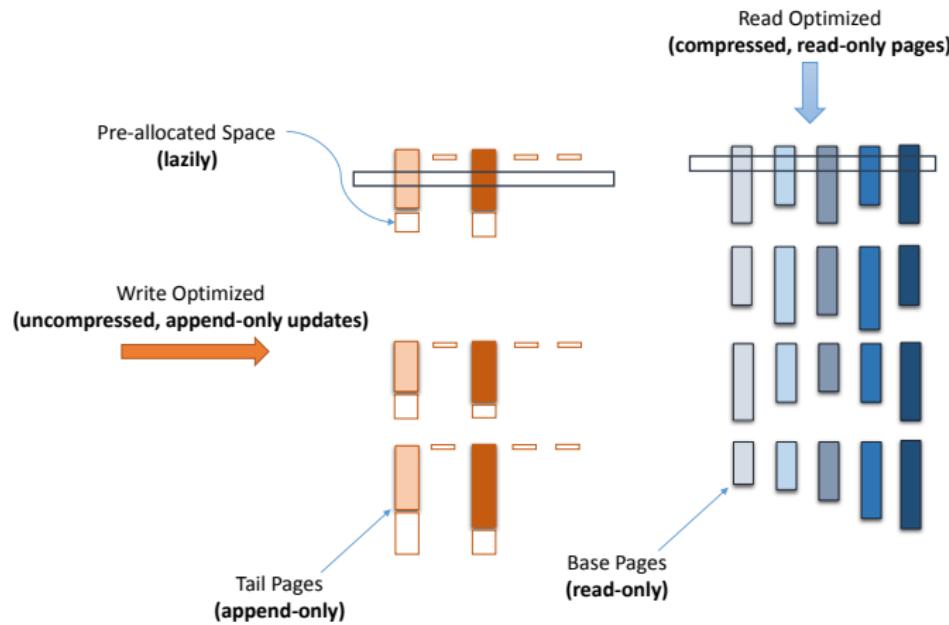
Recent updates for a range of records are clustered in their **tails pages** (transforming costly point updates into an amortized analytical-like query)

L-Store: Detailed Design



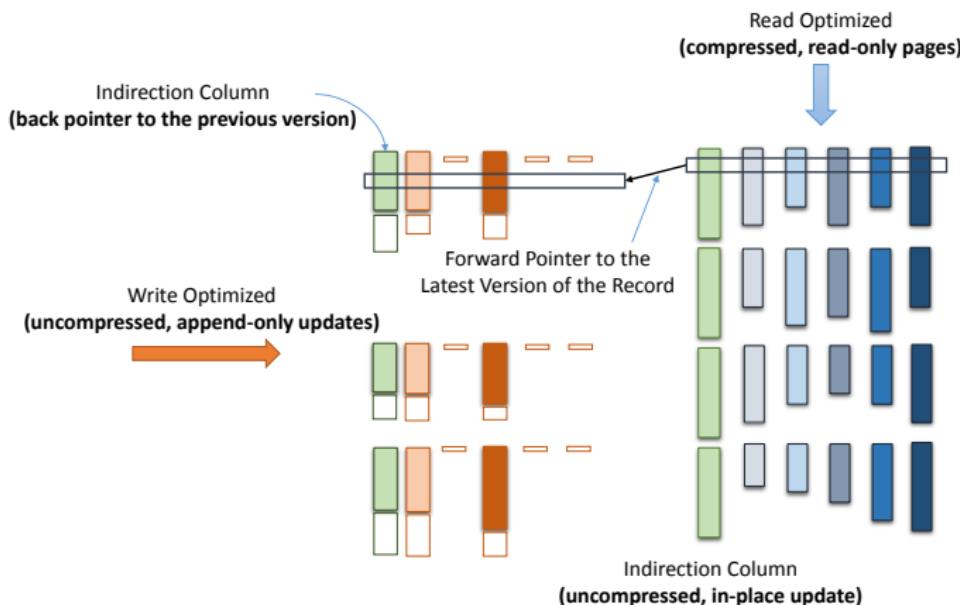
Recent updates for a range of records are clustered in their **tails pages** (transforming costly point updates into an amortized analytical-like query)

L-Store: Detailed Design



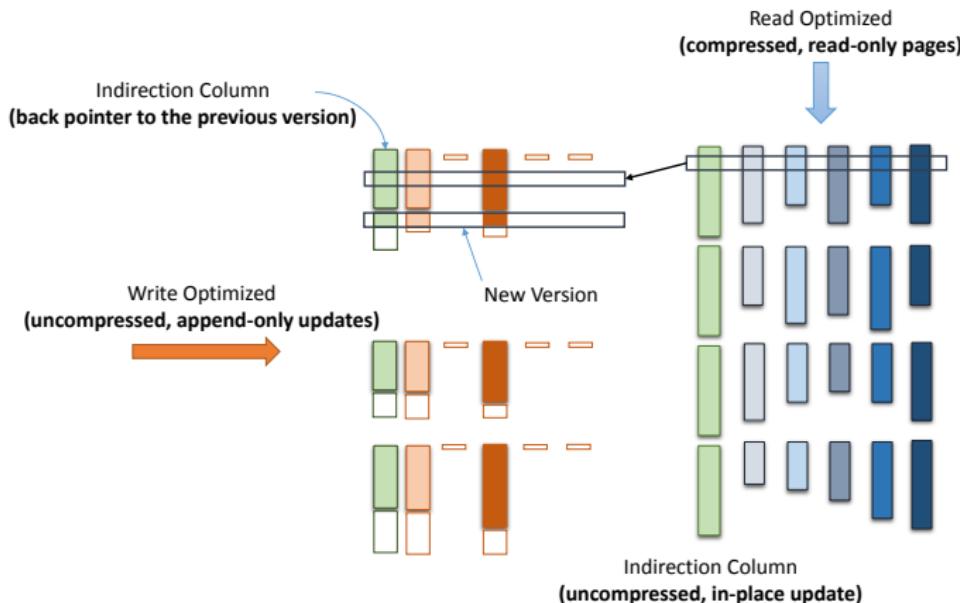
Recent updates are strictly appended, uncompressed in the pre-allocated space
(eliminating the read/write contention)

