# Tabular: Efficiently Building Efficient Indexes

Ziyi Yan, Mohammed Farouk Drira, Tianxun Hu, Tianzheng Wang



#### Concurrent data structure (index) is hard to build

- CPU/OS building blocks are hard to use
  - Synchronization primitives: latches, atomics
  - Locking protocol: pessimistic, optimistic, latch-free
  - Error-prone programming experience: deadlocks, memory leaks
- Need expertise in many neighboring areas
  - Memory management/Storage Devices/Networking
- Transaction semantics are easy to use
  - A transaction only need single-threaded logic
  - Concurrency control protocol guarantee ACID => Concurrent safety!
- DBMS engines have concurrency control but only for user transaction

Re-use concurrency control to implement concurrent indexes





Start a new transaction

```
txn = BeginTransaction()
```



Start a new transaction

```
txn = BeginTransaction()
```

Insert a new record



Start a new transaction

```
txn = BeginTransaction()
```

Insert a new record

Read the record by rid

```
record = txn.Read(table, rid)
```



Start a new transaction

```
txn = BeginTransaction()
```

Insert a new record

Read the record by rid

```
record = txn.Read(table, rid)
```

Update the record by rid

```
record.name = "tabular"
txn.Update(table, rid, record)
```



Start a new transaction

```
txn = BeginTransaction()
```

Insert a new record

• Read the record by rid

```
record = txn.Read(table, rid)
```

Update the record by rid

```
record.name = "tabular"
txn.Update(table, rid, record)
```

Commit the transaction

```
txn.Commit()
```



Start a new transaction

```
txn = BeginTransaction()
```

Insert a new record

• Read the record by rid record = txn.Read(table, rid)

• Update the record by rid

```
record.name = "tabular"
txn.Update(table, rid, record)
```

Commit the transaction

```
txn.Commit()
```

```
int COUNT;
void foo() {
   COUNT++;
}
```



Start a new transaction

```
txn = BeginTransaction()
```

Insert a new record

Read the record by rid record = txn.Read(table, rid)

Update the record by rid

```
record.name = "tabular"
txn.Update(table, rid, record)
```

Commit the transaction

```
txn.Commit()
```

```
int COUNT;

void foo() {
    COUNT++;
}

mutex.lock();
    COUNT++;
    mutex.unlock();
}
```



Start a new transaction

```
txn = BeginTransaction()
```

Insert a new record

- Read the record by rid record = txn.Read(table, rid)
- Update the record by rid

```
record.name = "tabular"
txn.Update(table, rid, record)
```

Commit the transaction

```
txn.Commit()
```

```
std::mutex mutex;
int COUNT;
                 void foo_lock() {
void foo() {
                   mutex.lock();
  COUNT++;
                   COUNT++;
                   mutex.unlock();
void foo tabular() {
  cnt = txn.Read(tbl, cnt id);
  cnt++;
  txn.Update(tbl, cnt_id, cnt);
```

Start a new transaction

```
txn = BeginTransaction()
```

Insert a new record

- Read the record by rid record = txn.Read(table, rid)
- Update the record by rid

```
record.name = "tabular"
txn.Update(table, rid, record)
```

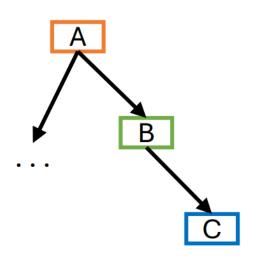
Commit the transaction

```
txn.Commit()
```

```
std::mutex mutex;
int COUNT;
                 void foo lock() {
void foo() {
                   mutex.lock();
  COUNT++;
                   COUNT++;
                   mutex.unlock();
void foo tabular() {
  txn = BeginTransaction()
  cnt = txn.Read(tbl, cnt id);
  cnt++;
  txn.Update(tbl, cnt_id, cnt);
  txn.Commit()
```

#### B+-tree on transactional tables

- Each node is stored as a record on table
- Each field in struct definition of BTreeNode becomes a column in schema definition
- Nodes are "connected" via
   RIDs instead of pointers



