Reducing the Storage Overhead of Main-Memory OLTP Databases with **Hybrid Indexes**

Huanchen Zhang

David G. Andersen, Andrew Pavlo, Michael Kaminsky, Lin Ma, Rui Shen

PARALLEL DATA LABORATORY

Carnegie Mellon University



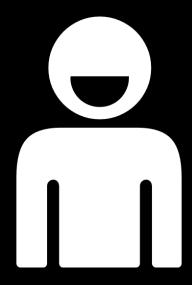




Part I Initial Exploration of Hybrid Indexes [SIGMOD'16]







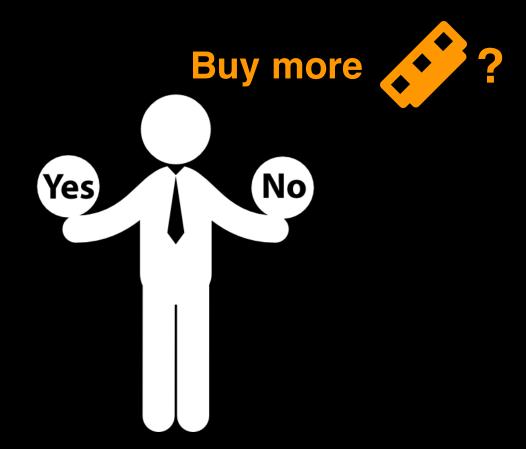




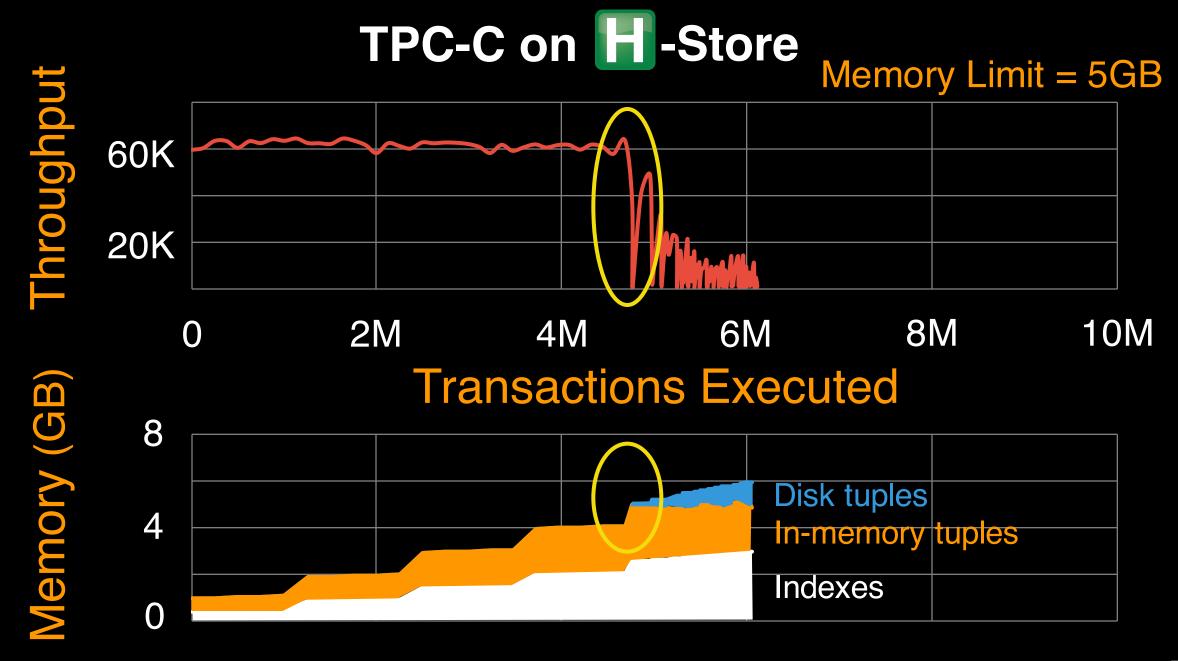








You are running out of memory





The better way: Use memory more efficiently



Indexes are LARGE

Hybrid Index % space for index Benchmark **58%** → **34%** TPC-C **55%** → **41%** Voter 34% -> 18% Articles

Our Contributions [SIGMOD'16]

- 1) The hybrid index architecture
- 2 The Dual-Stage Transformation
- (3) Applied to 4 index structures
 - B+tree Skip List
 - Masstree Adaptive Radix Tree (ART)

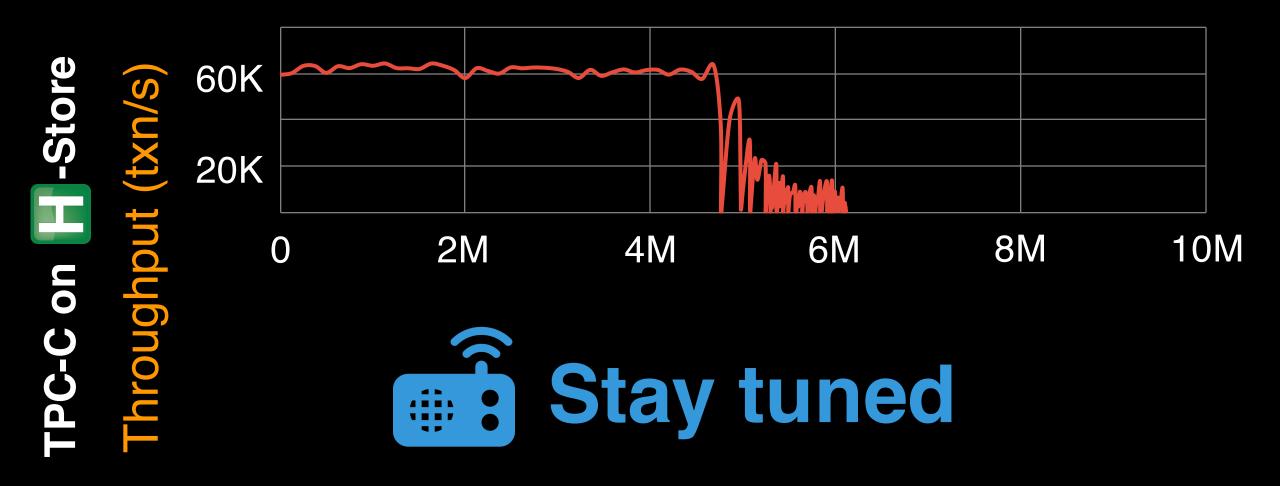
Performance

Space



30 – 70%

Did we solve this problem?



Transactions Executed

How do hybrid indexes achieve memory savings?