L-Store: Detailed Design



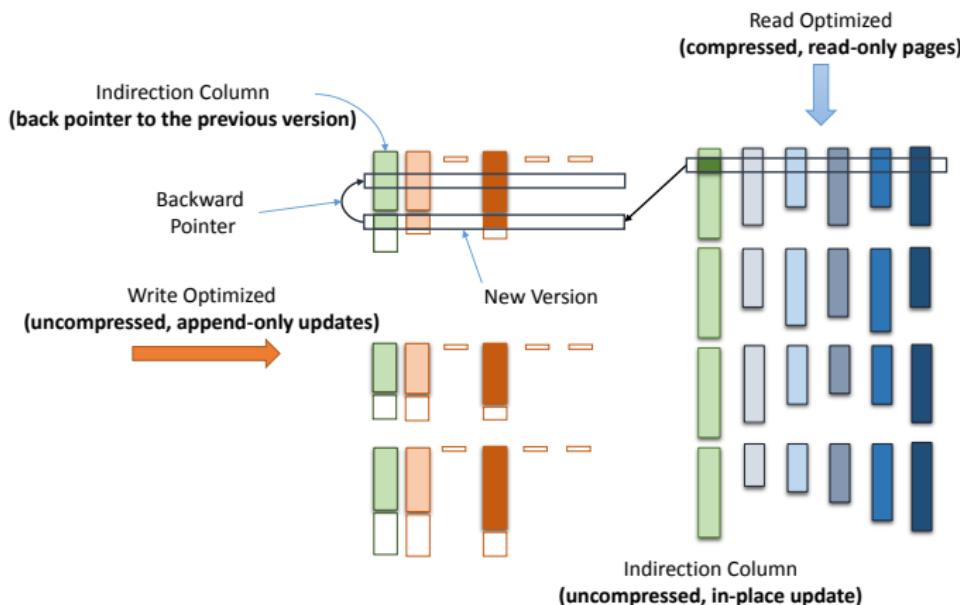
Achieving (at most) 2-hop access to the latest version of any record
(avoiding read performance deterioration for point queries)

L-Store: Detailed Design



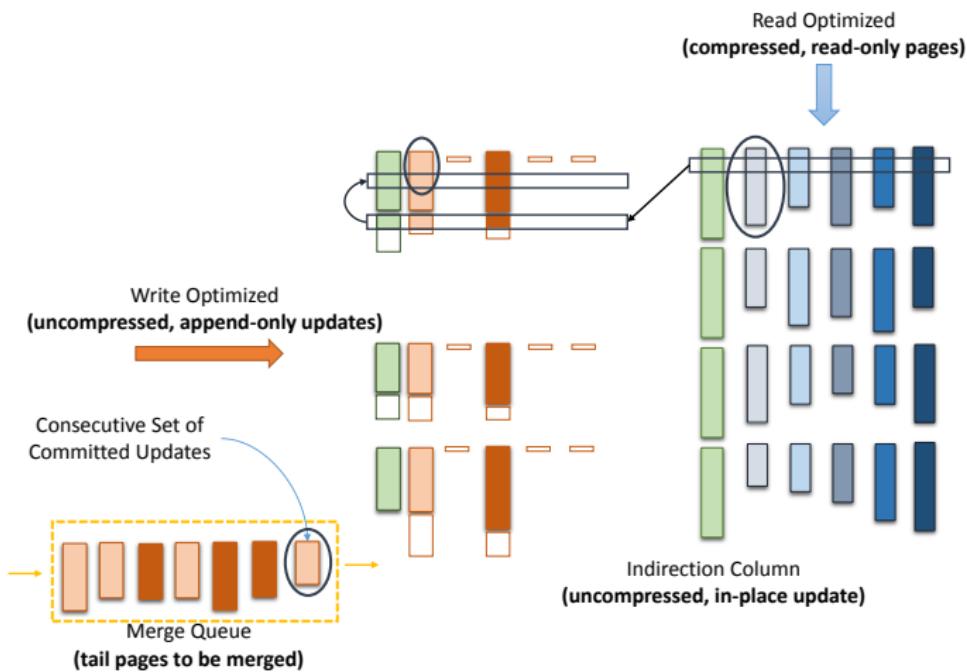
Achieving (at most) 2-hop access to the latest version of any record
(avoiding read performance deterioration for point queries)

L-Store: Detailed Design



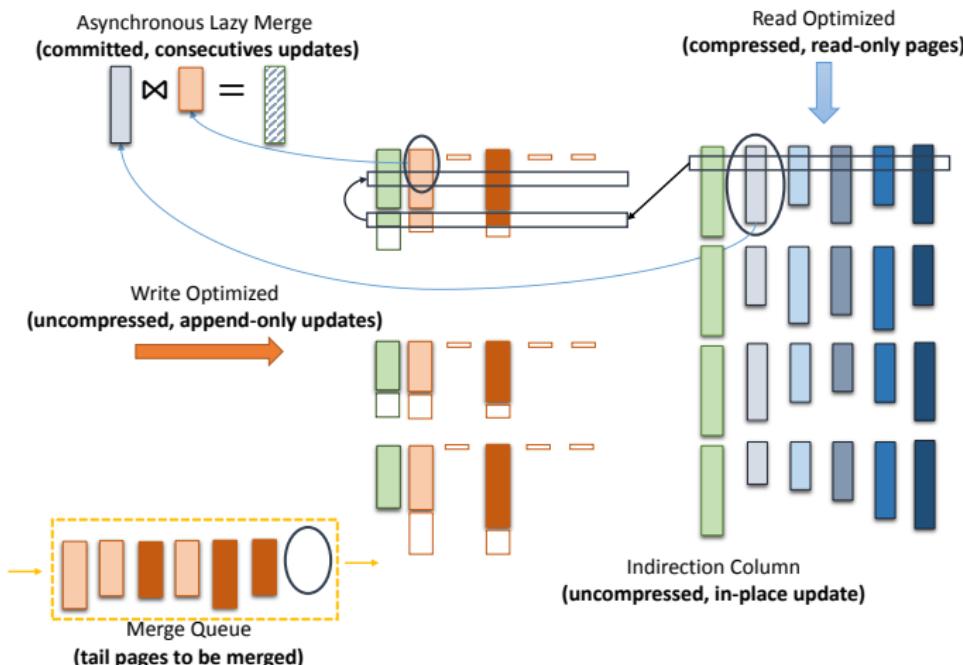
Achieving (at most) 2-hop access to the latest version of any record
(avoiding read performance deterioration for point queries)

