

Outline

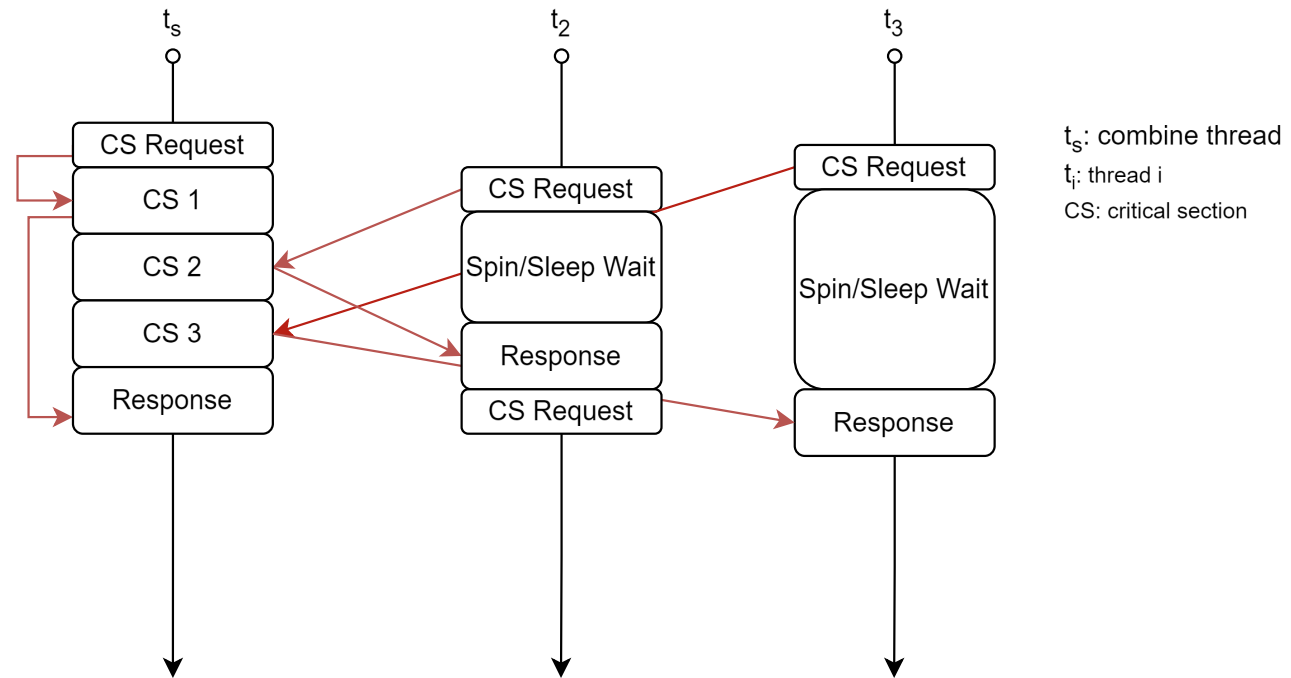
1. Recap: Delegation Styled Lock	1	4.4. Mutex	27
1.1. Imbalanced Workload	3	5. Code Change	28
1.2. Scheduler Subversion	5	5.1. Shared Counter	29
2. Solutions to Lock Usage Fairness	6		
2.1. Banning	8		
2.2. Naive Priority Queue	9		
2.3. Priority Queue (CFS like)	10		
3. Implementation	12		
3.1. Flat Combining	13		
3.2. CCSynch	17		
3.3. Flat Combining with Banning	19		
3.4. CCSynch with Banning	20		
3.5. FC-PQ	21		
3.6. FC-Skiplist	22		
4. Profiling Result	23		
4.1. Flat Combining	24		
4.2. CCSynch	25		
4.3. FC-PQ-BHeap	26		

1. Recap: Delegation Styled Lock

Thread publish their critical section to a job queue, and one server thread (combiner) will execute the job to prevent control flow switching.

