



Lock-Free Programming: Pro Tips

Jean-Philippe BEMPEL
Performance Architect

@jpbempel
<http://jpbempel.blogspot.com>

Agenda

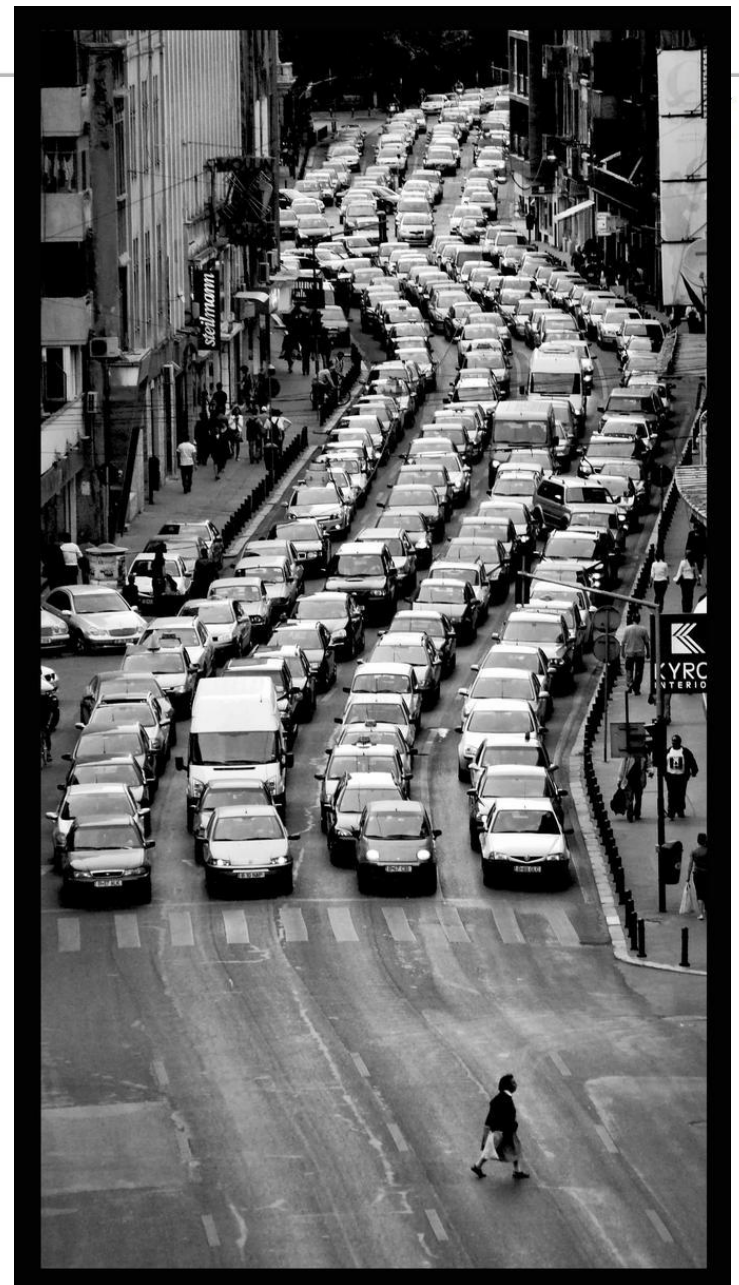
- Measuring Contention
- Lock Striping
- Compare-And-Swap
- Introduction to Java Memory Model
- Disruptor & RingBuffer
- Spinning
- Ticketing: OrderedScheduler

Immutability

Contention

Contention

- Two or more thread competing to acquire a lock
- Thread parked when waiting for a lock
- Number one reason we want to avoid lock



Measure, don't guess!

Kirk Pepperdine & Jack Shirazi

Measure, don't premature!