0— Static

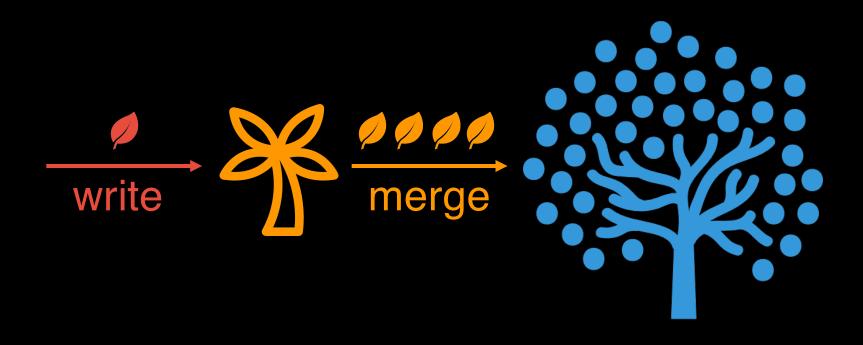
Hybrid Index: a dual-stage architecture





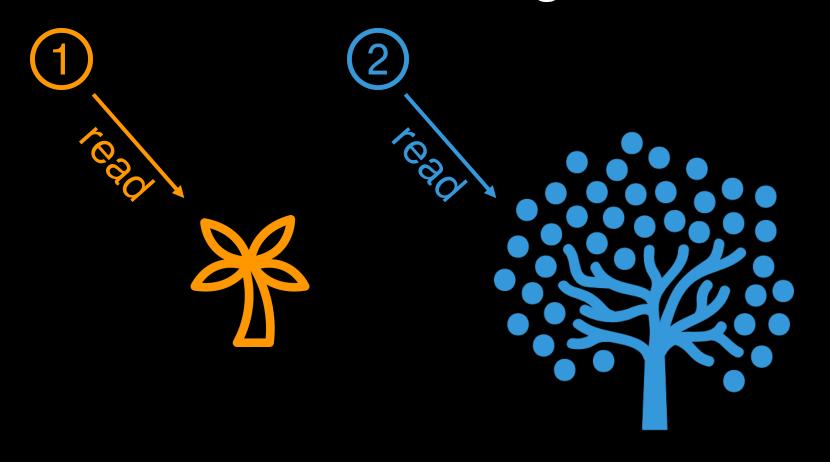
dynamic stage

Inserts are batched in the dynamic stage



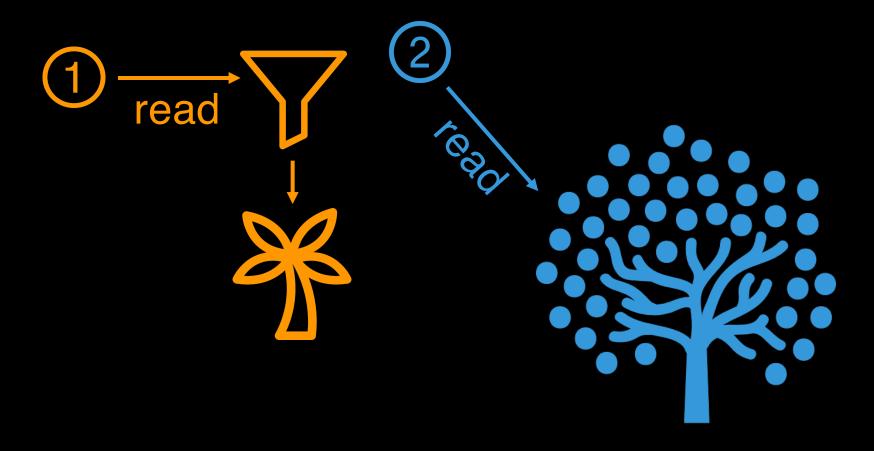
dynamic stage

Reads search the stages in order

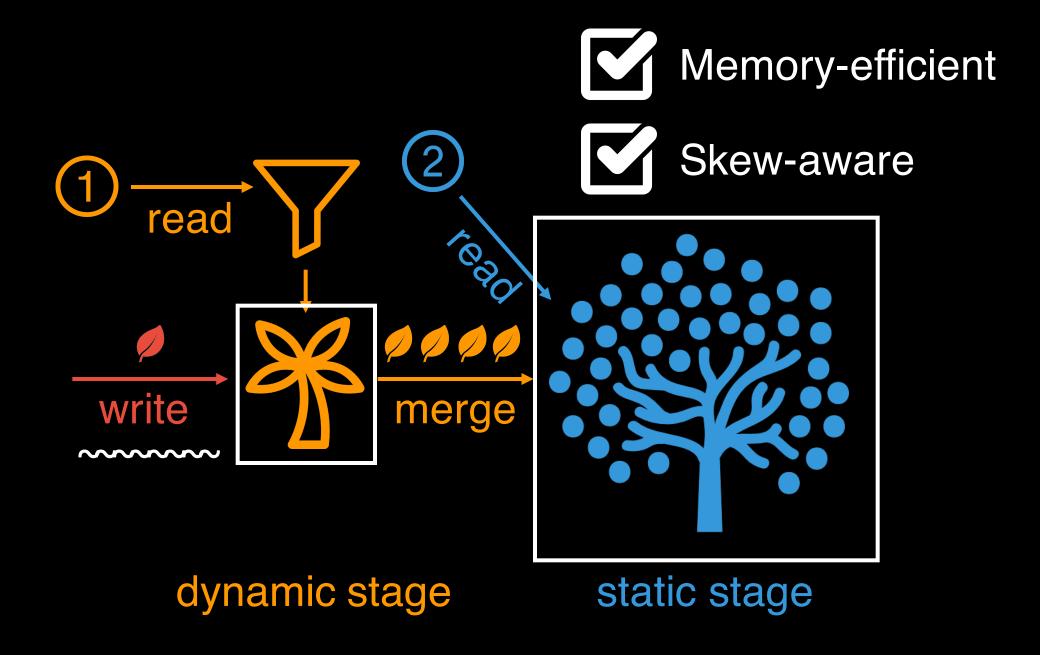


dynamic stage

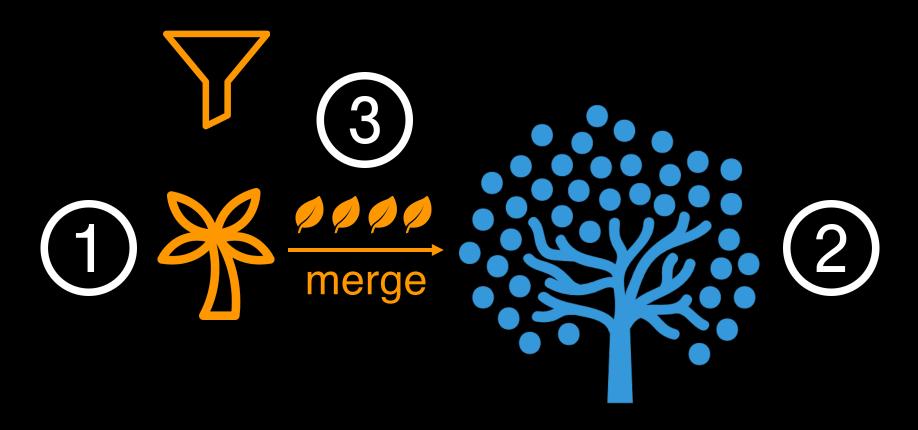
A Bloom filter improves read performance



