

Polaris

Enabling Transaction Priority in Optimistic
Concurrency Control

Rahul Kushwaha, Jan 28, 2024

Polaris: Introduction

- Optimistic Concurrency Control protocol which supports multiple levels of priority.
- Transactions with same priority are fully optimistic.
- Prioritization is accomplished via Reservation.
- Benefits
 - Significantly lower p999 tail latency.
 - Higher throughput for high contention workloads.

Optimistic Concurrency Control

- 3 Phases of a transaction: Read, Validation, and Write.

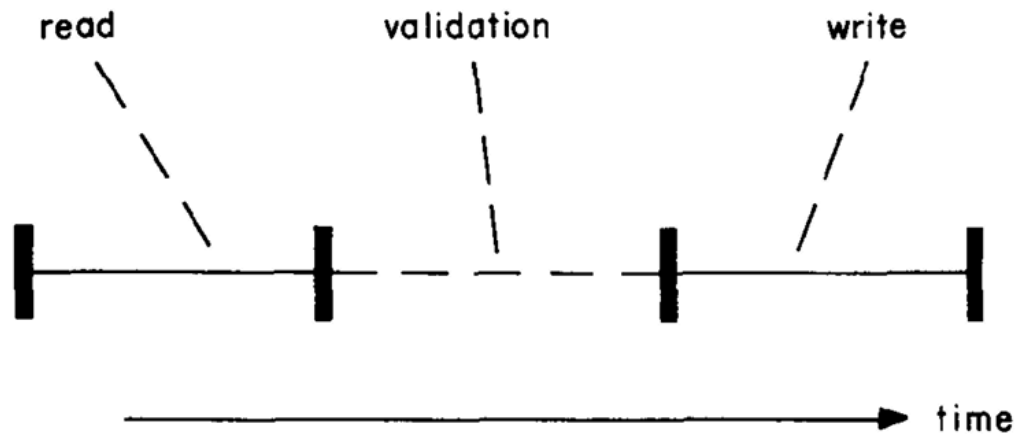


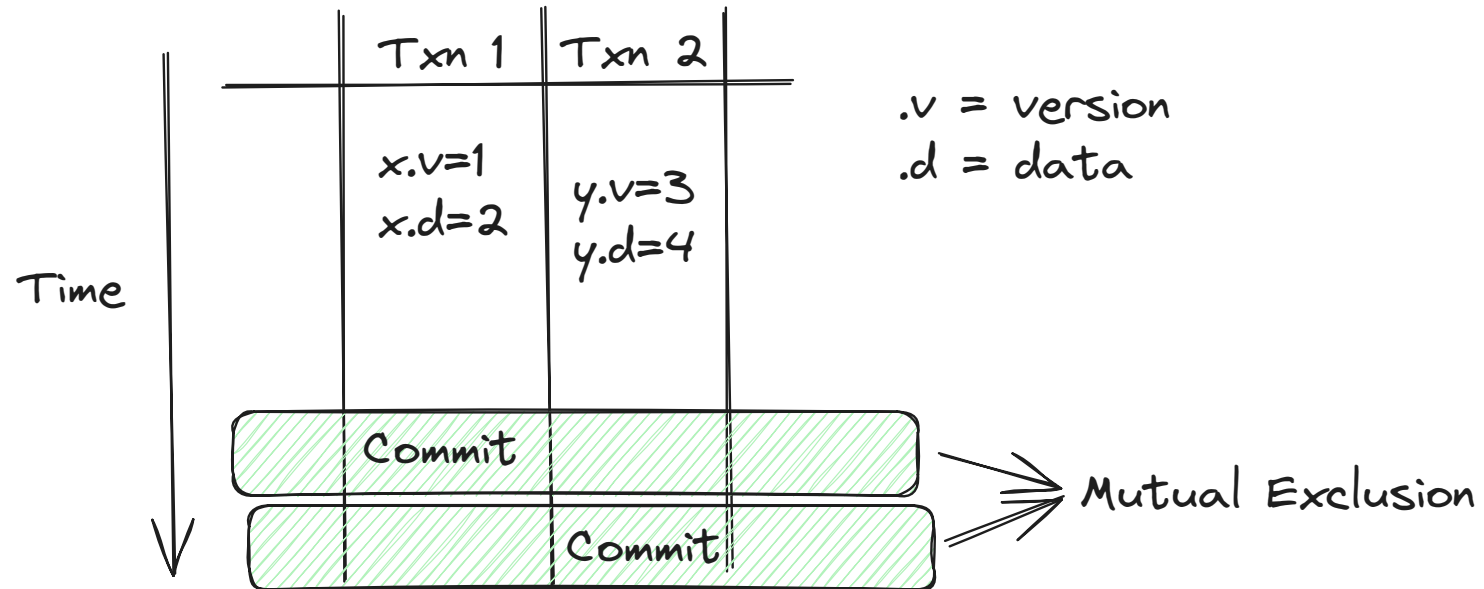
Fig. 1. The three phases of a transaction.