Measuring Contention

Synchronized blocks:

- Profilers (YourKit, JProfiler, ZVision)
- JVMTI native agent
- Results may be difficult to exploit

Measuring Contention: JProfiler

Local attach [2] - JProfiler 8.1.3

Session View Profiling Window Help

Start Center Detach Save Snapshot Session Settings Start Recordings Stop Recordings Start Tracking Run GC Add Bookmark Export View Settings Help Record Monitors Freeze View Show In Heap Walker Show In Graph

Session Profiling View specific

Live memory Heap walker CPU views Threads Monitors & locks

Current Locking Graph Current Monitors Locking History Graph Monitor History Monitor Usage Statistics Telemetries

Show monitors: Waiting and blocking Threshold in ms: 0 Filter by All text fields Reset Filters

Time	Duration	Type	Monitor ID	Monitor Class	Waiting Thread	Owning Thread
0:03.344 [Feb 26, 2015 ...]	8976 µs	Blocked	1	com.ullink.ultools.locks.SyncContention\$MyLock	main [main]	Thread-0 [main]
0:03.353 [Feb 26, 2015 ...]	1940 ms	Blocked	1	com.ullink.ultools.locks.SyncContention\$MyLock	Thread-0 [main]	main [main]
0:05.293 [Feb 26, 2015 ...]	1007 µs	Blocked	1	com.ullink.ultools.locks.SyncContention\$MyLock	main [main]	Thread-0 [main]
0:05.294 [Feb 26, 2015 ...]	9980 µs	Blocked	1	com.ullink.ultools.locks.SyncContention\$MyLock	Thread-0 [main]	main [main]
0:05.304 [Feb 26, 2015 ...]	1009 µs	Blocked	1	com.ullink.ultools.locks.SyncContention\$MyLock	main [main]	Thread-0 [main]
0:05.305 [Feb 26, 2015 ...]	340 ms	Blocked	1	com.ullink.ultools.locks.SyncContention\$MyLock	Thread-0 [main]	main [main]
0:05.646 [Feb 26, 2015 ...]	160 ms	Blocked	1	com.ullink.ultools.locks.SyncContention\$MyLock	Thread-0 [main]	main [main]
0:05.806 [Feb 26, 2015 ...]	8984 µs	Blocked	1	com.ullink.ultools.locks.SyncContention\$MyLock	main [main]	Thread-0 [main]
0:05.815 [Feb 26, 2015 ...]	29 ms	Blocked	1	com.ullink.ultools.locks.SyncContention\$MyLock	Thread-0 [main]	main [main]
0:05.845 [Feb 26, 2015 ...]	1062 µs	Blocked	1	com.ullink.ultools.locks.SyncContention\$MyLock	main [main]	Thread-0 [main]
0:05.846 [Feb 26, 2015 ...]	160 ms	Blocked	1	com.ullink.ultools.locks.SyncContention\$MyLock	Thread-0 [main]	main [main]
0:06.006 [Feb 26, 2015 ...]	1907 µs	Blocked	1	com.ullink.ultools.locks.SyncContention\$MyLock	main [main]	Thread-0 [main]
0:06.008 [Feb 26, 2015 ...]	160 ms	Blocked	1	com.ullink.ultools.locks.SyncContention\$MyLock	Thread-0 [main]	main [main]
Total:	28 s					

Recording thresholds: 1 000 µs blocking / 100 000 µs waiting [\[Change\]](#)

Filtered stack trace for waiting thread:
<no stack trace was recorded>

Filtered stack trace for owning thread:
<no stack trace was recorded>

unlicensed copy for evaluation purposes, 10 days remaining 0 active recordings Auto-update 2 s VM #3 04:57 Profiling

Measuring Contention: ZVision



Azul ZVision

User: root Host: archi-srv.ullink.lan / 11098 Version: 1.7.0-zing_14.11.0.0-b12

Domain: Group label: Application label:

Uptime: 00:11:38

[Overview](#) | [Azul Support](#) | [Threads](#) | [CPU](#) | [Memory](#) | [Compilers](#) | [Applications](#)

[List](#) | [Stack trace](#) | [Deadlocks](#) | [Contention](#)

Lock Statistics of [sun.misc.Launcher\\$AppClassLoader](#)