dynamic stage

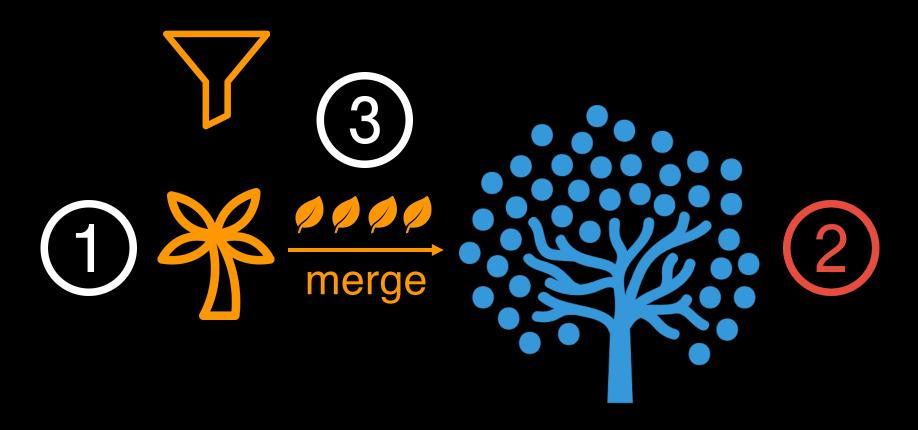


The Dual-Stage Transformation



dynamic stage

The Dual-Stage Transformation



dynamic stage

The Dynamic-to-Static Rules





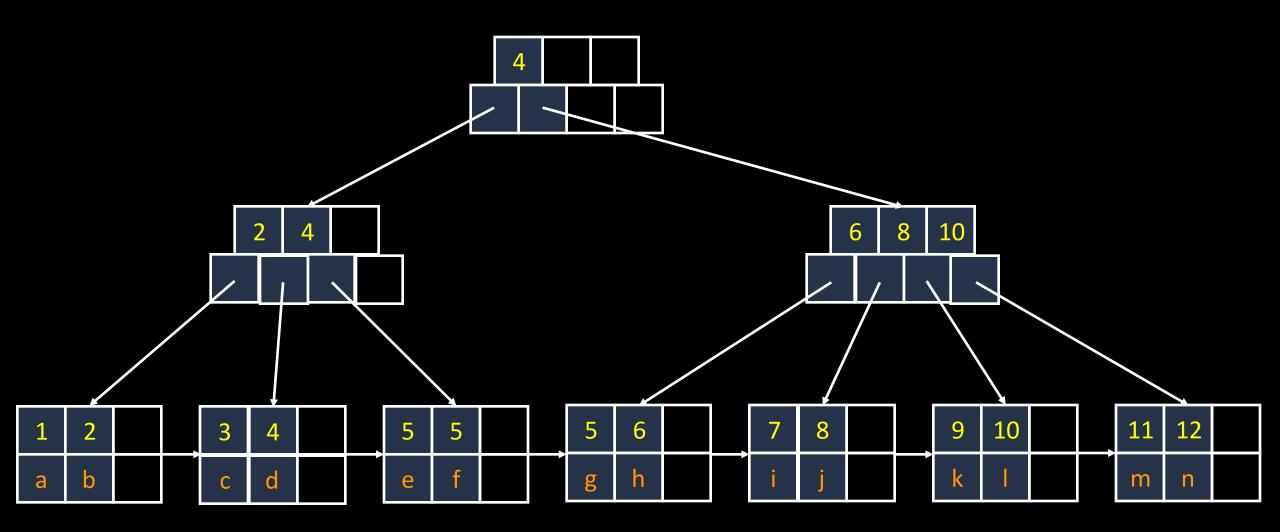


The Dynamic-to-Static Rules

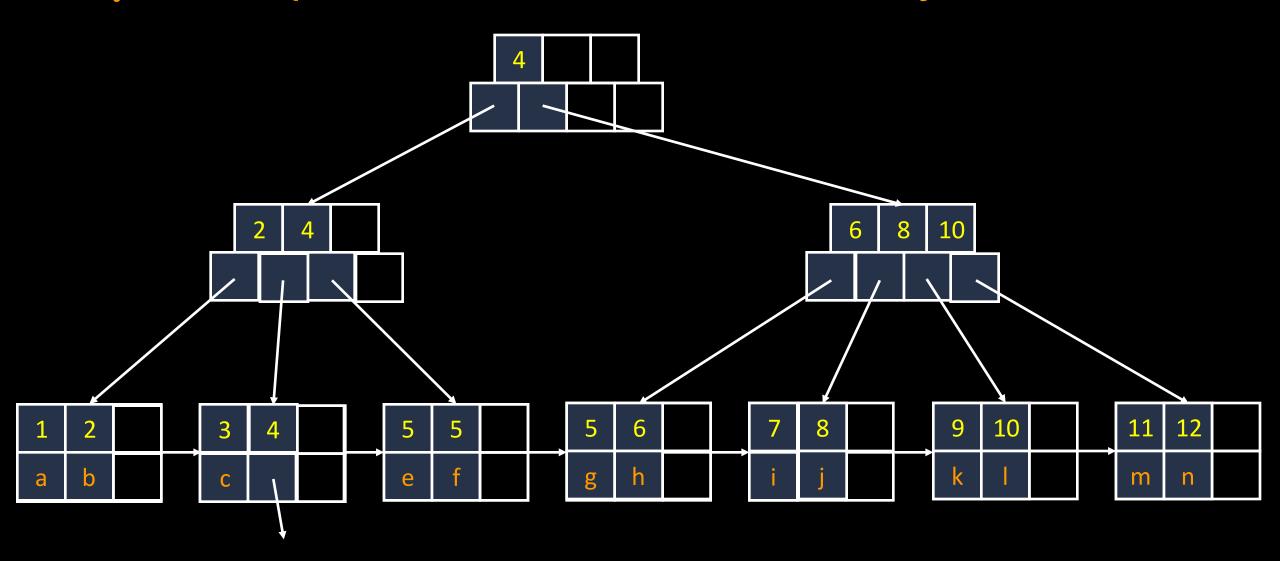




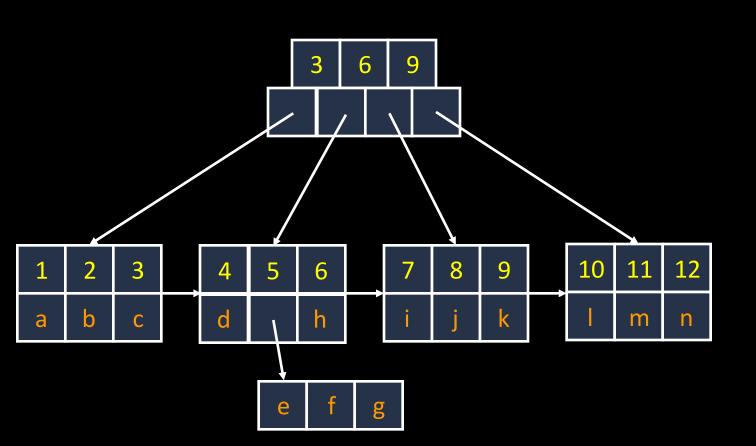




Compaction: minimize # of memory blocks

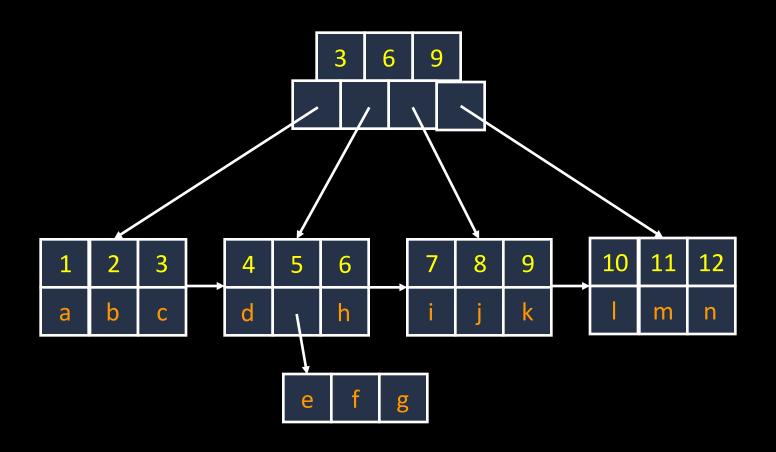


Compaction: minimize # of memory blocks



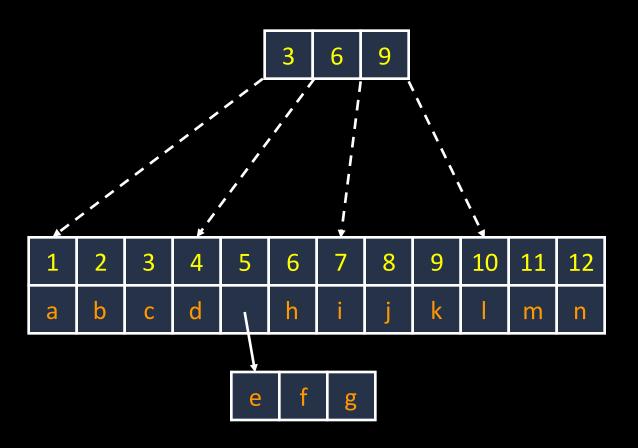


Reduction: minimize structural overhead