L-Store: Contention-free Merge



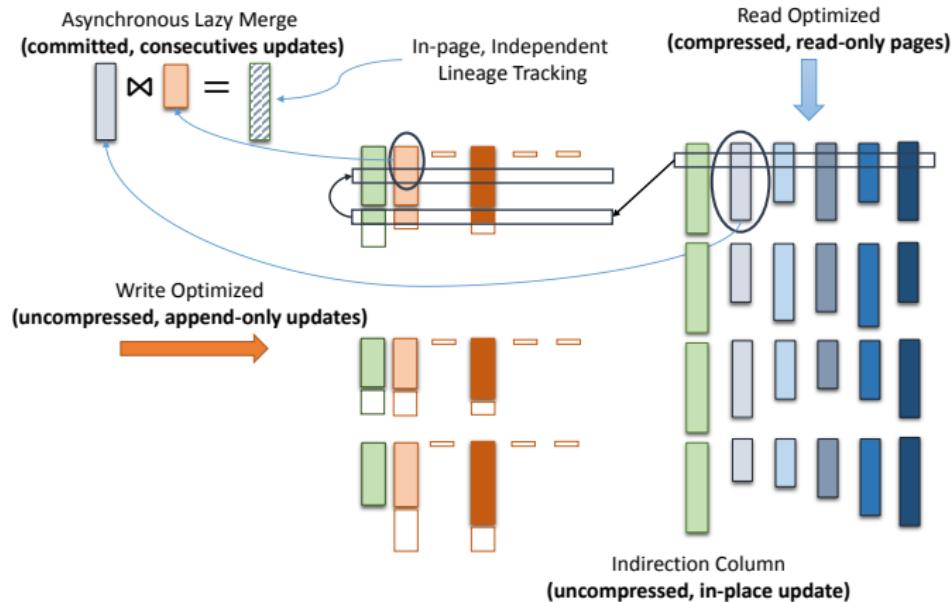
Contention-free merging of only stable data: read-only and committed data
(no need to block on-going and new transactions)

L-Store: Contention-free Merge



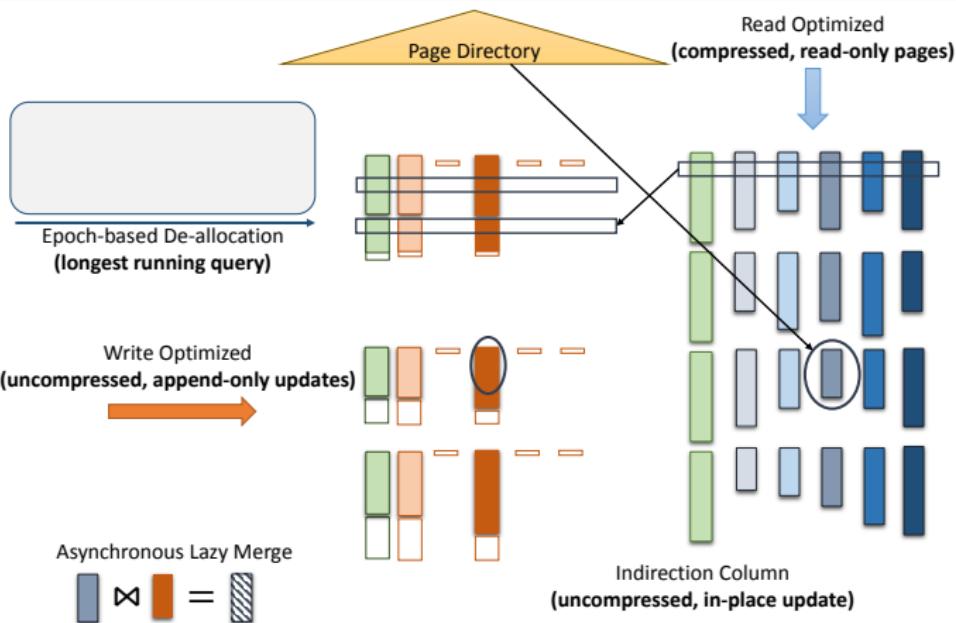
Lazy independent merging of **base pages** with their corresponding **tail pages**
(resembling a local left outer-join of the base and tail pages)

L-Store: Contention-free Merge



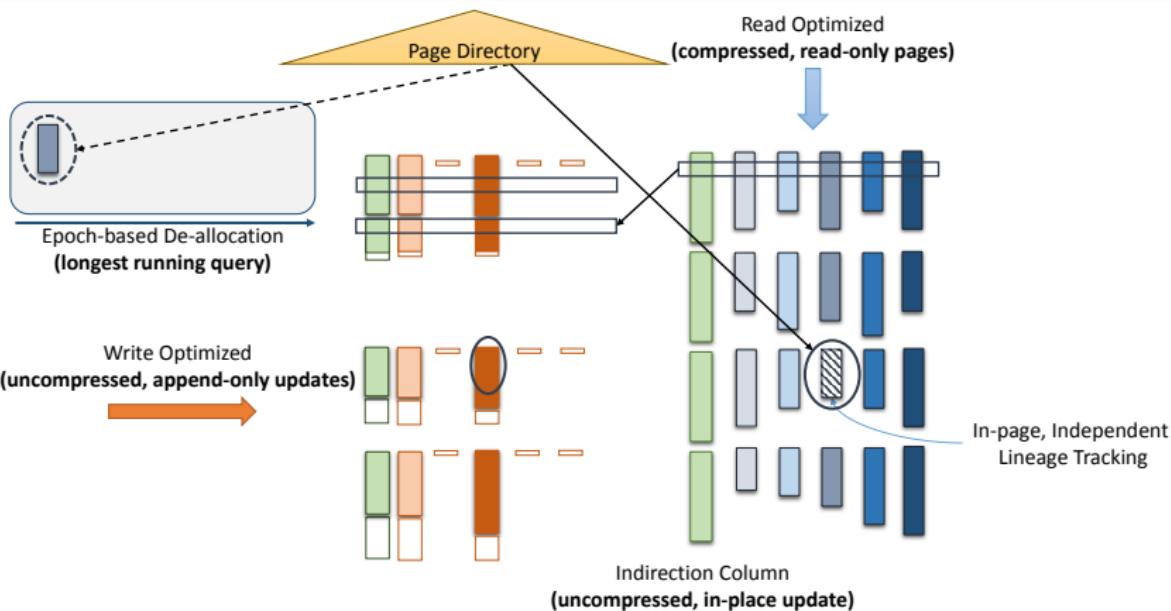
Independently tracking the lineage information within every page
(no need to coordinate merges among different columns of the same records)

