Tabular: Efficiently Building Efficient Indexes

Ziyi Yan Mohamed Farouk Drira Tianxun Hu Tianzheng Wang





Building DBMSs: a Black Art

- Hand-crafted data structures
 - Hard to code up, hard to debug
 - Pretend correctness until it crashes

Know the basics



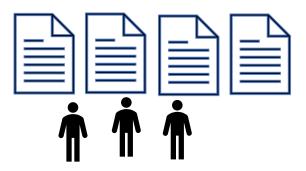


Know the hardware





+ DBMS research



+ Neighbouring areas



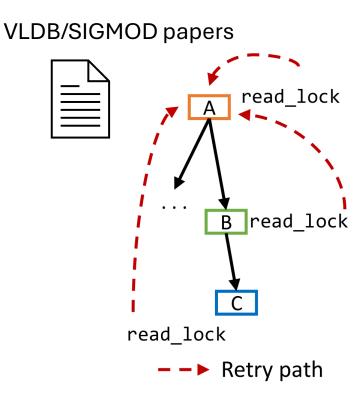


Need to know multiple areas very well; rare



Let's hand-craft a concurrent B+-tree

Memory-friendly layout + optimistic lock coupling



```
1 bool BTree::Insert(Key k, Value v):
 2 restart:
                             Epoch manager
 3 epoch enter() -
                            implementation...
 4 retry = false
     n = root
    ver = n.read lock(retry)
    if retry or n != root: goto restart
    while n is inner:
                                               PPoPP/SIGMOD/VLDB papers
      next = n.children[findChild(n, k)]
      n.verify read(ver, retry)
      if retry: goto restart
      ver next = next.read lock(retry)
      if retry: goto restart
      if next is inner node:
                                     Optimistic lock
      else: ...
                                    implementation...
16 ...
17 epoch exit()
```

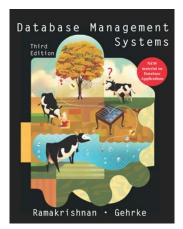
PPoPP/SOSP papers

Complex stuff!



Building DBMSs: a Black Art

Know the basics





Know the hardware





But many more know how to use a DBMS!

+ DBMS internal basics





+ Neighbouring areas

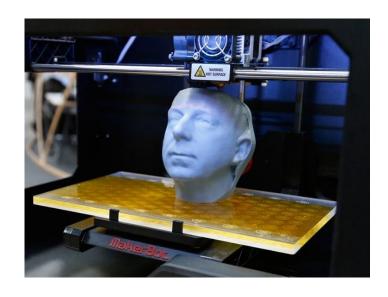




- Low-level programming
 - Hard to code up, hard to debug
 - Pretend correctness until it crashes



DB Transactions vs. Manual Parallel Programming



- Transactions
- Tables and records
- Transparent write-ahead-logging
- Buffered tables



- Locks/lock-free algorithms
- Explicit memory management
- Manual persistence, I/O
- Manual caching solutions

Can we do the same for parallel programming?

* To Lock, Swap or Elide: On the Interplay of Hardware Transactional Memory and Lock-free Indexing, VLDB 2015



Modelling Data Structures as Relational Tables

struct BTreeNode { int n_keys; bool is_leaf; int lock; BTreeNode *right_child; KVPair kvs[16]; }:

A "B-tree node table"

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

- Struct/class → Table
 - Member variables → Table columns
 - Defined by a schema, just like a DB app
- Instances of struct/class → Table records
- - An OLTP engine takes care of concurrency and persistence



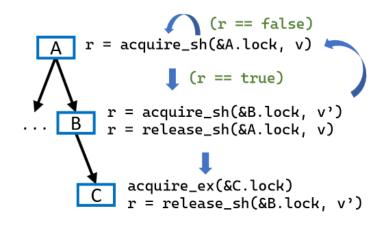
Modelling Data Structures as Relational Tables

Before

Allocation:

```
struct *mynode = (Node *)malloc(sizeof(Node));
```

Concurrency:





After

Allocation:

```
struct *mynode = table.insert(...)
```

Concurrency:

```
Transaction {
  table.read(A);
  table.read(B);
  table.read(C);
  ...
  table.update(C);
}
```

```
Single-
threaded logic

Just need
Cowbook
```

"Wasn't this tried and failed?"