(a) Logical view

```
int capacity;
int n_keys;
bool is_leaf;
int lock;
BTreeNode *right_child;
KVPair kvs[16];
};
```

(b) C/C++ struct/class definition

RID	n_keys	is_leaf	right_child	kvs	
0	10	False	1	• • •	
1	8	False	2	• • •	
2	4	True	INVALID_RID	• • •	
• •					

(c) Table storing BTreeNode objects



#### Single-threaded code

```
1 def BTree::lookup(Key k, Value *out_v):
2    n = root
3    while n is inner:
4     next = n.findChild(n, k)
5    n = next
6    *out_v = n.findValue(k)
```



#### Single-threaded code

```
1 def BTree::lookup(Key k, Value *out_v):
2    n = root
3    while n is inner:
4     next = n.findChild(n, k)
5    n = next
6    *out v = n.findValue(k)
```

#### Hand-crafted code using Optimistic Lock Coupling (OLC)

```
1 def BTree::lookup(Key k, Value *out v):
 Epoch manager
                       2 restart:
implementation ...
                       3 epoch enter()
                                                  Optimistic lock
                       4 retry = false
                                                 implementation...
         read_lock()
                           n = root
                           ver = n.read lock(retry)
                           if retry or n != root: goto restart
                           while n is inner:
                             next = n.children[findChild(n, k)]
                      10
                             n.verify read(ver, retry)
          read_lock() 11
                             if retry: goto restart
                             ver next = next.read lock(retry)
                             if retry: goto restart
                             n = next
                      15
                           *out v = n.findValue(k)
        Retry path
                      16 epoch exit()
```

#### Single-threaded code

```
1 def BTree::lookup(Key k, Value *out_v):
2    n = root
3    while n is inner:
4        next = n.findChild(n, k)
5        n = next
6    *out_v = n.findValue(k)
```

#### Tabular-based code

```
1 def BTree::lookup(Key k, Value *out_v):
2    t = BeginTransaction()
3    n = t.Read(table, root_id)
4    while n is inner:
5        next_id = n.findChild(n, k)
6        n = t.Read(table, next_id)
7    *out_v = n.findValue(k)
8    t.Commit()
```

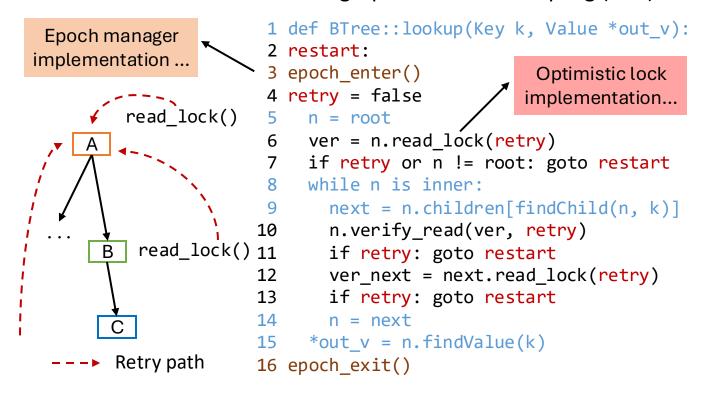
#### Hand-crafted code using Optimistic Lock Coupling (OLC)

```
1 def BTree::lookup(Key k, Value *out v):
 Epoch manager
                       2 restart:
implementation ...
                       3 epoch enter()
                                                  Optimistic lock
                       4 retry = false
                                                 implementation...
        read_lock()
                           n = root
                           ver = n.read lock(retry)
                           if retry or n != root: goto restart
                           while n is inner:
                             next = n.children[findChild(n, k)]
                      10
                             n.verify read(ver, retry)
          read_lock() 11
                             if retry: goto restart
                             ver next = next.read lock(retry)
                             if retry: goto restart
                             n = next
                      15
                           *out v = n.findValue(k)
        Retry path
                      16 epoch exit()
```