L-Store: Epoch-based Contention-free De-allocation



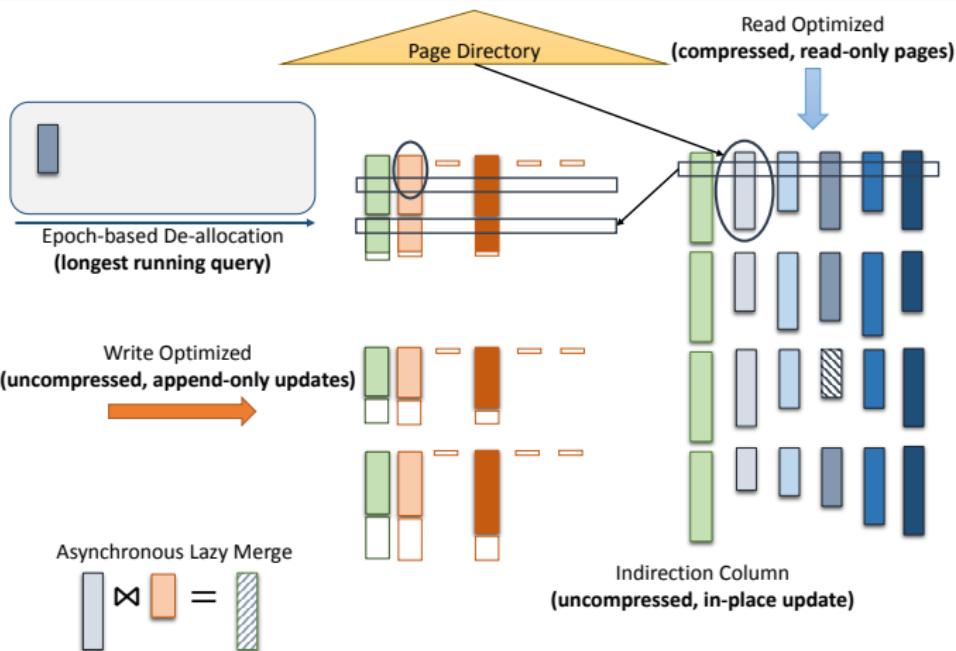
Contention-free page de-allocation using an epoch-based approach
(no need to drain the ongoing transactions)

L-Store: Epoch-based Contention-free De-allocation



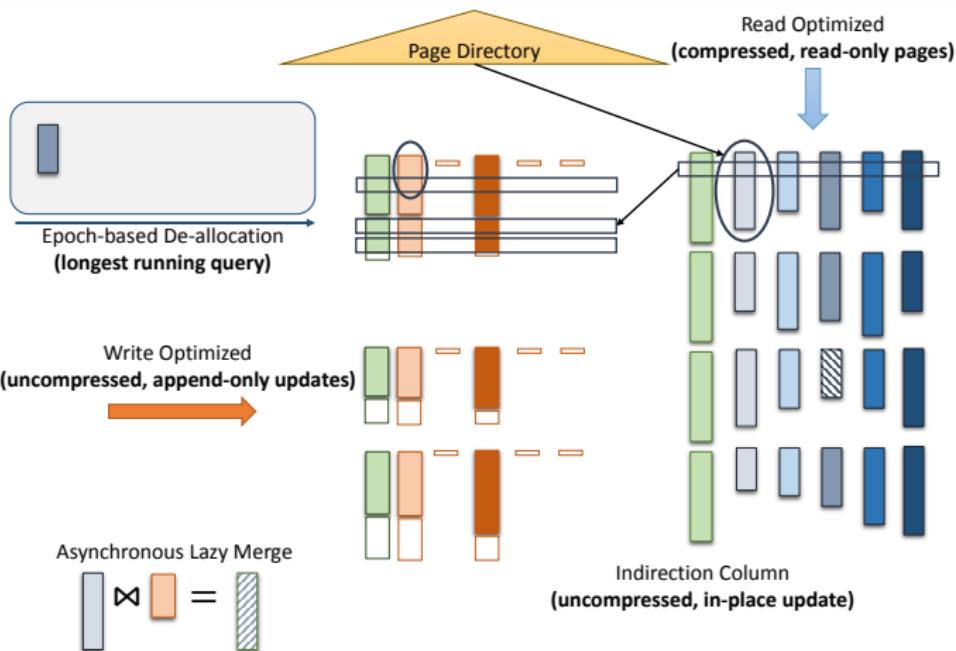
Contention-free page de-allocation using an epoch-based approach
(no need to drain the ongoing transactions)

L-Store: Epoch-based Contention-free De-allocation



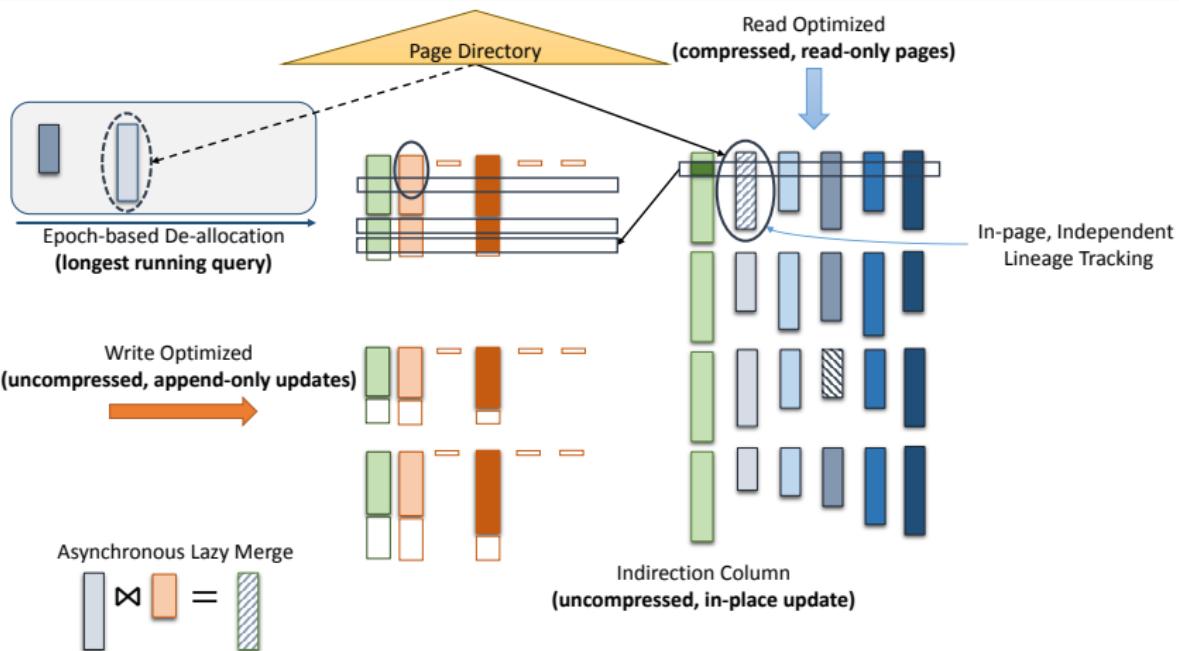
Contention-free page de-allocation using an epoch-based approach
(no need to drain the ongoing transactions)