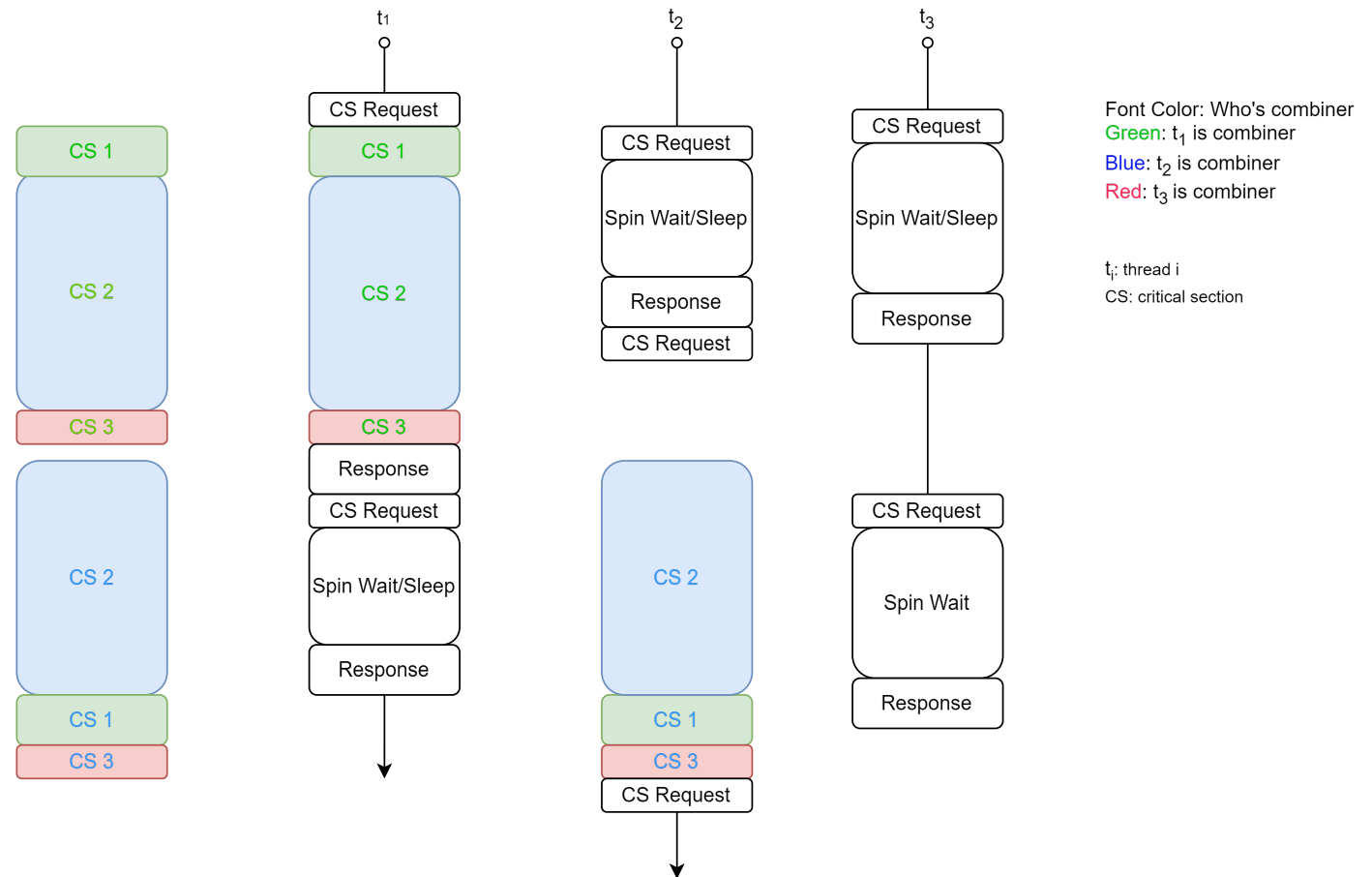
Two aspects:

1. How to elect the combiner thread
2. How to schedule the job



1.1. Imbalanced Workload

t_2 is occupying
more lock-
usage than t_1
and t_3



1.1. Imbalanced Workload

1.1.1. Lock Usage Fairness

- t_2 uses the lock longer than t_1 and t_3 because it has a longer critical section.
- t_1 can leave the lock earlier if using a normal lock, since when it acquires the lock, the lock is uncontended. However, now it needs to help other threads to execute the critical section.

