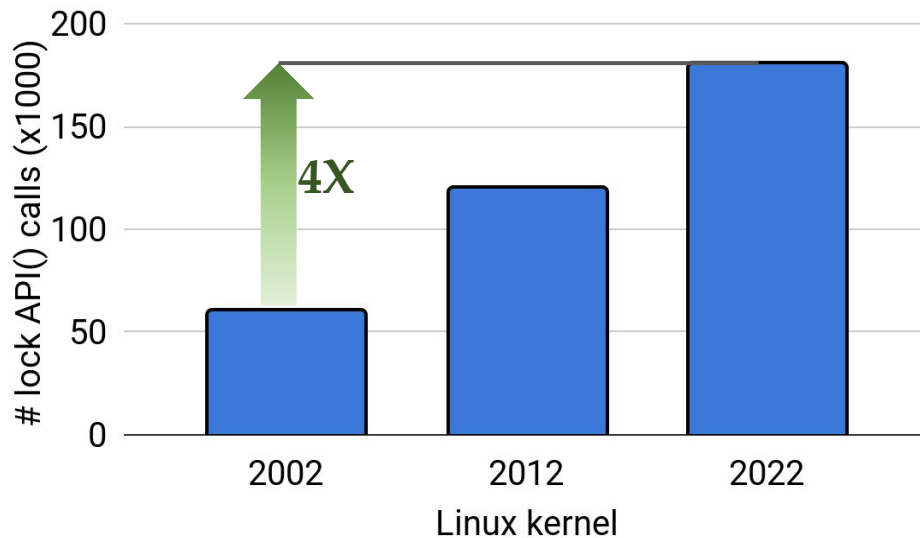


# Ship your Critical Section Not Your Data: Enabling Transparent Delegation with TCLocks

**Vishal Gupta**   Kumar Kartikeya Dwivedi   Yugesh Kothari  
Yueyang Pan   Diyu Zhou   Sanidhya Kashyap



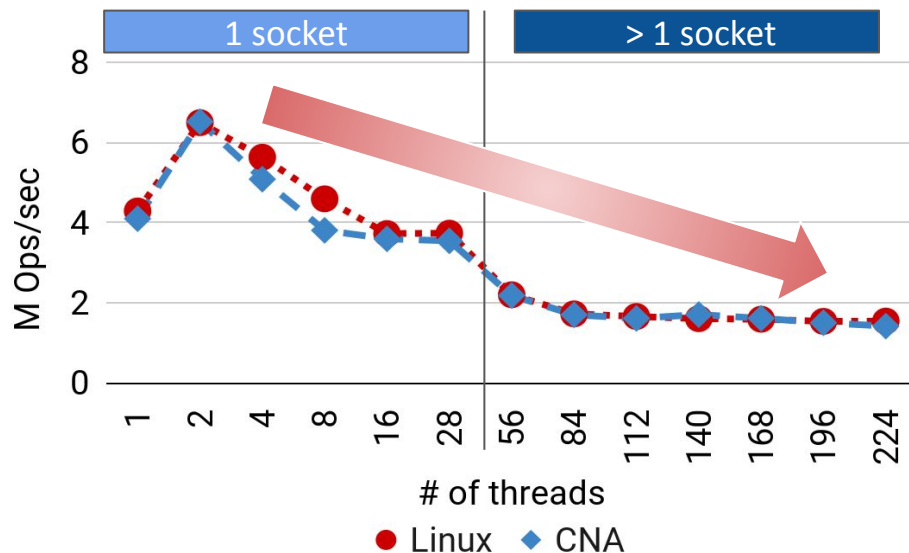
# Locks: **MOST WIDELY** used mechanism



More locks are in use to improve OS scalability

