

Kvell+: Snapshot Isolation without Snapshots

Baptiste Lepers
Oana Balmau
Karan Gupta
Willy Zwaenepoel

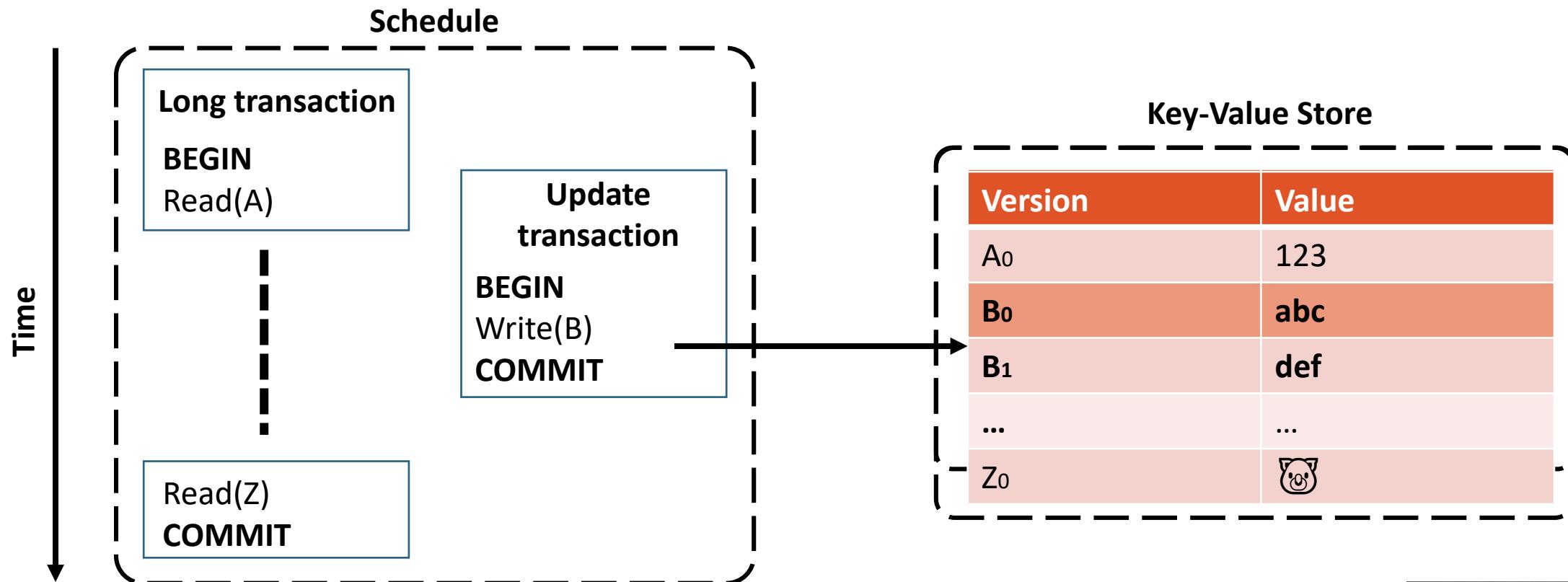


THE UNIVERSITY OF
SYDNEY

NUTANIX™



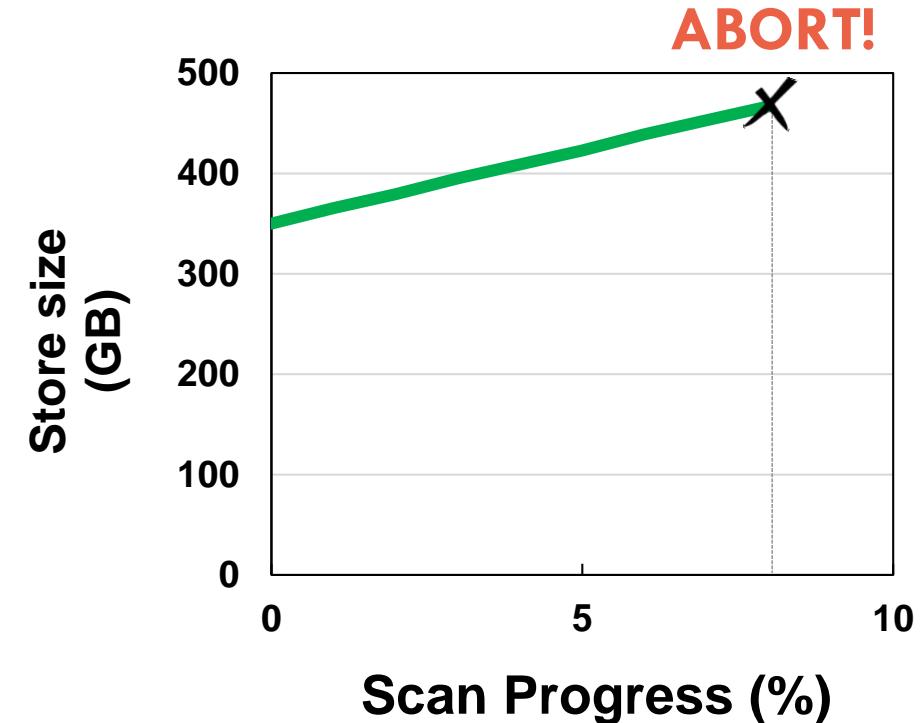
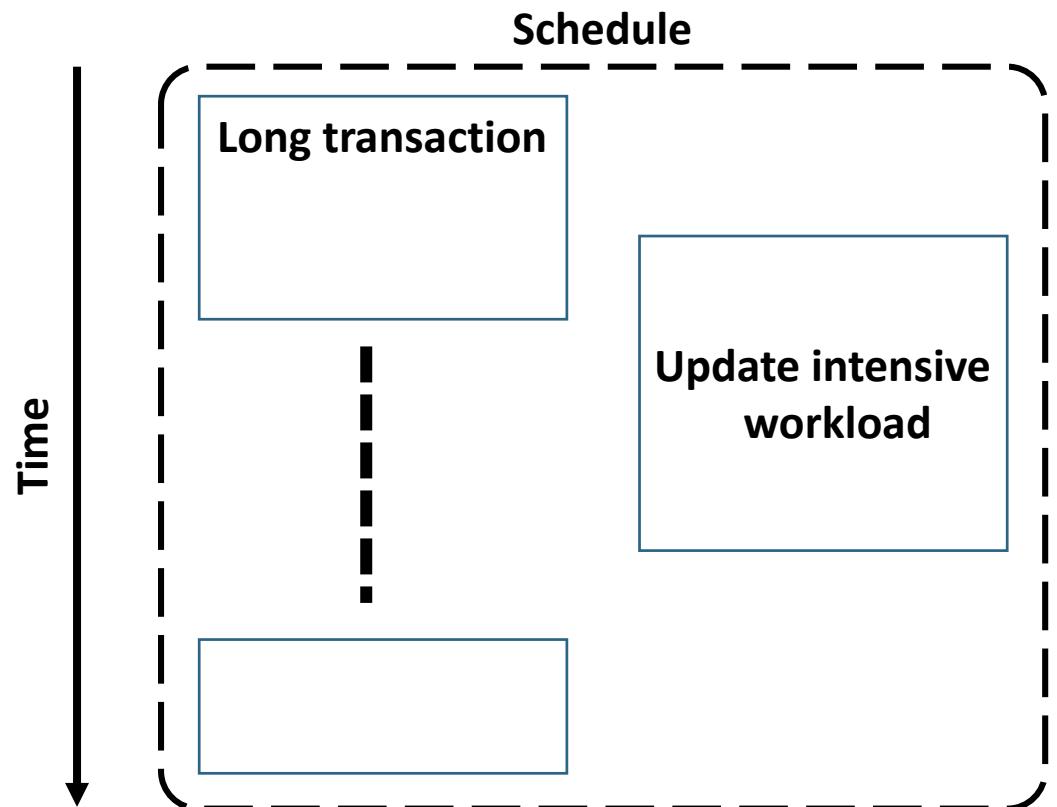
Snapshot Isolation



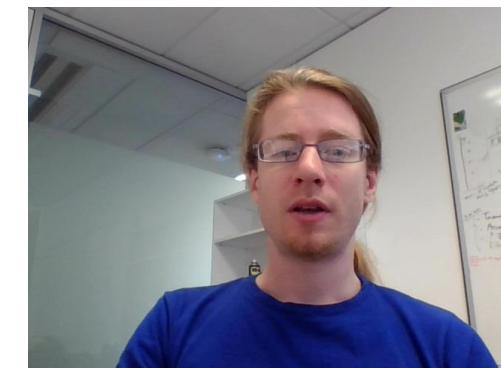
Versions are kept if they belong to a snapshot



Old versions → Space amplification

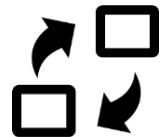


Long transactions cause massive space amplification
Store may run out of space



Do we really need to keep versions?