Contention Tree

- 1. 0.00% 0ms 1 [com.ullink.ultools.tracefileng.TraceFileConsumer.run](#) (TraceFileConsumer.java:105, bci=69)
 - 2. - 0ms 1 [java.util.concurrent.ThreadPoolExecutor.runWorker](#) (ThreadPoolExecutor.java:1145, bci=95)
 - 3. - 0ms 1 [java.util.concurrent.ThreadPoolExecutor\\$Worker.run](#) (ThreadPoolExecutor.java:615, bci=5)
 - 4. - 0ms 1 [java.lang.Thread.run](#) (Thread.java:745, bci=11)
- 1. 0.00% 0ms 1 [com.sun.jna.Native.initIDs](#) (native method)
 - 2. - 0ms 1 [com.sun.jna.Native](#) (Native.java:135, bci=91)
 - 3. - 0ms 1 [com.ullink.ultools.sys.CpuAffinity\\$LinuxCpuAffinity\\$LinuxCLibrary](#) (CpuAffinity.java:605, bci=5)
 - 4. - 0ms 1 [com.ullink.ultools.sys.CpuAffinity\\$LinuxCpuAffinity.getCurrentThreadId](#) (CpuAffinity.java:690, bci=0)
 - 5. - 0ms 1 [com.ullink.ultools.sys.CpuAffinity\\$LinuxCpuAffinity.getCpuBinding](#) (CpuAffinity.java:620, bci=1)
 - 1. 0.00% 0ms 1 [com.ullink.ultools.jvm.impl.HotspotGCMonitor.run](#) (HotspotGCMonitor.java:77, bci=89)
 - 2. - 0ms 1 [com.ullink.ultools.jvm.MXBeanMonitor.reportHighFrequency](#) (MXBeanMonitor.java:197, bci=11)
 - 3. - 0ms 1 [com.ullink.ultools.jvm.MXBeanMonitor.run](#) (MXBeanMonitor.java:165, bci=17)
 - 4. - 0ms 1 [java.lang.Thread.run](#) (Thread.java:745, bci=11)
 - 1. 0.00% 0ms 1 [com.ullink.ulbridge2.ULBridge.registerMBeans](#) (ULBridge.java:3452, bci=53)
 - 2. - 0ms 1 [com.ullink.ulbridge2.ULBridge.access\\$4700](#) (ULBridge.java:208, bci=1)
 - 3. - 0ms 1 [com.ullink.ulbridge2.ULBridge\\$PromoteInitFinalization.promote](#) (ULBridge.java:2494, bci=4)
 - 4. - 0ms 1 [com.ullink.ultools.failover.impl.HotFailoverManager.promote](#) (HotFailoverManager.java:616, bci=30)
 - 5. - 0ms 1 [com.ullink.ulbridge2.failover.BridgeHotFailoverManager.promote](#) (BridgeHotFailoverManager.java:51, bci=7)

Raw Counters

klass_sma_successes	0
klass_sma_failures	0
contended_lock_attempts	4
contended_max_time	0
contended_total_time	1
wait_count	0
wait_max_time	0
wait_total_time	0

Measuring Contention

`java.util.concurrent.Lock:`

- JVM cannot help us here
- JDK classes (lib), regular code
- JProfiler can measure them
- `j.u.c` classes modification + bootclasspath (jucprofiler)

Measuring Contention: JProfiler

Local attach [2] - JProfiler 8.1.3

Session View Profiling Window Help

Start Center Detach Save Snapshot Session Settings Start Recordings Stop Recordings Start Tracking Run GC Add Bookmark Export View Settings Help Record Monitors Freeze View Show In Heap Walker Show In Graph

Session Profiling View specific

Live memory Heap walker CPU views Threads Monitors & locks

Current Locking Graph Current Monitors Locking History Graph Monitor History Monitor Usage Statistics Telemetries

Show monitors: **Waiting and blocking** Threshold in ms: **0** Filter by **All text fields** **Reset Filters**