L-Store: Epoch-based Contention-free De-allocation



Contention-free page de-allocation using an epoch-based approach
(no need to drain the ongoing transactions)

L-Store: Epoch-based Contention-free De-allocation



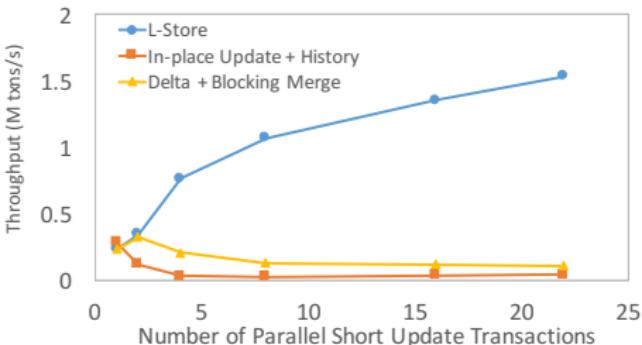
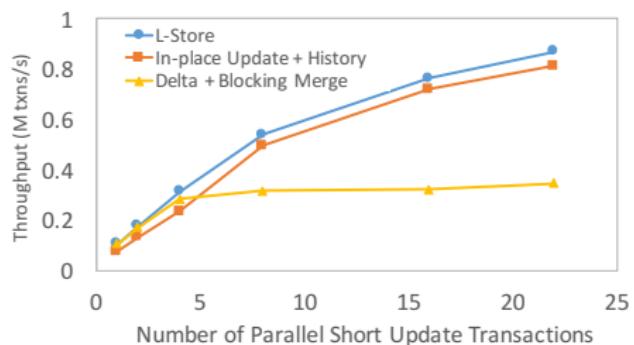
Contention-free page de-allocation using an epoch-based approach
(no need to drain the ongoing transactions)

Experimental Analysis

Experimental Settings

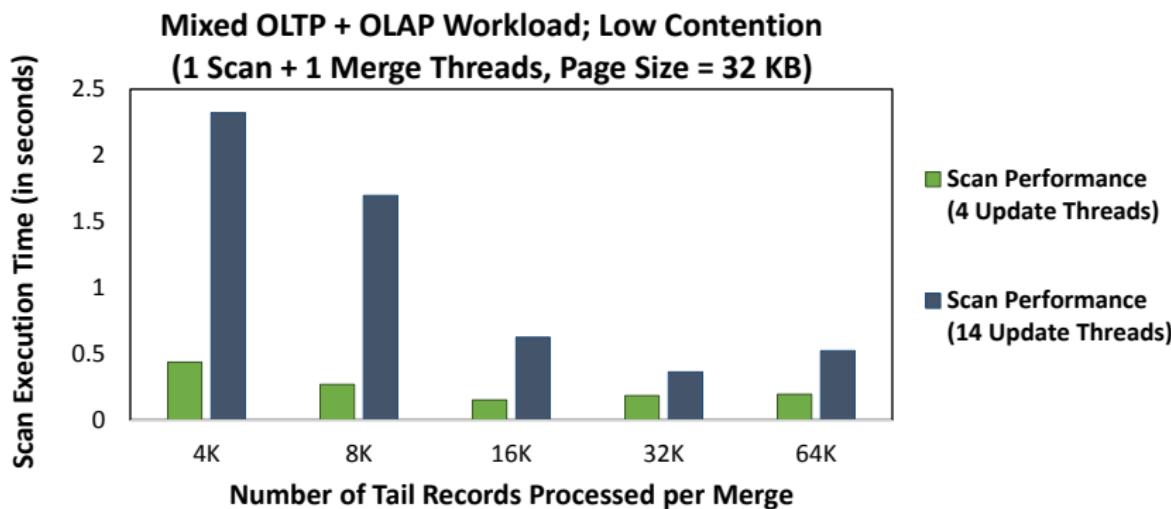
- Hardware:
 - 2 × 6-core Intel(R) Xeon(R) CPU E5-2430 @ 2.20GHz, 64GB, 15 MB L3 cache
- Workload: Extended Microsoft Hekaton Benchmark
 - Comparison with *In-place Update + History* and *Delta + Blocking Merge*
 - Effect of varying contention levels
 - Effect of varying the read/write ratio of short update transactions
 - Effect of merge frequency on scan
 - Effect of varying the number of short update vs. long read-only transactions
 - Effect of varying L-Store data layouts (row vs. columnar)
 - Effect of varying the percentage of columns read in point queries
 - Comparison with log-structured storage architecture (*LevelDB*)