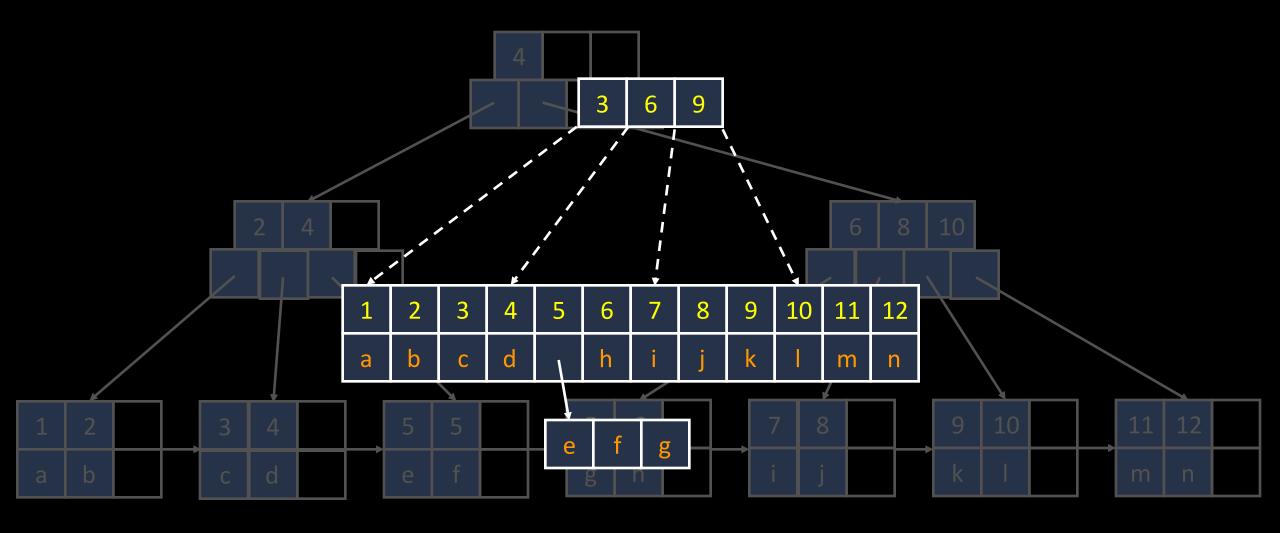




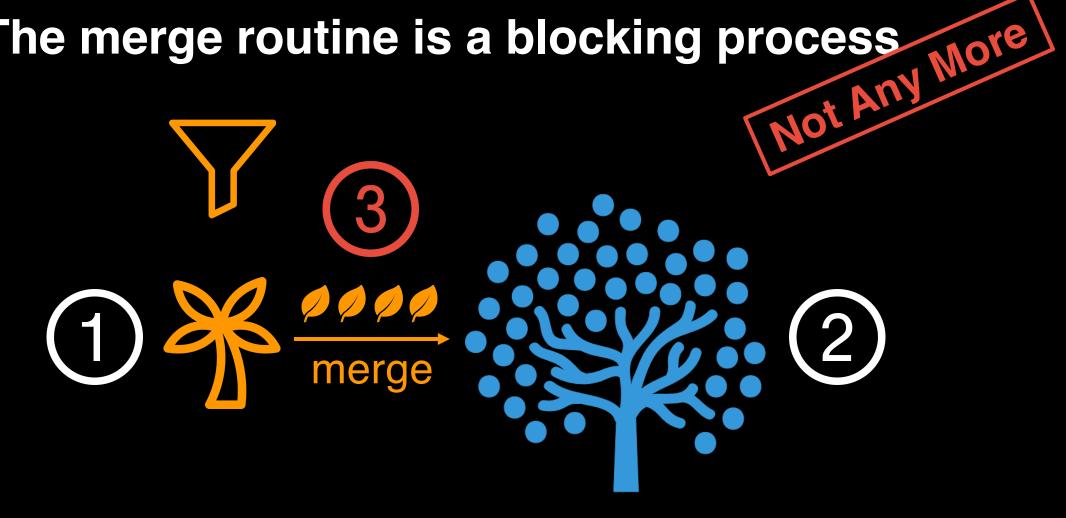
Reduction: minimize structural overhead



Reduction: minimize structural overhead

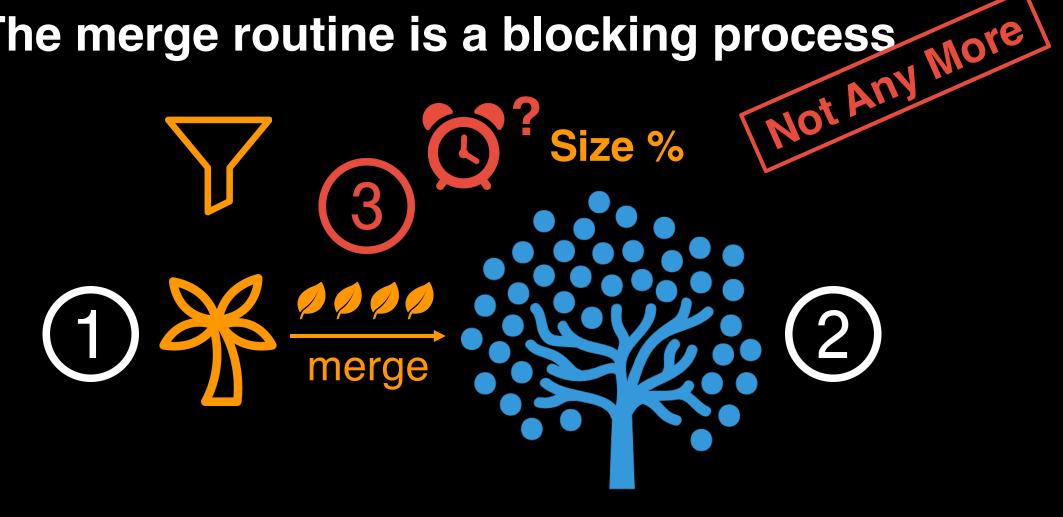


The merge routine is a blocking process



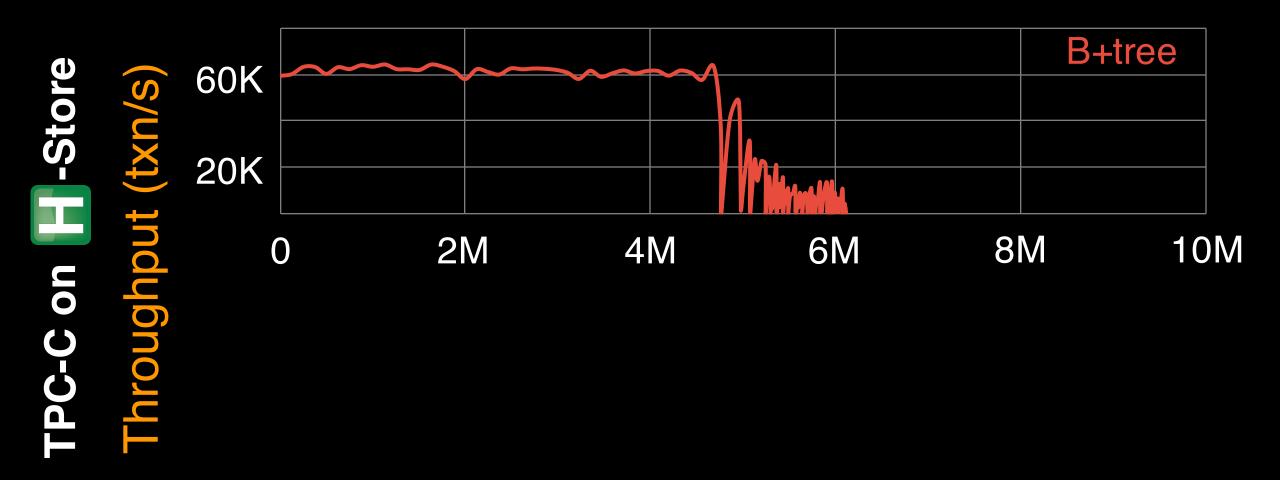
dynamic stage

The merge routine is a blocking process

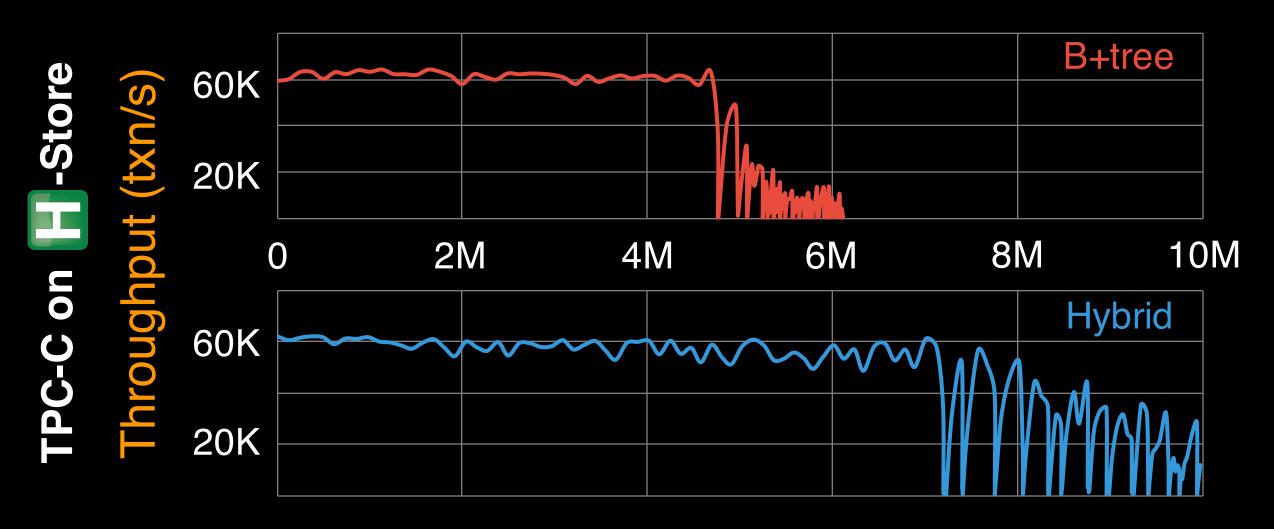


dynamic stage

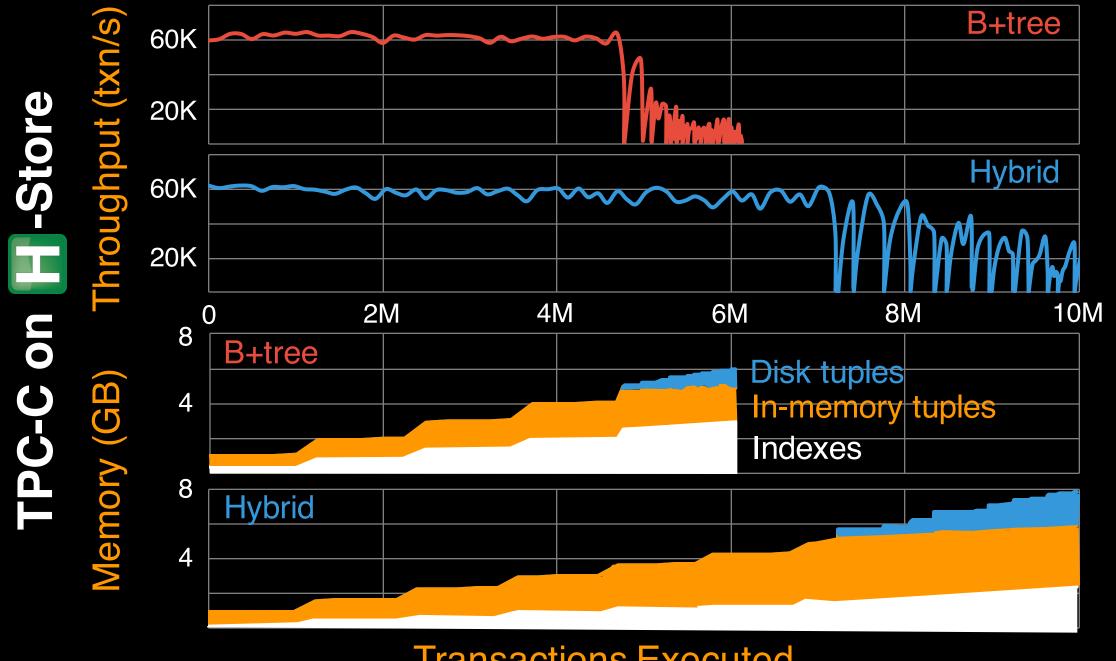
Did we solve this problem?



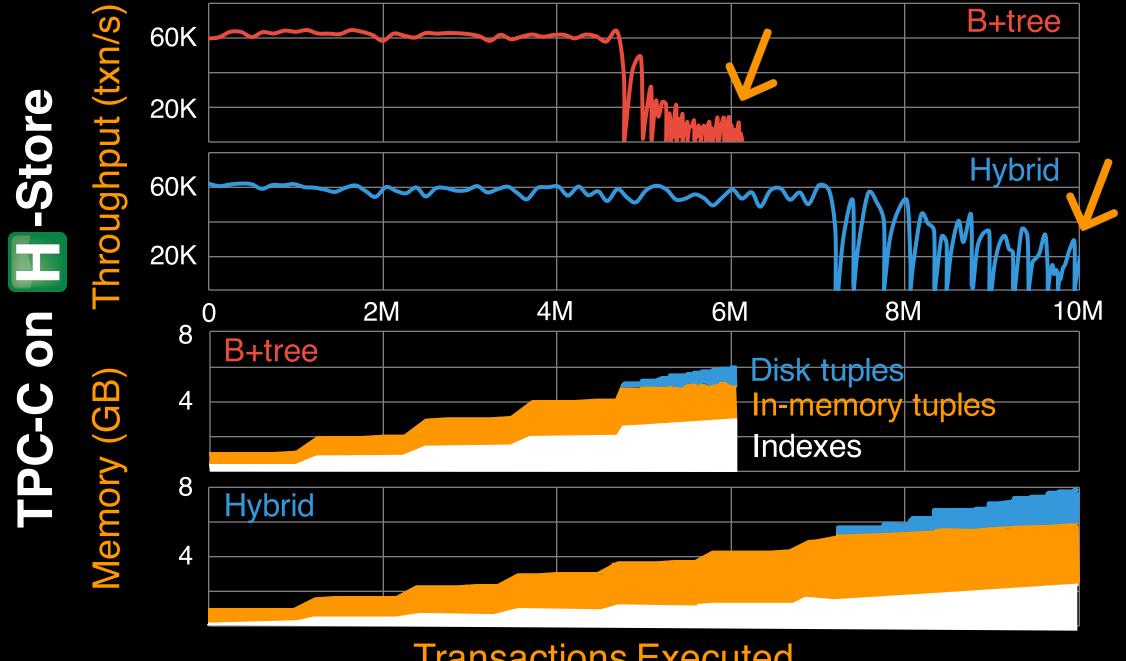
Yes, we improved the DBMS's capacity!