Time	Duration	Type	Monitor ID	Monitor Class	Waiting Thread	Owning Thread
983:52.472 [Feb 26, 201...	49 ms	Blocked	1	java.util.concurrent.locks.ReentrantReadWriteLock\$Nonfair...	main [main]	Thread-0 [main]
983:52.522 [Feb 26, 201...	50 ms	Blocked	1	java.util.concurrent.locks.ReentrantReadWriteLock\$Nonfair...	main [main]	Thread-0 [main]
983:52.572 [Feb 26, 201...	50 ms	Blocked	1	java.util.concurrent.locks.ReentrantReadWriteLock\$Nonfair...	main [main]	Thread-0 [main]
983:52.622 [Feb 26, 201...	50 ms	Blocked	1	java.util.concurrent.locks.ReentrantReadWriteLock\$Nonfair...	main [main]	Thread-0 [main]
983:52.672 [Feb 26, 201...	49 ms	Blocked	1	java.util.concurrent.locks.ReentrantReadWriteLock\$Nonfair...	main [main]	Thread-0 [main]
983:52.722 [Feb 26, 201...	49 ms	Blocked	1	java.util.concurrent.locks.ReentrantReadWriteLock\$Nonfair...	main [main]	Thread-0 [main]
983:52.772 [Feb 26, 201...	50 ms	Blocked	1	java.util.concurrent.locks.ReentrantReadWriteLock\$Nonfair...	main [main]	Thread-0 [main]
983:52.822 [Feb 26, 201...	49 ms	Blocked	1	java.util.concurrent.locks.ReentrantReadWriteLock\$Nonfair...	main [main]	Thread-0 [main]
983:52.872 [Feb 26, 201...	49 ms	Blocked	1	java.util.concurrent.locks.ReentrantReadWriteLock\$Nonfair...	main [main]	Thread-0 [main]
983:52.922 [Feb 26, 201...	50 ms	Blocked	1	java.util.concurrent.locks.ReentrantReadWriteLock\$Nonfair...	main [main]	Thread-0 [main]
983:52.972 [Feb 26, 201...	49 ms	Blocked	1	java.util.concurrent.locks.ReentrantReadWriteLock\$Nonfair...	main [main]	Thread-0 [main]
983:53.022 [Feb 26, 201...	49 ms	Blocked	1	java.util.concurrent.locks.ReentrantReadWriteLock\$Nonfair...	main [main]	Thread-0 [main]
983:53.072 [Feb 26, 201...	50 ms	Blocked	1	java.util.concurrent.locks.ReentrantReadWriteLock\$Nonfair...	main [main]	Thread-0 [main]
Total:	12 s					

Recording thresholds: **1 000 µs blocking / 100 000 µs waiting** [\[Change\]](#)

Filtered stack trace for waiting thread:

```
java.util.concurrent.locks.ReentrantReadWriteLock$WriteLock.lock()
com.ullink.ultools.locks.ProfilingReentrantLock.p(java.util.concurrent.locks.Lock)
com.ullink.ultools.locks.ProfilingReentrantLock.profReentrantReadWriteLock()
com.ullink.ultools.locks.ProfilingReentrantLock.main(java.lang.String[])
com.intelij.rt.execution.application.AppMain.main(java.lang.String[])
```

Filtered stack trace for owning thread:

```
java.lang.Thread.sleep(long)
com.ullink.ultools.locks.ProfilingReentrantLock$2.run()
```

unlicensed copy for evaluation purposes, 10 days remaining 0 active recordings Auto-update 2 s VM #2 985:05 Profiling

Measuring Contention

- Insertion of contention counters
 - Identify place where lock fail to be acquired
 - increment counter
- Identify Locks
 - Call stacks at construction
 - Logging counter status
- How to measure existing locks in your code
 - Modify JDK classes
 - Reintroduce in bootclasspath