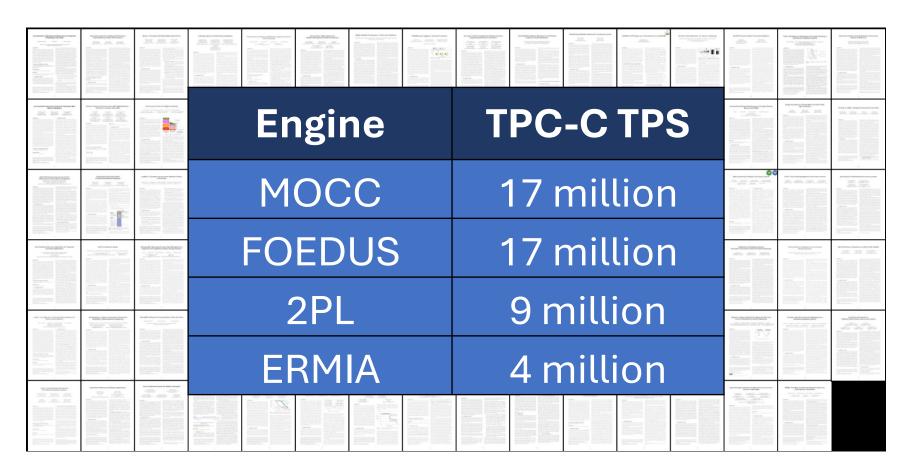


Past Failed Aspirations

- Object-on-DB / Object-oriented DBMSs
 - High overhead over full RDBMSs
- Hardware transactional memory (HTM)
 - No persistence
 - Hard to use "weird" aborts
- Software transactional memory (STM)
 - High overhead
 - No persistence
 - "Research toy"



It's different this time: in-memory OLTP since 2010+



On a 16-socket, 288-core HPE server

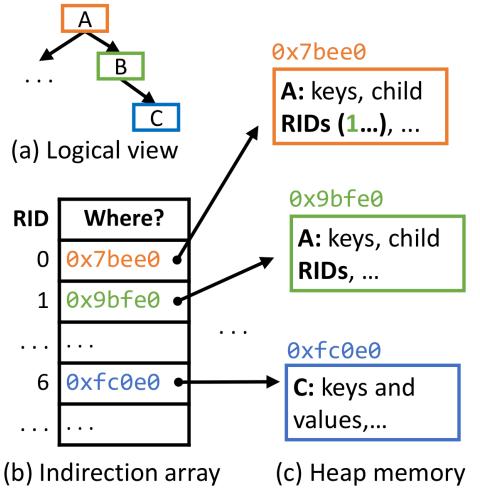
"Who need these?"

"Very few!"

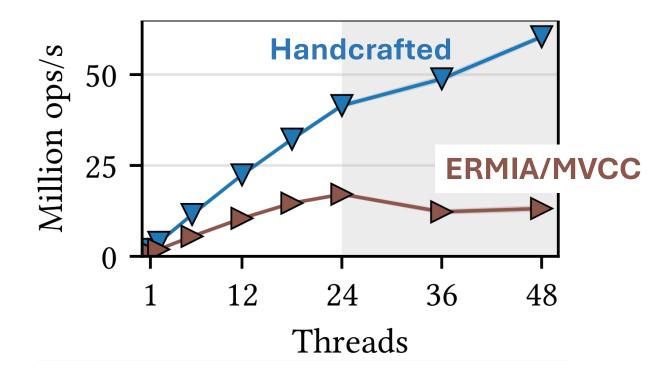
Much more headroom for objects-over-DB



First try: B+-tree over ERMIA*



Pointer chasing dominates: 80% slower



Multi-versioning considered harmful

* ERMIA: Fast Memory-Optimized Database System for Heterogeneous Workloads, SIGMOD 2016