# Performance: Micro-benchmark

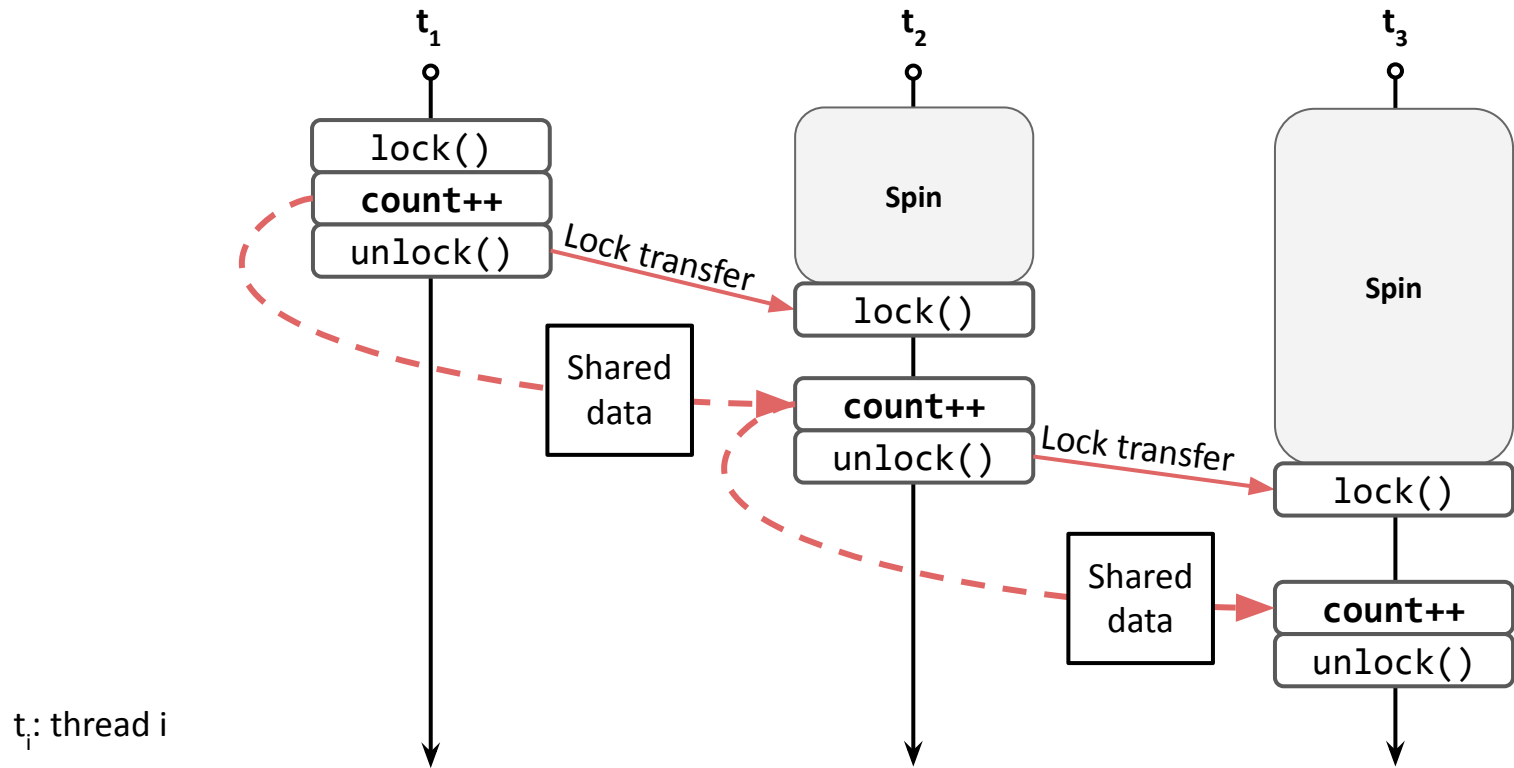
Benchmark: Each thread renames a file in a directory, serialized by a directory lock



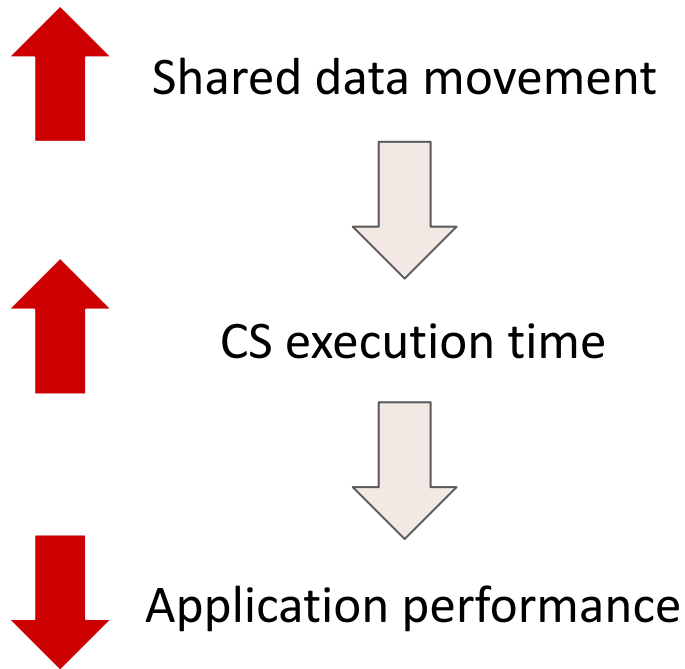
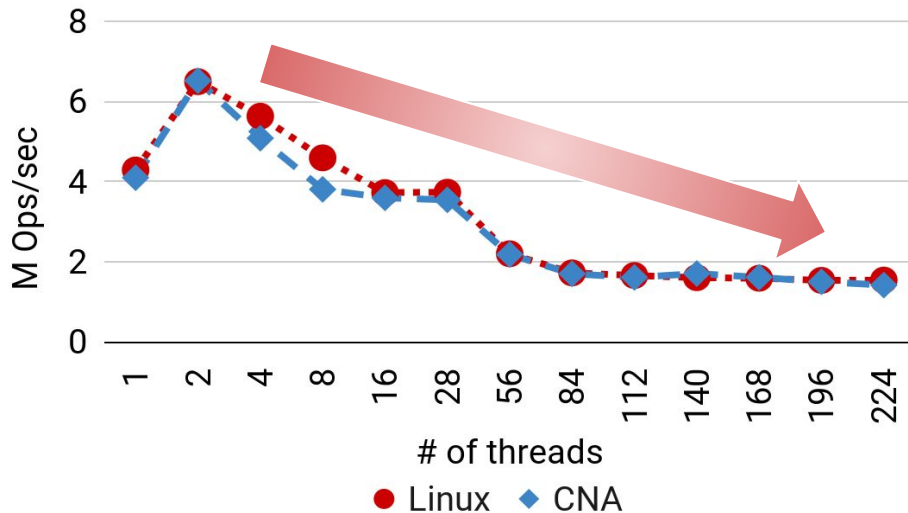
- Performance decreases with increasing core count
- NUMA-aware locks (CNA) follow a similar trend