Img Src: [1]

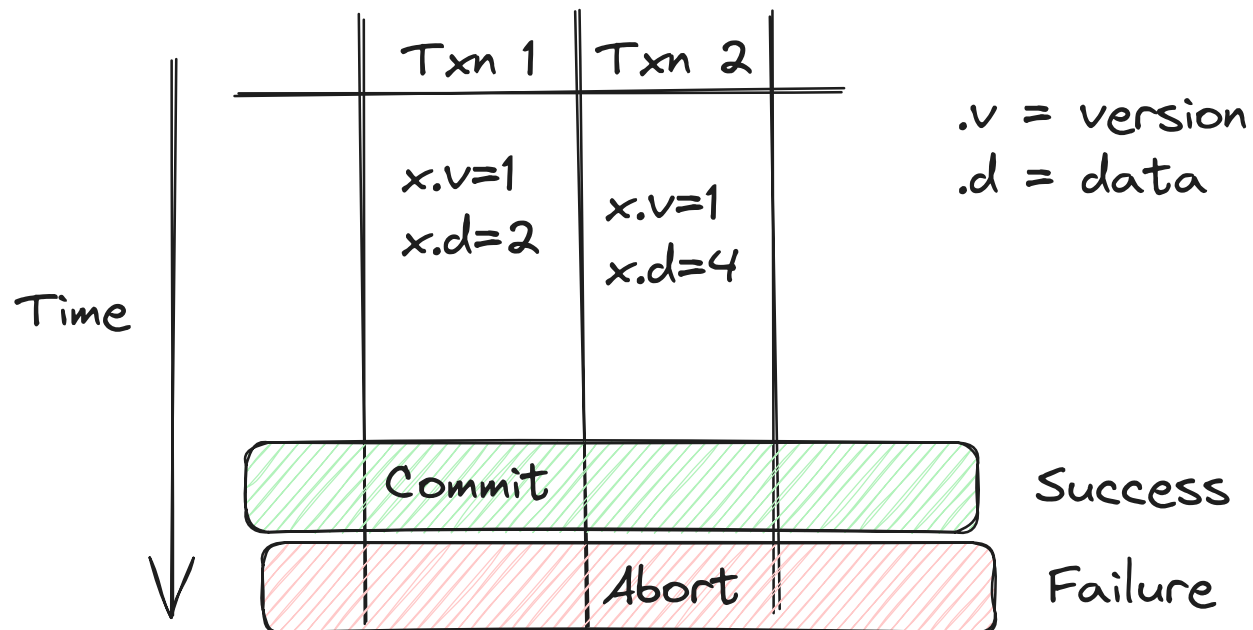
Optimistic Concurrency Control

- Read: Concurrent, multiple transactions can be executing in parallel in this phase.
 - Any mutation to the data is kept within the context of txn.
 - Read your own writes.
- Validation: Serial, global critical section.
- Write: Serial, global critical section.

Optimistic Concurrency Control



Optimistic Concurrency Control



Silo

- In-Memory database designed for modern multicore machines (high processor count & lots of main memory).
- Is a Serializable Database.
- Uses a variant of Optimistic Concurrency Control.
- Avoids centralized contention.
 - Example: Requires writes to shared memory only during commit phase.

