



Distributed Systems Part II

Solution to Exercise Sheet 10

1 Spin Locks

A read-write lock is a lock that allows either multiple processes to read some resource, or one process to write some resource.

- a) Write a simple read-write lock using only spinning, one shared integer and the CAS operation. Do not use local variables (it is ok to have variable within a method, but not outside).
- b) What is the problem with your lock?

Hint: what happens if a lot of processes access the lock repeatedly?

We now build a queue lock using only spinning, one shared integer, one local integer per process and the CAS operation.

- c) To prepare for this task, answer the following questions:
 - i) Head and tail of the queue have to be stored in the shared integer. What is the “head” and the “tail”, and how can they be stored in one integer?
Hint: could the head be a process id? Or is there a much easier solution?
 - ii) How could a process add itself to the queue?
Hint: you need the local integer of the process for this operation.
 - iii) When has a process acquired the lock?
 - iv) How does a process release the lock?
- d) Write down the lock using pseudo-code. Do not forget to initialize all variables.

Solution

- a) We use the shared integer `state` to indicate the state of the lock. The lock is free if `state` is 0. The lock is in write mode if `state` is -1. And it is in read-mode if `state` is n , with $n > 0$.

```
// the shared integer
int state = 0;

// acquire the lock for a read operation
read_lock(){
    while( true ){
        int value = state.read();
        if( value >= 0 ){
            if( state.CAS( value, value+1 ) == value ){
                // lock acquired
                return;
            }
        }
    }
}

// release the lock
read_unlock(){
    while( state.CAS( state, state-1 ) != state );
}

// acquire the lock for a write operation
write_lock(){
    while( true ){
        int value = state;
        if( state == 0 ){
            if( state.CAS( 0, -1 ) == 0 ){
                // lock acquired
                return;
            }
        }
    }
}

// release the lock
write_unlock(){
    // no need to test, no other process can call this at
    // the same time.
    state.CAS( -1, 0 );
}
```

- b) Starvation is a problem. Example: if many processes constantly acquire and release the read-lock, then the `state` variable always remains bigger than 0. If one process wants to acquire the write-lock, it will never get the chance.
- c) The basic idea behind this lock is a ticketing service as can be found in swiss post offices.
- i) The tail is the ticket which can be drawn by the next process. The head denotes the ticket which can acquire the lock. If we assume an integer consists of 32 bits, then we can use the first 16 bits for the head, and the last 16 bits for the tail.

- ii) The process reads the value of the tail, and then increments the tail. This should of course happen in a secure way, i.e. no two processes have the same ticket.
- iii) When its ticket equals the head.
- iv) The process increments the head by one.

d) // the shared integer containing head|tail
shared int queue = 0;

```

// the ticket of this process
int local = 0;

// acquire the lock
lock(){
    // 1. add this process to the queue
    local = add();
    // 2. wait until the lock is acquired
    while( head() != local() );
}

// add this process to the queue
int add(){
    while( true ){
        int value = queue.read();
        if( queue.CAS( value, value+1 ) == value ){
            return value & 0xFF;
        }
    }
}

// returns the current head of the queue
int head(){
    int value = state.read();
    return (value >>> 16) & 0xFF;
}

// releases the lock
unlock(){
    while( true ){
        int value = queue.read();
        int head = (value >>> 16) & 0xFF
        int tail = value & 0xFF
        int next = (head+1) << 16 | tail;
        if( queue.CAS( value, next ) == value ){
            return;
        }
    }
}

```

2 ALock2

Have a look at the source code below. It is a modified version of the ALock (slides 3/46 ff).

```

public class ALock2 implements Lock {
    int [] flags = {true, true, false, ..., false };
    AtomicInteger next = new AtomicInteger(0);
}

```

```

ThreadLocal<Integer> mySlot;

public void lock() {
    mySlot = next.getAndIncrement();
    while (!flags[mySlot % n]) {}
    flags[mySlot % n] = false;
}

public void unlock() {
    flags[(mySlot+2) % n] = true;
}
}

```

- a) What was the intention of the author of “ALock2”?
- b) Will ALock2 work properly? Why (not)?
- c) Give an idea how to repair ALock2.
Hint: don't bother about performance.

Solution

- a) The author wants that two processes can acquire the lock simultaneously.
- b) The lock is seriously flawed. An example shows how the lock will fail: Assume there are n processes, all processes try to acquire the lock. The first two processes (p_1, p_2) get the lock, the others have to wait. Process p_1 keeps the lock a very long time, while p_2 releases the lock almost immediately. Afterwards every second process (p_4, p_6, \dots) acquires and releases the lock. One half of all process are waiting on the lock (p_3, p_5, \dots), the others continues to work (p_4, p_6, \dots). If the working process now start to acquire the lock again, then they wait in slots that are already in use.
- c) A solution would be to increase the size of the array to at least $2 * n$ and further block the lock() method if a process holds the other lock for a (too) long time. This way, wrap-arounds are handled correctly.