Setup: 8-socket/224-core machine

# Traditional lock design: Large data movement



# Traditional lock design: Not ideal



# Delegation-style locks

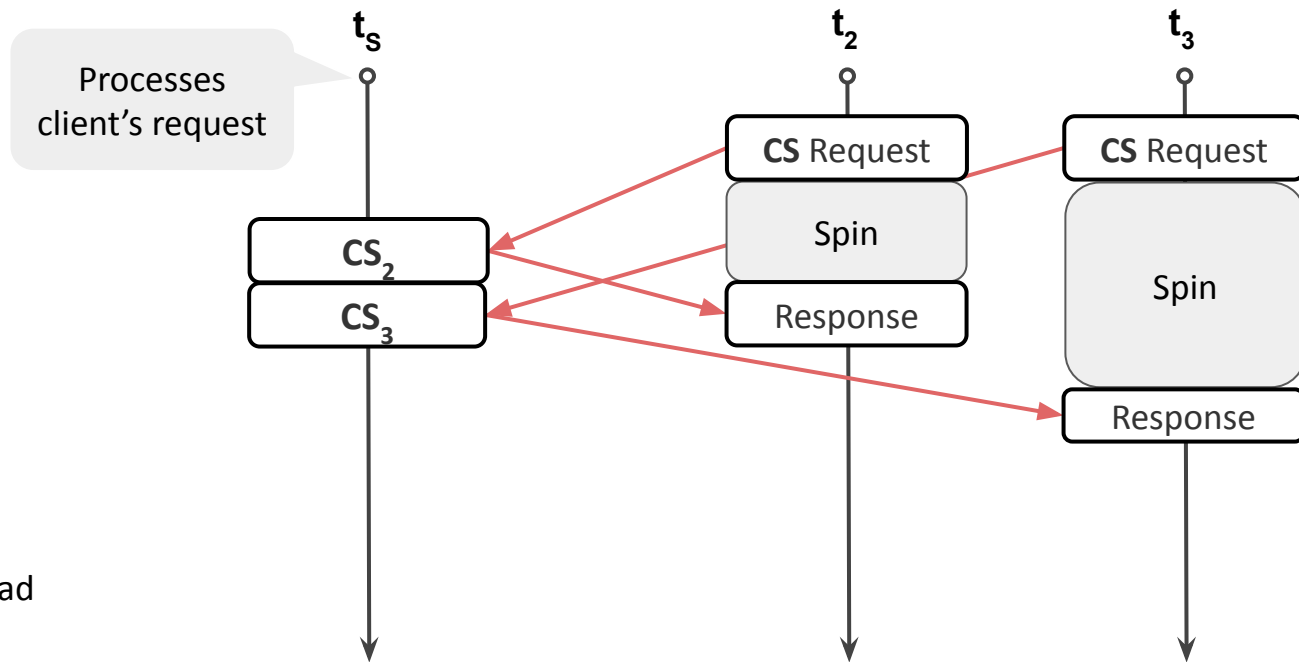
- Similar to a server-client model
  - Server: Lock holder
  - Client: Waits to acquire the lock
- Client ships its critical section request in the form of a function to the server thread

```
lock()  
count++  
unlock()
```



```
void incr_func() =  
    count++  
  
send_req_to_server(&incr_func)
```

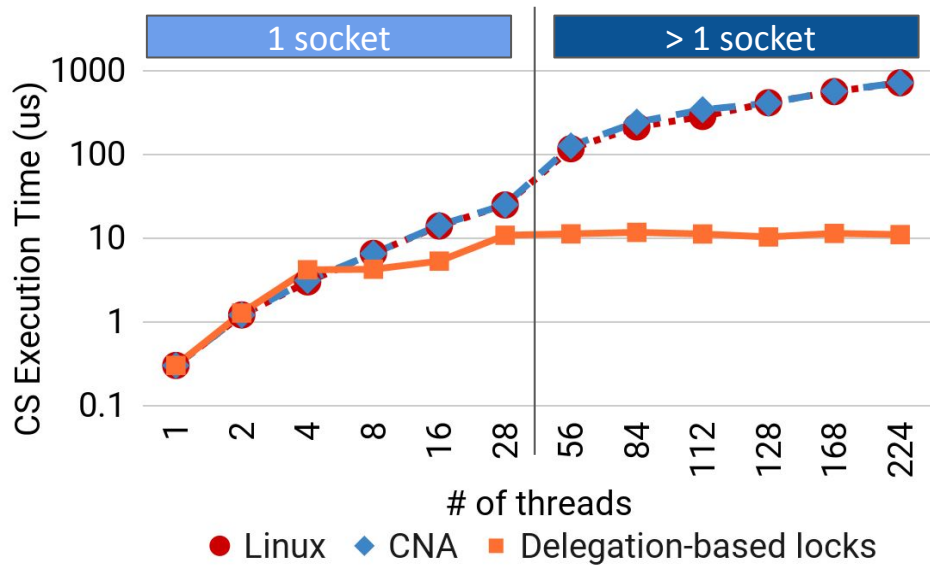
# Delegation-style locks



$t_s$ : server thread  
 $t_i$ : thread  $i$   
CS: critical section

# Delegation-style locks

Benchmark: Each thread renames a file in a directory, serialized by a directory lock