1.2. Scheduler Subversion

- Assuming the threads are over-subscribed¹, scheduler will pre-empt the threads.
 - Assuming instead of spin-waiting, threads are sleeping during waiting².
1. In a normal lock, Scheduler Subversion happens when threads holding shorter critical section are spinning “longer”, which consumes CPU time but not progressing, while threads that holding longer critical section are using the lock more proportionally. This subverts the fairness goal provided by the scheduler.
 2. In delegation styled lock, the combiner thread is helping the other threads. Since this delegation is transparent to the scheduler, scheduler will penalize the combiner thread because of its voluntary work.

Example

1. When t_1 is helping t_2 and t_3 , it will use more CPU time, but not doing its own job.
2. Scheduler is not aware of the delegation, so it will try to schedule t_1 less time because of its voluntary work.

¹This is actually valid assumption, since delegation styled lock is much more scalable than normal lock which enable the potential of very large number of threads

²Consider the case when we needs very large number of threads

Outline

1. Recap: Delegation Styled Lock	1	4.4. Mutex	27
1.1. Imbalanced Workload	3	5. Code Change	28
1.2. Scheduler Subversion	5	5.1. Shared Counter	29
2. Solutions to Lock Usage Fairness	6		
2.1. Banning	8		
2.2. Naive Priority Queue	9		
2.3. Priority Queue (CFS like)	10		
3. Implementation	12		
3.1. Flat Combining	13		
3.2. CCSynch	17		
3.3. Flat Combining with Banning	19		
3.4. CCSynch with Banning	20		
3.5. FC-PQ	21		
3.6. FC-Skiplist	22		
4. Profiling Result	23		
4.1. Flat Combining	24		
4.2. CCSynch	25		
4.3. FC-PQ-BHeap	26		

2. Solutions to Lock Usage Fairness

1. Banning
2. Priority Queue (CFS like)
3. Other Scheduling Mechanism

2.1. Banning

- Similar to U-SCL, we ban threads that executes long critical section.

2.1.1. Implementation

- Flat Combining with Banning (Section 3.3)
- CCSynch with Banning (Section 3.4)

2.1.2. Problem

- If there are threads that enters the lock sparsely, there may be chance that all current contending threads are banned, while the lock remains unused.

2.2. Naive Priority Queue

- Use a priority queue as the job queue (e.g. skip-list).

2.2.1. Procedure

- Combiner elected via an `AtomicBool`
- Priority queue is implemented via skip-list ([crossbeam-skiplist](#))
- To execute a critical section, thread post the request to a local node and pushes it into the skip-list.
- Combiner will pop the job from the skip-list.

2.2.2. Implementation

1. Section 3.6

2.2.3. Illustration

TODO

2.2.4. Problem

- Performance overhead of skip-list is really high.
- Dequeue is expensive, which waste combiner's CPU time (i.e. wasting potential lock usage).

2.3. Priority Queue (CFS like)

- I want to propose something that is easy to implement, which allows more general scheduling mechanism.

2.3.1. Motivation

- Combiner has exclusive control over the lock-usage statistics (why do we need distributed ordering).
- We may want other scheduling mechanism, e.g. EEVDF.
- Node can be reused

2.3.2. Idea

- Use something like an MPSC channel to publish the job.
- The combiner thread polls the channel to get the job, and re-order the job based on execution.

2.3.3. Illustration

TODO

2.3.4. Challenges

TODO!

1. Deadlock of a naive implementation (TODO).
 1. Workaround: TODO
2. How to elect the combiner thread (subversion problem)?
3. Publishing node can be expensive
 1. Caching?
4. When do combiner check the channel?