Effect of Varying Contention Levels



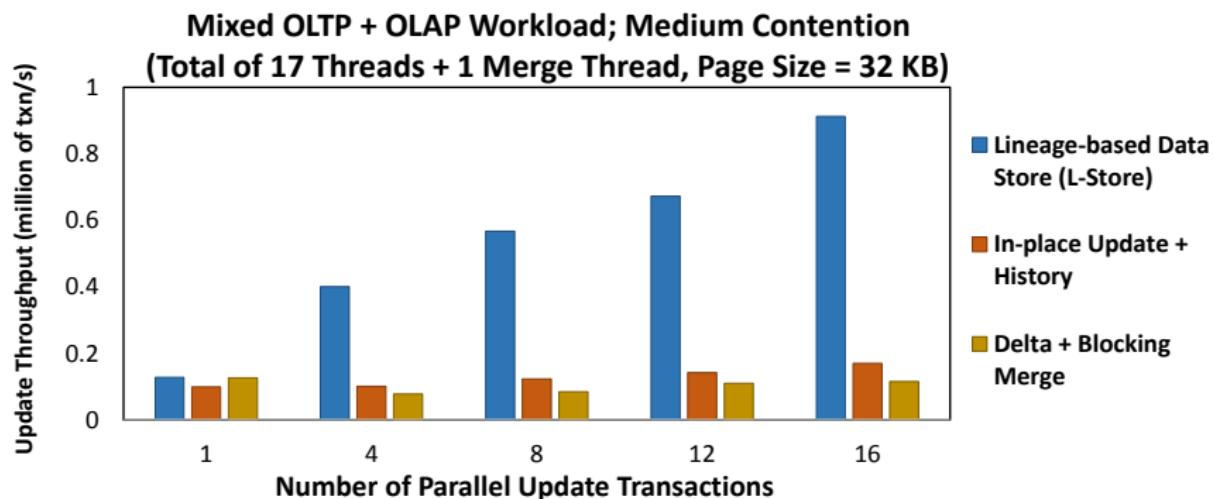
Achieving up to **40×** as increasing the update contention

Effect of Merge Frequency on Scan Performance



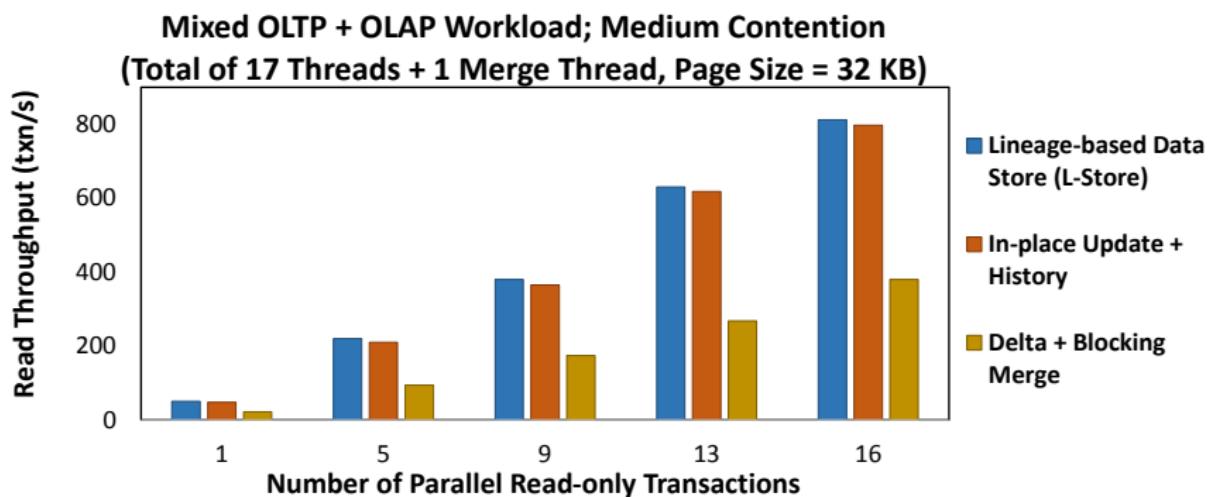
Merge process is essential in maintaining efficient scan performance

Effect of Mixed Workloads: Update Performance



Eliminating latching & locking results in a substantial performance improvement

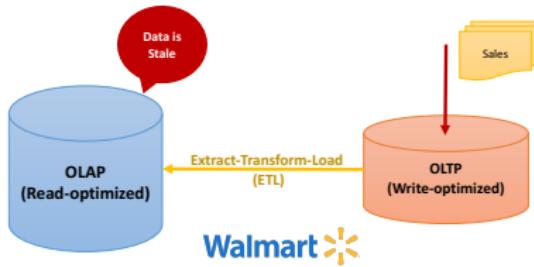
Effect of Mixed Workloads: Read Performance



Coping with tens of update threads with a single merge thread

- 1 Data Management at Microscale
- 2 Data Management at Microscale
- 3 Data Velocity: Index Maintenance
- 4 Data Volume: MVCC Concurrency
- 5 Data Volume: Coordination-free Concurrency
- 6 Combining Volume & Velocity: Lineage-based Storage Architecture
- 7 Data at Macroscale: Decentralized & Democratic Data Platform
- 8 Conclusions
- 9 References