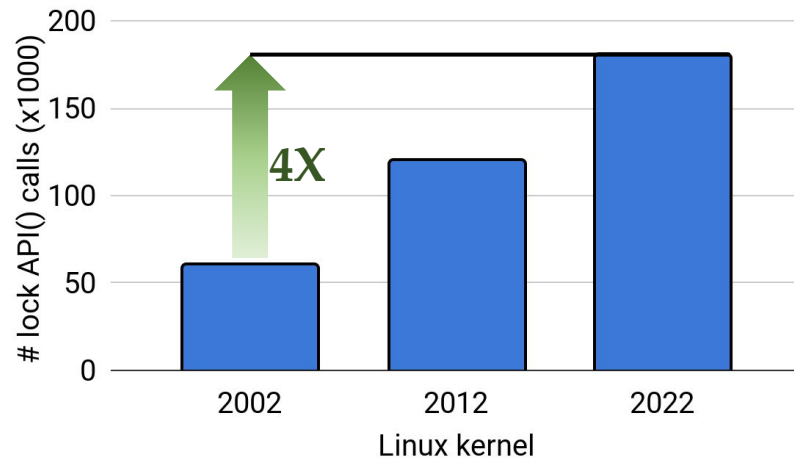
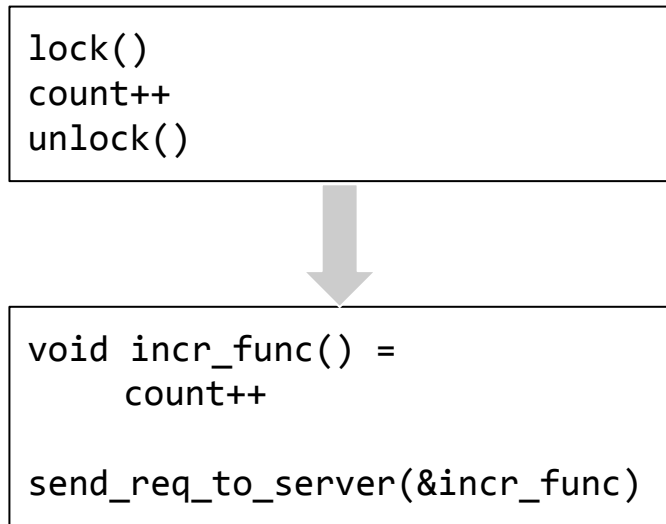
CS execution time similar with increasing core count

- Minimal shared data movement

Setup: 8-socket/224-core machine



# Delegation locks require code rewrite



Existing delegation-based design is impractical for Linux

# TCLocks: Goals

- **Transparency**
  - Use standard lock/unlock APIs without rewriting applications
- **Delegation**
  - Minimal shared data movement

**Transparent delegation**

# Agenda

- Motivation
- TCLock Design
- TCLock in Linux
- Evaluation
- Discussion

# How to achieve transparent delegation?

- **How to capture the thread's context?**
  - Without application rewrite
- **Where to capture the thread's context?**
  - Such that only critical section is captured
- **Does the waiter's thread modify its context?**
  - While the server is executing waiter's critical section