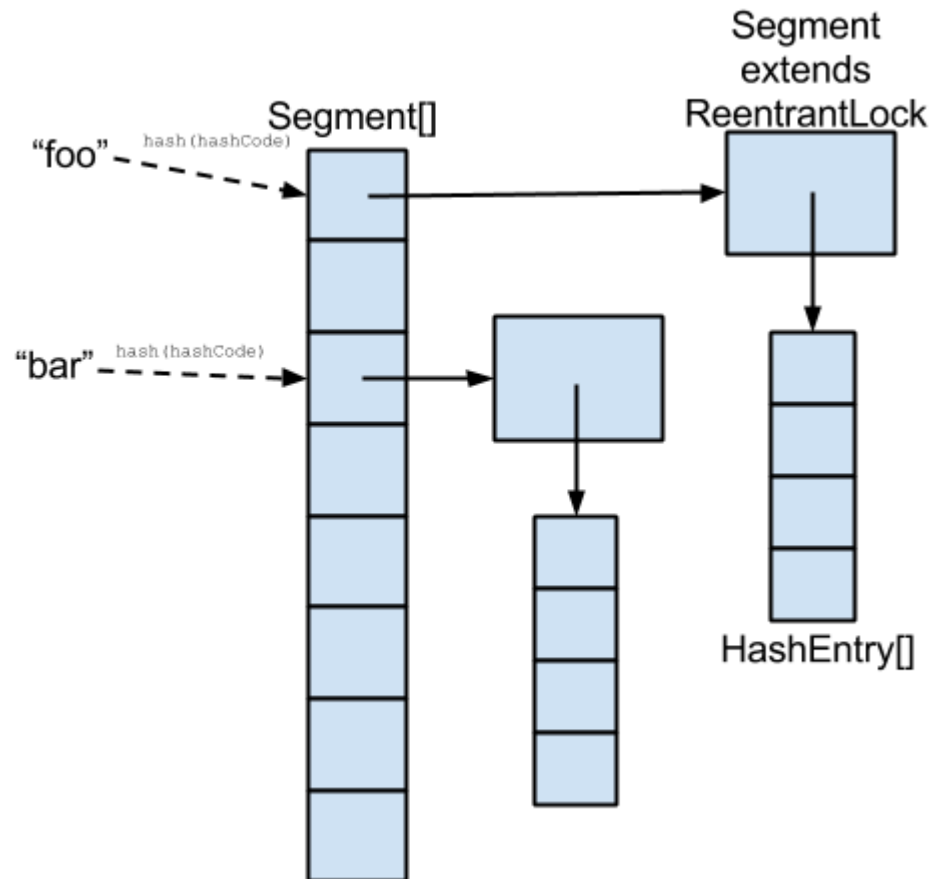
Lock striping

Lock striping

- Reduce contention by distributing it
- Not remove locks, instead adding more
- Good partitioning is key to be effective (like HashMap)

Lock striping

Best example in JDK: ConcurrentHashMap



Lock striping

- Relatively easy to implement
- Can be very effective as long as good partitioning
- Can be tuned (number of partition) regarding the contention/concurrency level

Compare-And-Swap

Compare-And-Swap

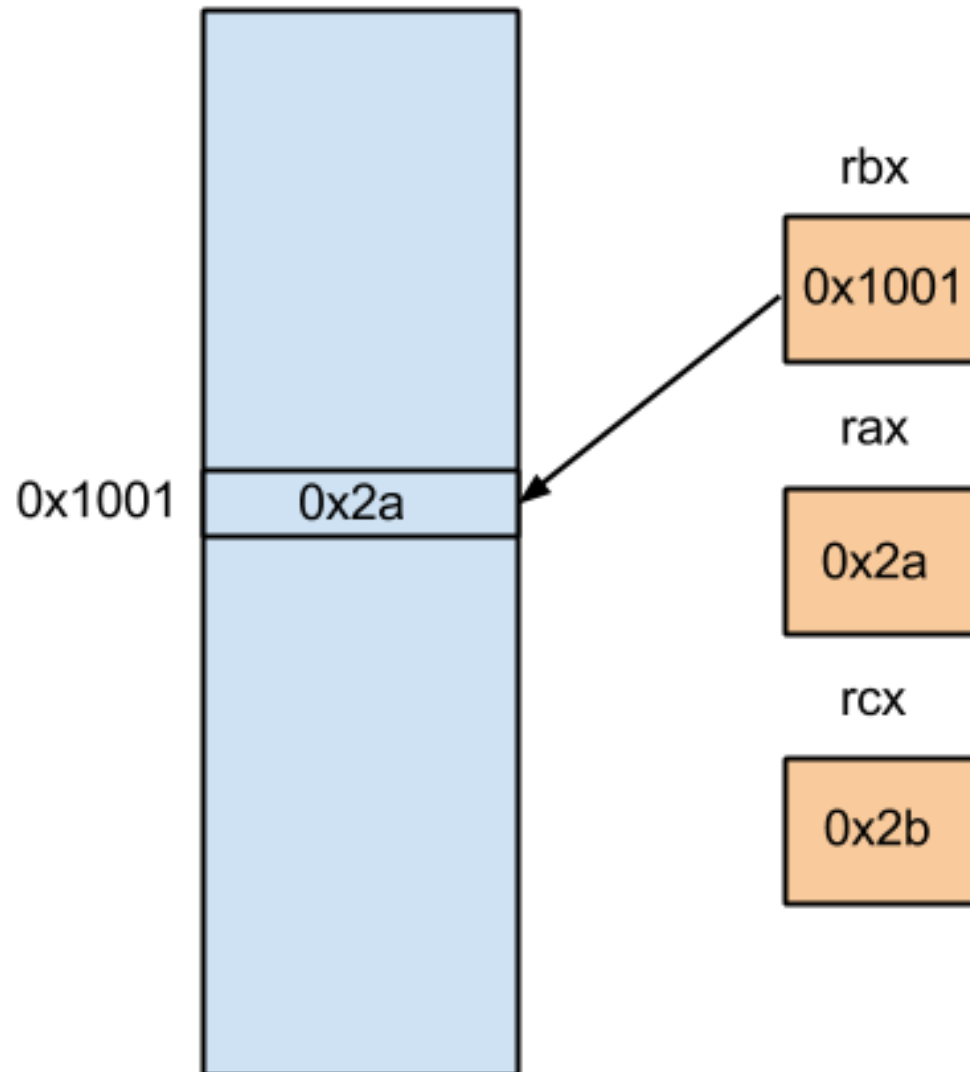
- Basic primitive for any lock-free algorithm
- Used to implement any locks or synchronization primitives
- Handled directly by the CPU (instructions)