Unfortunately FIFO (first in, first out) is still not guaranteed. In a second step one could make the unlock method more intelligent: instead of jumping two slots, the method searches for the oldest slot waiting for a lock. To simplify this search, the boolean array is replaced by an enum (or integer) array holding four states: unused, lockable, working, and finished. We can use a CAS operation to protect the unlock method against race conditions (two process may invoke the method concurrently).

```

public class ALock2 implements Lock{

    public enum State{UNUSED, LOCKABLE, WORKING, FINISHED};

    // >= 2n elements: 2 lockable, >= (n-2) unused, n finished
    State[] flags = {LOCKABLE, LOCKABLE, UNUSED, ..., UNUSED, FINISHED, ...};

    AtomicInteger next = new AtomicInteger ( 0 ) ;
    ThreadLocal<Integer> mySlot;

    public void lock(){
        // spin until slot becomes lockable
        mySlot = next.getAndIncrement();
        while(flags[mySlot%n] != LOCKABLE){}
    }
}

```

```

// mark as working
flags[mySlot%n] = WORKING;

// wait for other lock if its process is too slow (larger array helps here)
// & mark the slot as unused to support wrap-arounds
while(flags[(mySlot - flags.size() + n)%n] != FINISHED){}
flags[(mySlot - flags.size() + n)%n] = UNUSED;
}

public void unlock(){
    flags[mySlot] = FINISHED; // mark my slot as finished...

    // set next unused slot to lockable
    int index = mySlot + 1;
    while(true){
        if(flags[index%n] != UNUSED)
            index++;
        else if(flags[index%n].CAS(UNUSED, LOCKABLE) == UNUSED)
            break; // next unused slot became lockable
    }
}
}

```

3 MCS Queue Lock

See slides 3/65 ff.

- a) A developer suggests to add an **abort** flag to each node: if a process no longer wants to wait it sets this **abort** flag to **true**. If a process unlocks the lock, it may see the **abort** flag of the next node, jump over the aborted node, and check the successor's successor node. Modify the basic algorithm to support aborts.
Optional: sketch a proof for your answer.
Hint: Be aware of race-conditions!
- b) Assuming many processes may abort concurrently, does your answer from a) still work? Explain why. If it does not work: modify your algorithm to allow concurrent aborts.
Optional: sketch a proof for your answer.
- c) Instead of a **locked** and an **aborted** flag one could use an integer, and modify the integer with the CAS operation. What do you think about this idea? How is the algorithm affected? How is performance affected?
- d) The CLH lock (slide 3/56) is basically the same as an MCS lock. Conceptually the only difference is, that a process spins on the **locked** field of the predecessor node, not on its own node. What could be an advantage of CLH over MCS and what could be a disadvantage?

Solution

- a) There is more than one solution, but we can solve this problem without using RMW registers or other locks. It is important to set and read the flags in the right order: The **unlock** method first sets **locked**, then reads **aborted**. The **abort** method on the other hand first sets **aborted**, then reads **locked**. This way if **unlock** and **abort** run in parallel, one of them must already have written its flag before the other can read it. In the worst case **unlock** is called twice for some process, but that is not a problem. Unlocking an already unlocked lock results in no action.

```

public void unlock(){
    if( ... missing successor ... )
        ... wait for missing successor

    qnode.next.locked = false;
    if( qnode.next.aborted ){
        if( ... qnode.next misses successor ... ){
            if( ... really no successor ... )
                return;
        }
        else{
            ... wait for missing successor ...
        }
        qnode.next.next.locked = false;
    }
}

public void abort(){
    qnode.aborted = true;
    if( !qnode.locked ){
        unlock();
    }
}

```

- b) The solution of a) does not yet work for concurrent aborts. Making the `unlock` method recursive will help.

```

public void unlock(){
    unlock( qnode );
}
private void unlock( QNode qnode ){
    // as before...
    if( ... missing successor ... )
        ... wait for missing successor

    qnode.next.locked = false;
    if( qnode.next.aborted ){

        // wait for successor of qnode.next
        if( ... ){ ... } else{ ... }

        unlock( qnode.next );
    }
}

```

- c) There are four combinations of values the `locked` and `aborted` flag can have. We can easily encode these combinations in an integer. We would not need too worry about the order in which we read and write to the flags, as we could do this atomically. So the algorithm would get easier. We could also ensure that `unlock` is called only once. Depending on the benchmark this could increase the performance. On the other hand, a `CAS` operation is quite expensive and could decrease performance.
- d) - There could be problems with caches: spinning on a value that “belongs” to another process can introduce additional load on the bus, and thus slow down the entire application.

- + The implementation is much easier: when releasing the lock one has only to set its own `locked` flag to `false`.
- + Also aborting is easier: a blocked process could read the state of its predecessor. If the predecessor is aborted, then the successor can just remove the node from the queue, and continue reading values from its predecessor's predecessor.