2.3.5. Implementation

- Section 3.5

Outline

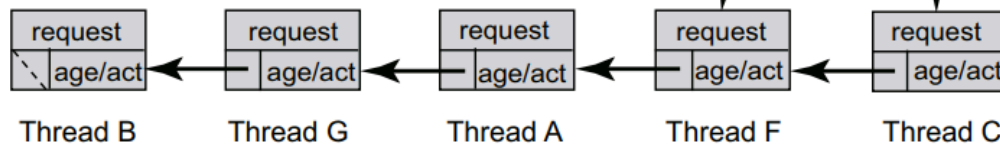
1. Recap: Delegation Styled Lock	1	4.4. Mutex	27
1.1. Imbalanced Workload	3	5. Code Change	28
1.2. Scheduler Subversion	5	5.1. Shared Counter	29
2. Solutions to Lock Usage Fairness	6		
2.1. Banning	8		
2.2. Naive Priority Queue	9		
2.3. Priority Queue (CFS like)	10		
3. Implementation	12		
3.1. Flat Combining	13		
3.2. CCSynch	17		
3.3. Flat Combining with Banning	19		
3.4. CCSynch with Banning	20		
3.5. FC-PQ	21		
3.6. FC-Skiplist	22		
4. Profiling Result	23		
4.1. Flat Combining	24		
4.2. CCSynch	25		
4.3. FC-PQ-BHeap	26		

3.1. Flat Combining

A singly linked list of nodes (belongs to each threads) are used to publish the job.

④ infrequently, new records are CASed by threads to head of list, and old ones are removed by combiner

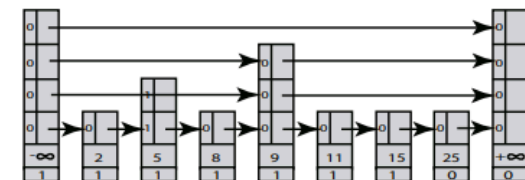
① thread writes request and spins on local record



publication list

③ combiner traverses list, performs scanCombineApply()