Long transactions are usually analytics



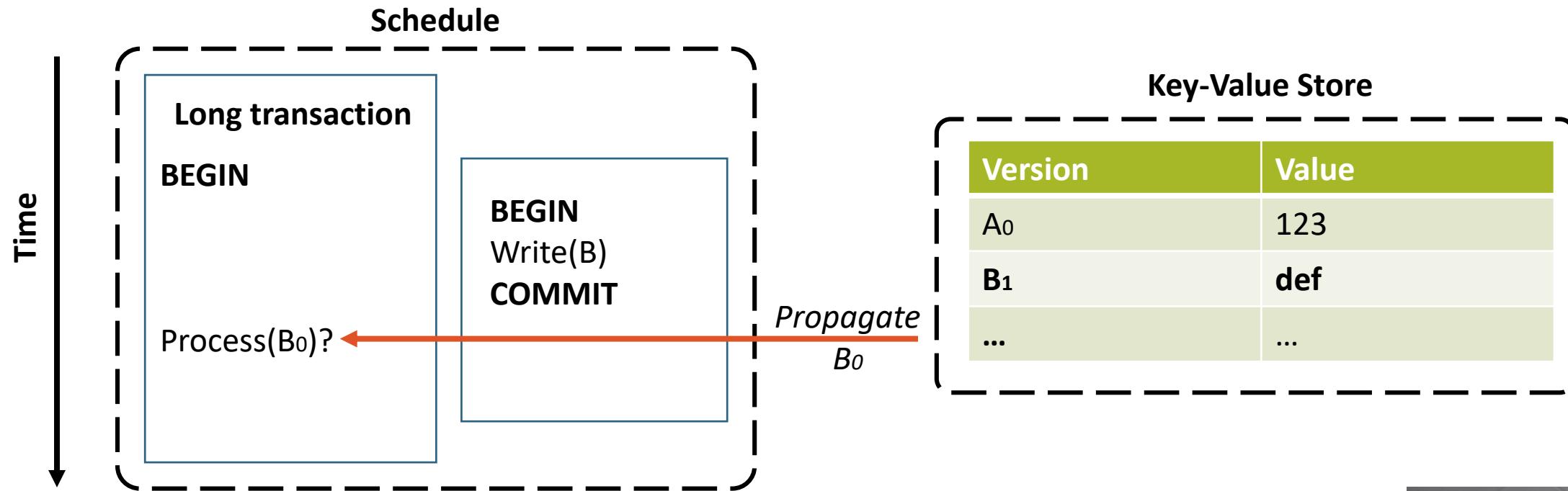
Oblivious to the order of reads



Only read items once



Solution: propagate & delete old versions



Idea: get a chance to process old versions before they are deleted
Solves space amplification



New class of analytics

OLCP: OnLine *Commutative* Processing



Scan portions of the store



Oblivious to the order of reads



Only read items once



OLCP interface

OLCP: OnLine *Commutative* Processing

`olcp_query(map, payload, scan ranges, point ranges)`

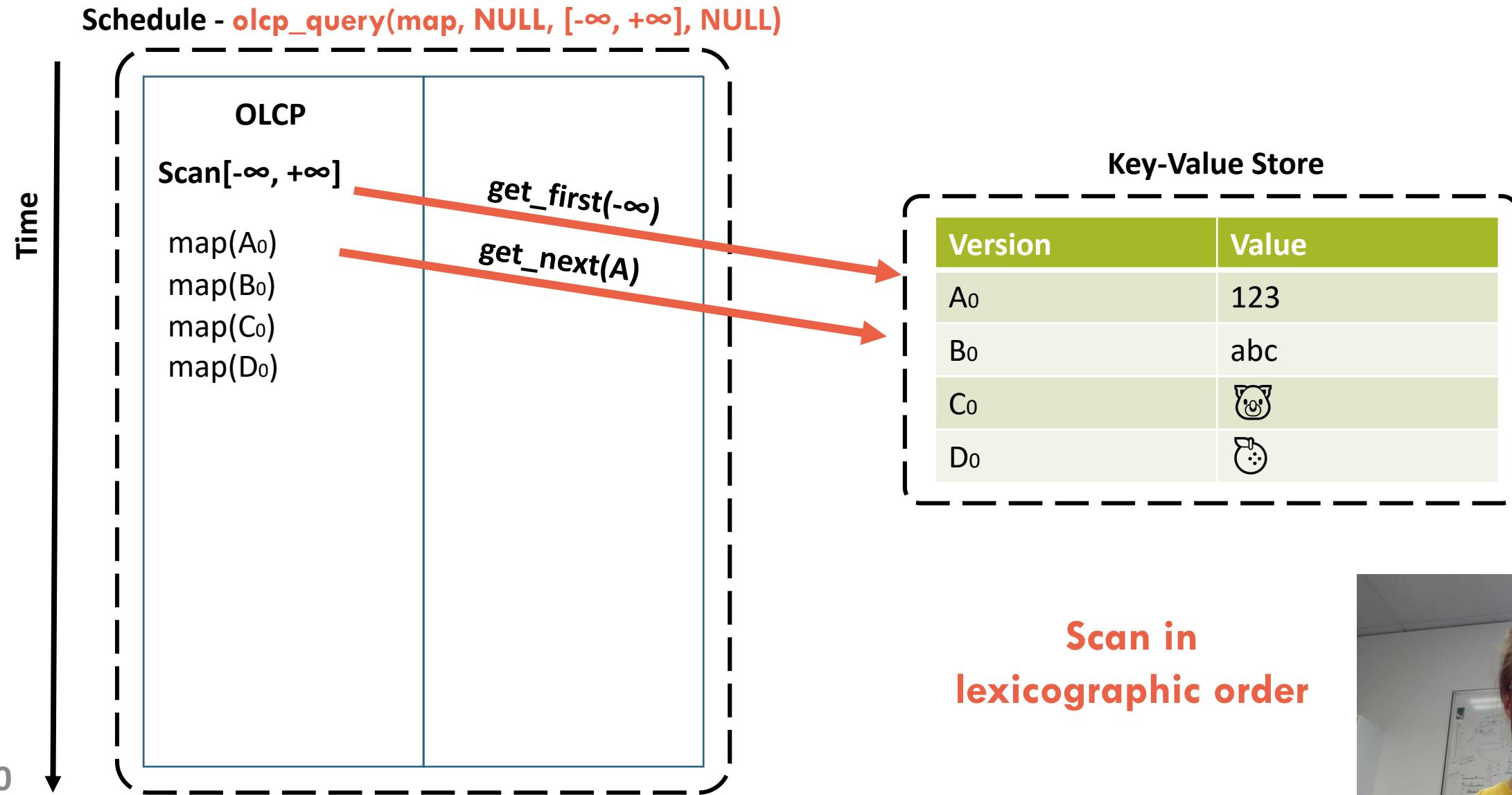
No versioning (but propagations)
For commutative operations