#### Single-threaded code

```
1 def BTree::lookup(Key k, Value *out v):
   n = root
   while n is inner:
     next = n.findChild(n, k)
     n = next
    *out v = n.findValue(k)
         Tabular-based code
1 def BTree::lookup(Key k, Value *out v):
   t = BeginTransaction()
   n = t.Read(table, root id)
   while n is inner:
     next id = n.findChild(n, k)
     n = t.Read(table, next id)
   *out v = n.findValue(k)
   t.Commit()
```

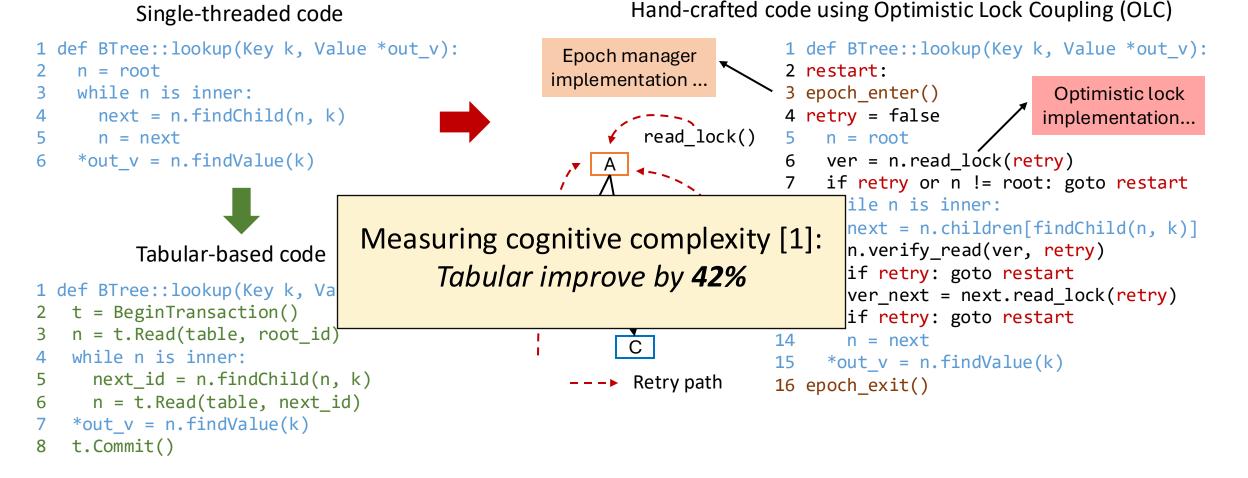
Hand-crafted code using Optimistic Lock Coupling (OLC)



Easy to develop, debug and maintain

Lots of branches, massive mental overhead





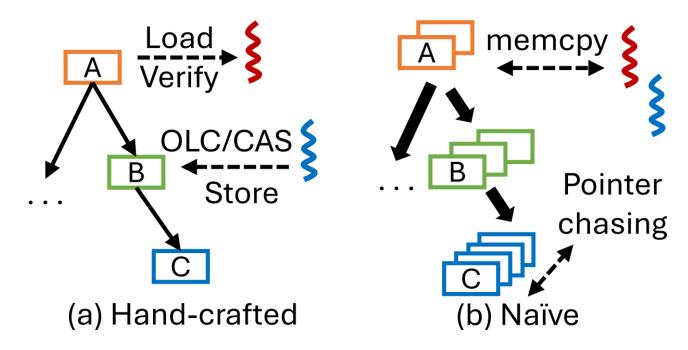
Easy to develop, debug and maintain

Lots of branches, massive mental overhead

[1]: G Ann Campbell. 2018. Cognitive complexity: An overview and evaluation. In Proceedings of the 2018 international conference on technical debt. 57–58.