Silo

Details

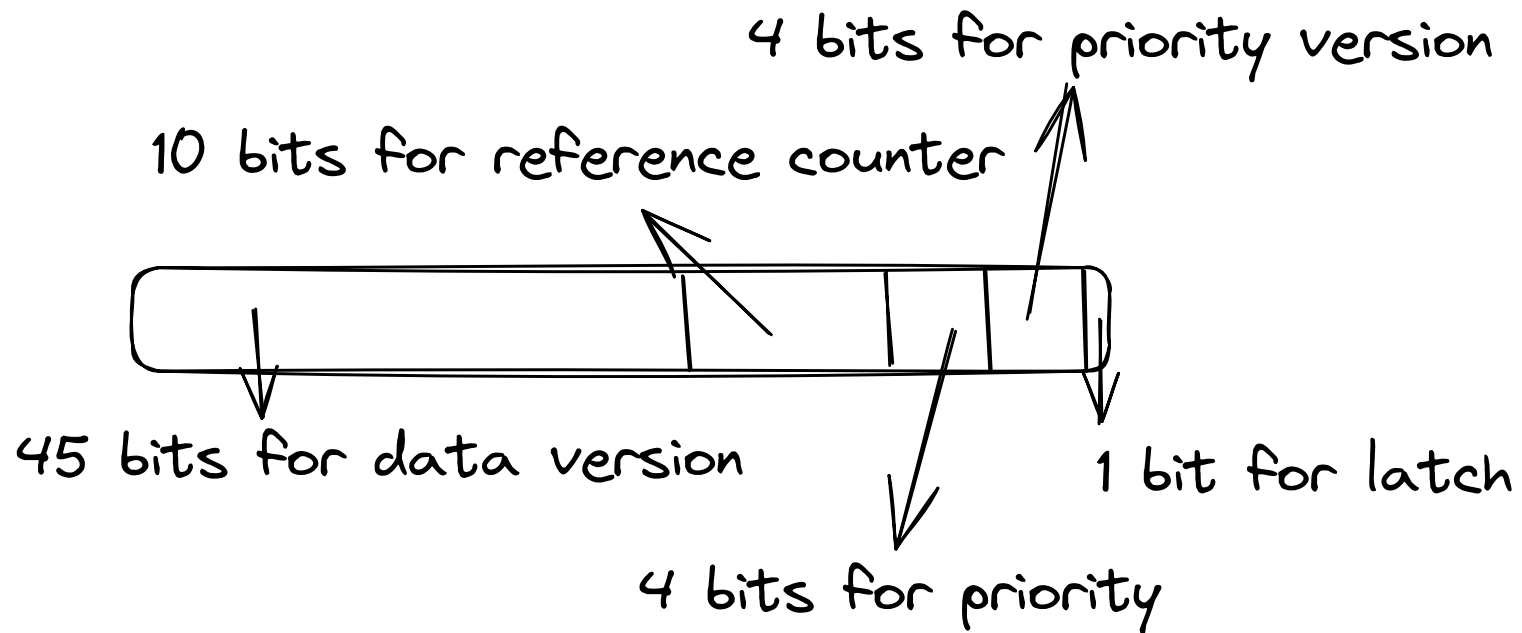
- Is based on time periods called epochs.
- Epochs form serialization points.
- A dedicated thread is responsible for periodically incrementing the epoch number(Global).
- All worker threads access Global Epoch Number during commit.

Silo

- Per-Record field, TransactionID (TID)
 - Data Version
 - Latch

Polaris TID

- Each record has a TID field.



Algorithm 1: Record Access Protocol

Algorithm 1: Record Access Protocol

Data: transaction priority $tx.prio$, record r , read-set R , write-set W , access type is_write

```
1 do
2   | do
3   |   |  $tid = r.tid$  // atomic load
4   |   | while  $tid.latch == LOCKED$ 
5   |   |  $new\_tid, is\_reserved = try\_reserve(tid, tx.prio, is\_write)$ 
6   |   |  $r\_local\_copy = r.copy()$ 
7   |   | while  $!compare\_and\_swap(r.tid, tid, new\_tid)$ 
8   |   |  $R.add(r, is\_reserved, new\_tid)$ 
9   |   | if  $is\_write$  then
10  |   |   |  $W.add(r)$ 
11  |   | return  $r\_local\_copy$ 
```

Algorithm 2: Reservation Protocol

- Goal: Low-priority txn should not abort a high-priority txn.
- Example:
 - Two txns, [A -> high priority, B -> low priority]
 - If A has read the record but not committed, B should not be able to write to it, as that will cause A to abort.
 - If B has read the record but not committed, A must be able to ignore B and proceed to read/write/commit.