Compare-And-Swap

- Update atomically a memory location by another value if the previous value is the expected one
- instruction with 3 arguments:
 - memory address (rbx)
 - expected value (rax)
 - new value (rcx)

```
movabs rax,0x2a
movabs rcx,0x2b
lock cmpxchg QWORD PTR [rbx],rcx
```

Compare-And-Swap



Compare-And-Swap

- In Java for AtomicXXX classes:

```
boolean compareAndSet(long expect, long update)
```

- Memory address is the internal `value` field of the class

Compare-And-Swap: AtomicLong

- Atomic increment with CAS

[JDK7] getAndIncrement():

```
while (true) {
    long current = get();
    long next = current + 1;
    if (compareAndSet(current, next))
        return current;
}
```

[JDK8] getAndIncrement():

```
return unsafe.getAndAddLong(this, valueOffset, 1L);
```

intrinsicified to:

```
movabs rsi, 0x1
lock xadd QWORD PTR [rdx+0x10], rsi
```

Compare-And-Swap: Lock implementation

ReentrantLock is implemented with a CAS:

```
volatile int state;
```

```
lock()  
    compareAndSet(0, 1);
```

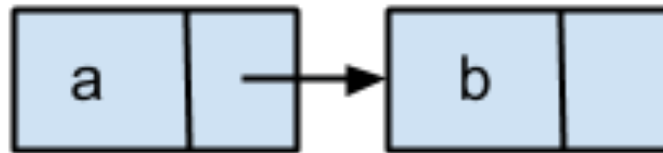
if CAS fails => lock already acquired

```
unlock()  
    setState(0)
```

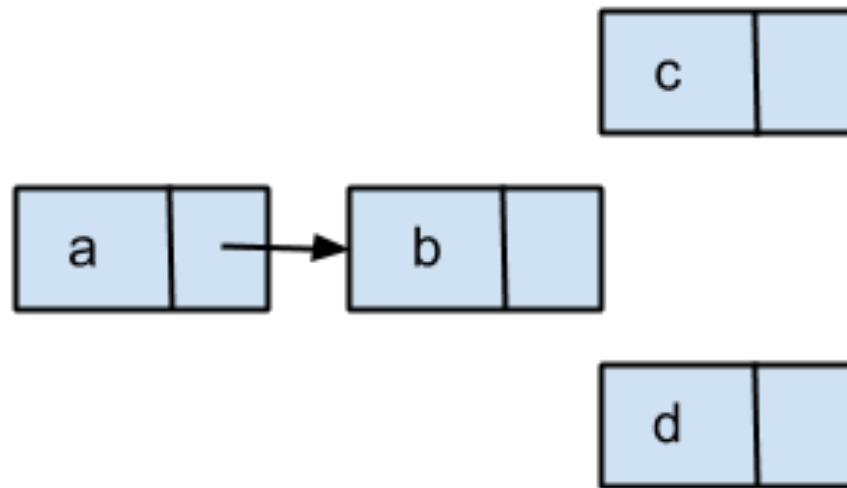

Compare-And-Swap: ConcurrentLinkedQueue

- Simplest lock-free algorithm
- Use CAS to update the next pointer into a linked list
- if CAS fails, means concurrent update happened
- Read new value, go to next item and retry CAS