### Naïve objects-on-DB performance is bad

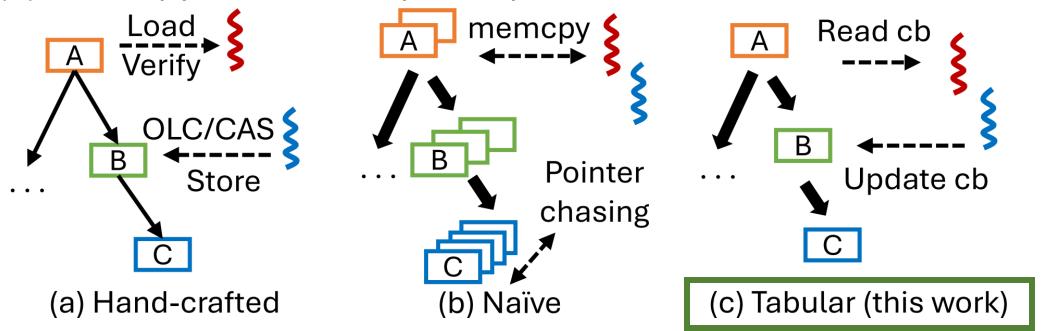
- Traditional concurrency control (CC) scheme incurs unnecessary overhead
  - Optimistic CC (OCC) needs memory copying to transaction-local memory
  - Multi-versioned CC (MVCC) stalls at pointer-chasing while accessing records





# Naïve objects-on-DB performance is bad

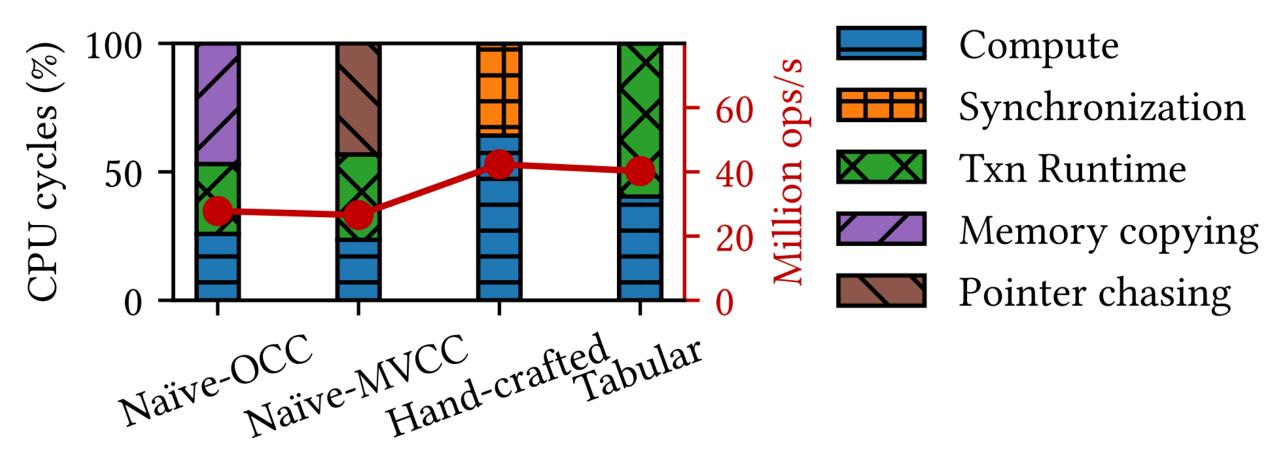
- Traditional concurrency control (CC) scheme incurs unnecessary overhead
  - Optimistic CC (OCC) needs memory copying to transaction-local memory
  - Multi-versioned CC (MVCC) stalls at pointer-chasing while accessing records
- Tabular (this work) adopts: (1) single-versioned OCC for direct record access and (2) zero-copy reads and in-place updates via callback functions