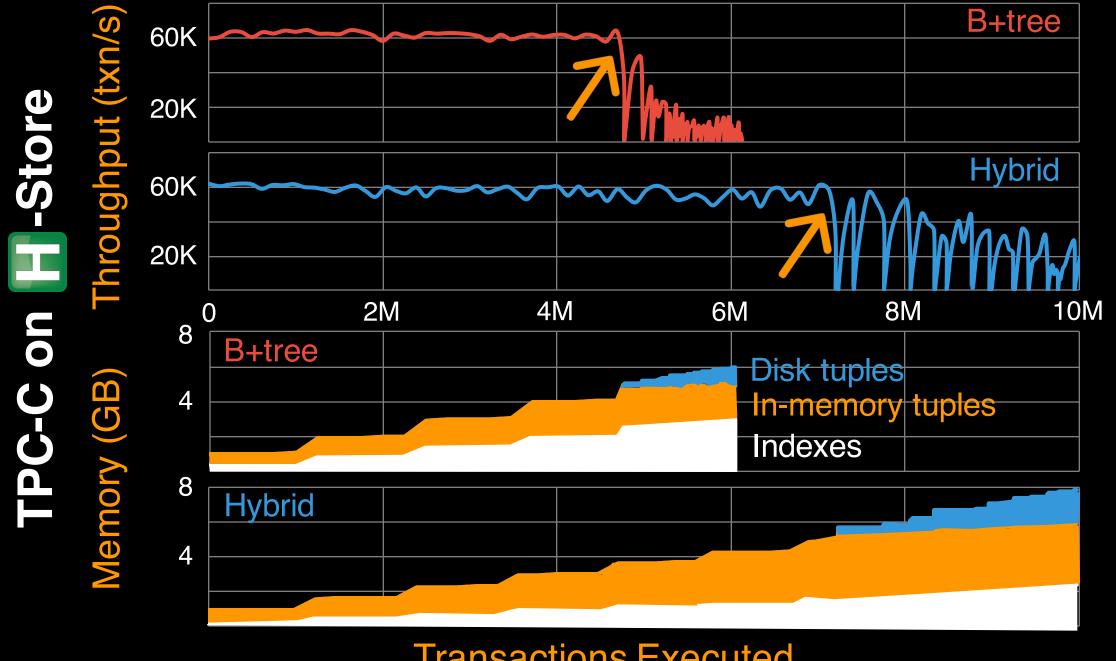
Transactions Executed



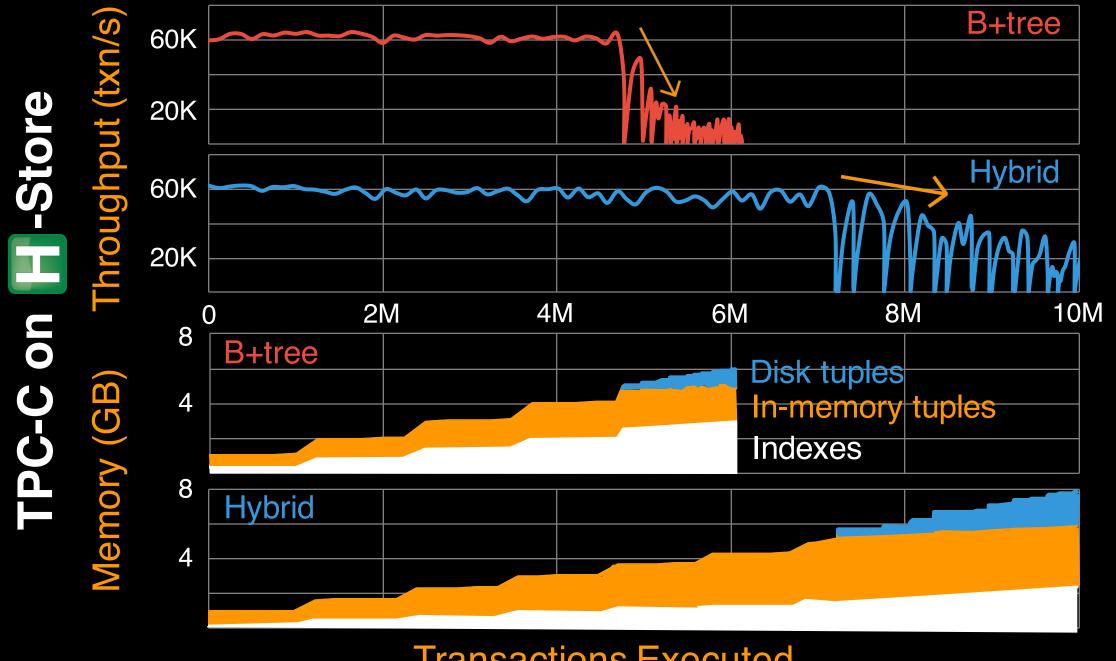
Transactions Executed



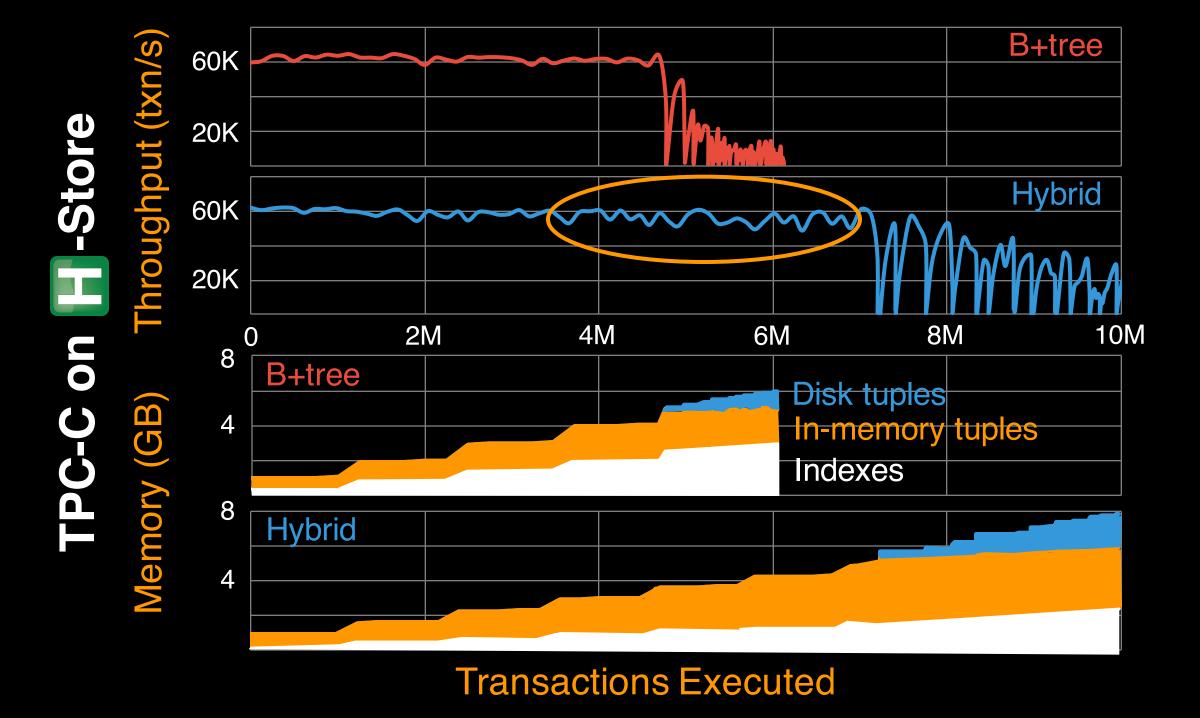
Transactions Executed

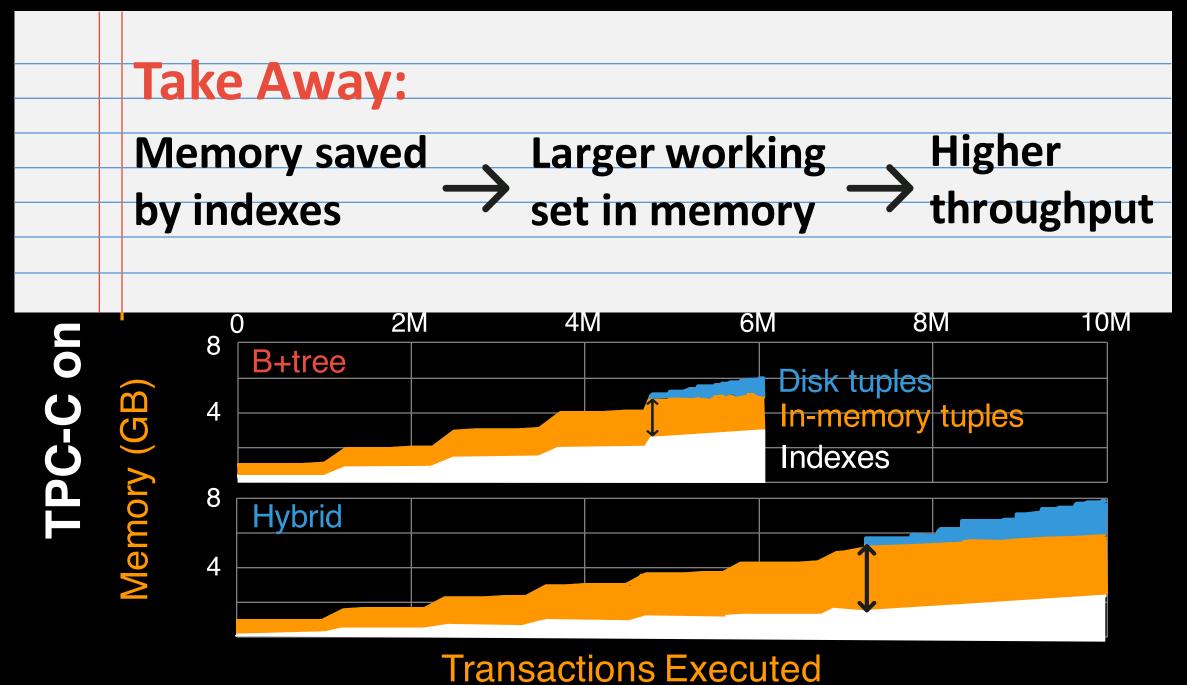


Transactions Executed



Transactions Executed





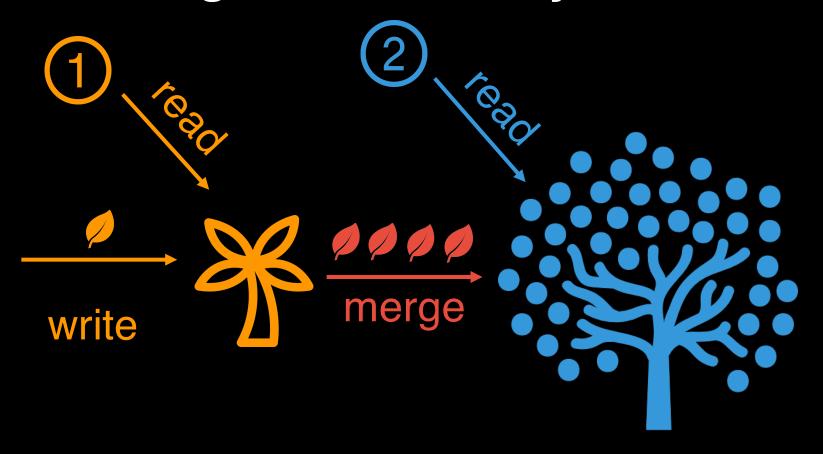
Part I Recap

- (1) The hybrid index architecture GENERAL
- (2) The Dual-Stage Transformation PRACTICAL
- (3) Applied to 4 index structures USEFUL
 - B+tree Skip List
 - Masstree Adaptive Radix Tree (ART)

Part II Concurrent hybrid indexes with non-blocking merge

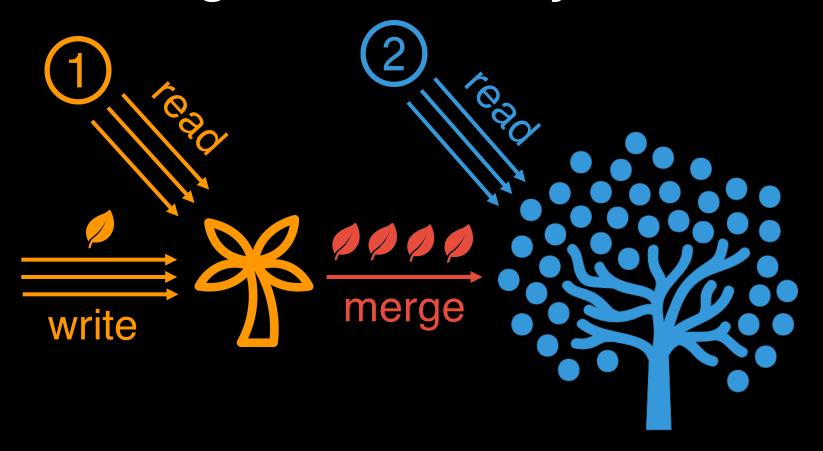


Building Concurrent Hybrid Index?



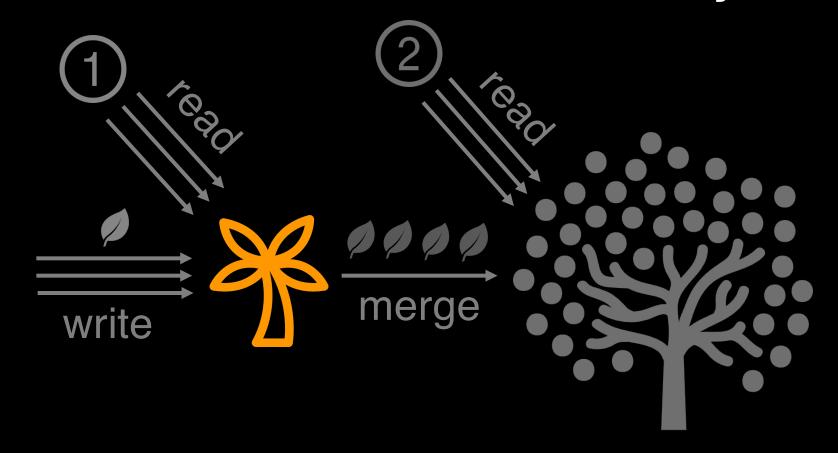
dynamic stage

Building Concurrent Hybrid Index?