Recap: Data Management Challenges at Microscale



OLTP and OLAP data are isolated at microscale

Recap: Data Management Challenges at Microscale



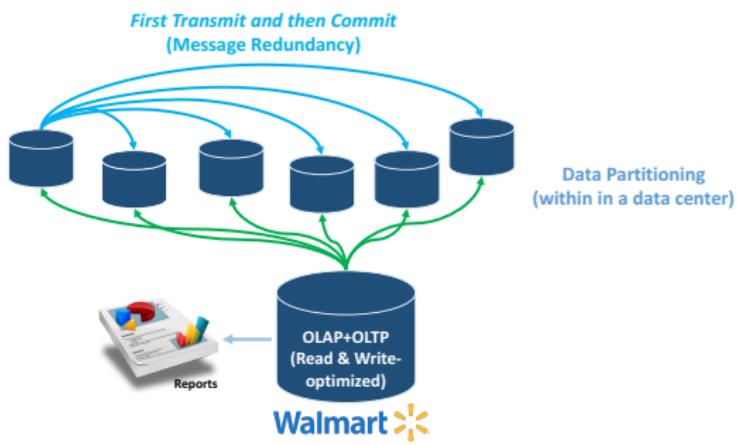
First step is to unify OLTP and OLAP

Platform Scaling: Data Partitioning



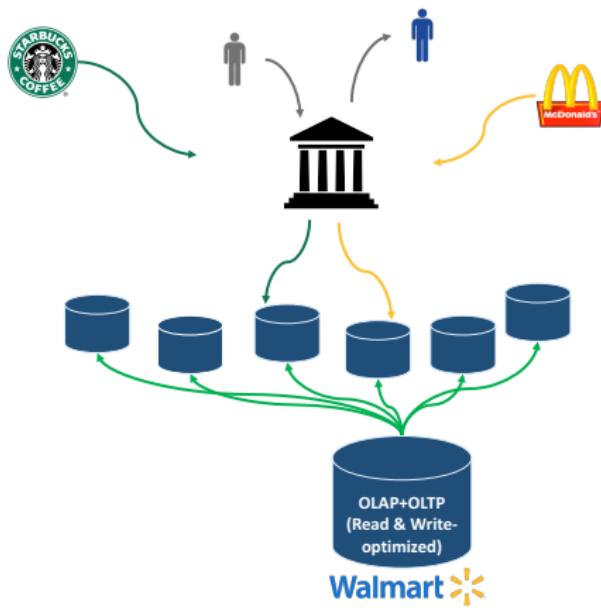
Moving towards distributed environment

Platform Scaling: Non-blocking Agreement Protocols



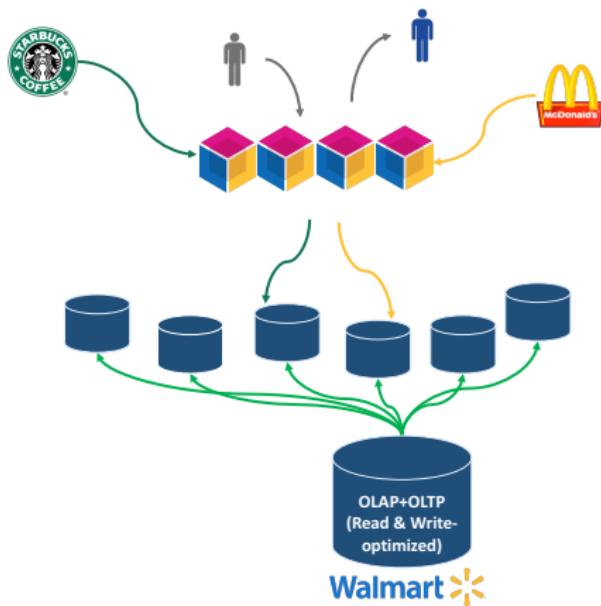
Message redundancy vs. latency trade-offs [EasyCommit, EDBT'18]

Central Control: Data Gate Keeper



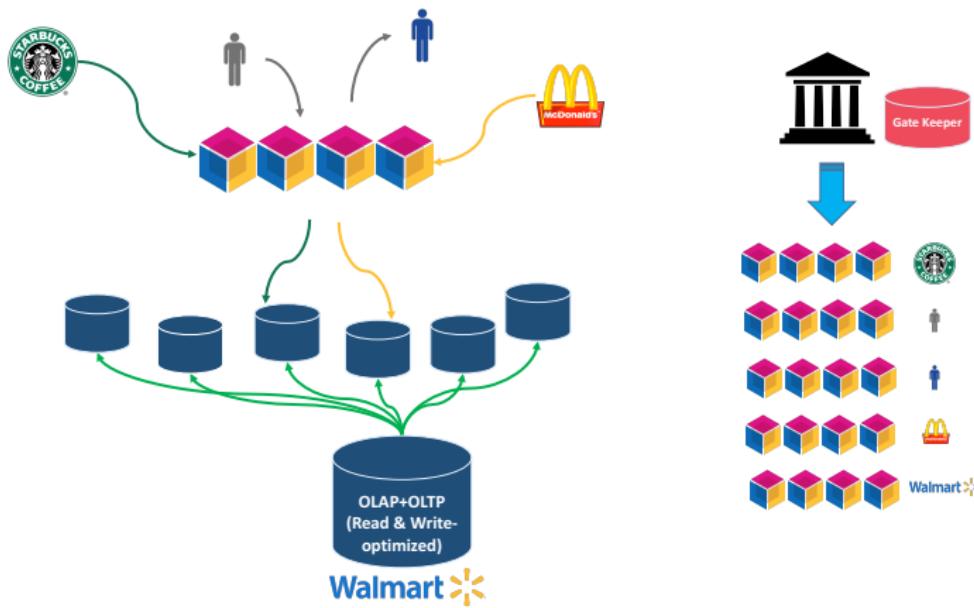
Conform to trusting the central authority and governance

Decentralized Control: Removing Data Barrier



Seek trust in *decentralized* and *democratic* governance [PoE (under submission)]

Democratic Control: Removing Trust Barrier



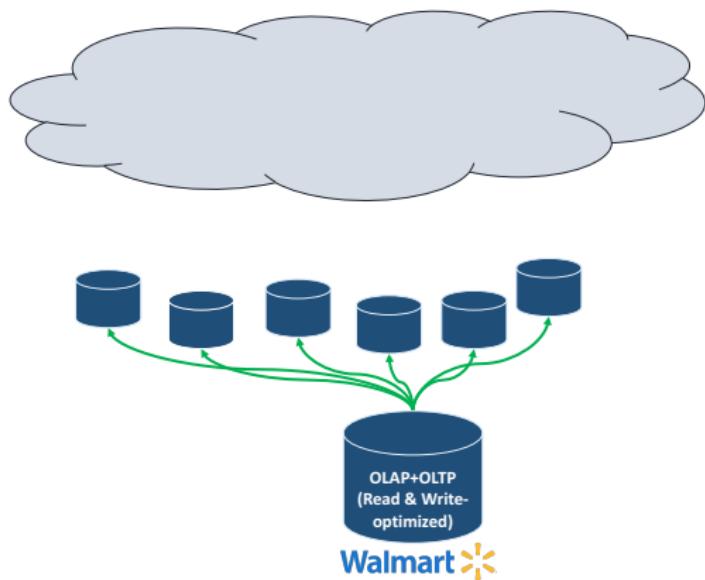
Seek trust in *decentralized* and *democratic* governance [PoE (under submission)]