Algorithm 2: Reservation Protocol

Algorithm 2: Reservation Protocol

Data: transaction priority $tx.prio$, a copy of record TID tid , access type is_write

```
1 function try_reserve( $tid, tx.prio, is\_write$ ):
2    $new\_tid = tid$ 
3   if  $tid.prio == tx.prio$  then
4     /* reserve with same-priority transactions
5      $new\_tid.ref\_cnt++$ 
6      $is\_reserved = true$ 
7   else if  $tid.prio < tx.prio$  then
8     /* preempt from low-priority transactions
9      $new\_tid.prio = tx.prio$ 
10     $new\_tid.ref\_cnt = 1$ 
11     $is\_reserved = true$ 
12  else
13    /* reserved by high-priority transactions
14    if  $is\_write$  then
15       $ABORT()$ 
16     $is\_reserved = false$ 
17  return  $new\_tid, is\_reserved$ 
```

Algorithm 3: Commit Protocol(Pt-1)

Algorithm 3: Commit Protocol

Data: transaction priority $tx.prio$, read-set R , write-set W

```
1 for  $r$  in  $sorted(W)$  do
2   do
3      $tid = r.tid$  // atomic load
4     if  $tid.prio > tx.prio$  or  $tid.latch == LOCKED$  then
5       ABORT()
6        $locked\_tid = tid$ 
7        $locked\_tid.latch = LOCKED$ 
8     while ! $compare\_and\_swap(r.tid, tid, locked\_tid)$ 
9   for  $r, is\_reserved, tid$  in  $R$  do
10     $curr\_tid = r.tid$  // atomic load
11    if  $curr\_tid.latch == LOCKED$  and  $r$  not in  $W$  then
12      ABORT() // locked by another transaction
13    if  $curr\_tid.data\_ver != tid.data\_ver$  then
14      ABORT() // data has been updated
/* validation pass; transaction can commit
```

Algorithm 3: Commit Protocol(Pt-2)

```
15 new_data_ver = W.max_data_ver() + 1
16 for r in W do
17     r.install_write()
18     tid = r.tid // atomic load
19     new_tid = cleanup_write(tid, new_data_ver)
20     r.tid = new_tid // atomic store
21 for r, is_reserved, old_tid in R do
22     if r not in W and is_reserved then
23         do
24             tid = r.tid // atomic load
25             new_tid = cleanup_read(tid, tx.prio, old_tid.prio_ver)
26             if new_tid == tid then
27                 break // no cleanup needed
28             while !compare_and_swap(r.tid, tid, new_tid)
```

Algorithm 4: Write CleanUp

```
12 function cleanup_write(tid, new_data_ver):  
13     new_tid.data_ver = new_data_ver  
14     new_tid.latch = UNLOCKED  
15     new_tid.prio = 0  
16     new_tid.prio_ver = tid.prio_ver + 1  
17     new_tid.ref_cnt = 0  
18     return new_tid
```