### Naïve objects-on-DB cycle breakdowns

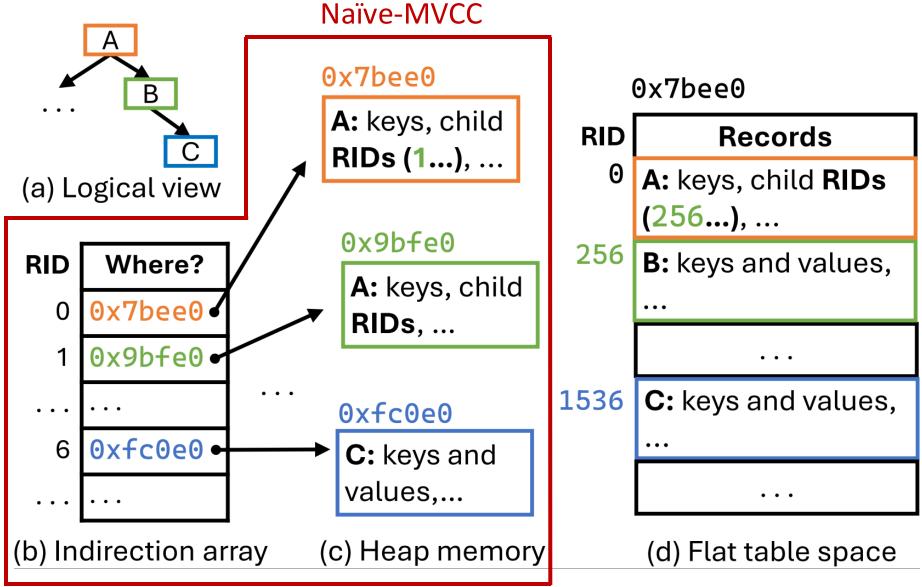
- Unnecessary overhead: Memory copying, Pointer chasing
- They took nearly 50% of cycles





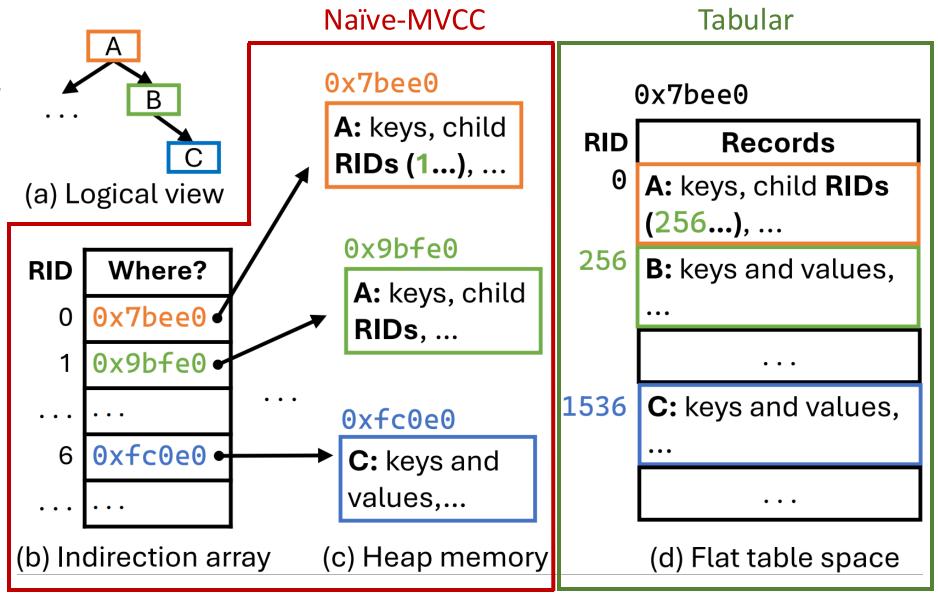
#### Tabular: Single-Versioned, Direct Record Access

- RID is the offset from base pointer of flat memory
- Records are stored directly on the table instead of referenced via pointers
- It saves ~19-29% cycles in B+-tree



#### Tabular: Single-Versioned, Direct Record Access

- RID is the offset from base pointer of flat memory
- Records are stored directly on the table instead of referenced via pointers
- It saves ~19-29% cycles in B+-tree



#### Tabular: Copy-based vs Callback-based interface

 OCC requires memory copying from tables to transactions

```
# copy-based
node = txn.Read(table, rid)
is_full = node.num_entries == MAX_ENTRIES

# callback-based
var is_full: bool
def callback(record):
    is_full = node.num_entries == MAX_ENTRIES
txn.ReadCallback(table, rid, callback)
```