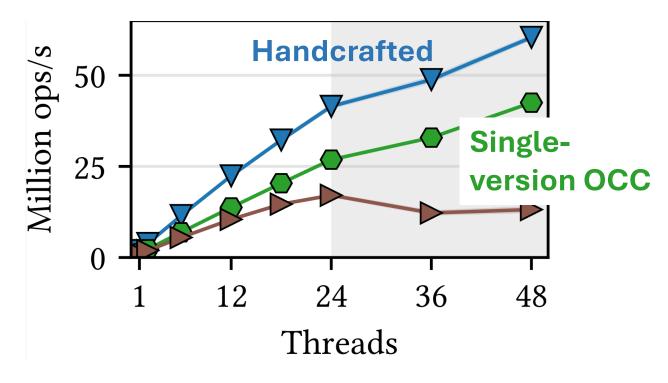


One more time: B+-tree over Single-Versioned OCC

Just plain flat space

0x7bee0 **RID** Records A: keys, child RIDs (256...), ... 256 **B:** keys and values, ... memcpy 1536 C: keys and values, ... Local copy of B ...other data... **Transaction context**

~50% cycles on memcpy, still a ~30% gap:



Memory copy considered harmful



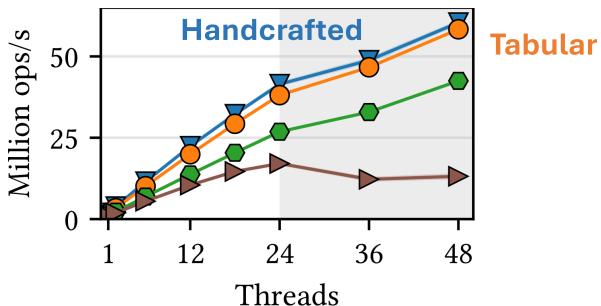
Tabular: Single-Version + (Near) Zero Memory Copy

- Flat table space to avoid versioning/pointer chasing overhead
- An old trick ship the function, not data to remove unnecessary memcpy

Before:

```
Node n = table.read(rid);
k, v = search(&n)...

After:
k, v = table.read(rid, search);
```

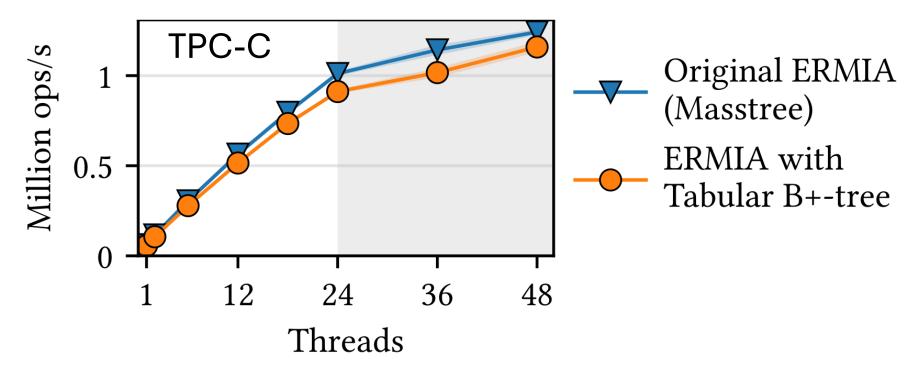


Applies to writes, too
Only need simple tweaks in OCC protocol to make this work



As a Drop-in Replacement

Masstree in ERMIA → Tabular B+-tree



Over 95% of hand-crafted



Summary

- Building a DBMS is hard
 - Have to know much more beyond DBMS itself: Parallel programming, hardware...

- Tabular: Objects-on-DB via modern OLTP techniques
 - Single-versioning + OCC + zero-copy transactions
 - Database concurrency control for data structures
 - Various benefits easier programming, debugging, migration, transparent persistence...

More in our paper and code repo:

https://github.com/sfu-dis/tabular