② thread acquires lock, becomes combiner, updates count



sequential data structure

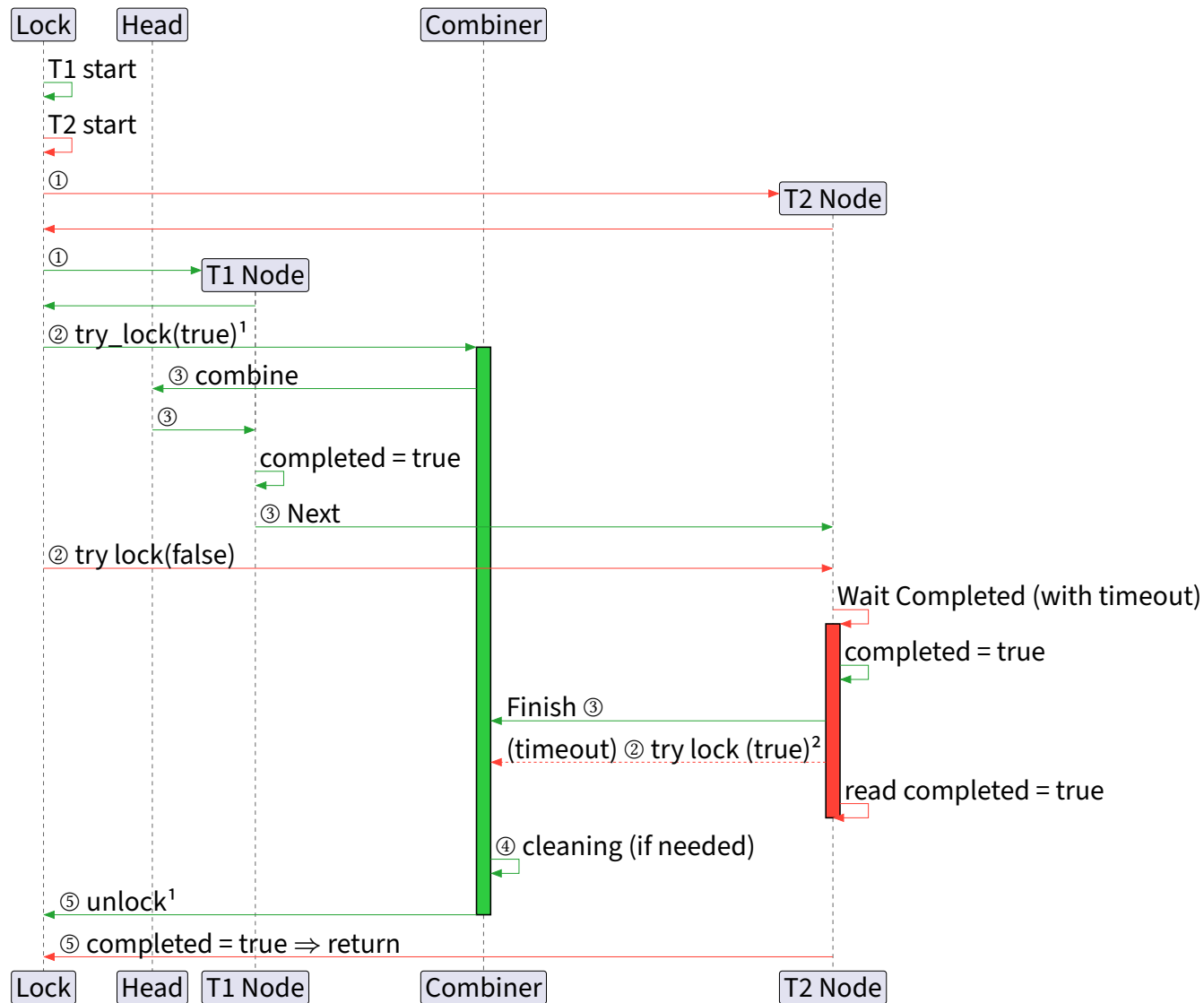
3.1. Flat Combining

3.1.1. Illustration 1

3.1. Flat Combining

¹A “lock” inside the `fc_lock` that is used to provide mutual exclusion for combiner