#### Tabular: Zero-Copy Record Access

- OCC requires memory copying from tables to transactions
- Observation: Record access can be re-tried cheaper compared to memory copy
- User-defined callback 9
   executes directly on the 10
   records 11
- It saves ~20-31% cycles in B+-tree

```
1 def Transaction::ReadCallback(table, rid, callback):
    # Take a reference (pointer) to the record
    Record & record = table.GetRecord(rid);
4 retry:
    tid = record.get_consistent_tid()
    callback(record) # Execute the callback in-place
6
    # Verify the record did not change
8
    tid_now = record.get_consistent_tid()
9
    if (tid_now == tid):
      read_set.Add(record, tid);
11
    else:
12
13
      goto retry
```

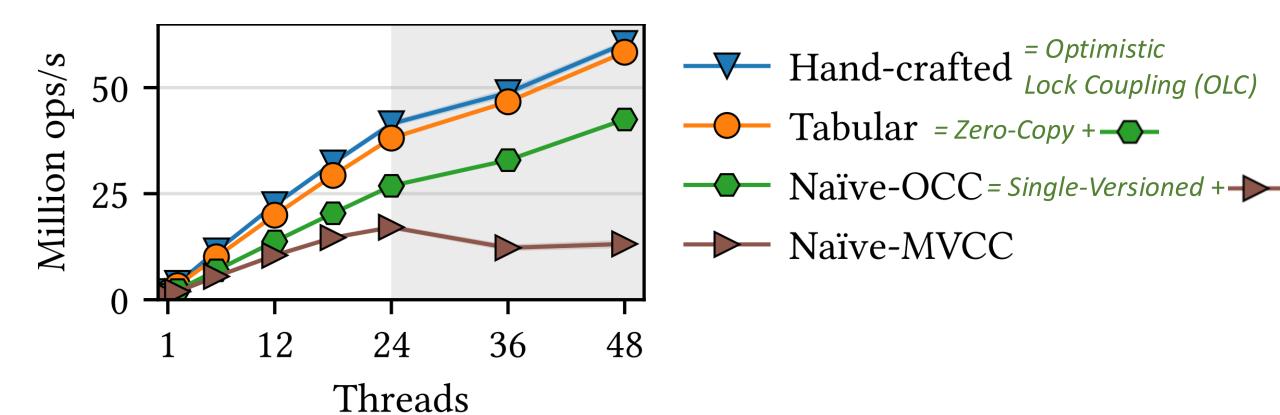
#### Tabular: Zero-Copy Record Access

- OCC requires memory copying from tables to transactions
- Observation: Record access can be re-tried cheaper compared to memory copy
- User-defined callback 9
   executes directly on the 10
   records 11
- It saves ~20-31% cycles in B+-tree

```
1 def Transaction::ReadCallback(table, rid, callback):
    # Take a reference (pointer) to the record
    Record & record = table.GetRecord(rid);
4 retry:
5
    tid = record.get_consistent_tid()
    callback(record) # Execute the callback in-place
6
    # Verify the record did not change
8
    tid_now = record.get_consistent_tid()
9
    if (tid_now == tid):
      read_set.Add(record, tid);
11
    else:
12
13
      goto retry
```

#### Tabular achieves competitive performance

- Workload: 100M of 8B-key B+-tree running 50% lookups and 50% updates
- Tabular performs at ~90% of Hand-crafted (Upper bound)





# Tabular: ACI(D) Guarantee



### Tabular: ACI(D) Guarantee

#### Correctness

- Achieve serializability the same way as OCC
- Need some tweaks in commit protocol when UpdateCallback() is involved

```
def UpdateCallback(table, rid, cb):
    write_set.Add(table, rid, cb)
```

```
1 def Transaction::Commit(table, rid, callback):
     # Sort out-of-place and lock write set
     sorted_ws = sort(write_set)
     foreach w in sorted_ws:
      lock w
 6
     ... fences and get commit epoch ...
 8
     # Validate reads - exactly the same as OCC
    foreach r in read set:
10
      if r.tid != r.record.get_tid() and r is not locked by me:
11
        rollback and return
12
13
     ... iterate read/write sets to get new commit_tid ...
14
15
     # Apply write callbacks and updates
     foreach w in write_set:
17
      if w.is_callback:
18
        w.callback() # invoke the callback
19
      else: # "normal" materialized update
21
        memcpy(w.record.data, w.data, w.size)
      w.tid = commit_tid
23
     ... other remaining operations ...
```