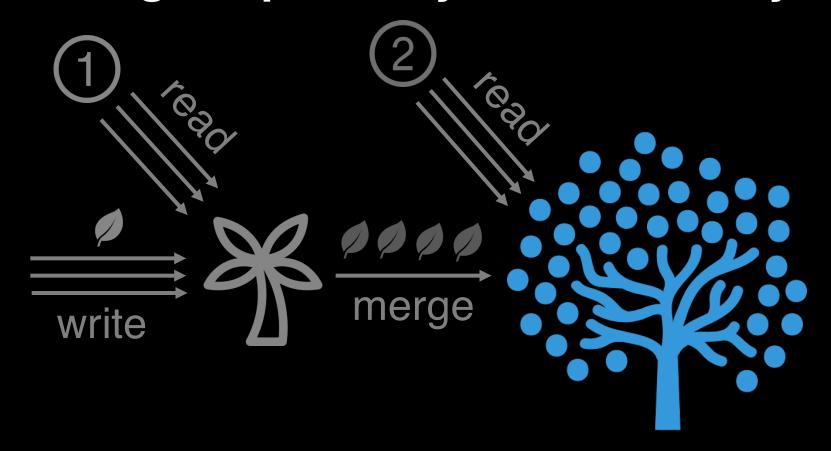
dynamic stage

Use concurrent data structures for dynamic-stage



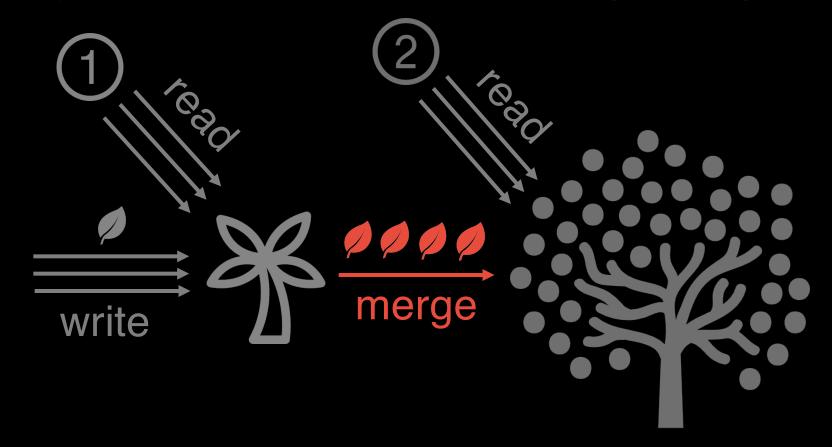
dynamic stage

Static-stage is perfectly concurrent by default



dynamic stage

Challenge: efficient non-blocking merge algorithm



dynamic stage

Merge Algorithm Requirements

1

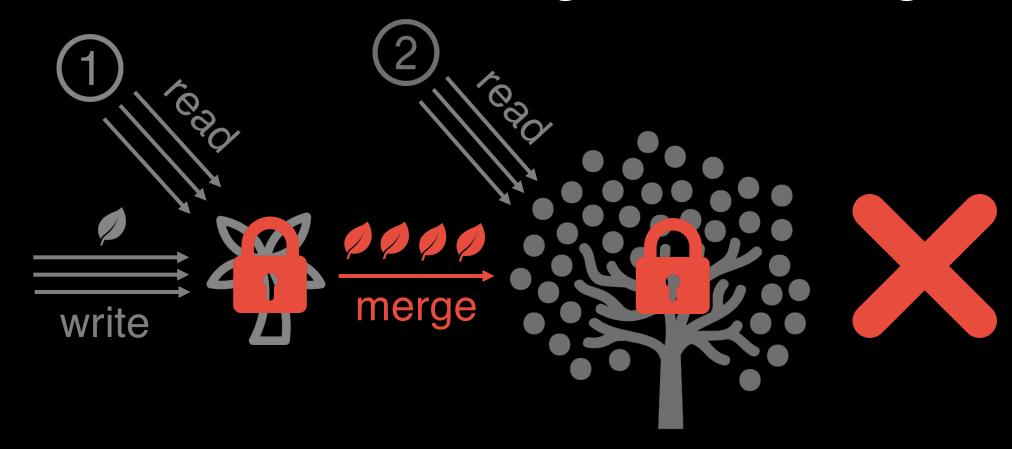
Non-blocking

- All existing items are accessible during merge
- New items can still enter
- 2

Efficient

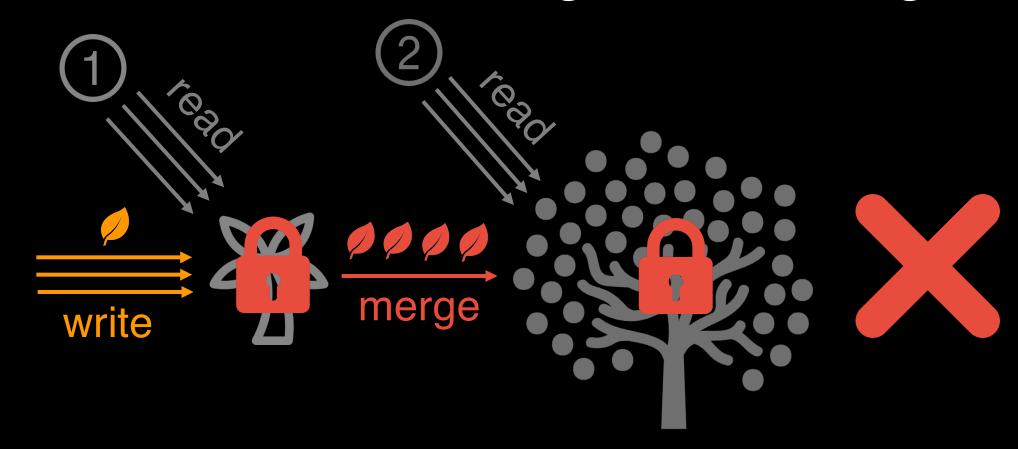
- Fast
- Bounded temporary memory use

Naïve Solution 1: Coarse-grained Locking



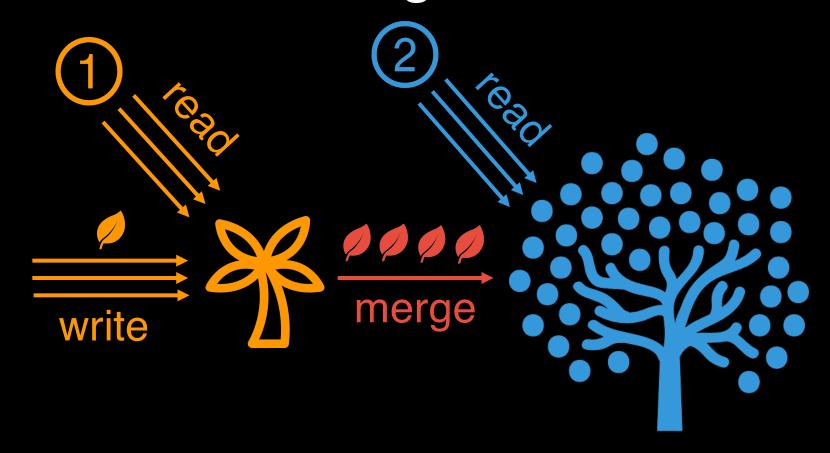
dynamic stage

Naïve Solution 1: Coarse-grained Locking



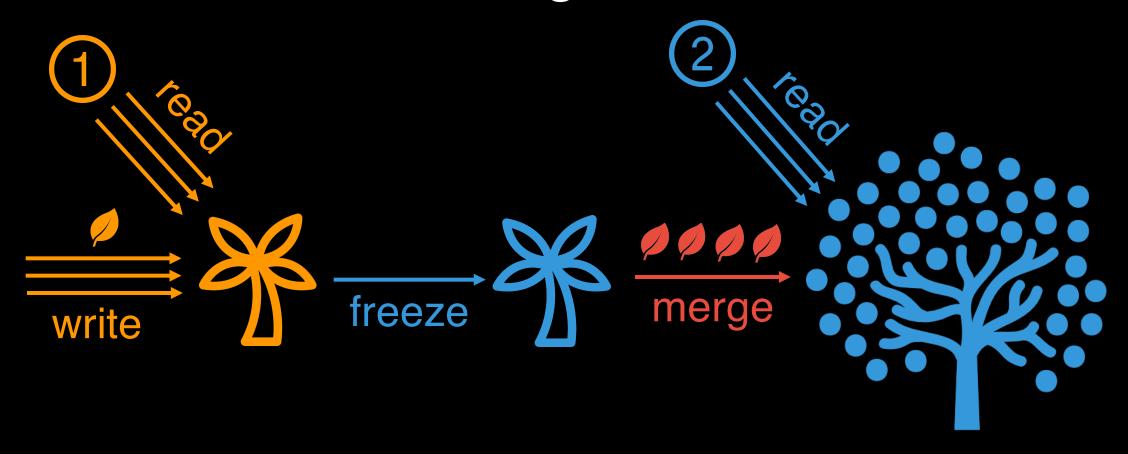
dynamic stage

The intermediate stage unblocks write traffic



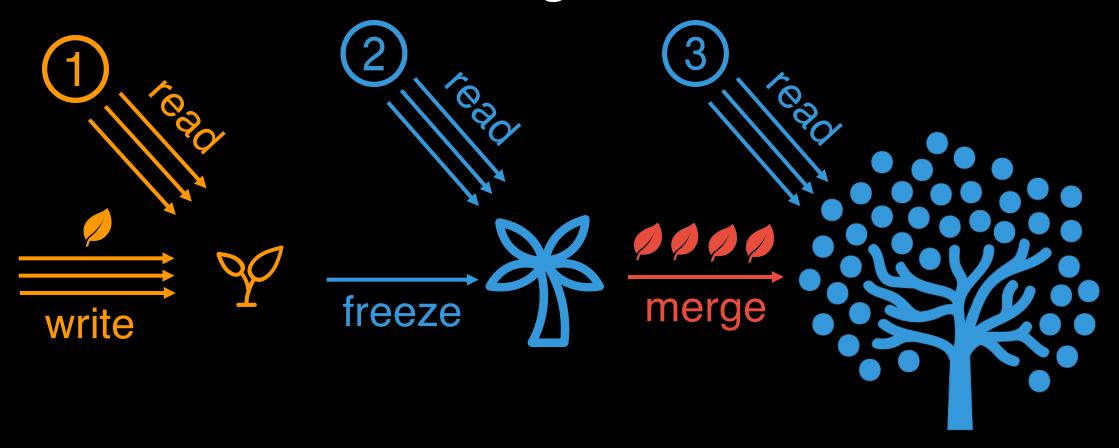
dynamic stage

The intermediate stage unblocks write traffic



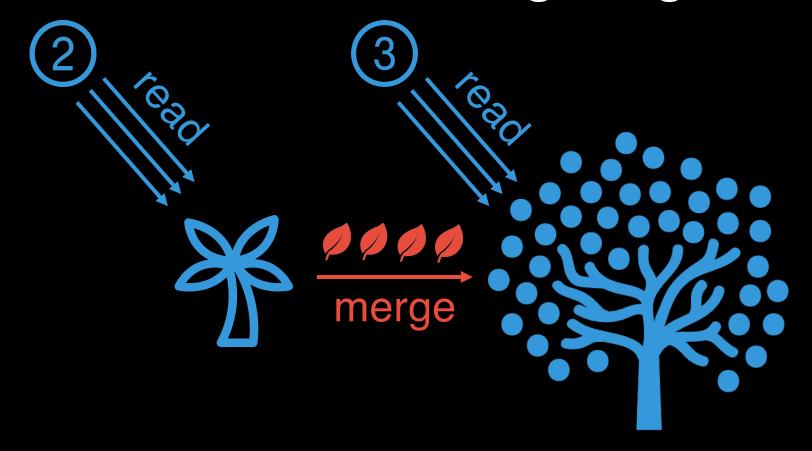
dynamic stage Intermediate stage

The intermediate stage unblocks write traffic



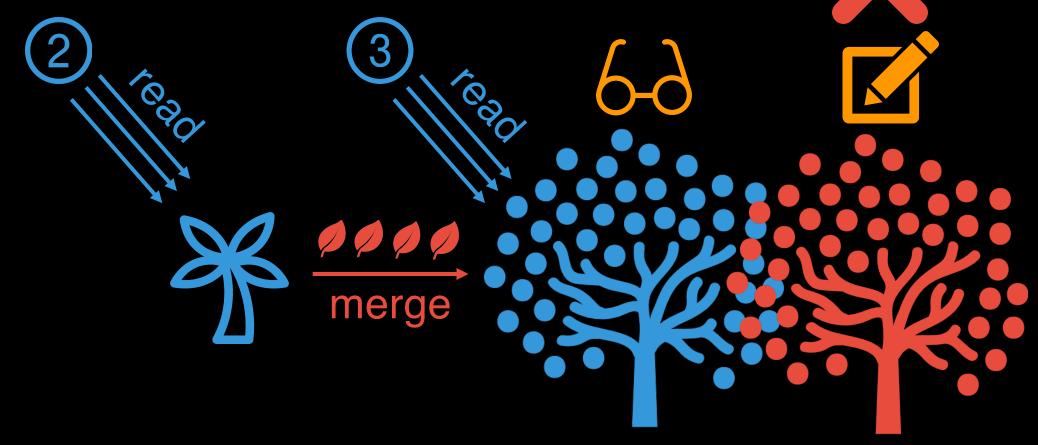
dynamic stage Intermediate stage

How do we unblock reads during merge?