Algorithm 4: Read CleanUp

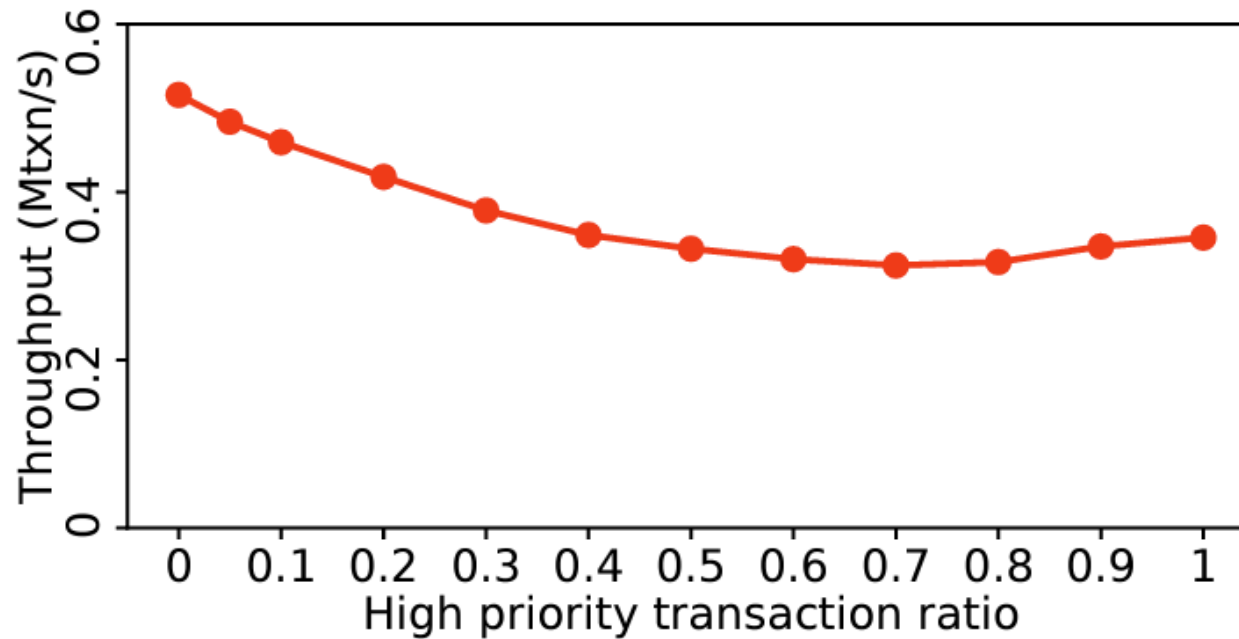
```
1 function cleanup_read(tid, tx.prio, prio_ver):
2   if tid.latch == LOCKED then
3     return tid // no cleanup needed
4   if tid.prio != tx.prio or tid.prio_ver != prio_ver then
5     return tid // no cleanup needed
6   new_tid = tid
7   new_tid.ref_cnt--
8   if new_tid.ref_cnt == 0 then
9     new_tid.prio = 0
10    new_tid.prio_ver++
11  return new_tid
```

Priority Assignment

$$p = \begin{cases} p_0, & \text{if } abort_cnt < t \\ p_0 + \lfloor (abort_cnt - t)/s \rfloor, & \text{otherwise} \end{cases}$$