## Tabular: ACI(D) Guarantee

#### Correctness

- Achieve serializability the same way as OCC
- Need some tweaks in commit protocol when UpdateCallback() is involved

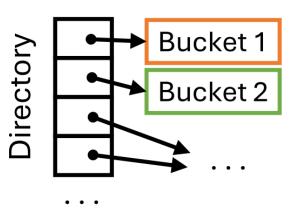
```
def UpdateCallback(table, rid, cb):
    write_set.Add(table, rid, cb)
```

- Transparent durability
  - Optionally via a per-table knob
  - Thread-local distributed logging and recovery

```
1 def Transaction::Commit(table, rid, callback):
     # Sort out-of-place and lock write set
     sorted_ws = sort(write_set)
     foreach w in sorted_ws:
      lock w
 6
     ... fences and get commit epoch ...
 8
     # Validate reads - exactly the same as OCC
    foreach r in read set:
10
      if r.tid != r.record.get_tid() and r is not locked by me:
11
12
        rollback and return
13
     ... iterate read/write sets to get new commit_tid ...
14
15
     # Apply write callbacks and updates
     foreach w in write_set:
17
18
      if w.is_callback:
        w.callback() # invoke the callback
19
      else: # "normal" materialized update
20
21
        memcpy(w.record.data, w.data, w.size)
      w.tid = commit_tid
23
     ... other remaining operations ...
```

### Case study: Tabular-based Hash-table

- Buckets are fixed-size records in a table
- Directory is stored in a separate table
- Multiple tables
   as one data
   structure is
   supported
   trivally



(a) Logical view

<pre>struct Bucket {</pre>						
<pre>int n_keys;</pre>						
<pre>int local_depth;</pre>						
<pre>int lock;</pre>						
<pre>KVPair kvs[16];</pre>						
<b>}</b> ;						

(b) Bucket definition

<pre>struct Directory {</pre>
<pre>int global_depth;</pre>
<pre>Bucket *buckets[];</pre>
<pre>int lock;</pre>
<b>}</b> ;