Versioned items
For ordered operations

Called once on every item of the scan ranges: `map(item, payload)`

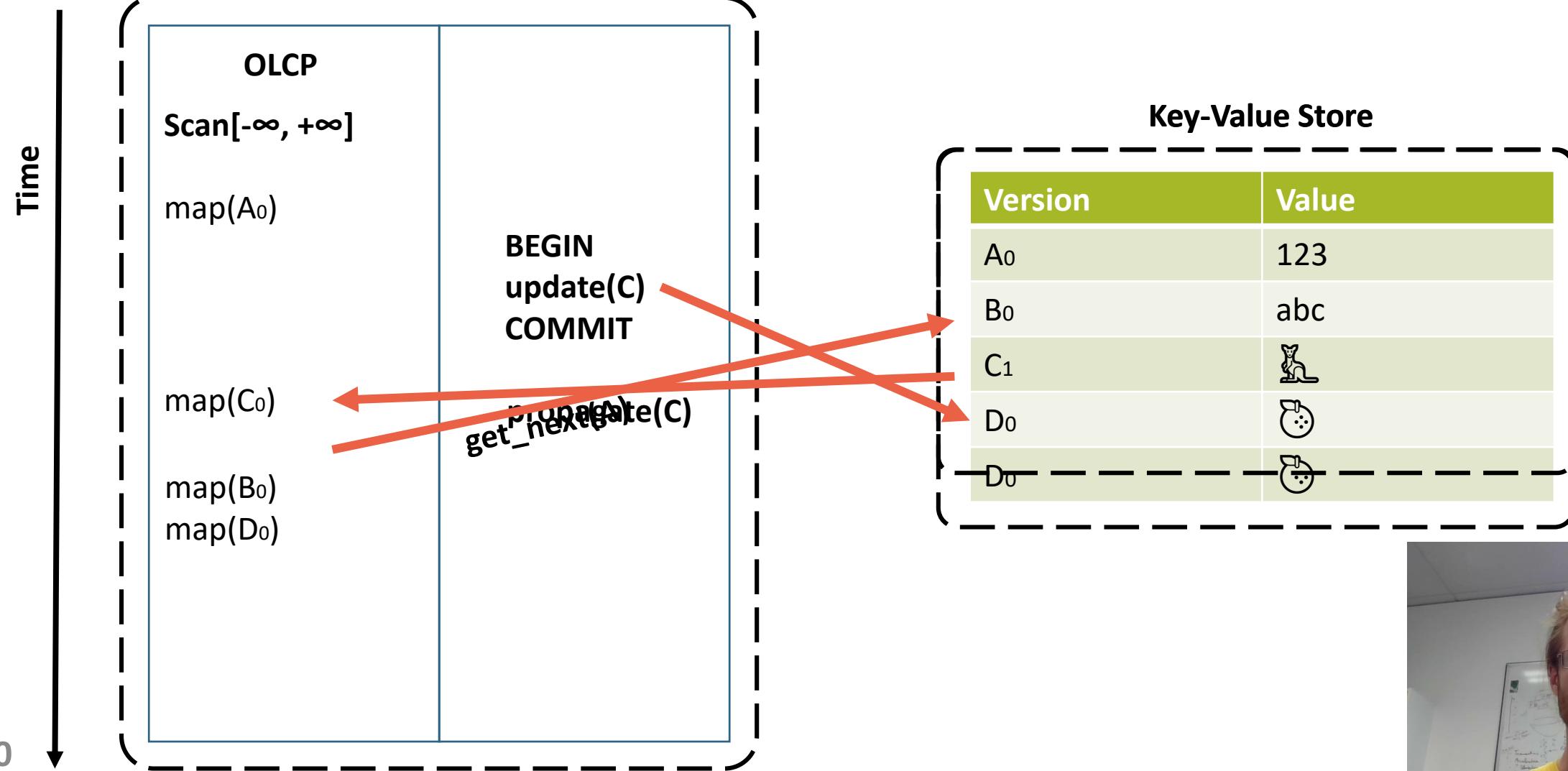


OLCP example, no concurrent transaction



OLCP example, concurrent transactions

Schedule - `olcp_query(map, NULL, [-∞, +∞], NULL)`



How to call map() exactly once?