Intermediate stage

Naïve Solution 2: Full Copy-on-write

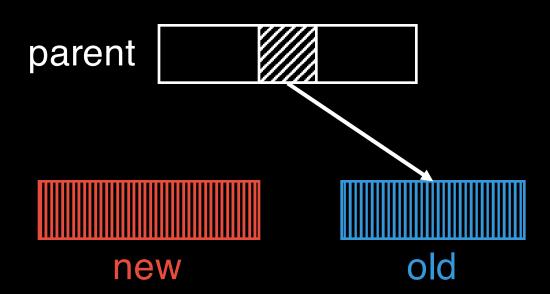


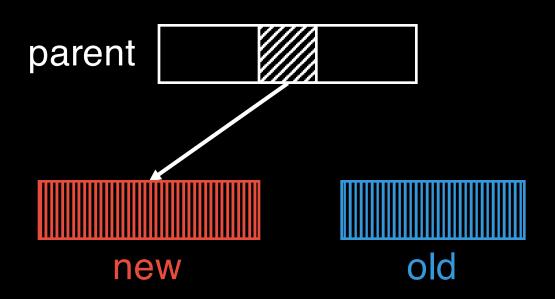
Intermediate stage

Key Observation

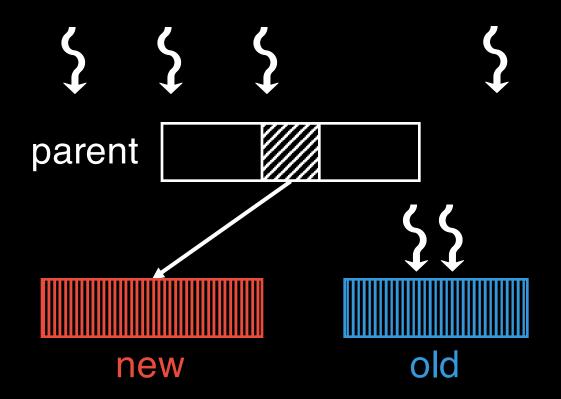
Merged-in items in the static-stage will NOT be accessed until the intermediate-stage is deleted



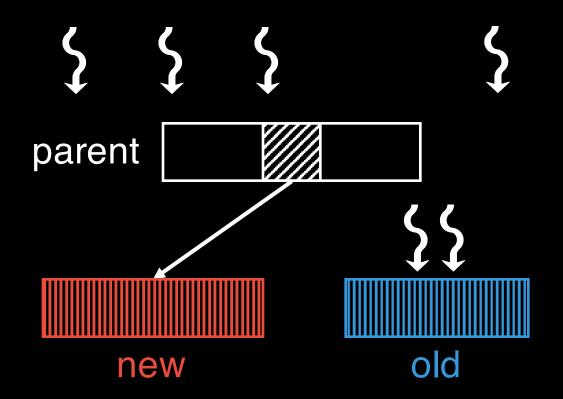




When can we safely reclaim the garbage?

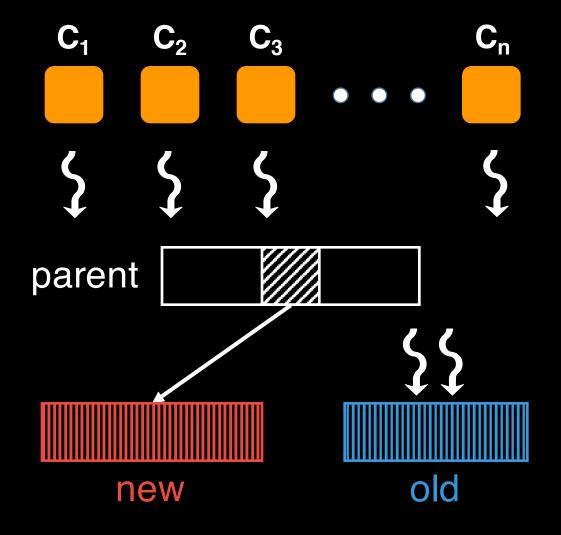


When can we safely reclaim the garbage?



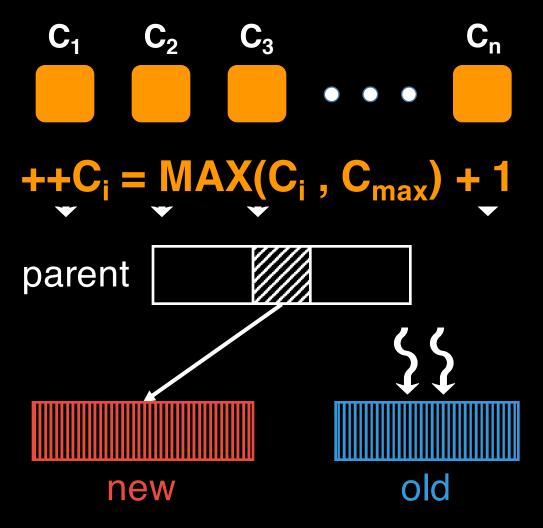
When no thread still holds a reference to it!

Thread-local counters

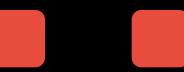


When no thread still holds a reference to it!







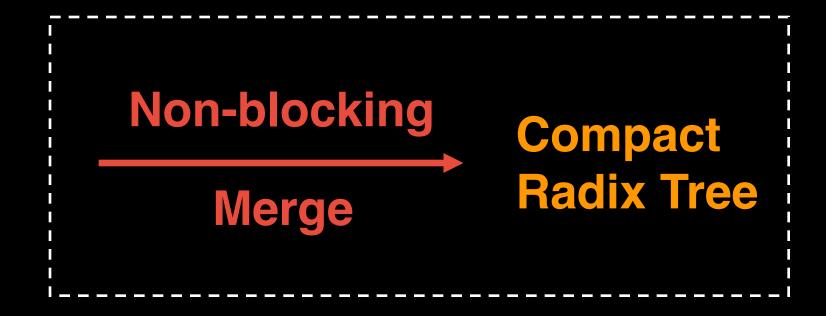


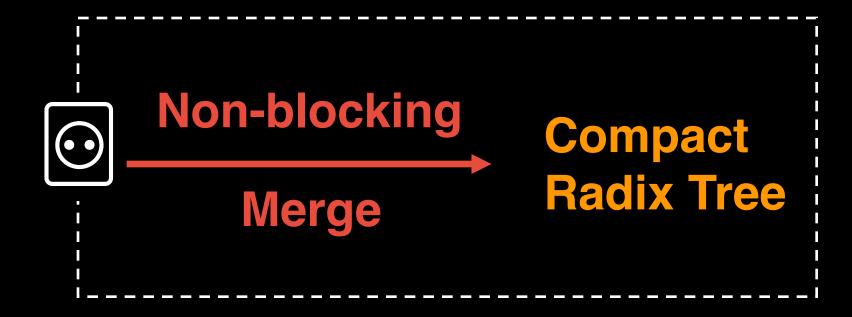
W6e6 on dition:still onld sagateragetagit!

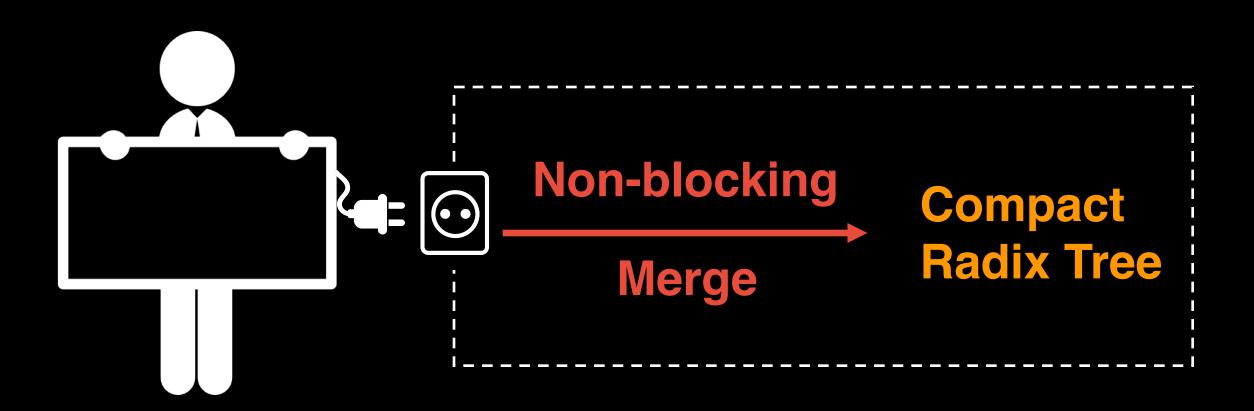
A Quick Recap of the Merge Algorithm

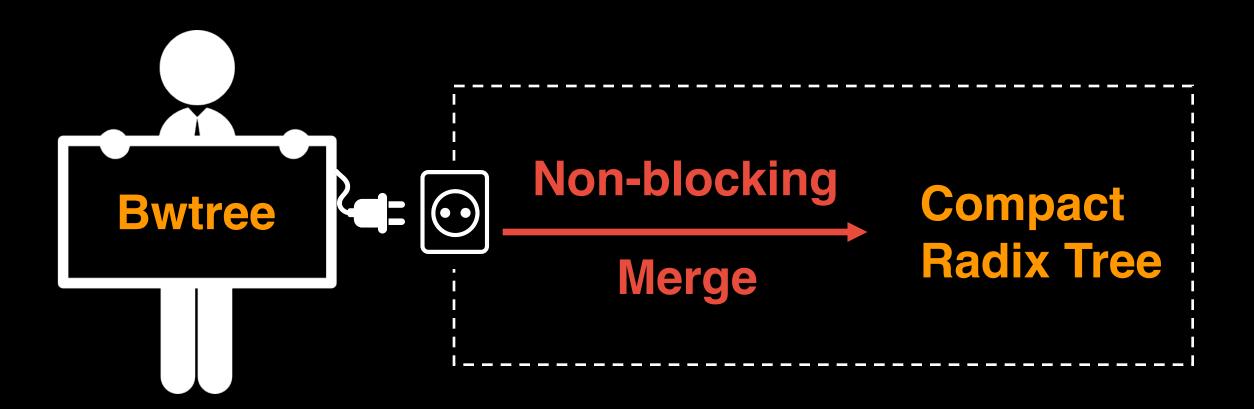
The intermediate stage separates writes from the merge process