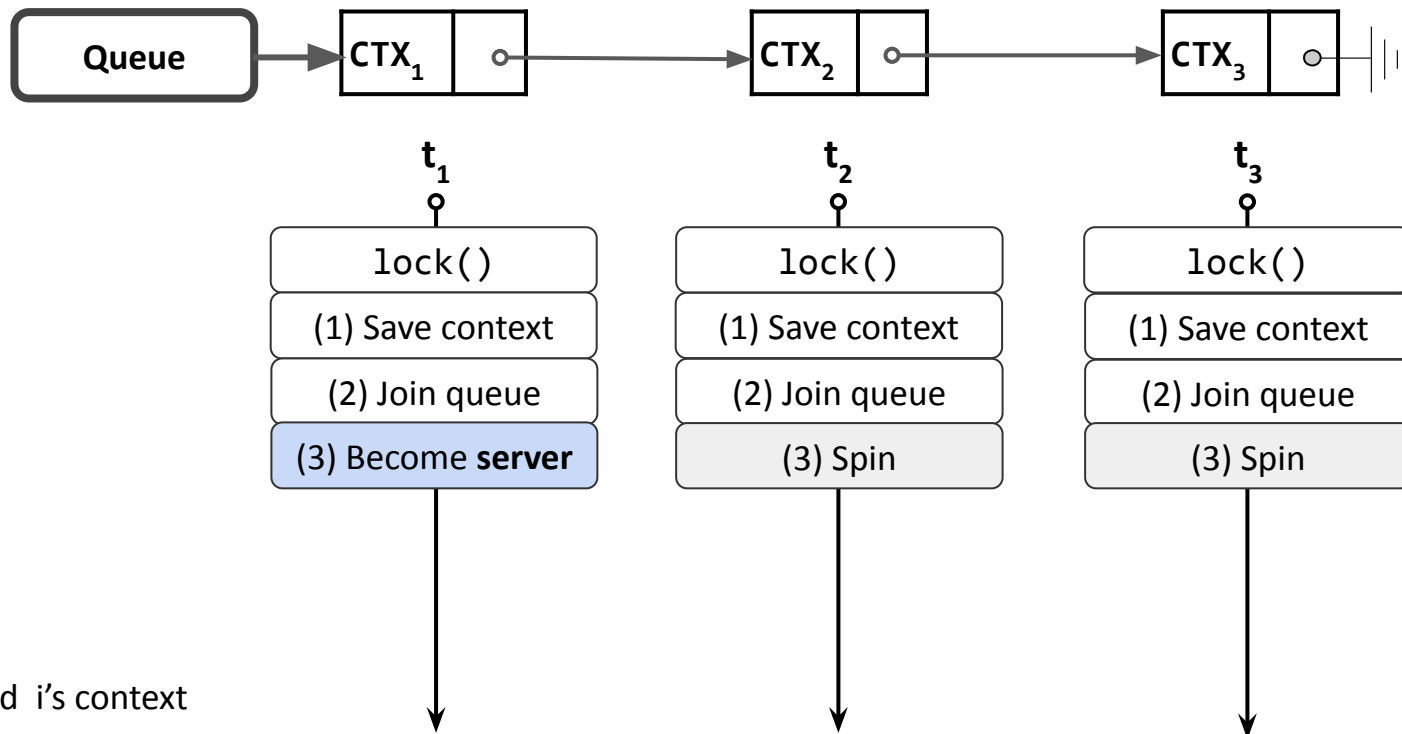
# Key idea: Transparent delegation

- **How to capture the thread's context?**
  - Instruction pointer + stack pointer + general-purpose registers
- **Where to capture the thread's context?**
  - Start and end of lock/unlock API
- **Does the waiter's thread modify its context?**
  - No, lock waiter busy waits to acquire the lock

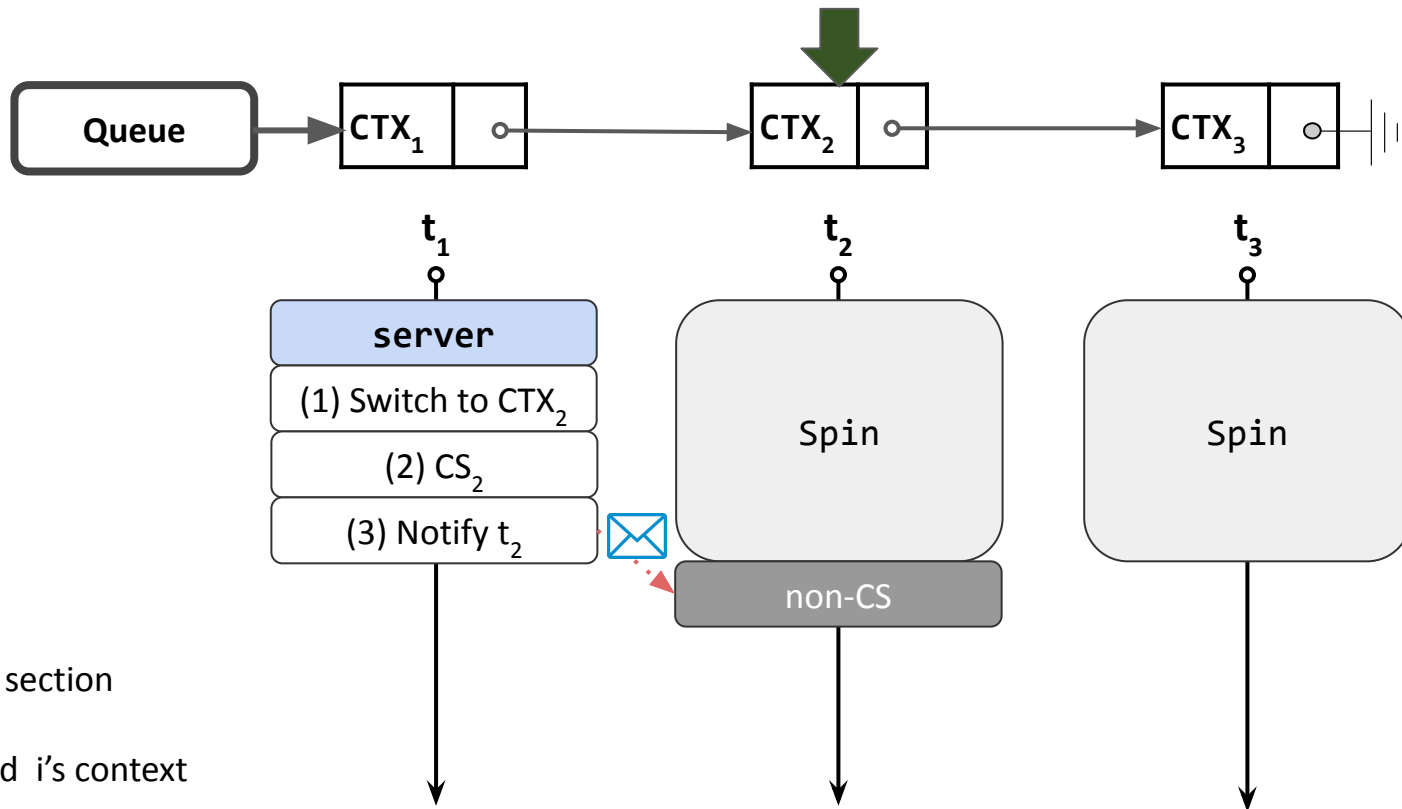
# TCLocks: Putting it all together

- Queue-based lock (Similar to qspinlock)
  - List of waiters maintained as a queue
  - Supports locking same lock in different contexts (Task, IRQs, NMI)
- Same lock/unlock API
- Server thread batches each waiters' request
- No dedicated server thread
  - Head of the queue becomes the server
  - The role is transferred to the next waiter after some threshold (Batch count)