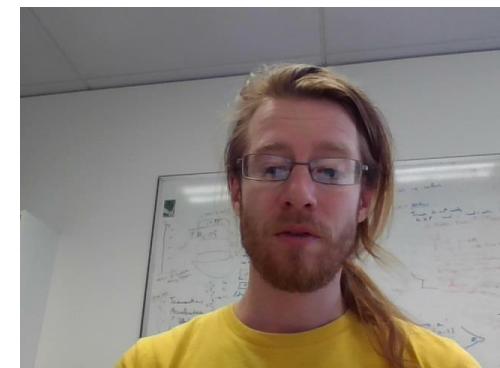


Snapshot timestamp

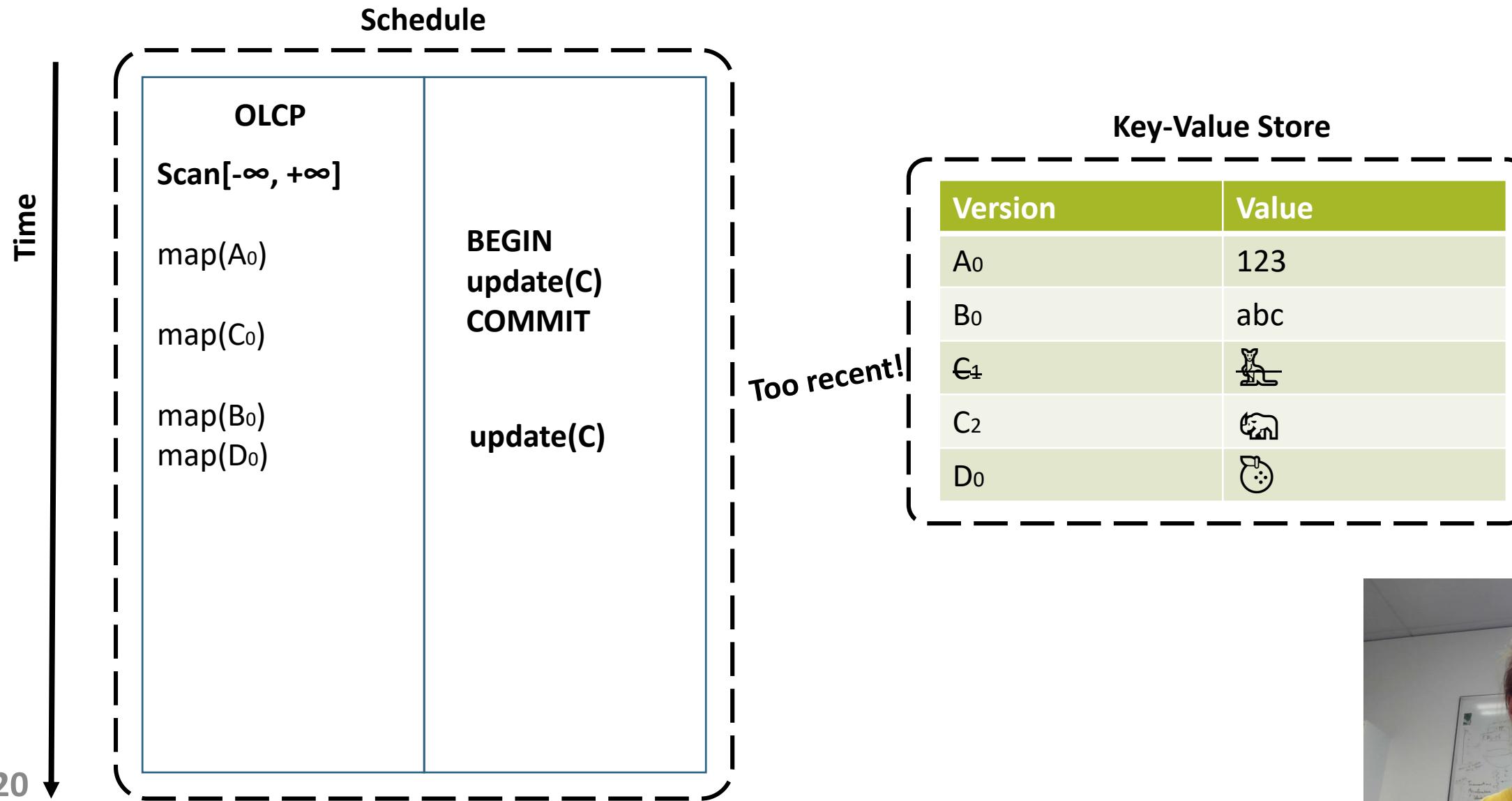


Last scanned item

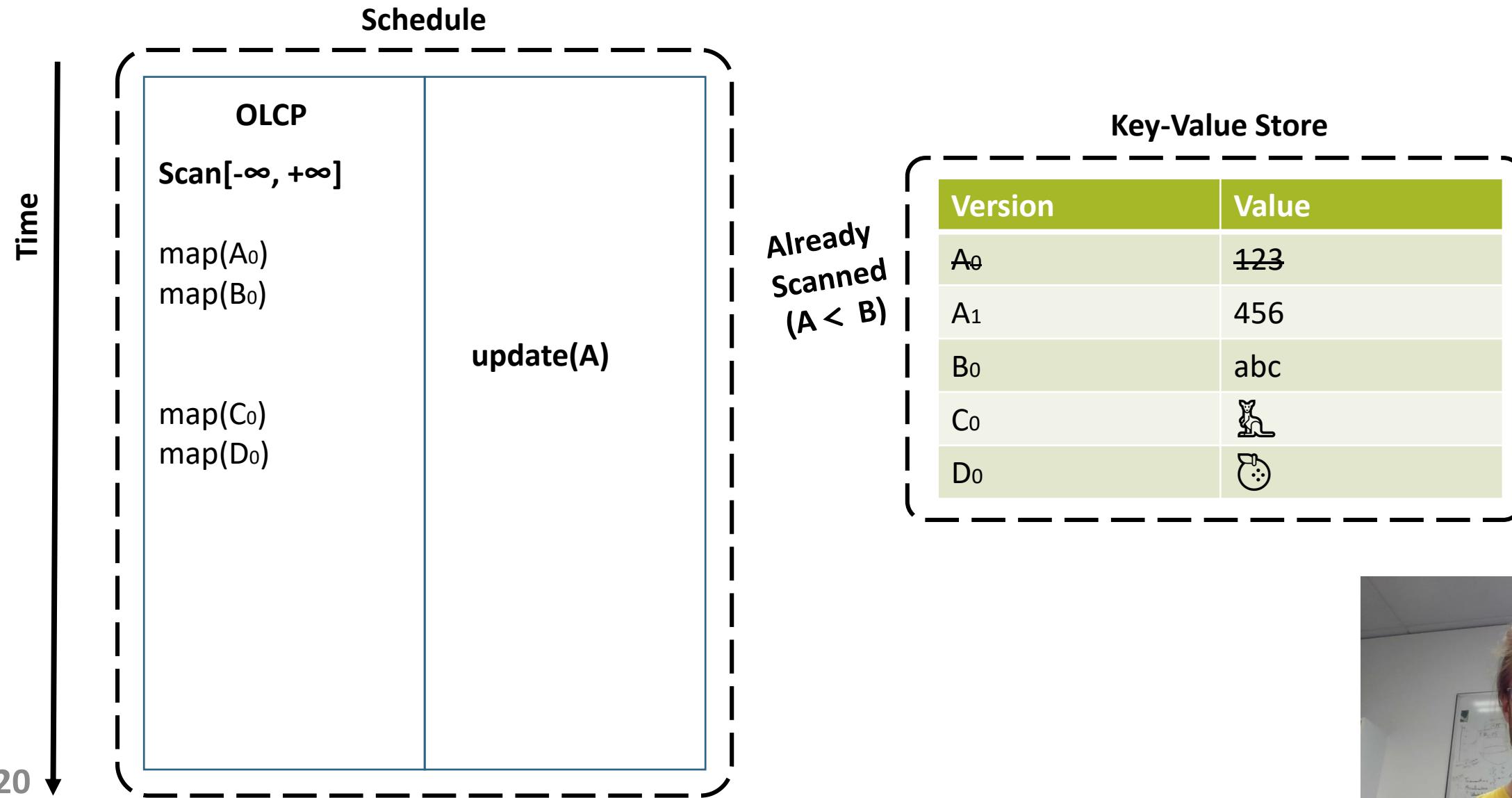
→ **Low memory overhead**



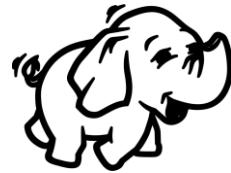
Snapshot timestamp usage



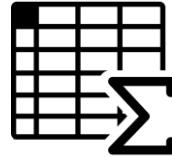
Last scanned item usage



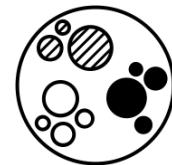
OLCP applicability



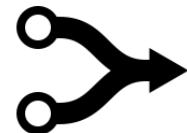
Map Reduce



SUM, COUNT



GROUP BY, CUBE, ROLLUP



Joins



OLCP example, TPC-H

SQL

```
select sum(l_quantity) as sum_qty
from lineitem
where l_shipdate <= '1998-09-04'
group by l_returnflag , l_linenumber
```

OLCP

```
map(item *i, payload *p) {
    if( i->l_shipdate <= "1998-09-04")
        string k = i->l_returnflag+"|"+i->l_linenumber;
        p->sum_qty[k] += i->l_quantity;
}
olcp_query(map, &p, [lineitems], NULL);
```

No space amplification with OLCP



Evaluation

Machine:

4 cores, 32GB RAM, Optane 905P drive (500K IOPS, 2GB/s)