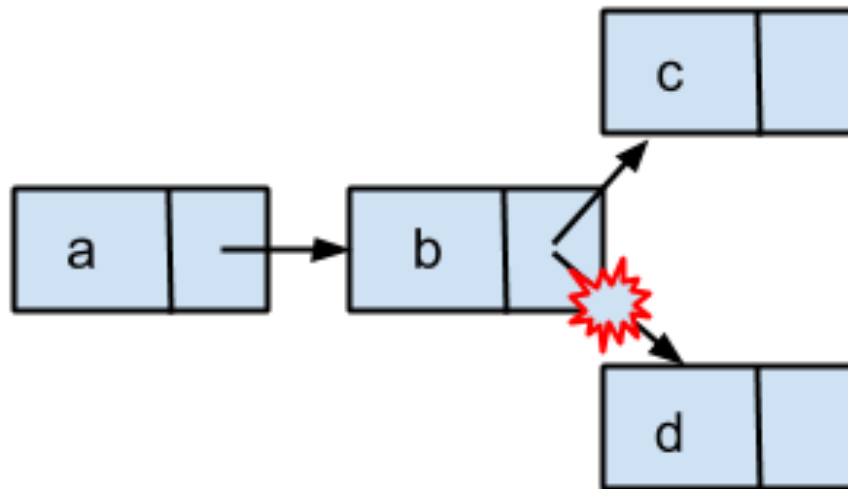
Compare-And-Swap: ConcurrentLinkedListQueue



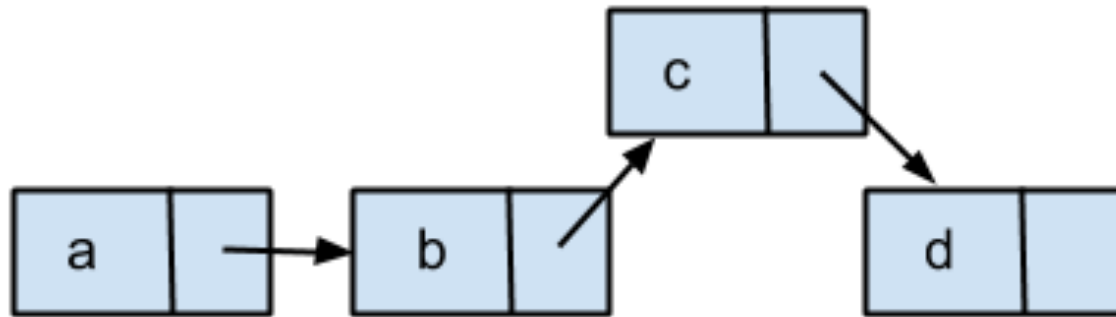
Compare-And-Swap: ConcurrentLinkedListQueue



Compare-And-Swap: ConcurrentLinkedListQueue



Compare-And-Swap: ConcurrentLinkedListQueue



Java Memory Model (introduction)

Memory Model

- First language having a well defined memory model:
Java JDK 5 (2004) with JSR 133
- C++ get a standard Memory Model in 2011 (C++11)
- Before that, some constructions may have
undefined/different behavior on different platform
(Double Check Locking)

Memory ordering

```
int a;
int b;
boolean enabled;
```

```
{
    a = 21;
    b = a * 2;
    enabled = true;
}
```

JIT Compiler



```
{
    enabled = true;
    a = 21;
    b = a * 2;
}
```


Memory ordering

```
int a;
int b;
boolean enabled;
```

Thread 1

```
{
    a = 21;
    b = a * 2;
    enabled = true;
}
```

Thread 2

```
{
    if (enabled)
    {
        int answer = b
        process(answer);
    }
}
```