3.1.2. Illustration 2

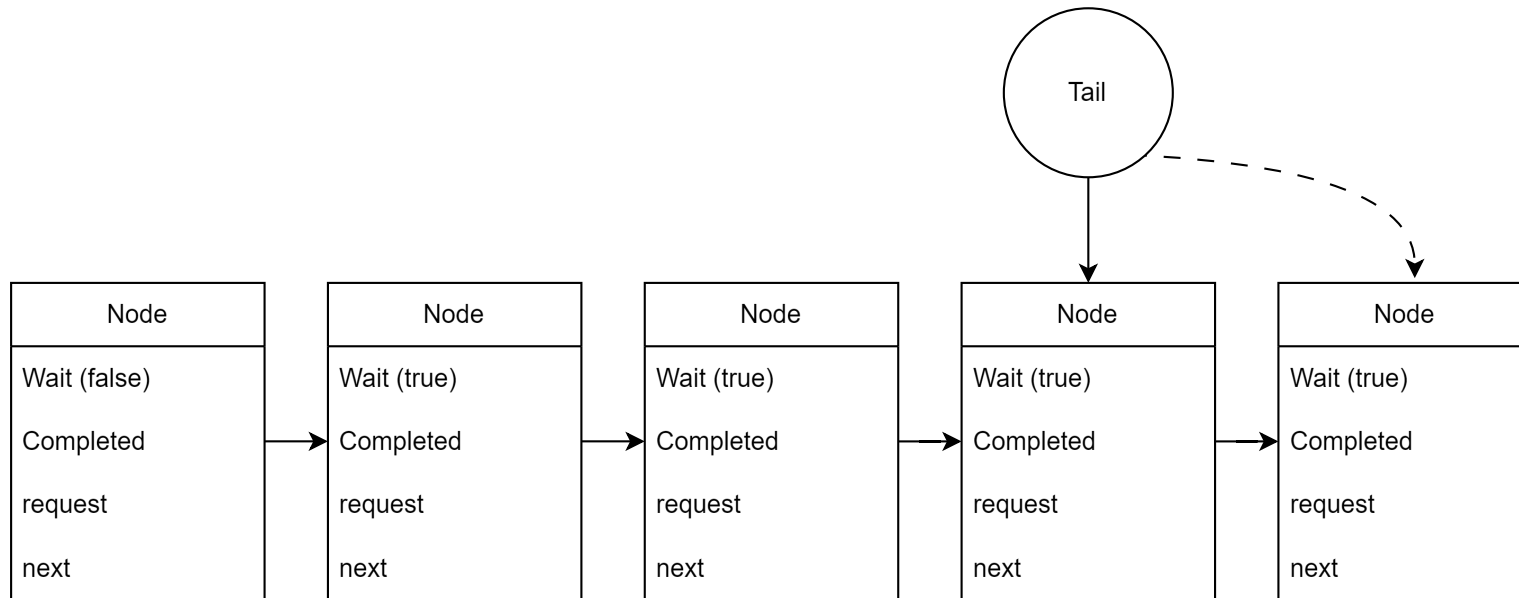


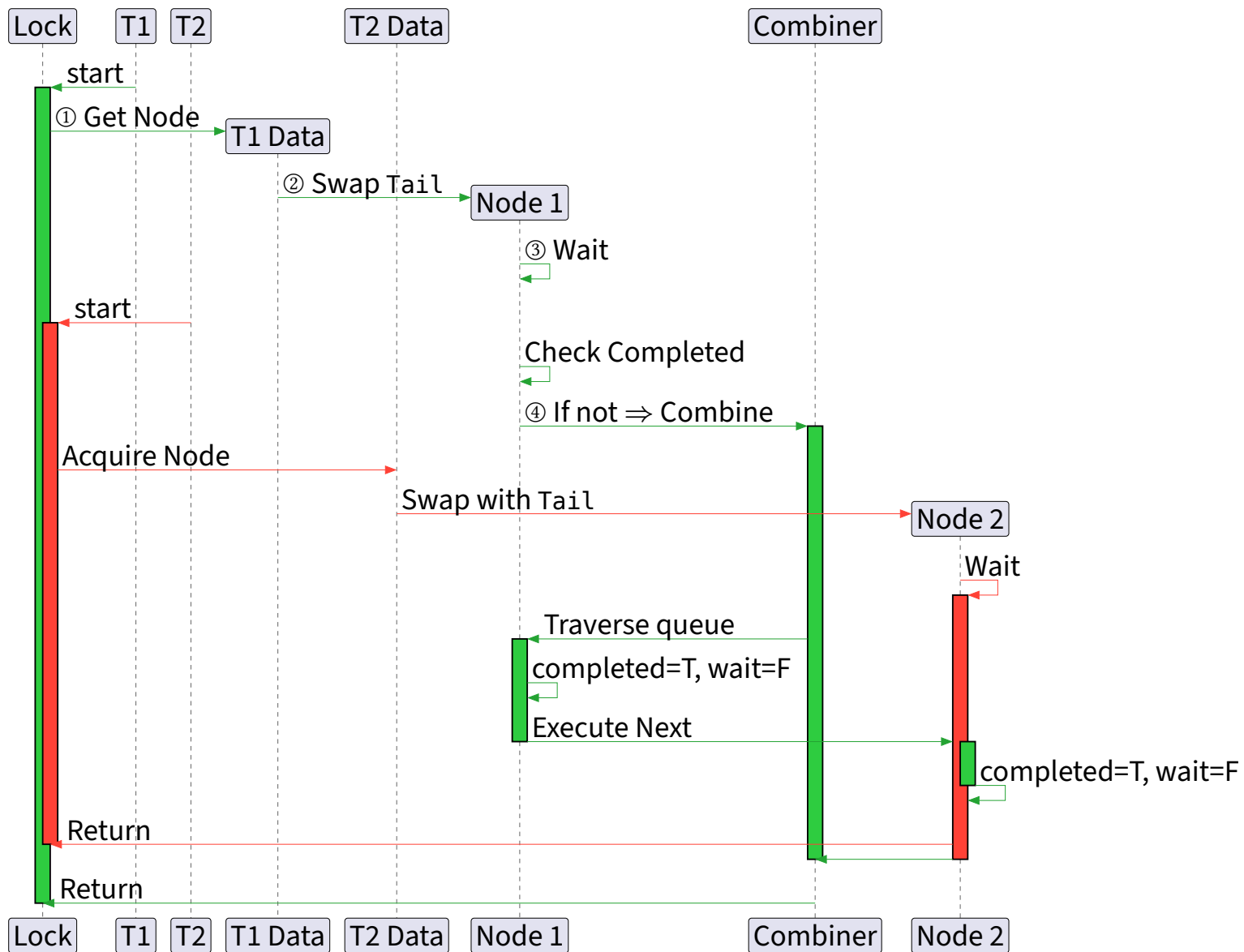
²Dashed Line means potentially executed but not in this case.