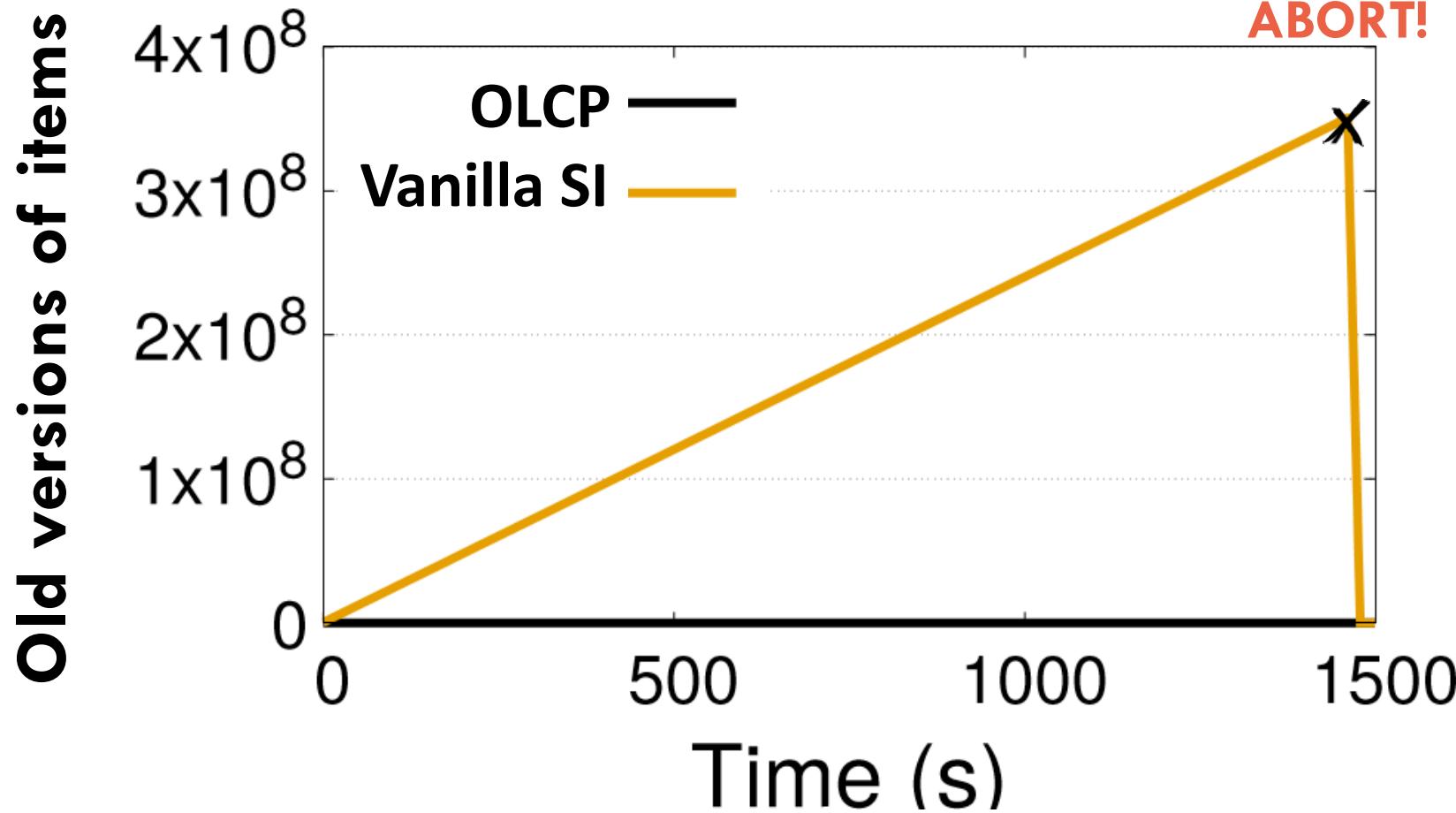
YCSB-T Benchmark:

1KB items, 100M elements (100GB)