Results: YCSB

Varying number of high priority & low priority txns



YCSB-A, $r = 50\%$, $w = 50\%$, $\theta = 0.99$

Results: YCSB

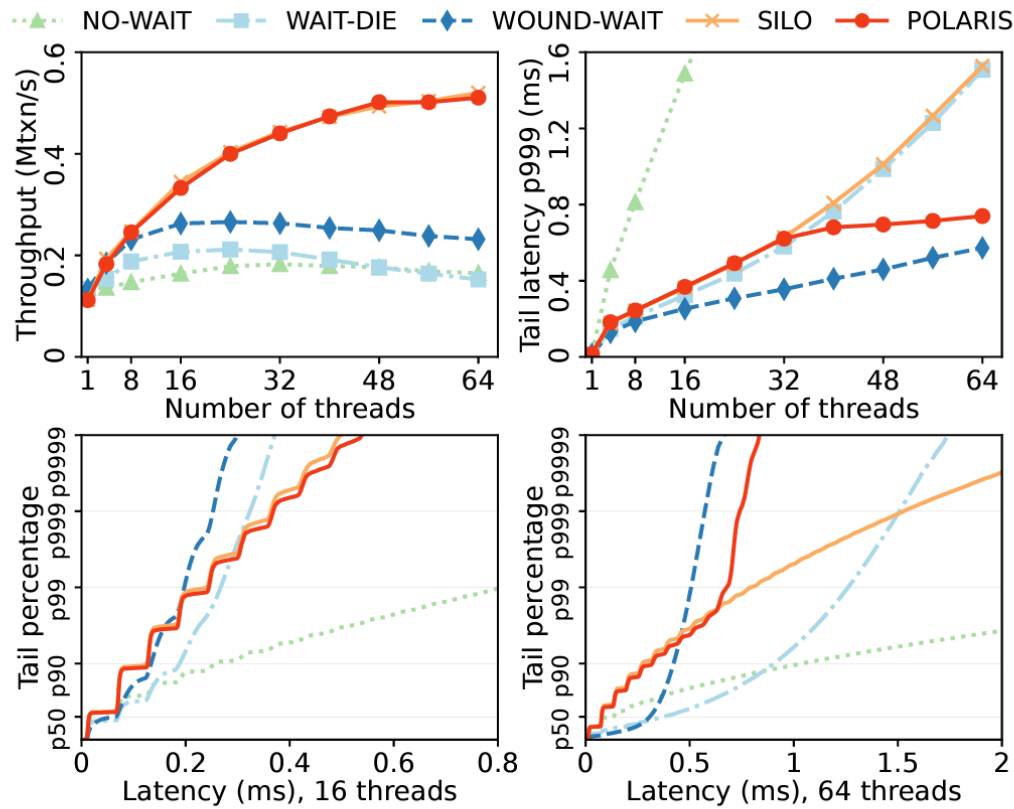
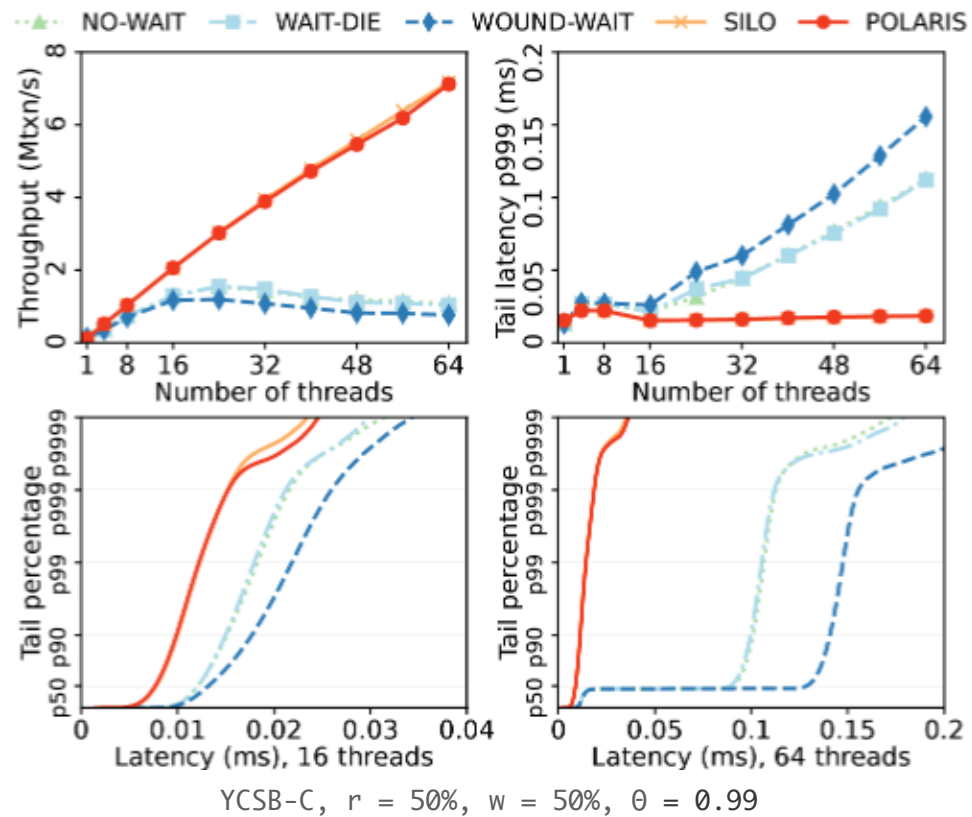
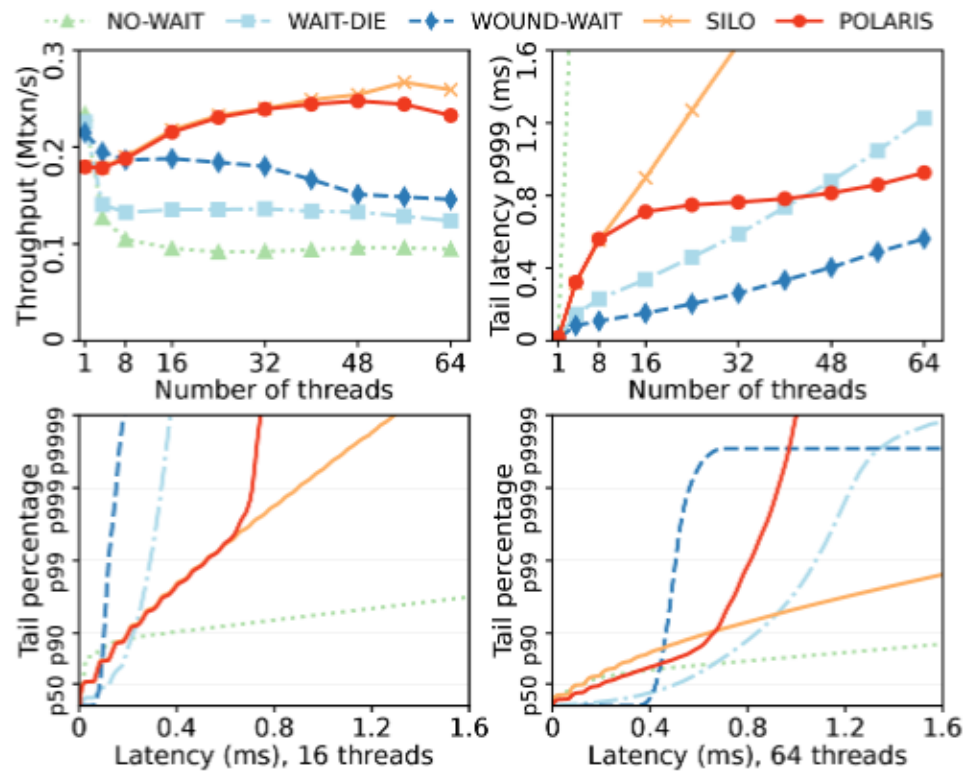