3.2. CCSynch

CCSynch maintains a FIFO queue of the job.

- ① Acquire a *next* node from `thread_local`
- ② Swap the `Tail` with the next node
- ③ Wait `wait`
- ④ if `!completed`, traverse the queue and execute jobs, set `completed` to `true`, and wait to `false`
- ⑤ when reaching combine limit, set next wait to `false`, and `completed` to `false`





3.3. Flat Combining with Banning

TODO!

3.4. CCSynch with Banning

TODO!

3.5. FC-PQ

TODO!

3.6. FC-Skiplist

TODO!

Outline

1. Recap: Delegation Styled Lock	1	4.4. Mutex	27
1.1. Imbalanced Workload	3	5. Code Change	28
1.2. Scheduler Subversion	5	5.1. Shared Counter	29
2. Solutions to Lock Usage Fairness	6		
2.1. Banning	8		
2.2. Naive Priority Queue	9		
2.3. Priority Queue (CFS like)	10		
3. Implementation	12		
3.1. Flat Combining	13		
3.2. CCSynch	17		
3.3. Flat Combining with Banning	19		
3.4. CCSynch with Banning	20		
3.5. FC-PQ	21		
3.6. FC-Skiplist	22		
4. Profiling Result	23		
4.1. Flat Combining	24		
4.2. CCSynch	25		
4.3. FC-PQ-BHeap	26		

4.1. Flat Combining

Function	CPU Time	Clockticks	Instructions Retired	CPI Rate	Module
lock	44.490s	1.11×10^{11}	1.30×10^9	85.489	dlock
bench code	4.682s	1.17×10^{10}	1.27×10^{10}	0.918	dlock
[vmlinux]	0.069s	1.49×10^8	6.27×10^7	2.368	vmlinux
thread::current	0.037s	9.68×10^7	1.98×10^7	4.889	dlock

Function	Retiring	Front-End Bound	Bad Speculation	Back-End Bound
lock	0.90%	0.10%	0.00%	99.10%
bench code	30.90%	0.30%	0.80%	68.00%
[vmlinux]	34.30%	57.20%	11.40%	0.00%
thread::current	87.80%	100.00%	0.00%	0.00%

4.2. CCSynch

4.3. FC-PQ-BHeap

4.4. Mutex

Outline

1. Recap: Delegation Styled Lock	1	4.4. Mutex	27
1.1. Imbalanced Workload	3	5. Code Change	28
1.2. Scheduler Subversion	5	5.1. Shared Counter	29
2. Solutions to Lock Usage Fairness	6		
2.1. Banning	8		
2.2. Naive Priority Queue	9		
2.3. Priority Queue (CFS like)	10		
3. Implementation	12		
3.1. Flat Combining	13		
3.2. CCSynch	17		
3.3. Flat Combining with Banning	19		
3.4. CCSynch with Banning	20		
3.5. FC-PQ	21		
3.6. FC-Skiplist	22		
4. Profiling Result	23		
4.1. Flat Combining	24		
4.2. CCSynch	25		
4.3. FC-PQ-BHeap	26		

5.1. Shared Counter

1. Remove blackbox for accessing the data
2. Change the blackbox position

```
1 - while black_box(loop_limit) > 0 {  
2 -     *data += 1;  
3 - }  
4 + while loop_limit > 0 {  
5 +     *black_box(&mut *data) += 1;  
6 +     loop_limit -= 1;  
7 + }
```

diff

5.1. Shared Counter

5.1.1. Reason for **blackbox**

1. `loop_limit` => the length of **Critical Section**.
2. Compiler will optimize the code to something like `*data += loop_limit;`, which will make varying the `loop_limit` not affecting the length of **Critical Section**.

5.1.2. Reason for the change

1. I want to mimic more access to the shared variable (hopefully something like `inc (rax)` in assembly).
2. The previous version contains too much overhead for doing the loop.