Update heavy + 1 Vanilla SI/OLCP scan



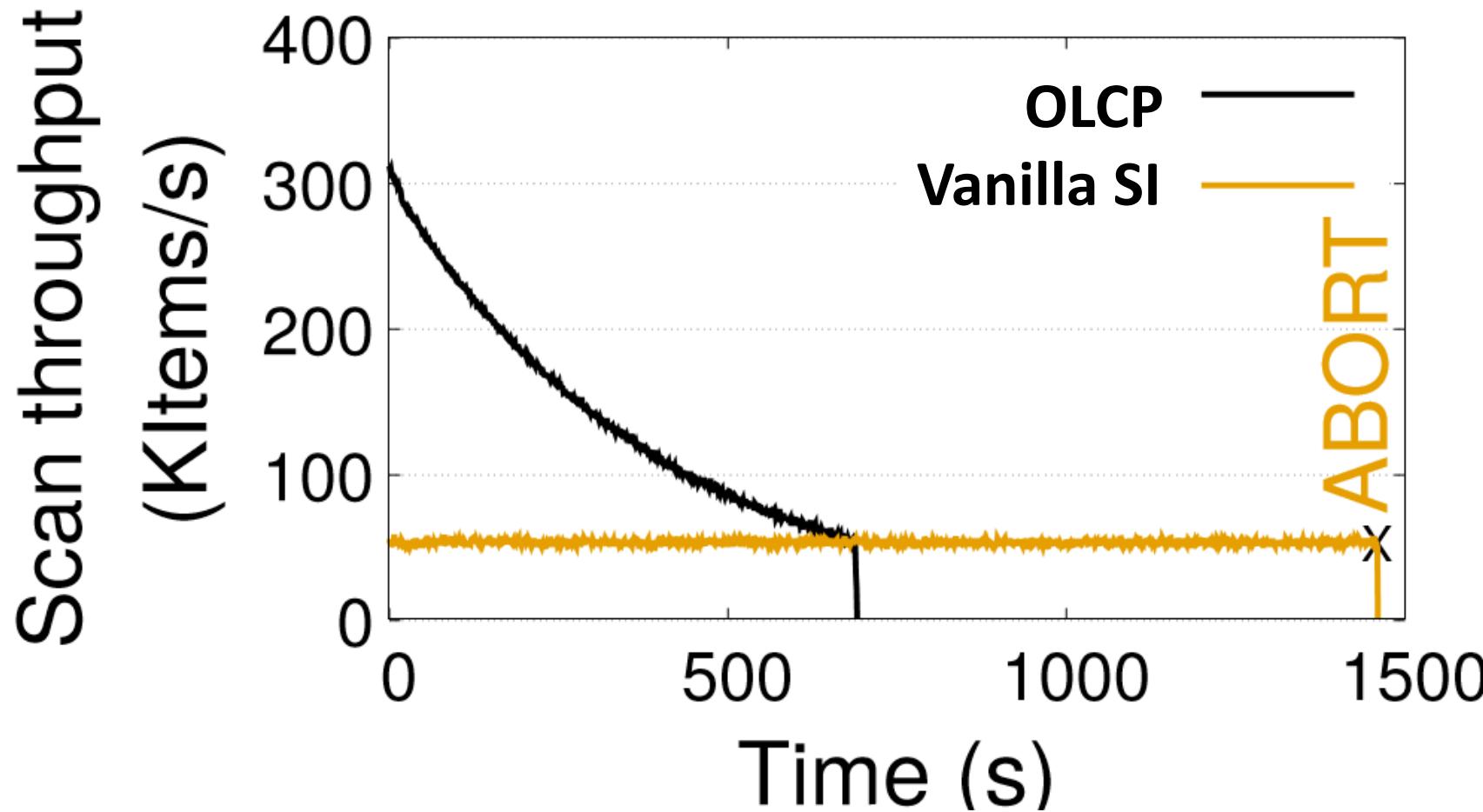
Evaluation – YCSB-T + 1 scan



Vanille SI scan aborts after 1500s
No space amplification with OLCP



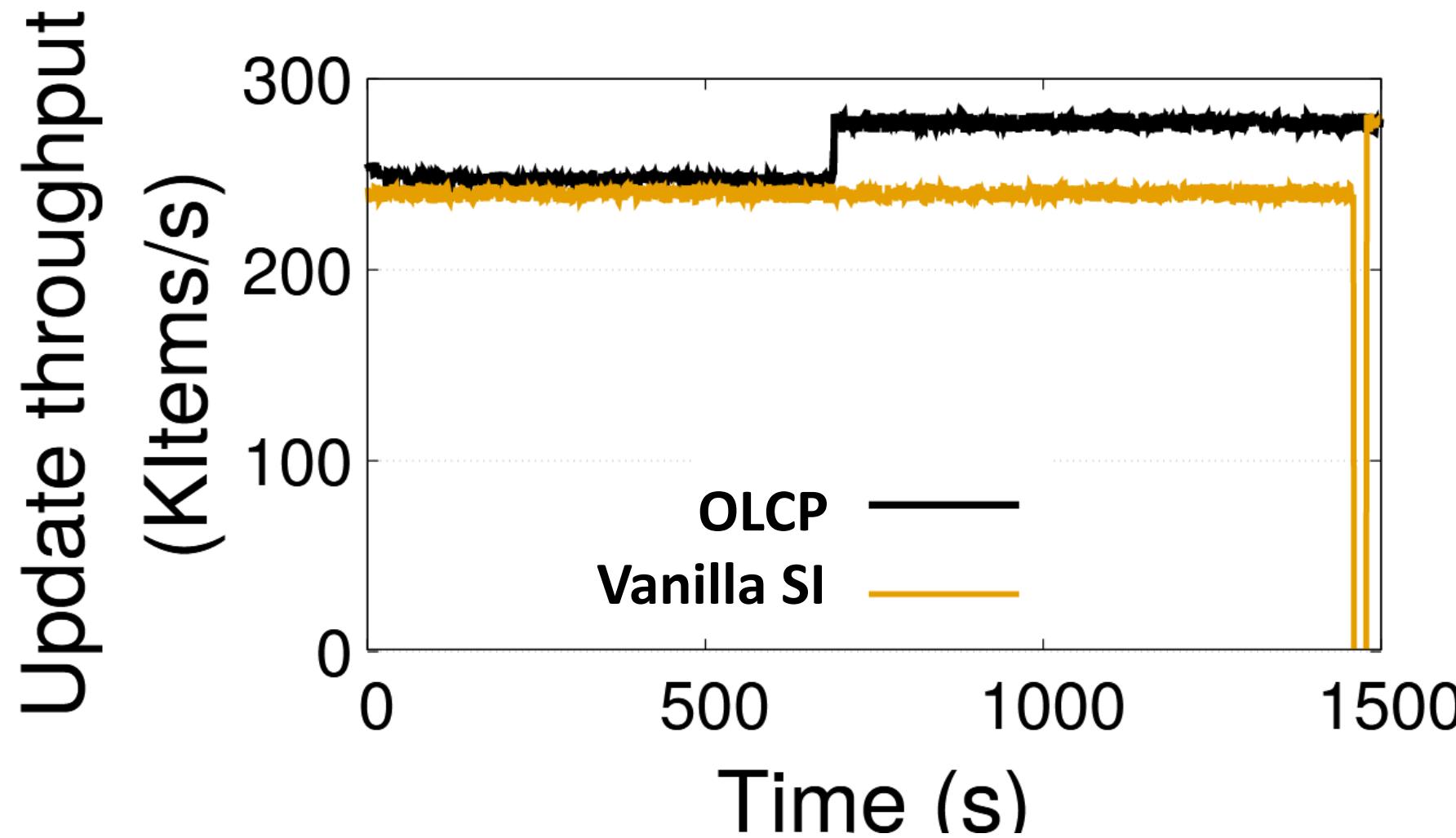
Evaluation – YCSB-T + 1 scan



OLCP scan is faster
(propagated items are mostly in memory)



Evaluation – YCSB-T + 1 scan



Propagations have little/no overhead



More in the paper



- How to do efficient propagations
- Concurrency issues
- TPC-H examples
- Detailed evaluation



Conclusion

- Snapshot Isolation induces space amplification
- OCLP Propagations → avoid space amplification
- Fits many analytics queries, is efficient



To kvell: to feel happy and proud

<https://github.com/BLepers/KVell>
baptiste.lepers@sydney.edu.au