# TCLocks in action: Phase 1



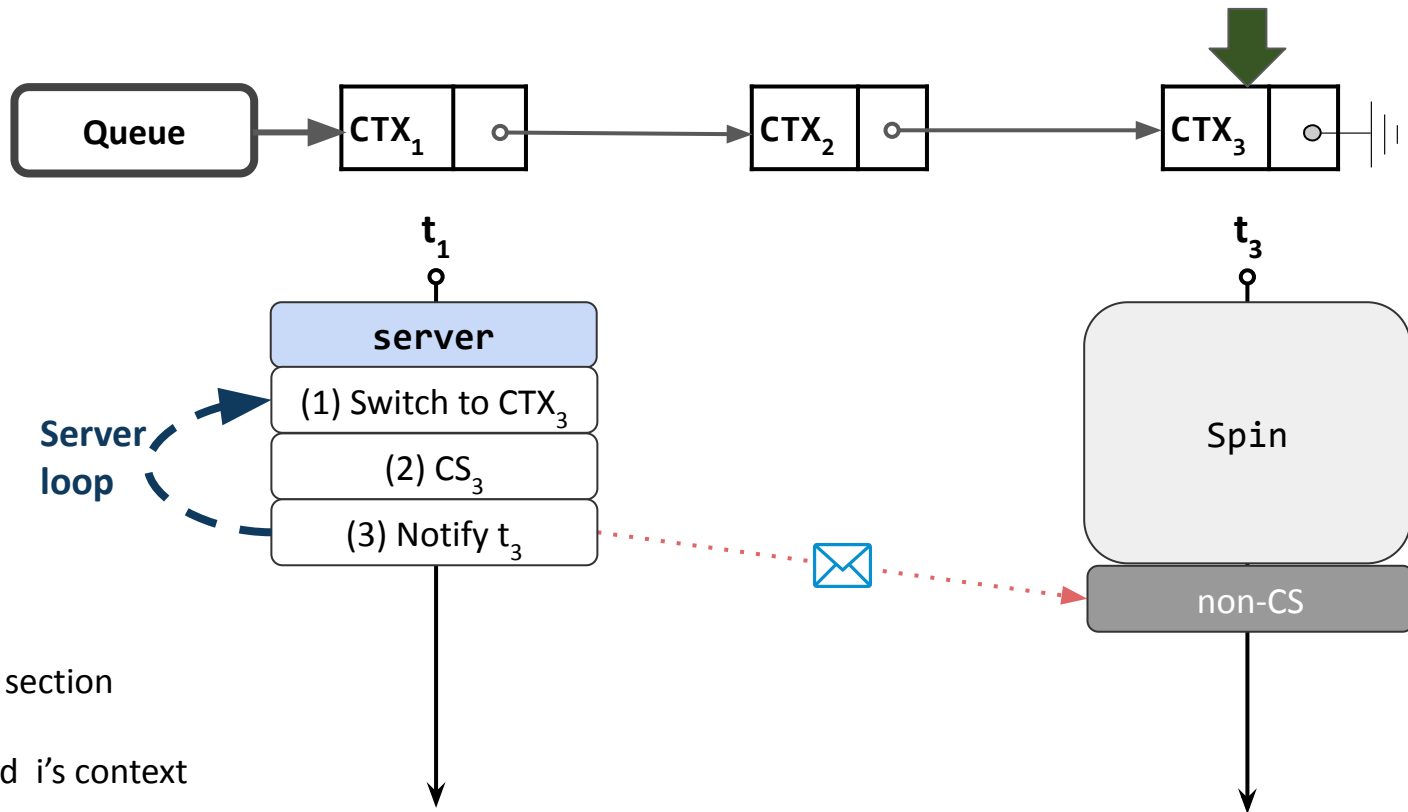
# TCLocks in action: Phase 2



CS: Critical section  
 $t_i$ : thread  $i$   
CTX <sub>$i$</sub> : thread  $i$ 's context



# TCLocks in action: Phase 2



# Agenda

- Motivation
- TCLock Design
- TCLock in Linux
- Evaluation
- Discussion

# TCLocks in Linux

- How to handle:
  - Waiter thread's state modification ?
  - per-CPU variables ?
  - Nested locking ?
  - Out-of-order unlocking ?
  - Mutex ?
  - Reader-Writer Semaphore ?

# TCLocks: Waiter's thread state modification

- Ideal scenario
  - Waiter's thread does not modify its context
- Reality
  - External events can modify waiter's context
    - Interrupts: Require stack access
    - Waiter's parking/wakeup mechanism: Require stack access
- **Solution: Ephemeral stack**
  - An empty piece of memory used only during critical section execution
  - Waiter's thread switches to Ephemeral stack after saving its context
  - This handles:
    - Interrupts on waiter's CPU
    - Waiter's thread parking/wakeup mechanism

# TCLocks: per-CPU Variables

- Kernel Assumption:
  - per-CPU variables are stable inside critical section
- With TLock
  - Critical section running on different CPU.
  - Different per-CPU variables are accessed.
  - Is this behavior correct ?

```
213 void mlock_drain_local(void)
214 {
215     struct folio_batch *fbatch;
216
217     local_lock(&mlock_fbatch.lock);
218     fbatch = this_cpu_ptr(&mlock_fbatch.fbatch);
219     if (folio_batch_count(fbatch))
220         mlock_folio_batch(fbatch);
221     local_unlock(&mlock_fbatch.lock);
222 }
```

- Yes, as long as it runs in a certain context
  - `!(irqs_disabled() || current->migration_disabled)` -> Run Combiner
  - Otherwise, fallback to `qspinlock`