Global-scale Reliable Platform over Unreliable Hardware



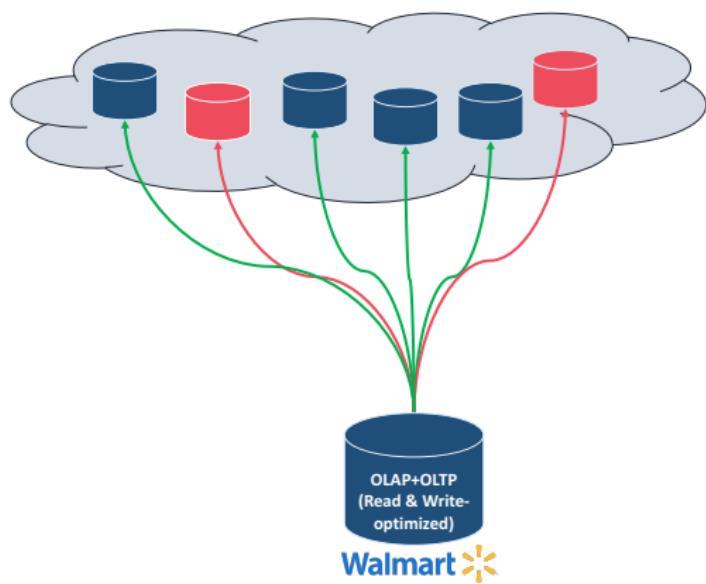
Self-managed infrastructure

Global-scale Reliable Platform over Unreliable Hardware



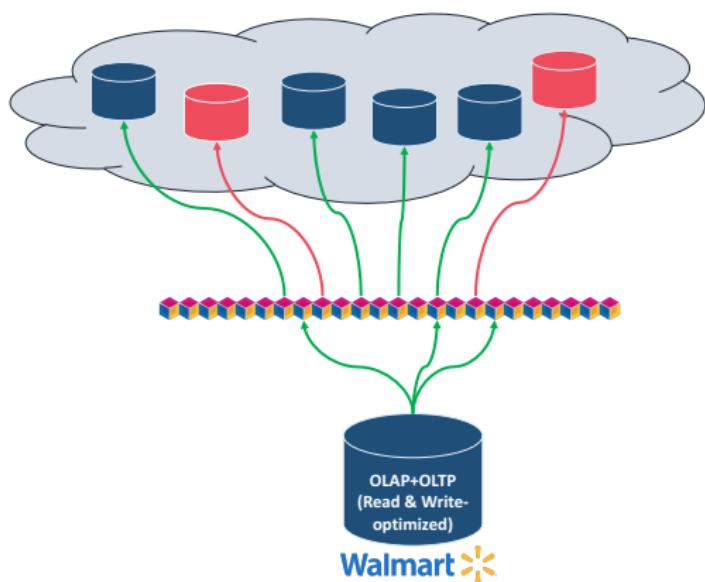
Cloud-managed infrastructure (trust the provider)

Global-scale Reliable Platform over Unreliable Hardware



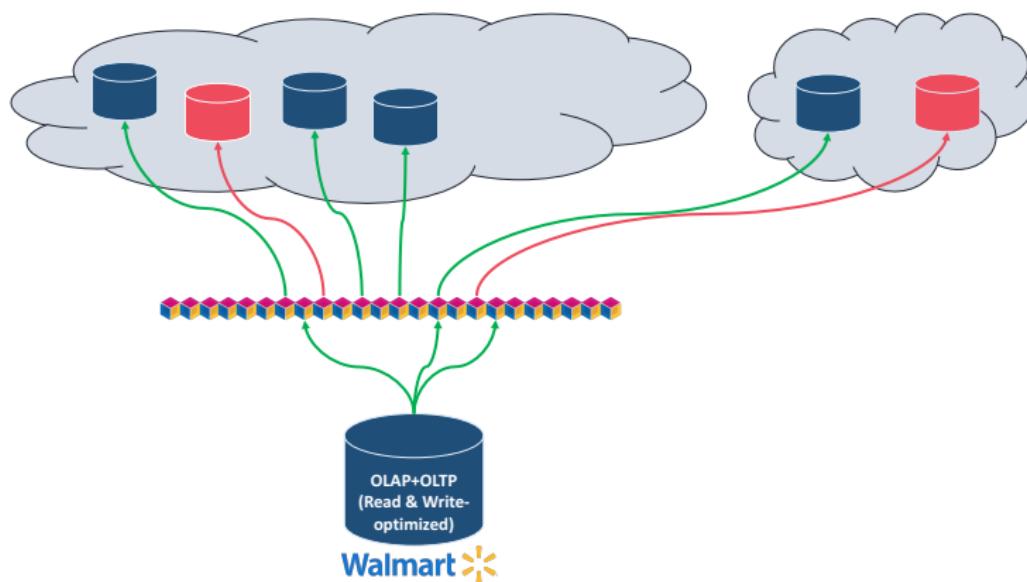
Cloud-managed infrastructure (trust the provider)

Global-scale Reliable Platform over Unreliable Hardware



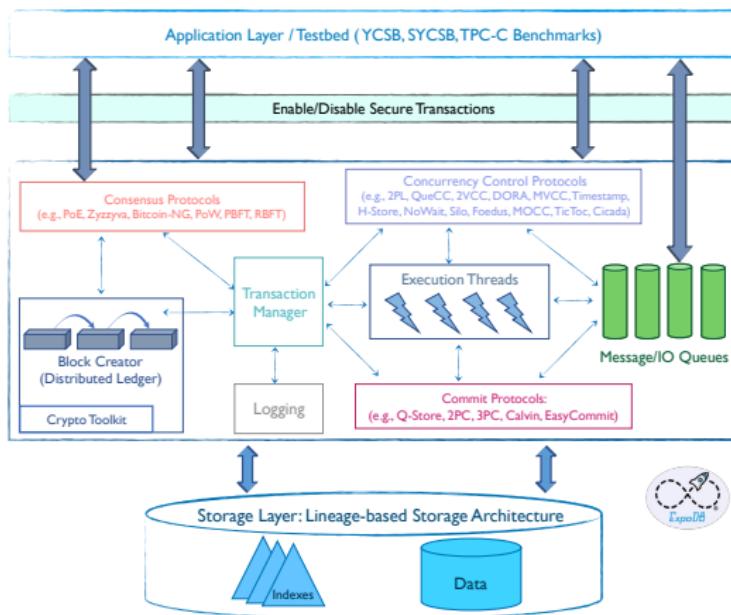
Light-weight, fault-tolerant, trusted middleware [Blockplane, (under submission)]

Global-scale Reliable Platform over Unreliable Hardware



Fault-tolerant protocols vs. consistency models [MultiBFT, GeoBFT (under submission)]

ExpoDB: Exploratory Data Platform Architecture



A decentralized & democratic platform to unify OLTP and OLAP

- 1 Data Management at Microscale
- 2 Data Management at Microscale
- 3 Data Velocity: Index Maintenance
- 4 Data Volume: MVCC Concurrency
- 5 Data Volume: Coordination-free Concurrency
- 6 Combining Volume & Velocity: Lineage-based Storage Architecture
- 7 Data at Macroscale: Decentralized & Democratic Data Platform
- 8 Conclusions
- 9 References

Contributions & Outlook



Questions?
Thank you!

Exploratory Systems Lab (ExpoLab)
Website: <https://msadoghi.github.io/>



Related Publications (Patents Omitted)