Fig. 3. Throughput and p999 tail latency over a spectrum of thread numbers; latency distribution in the cases of 16 and 64 threads (YCSB-A, $r = 50\%$, $w = 50\%$, $\theta = 0.99$).

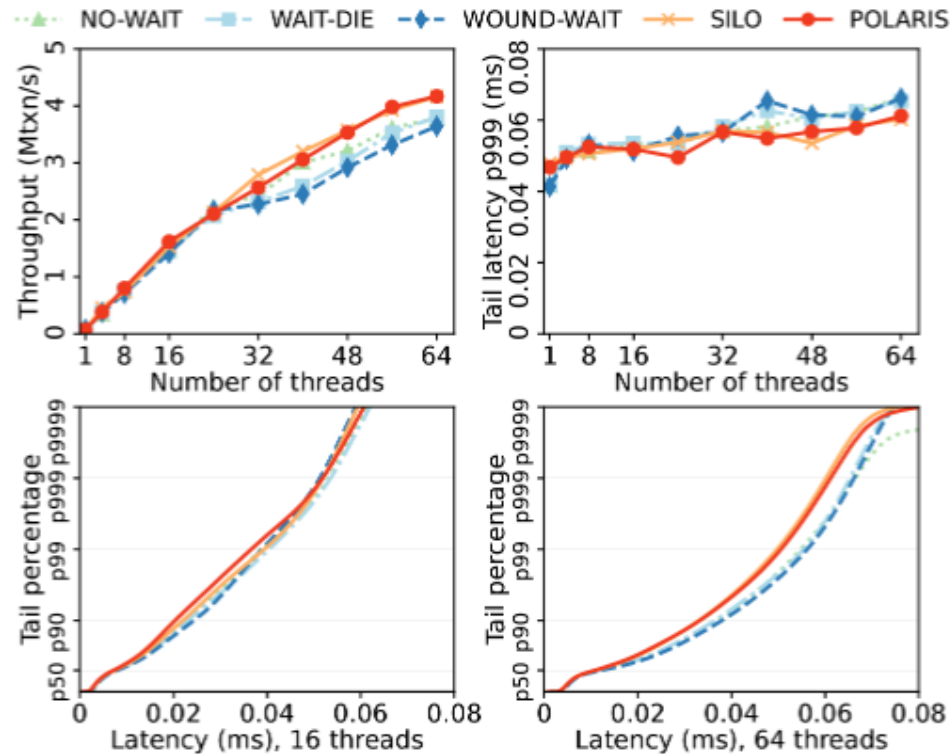
Results: YCSB



Results: TPC-C [High Contention]



Results: TPC-C [Low Contention]



Thank You!!

References

- [1] H.T.Kung and John T. Robinson. 1981. On Optimistic Methods for Concurrency Control. *ACM Trans. Database Syst.* 6, 2(jun1981), 213–226. <https://doi.org/10.1145/319566.319567>
- [2] Stephen Tu, Wenting Zheng, Eddie Kohler, Barbara Liskov, and Samuel Madden. 2013. Speedy Transactions in Multicore In-Memory Databases. In *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles (Farmington, Pennsylvania) (SOSP'13)*. Association for Computing Machinery, New York, NY, USA, 18–32. <https://doi.org/10.1145/2517349.2522713>