(c) Directory definition

RID	<pre>global_depth</pre>	buckets
0	0	•••
1	1	•••
2	2	•••
3	3	•••

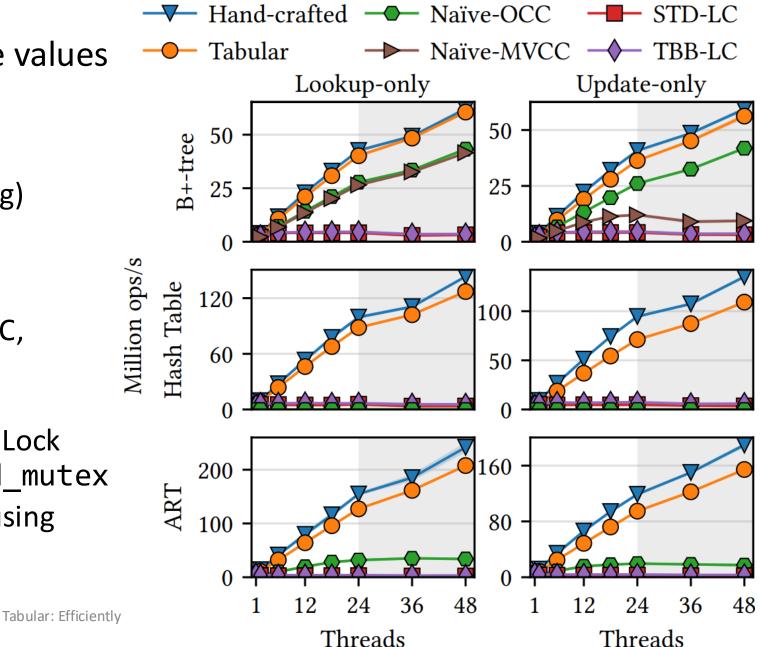
(d) Directory table

RID	n_keys	local_depth	kvs		
0	10	2	• • •		
1	8	2	• • •		
2	4	3	• • •		
• • •					

(e) Bucket table

#### Evaluation: YCSB-like index benchmarks

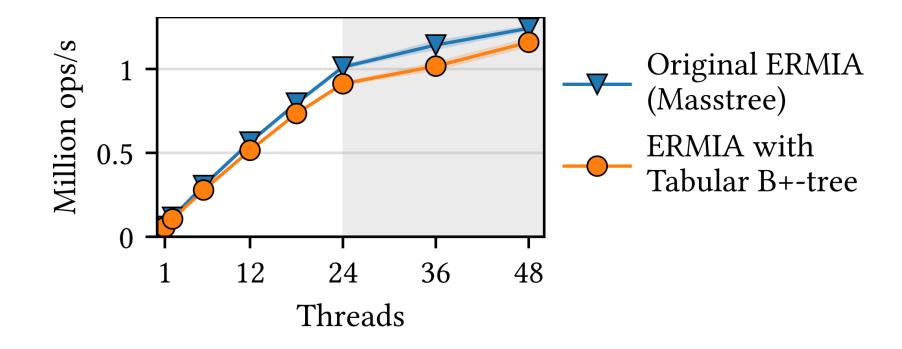
- 100M 8-byte keys and 8-byte values
- Index data structures
  - B+-tree
  - Hash-table (Extendible hashing)
  - Adaptive Radix Tree (ART)
- Variants
  - Hand-crafted, Naïve-MVCC, Naïve-OCC, Tabular are as aforementioned
  - STD-LC: Applying Pessimistic Lock Coupling using std::shared\_mutex
  - TBB-LC: Same as above but using tbb::spin\_rw\_mutex





#### Evaluation: TPC-C

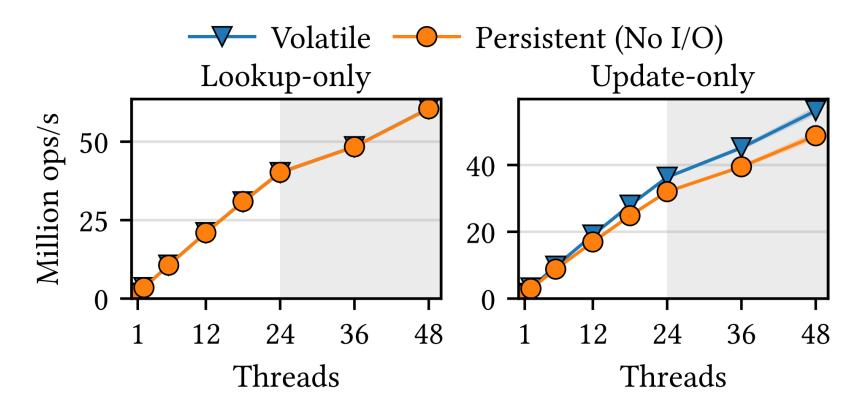
- ERMIA is a memory-resident OLTP engine
- Replace Masstree with Tabular-based B+-tree in ERMIA
- End-to-end performance is comparable to the original





#### Evaluation: Transparent Persistence

- In Persistent (No I/O), file operations are noops to show in-memory overhead
- Nearly zero overhead for lookups
- Less than 15% drop for updates





### Summary

- Tabular is a parallel programming library with transactional interfaces
- Data structures (not only indexes) can be modelled as tables to be transparently concurrent and persistent
- Tabular-based indexes deliver competitive performance compared to their hand-crafted counterparts with drastically lower programming complexity

• Code is available at <a href="https://github.com/sfu-dis/tabular">https://github.com/sfu-dis/tabular</a>