# TCLocks: Nested Locking

- Kernel Assumption:
  - Multiple different locks can nest with arbitrary depth.
  - Same lock can also nest in different execution contexts.
- With TCLock
  - Server thread can become a waiter thread for nested lock
- Solution similar to interrupt processing mechanism
  - Save server thread's context on the stack before calling the nested lock.
  - Restore the server thread's context when nested lock returns.

```
static void __d_move(struct dentry *dentry, struct dentry *target,
                    bool exchange)
{
    /* target is not a descendent of dentry->d_parent */
    spin_lock(&target->d_parent->d_lock);
    spin_lock_nested(&old_parent->d_lock, DENTRY_D_LOCK_NESTED);
} else {
    BUG_ON(p == dentry);
    spin_lock(&old_parent->d_lock);
    if (p != target)
        spin_lock_nested(&target->d_parent->d_lock,
                        DENTRY_D_LOCK_NESTED);
}
spin_lock_nested(&dentry->d_lock, 2);
spin_lock_nested(&target->d_lock, 3);
```

# TCLocks: Out-of-order Unlocking

- Kernel Assumption:

- Nested locks can be unlocked in any order

- With TLock

- Server thread returns to its own context in the unlock function.
- It can return before the lock it held is not unlocked

```
1670  /*
1671   * Splice contents of ipipe to opipe.
1672   */
1673  static int splice_pipe_to_pipe(struct pipe_inode_info *ipipe,
1674                                struct pipe_inode_info *opipe,
1675                                size_t len, unsigned int flags)

1699      pipe_double_lock(ipipe, opipe);

1792      pipe_unlock(ipipe);
1793      pipe_unlock(opipe);
```

- Solution: Use an array to track lock order

- Delay unlocking the out-of-order unlocked lock until the remaining locks are unlocked.

# TCLocks: Mutex

- Differences from Spinlock:
  - Server thread state is stored in `task_struct` instead of per-CPU variables.
- Rest is similar to mutex in the kernel:
  - Except, currently it doesn't support: `Mutex_lock_interruptible` / `mutex_lock_killable`.
  - It is handled same as `mutex_lock`.



# TCLocks: Reader-writer semaphore

- Phase-based reader-writer lock:
  - Reader phase allows all readers to proceed, while writers are waiting.
  - Writer phase combines all writers using a server thread, while readers are waiting.

# Agenda

- Motivation
- TCLock Design
- TCLock in Linux
- Evaluation
- Discussion

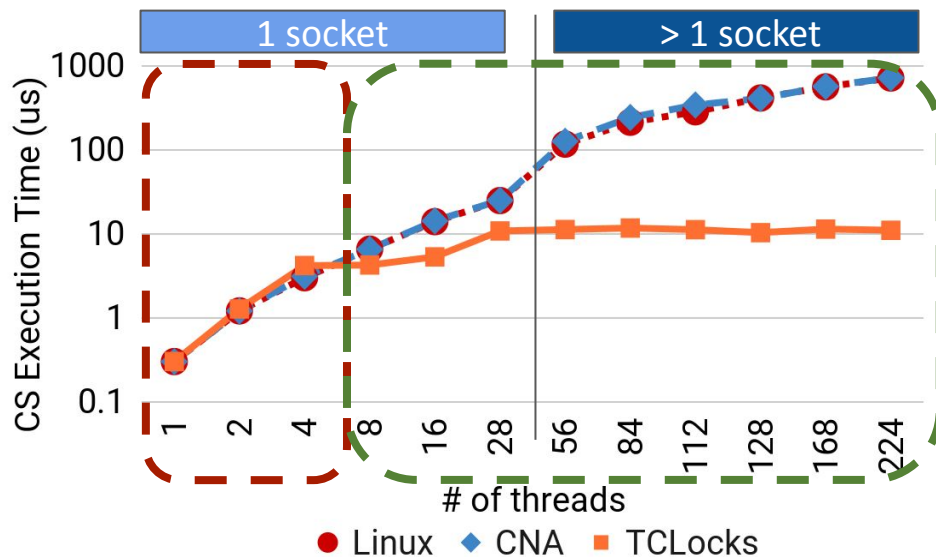
# TCLocks: Evaluation

- Does TCMutex reduce the time spent in critical section?
- Does TCMutex improve application performance?

Hardware:            8-socket/224-core Intel machine

# Evaluation: CS execution time

Benchmark: Each thread renames a file in a directory, serialized by a directory lock

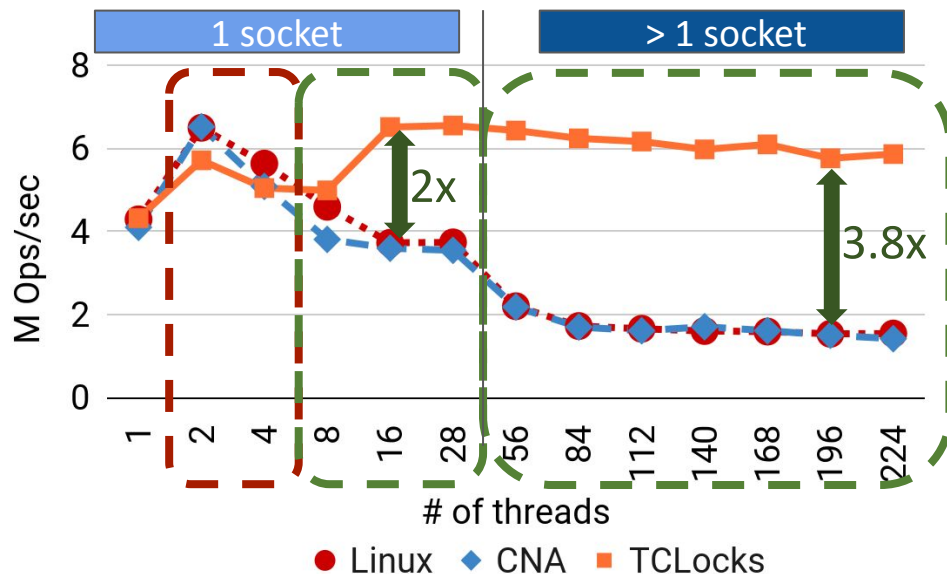


- **> 4 threads**
  - Minimal shared data movement
- **≤ 4 threads**
  - Context-switch overhead
  - Not enough batching