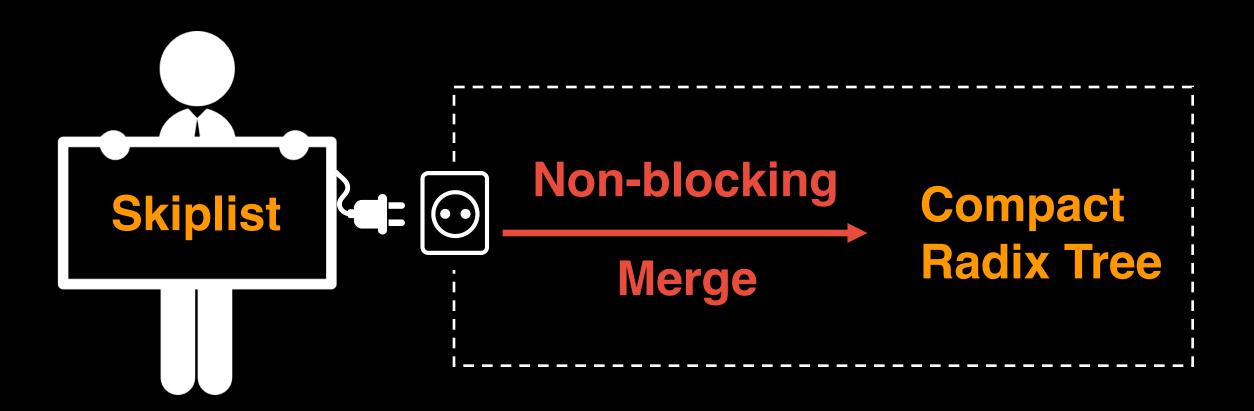
The incremental merge algorithm with rapid GC is non-blocking and space-efficient

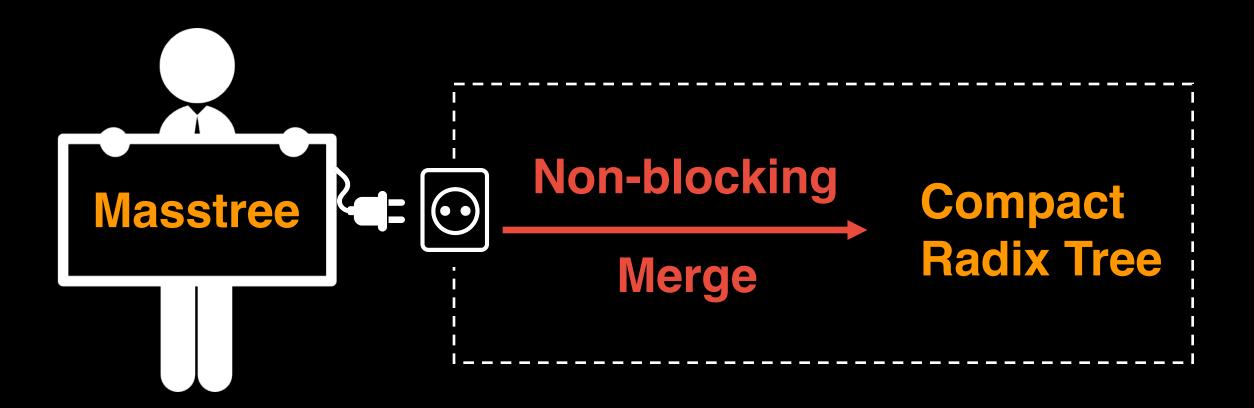












Part III Super-compact static-stage



Go "crazy" on space-efficiency

Succinct Data Structures

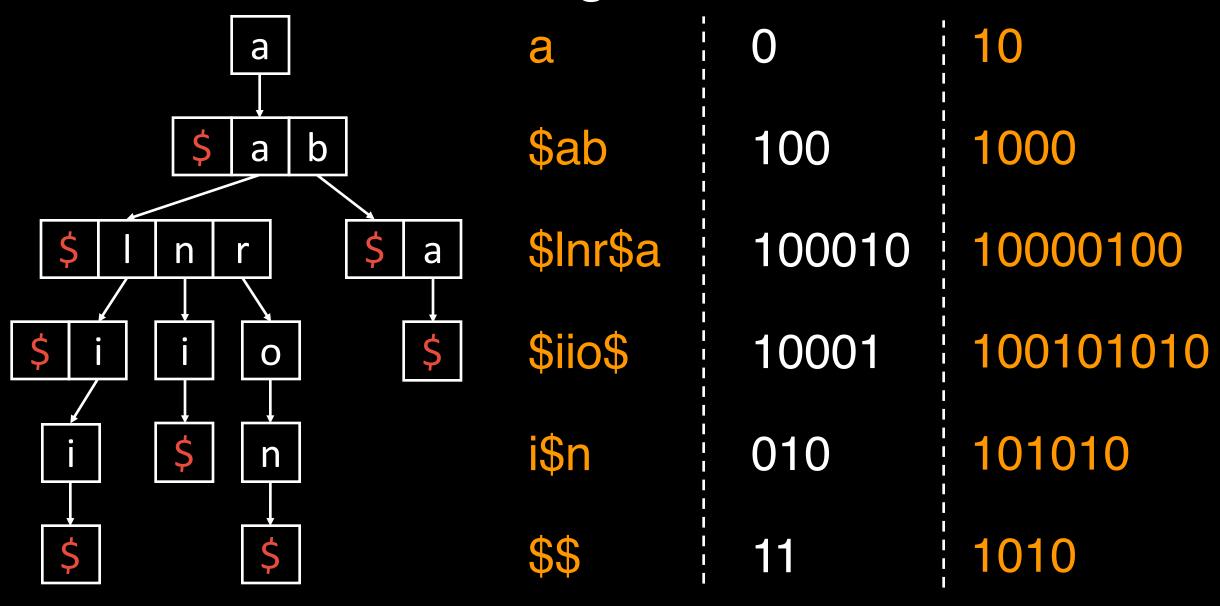
- Z + o(Z), where Z is the information-theoretic lower bound
- Still allow for efficient query operations

100011010000101...

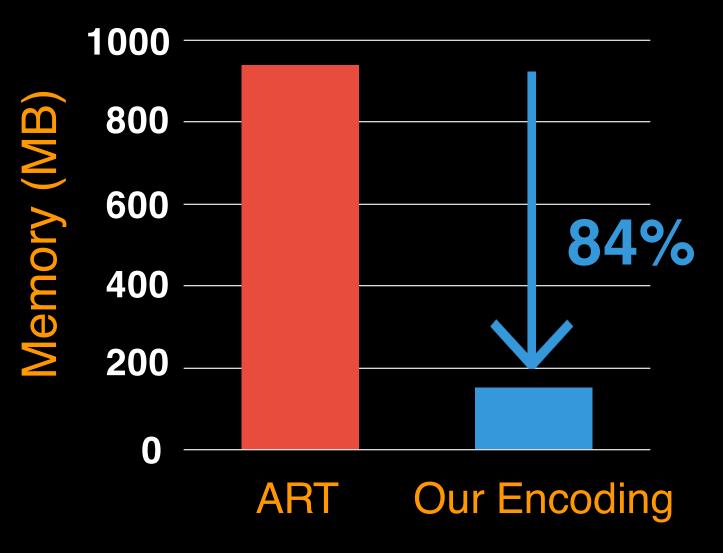
 $rank_1(x) = \#$ of 1's up to position x

 $select_1(x) = position of the x-occurrence of 1$

Encoding Radix Tree



Memory Savings with the New Encoding



50M email keys with average length = 20 bytes

The Takeaway Message

Hybrid indexes can save the precious memory resources with minimum performance penalty.

Toll-Free Hotline:



1-844-88-CMUDB

Back-up Slides

Latency (ms)

	B+tree	Hybrid
50%	10	10
99%	50	52
MAX	115	611

YCSB-based Microbenchmark Evaluation

Workload: insert, read/update(50/50)

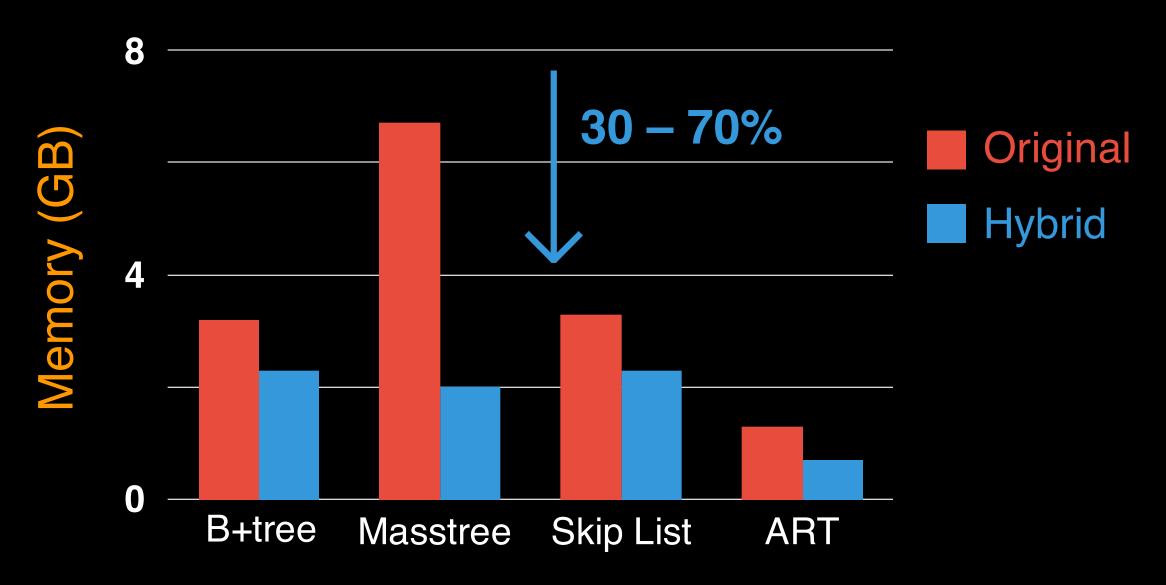
Key: email

Value: 64-bit unsigned integer (pointer)

Single thread

50M entries, 10M queries (Zipf distributed)

Hybrid index saves memory



Hybrid index provides comparable throughput