Setup: 8-socket/224-core machine

# Evaluation: Micro-benchmark

Benchmark: Each thread renames a file in a directory, serialized by a directory lock



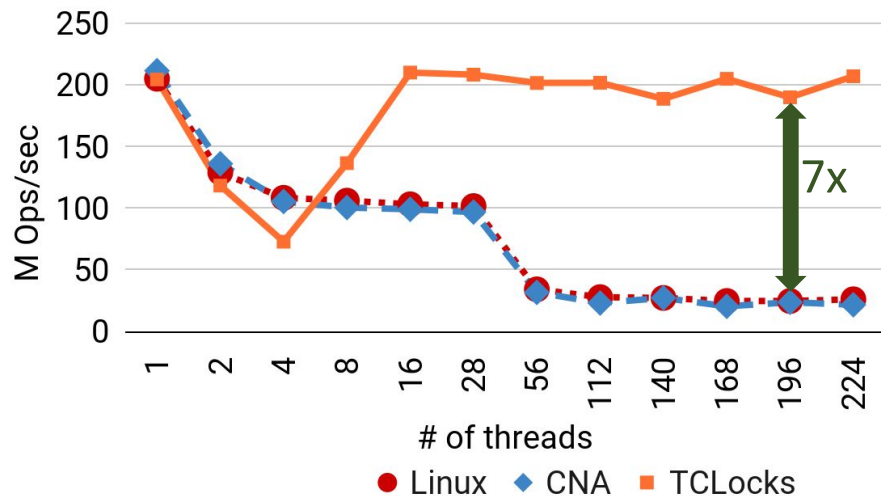
- **Within a socket:**
  - Minimal shared data movement
- **Across socket:**
  - NUMA-aware policy
- **2 - 4 cores:**
  - Context-switch overhead
  - Not enough batching

Setup: 8-socket/224-core machine

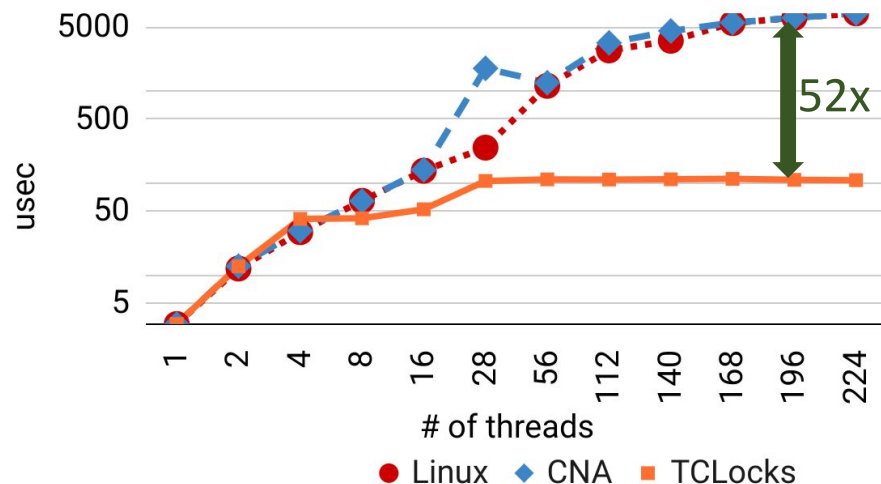
# Evaluation: Throughput and Latency

Benchmark: Each thread update an entry in hash-table, serialized by global spinlock

## Throughput



## 99%ile Latency (Lock + CS+Unlock)



TCLocks provides better throughput with lower 99% latency.

# Agenda

- Motivation
- TCLock Design
- TCLock in Linux
- Evaluation
- Discussion

# Discussion questions

- How to set the batching count ?
  - Throughput vs Latency
- How to handle performance regression at low contention (2-4 threads) ?
  - Switch between different lock mechanisms
  - TCLocks already uses qspinlock for certain contexts (IRQs disabled) and combining for others.
- How to handle CPU time accounting for server thread ?
  - Server thread might eat up the CPU time while executing other waiter's critical section.
  - Problem similar to CPU time accounting for interrupt processing.



# Discussion questions

- How to provide **current** macro correctly within and outside the critical section ?
  - Within a critical section, we need current of waiter's thread on server CPU.
  - Outside the critical section, we need current of server thread on server CPU.

# Conclusion

- Existing lock design:
  - Traditional lock design has more shared data movement
  - Delegation-based lock design requires application modification
- **TCLocks**: Provides transparent delegation
  - Capture thread's context at right time
- Key takeaway:
  - Applications can now use delegation-style locks without modification

<https://rs3lab.github.io/TCLocks/>

Thank you!

# Backup slides



# TCLocks: Pseudo-code

spin\_lock ():

node = get\_per\_cpu\_node()

save\_context\_on\_node(node)

join\_queue(node)

if( not head\_of\_queue() ):

    While node.wait is True:

        Continue

    restore\_context\_from\_node(node)

    Return

while True:

    qnext = get\_next\_thread()

    switch\_to(qnext)

    notify(qnext)

    If ( batch\_count\_exceeded() ); break

spin\_unlock():

    If (server\_context() ):

        switch\_to(server)

    else

        Lock = Unlocked