Memory ordering

```
int a;
int b;
volatile boolean enabled;
```

Thread 1

```
{
    a = 21;
    b = a * 2;
    enabled = true;
}
```

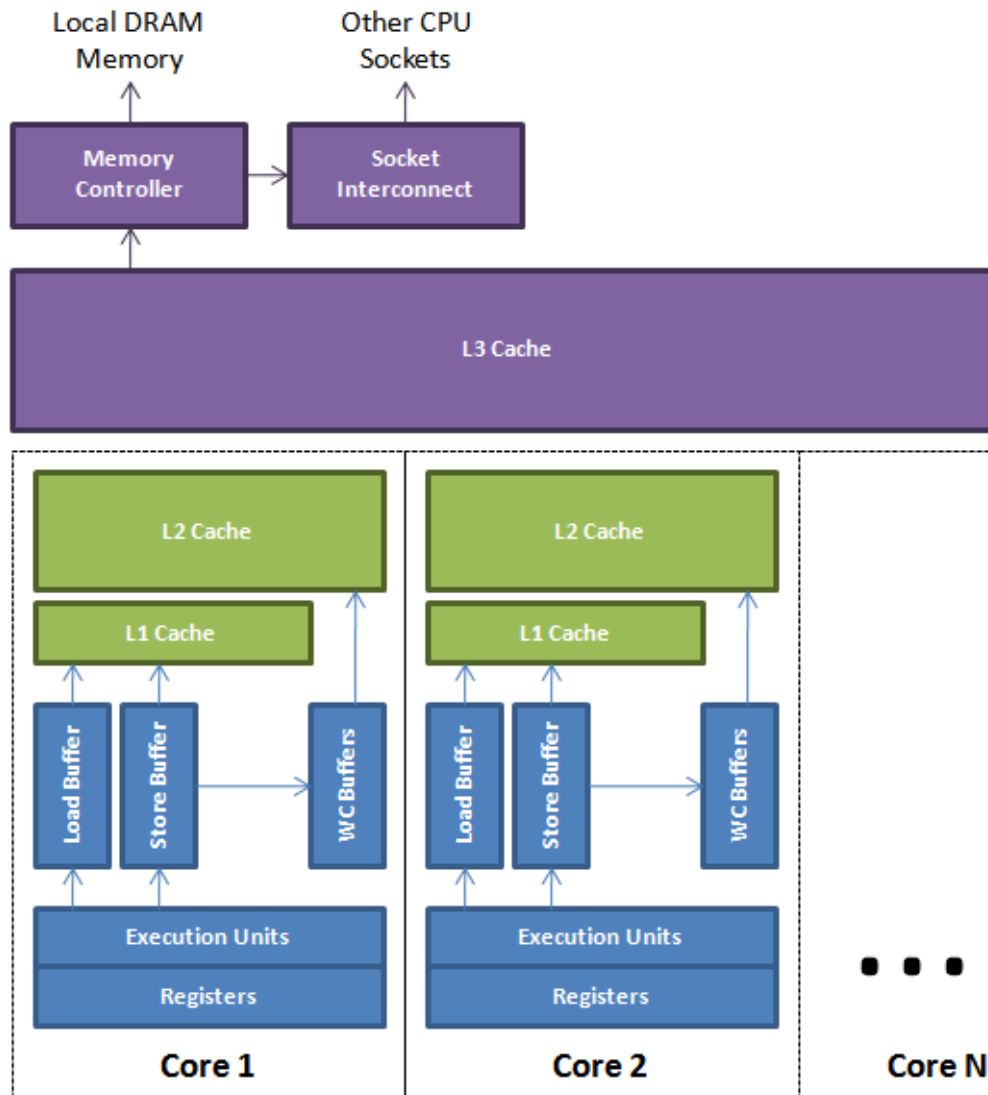
Thread 2

```
{
    if (enabled)
    {
        int answer = b;
        process(answer);
    }
}
```

Memory barriers

- Can be at 2 levels: Compiler & Hardware
- Depending on CPU architecture, barrier is not required
- on x86: Strong model, limited reordering

Memory barriers



Memory barriers: volatile

- volatile field implies memory barrier
- Compiler barrier: prevent reordering
- Hardware barrier: Ensure drain of the memory buffers
- on X86, only store barrier emits an hardware one

```
lock add DWORD PTR [rsp],0x0
```

Memory barriers: CAS

- CAS is also a memory barrier
- Compiler: recognized by JIT to prevent reordering
- Hardware: all lock instructions is a memory barrier

Memory barriers: synchronized

- Synchronized blocks have implicit memory barriers
- Entering block: Load memory barrier
- Exiting block: store memory barrier

Memory barriers: synchronized

```
synchronized (this)
{
    enabled = true;
    b = 21;
    a = b * 2;
}
```

Memory barriers: lazySet

- method from AtomicXXX classes
- Compiler only memory barrier
- Does not emit hardware store barrier
- Still guarantee non reordering (most important)
but not immediate effect for other thread

Disruptor & Ring Buffer

Disruptor

- LMAX library (incl. Martin Thompson)
- Not a new idea, circular buffers in Linux Kernel, Lamport
- Ported to Java

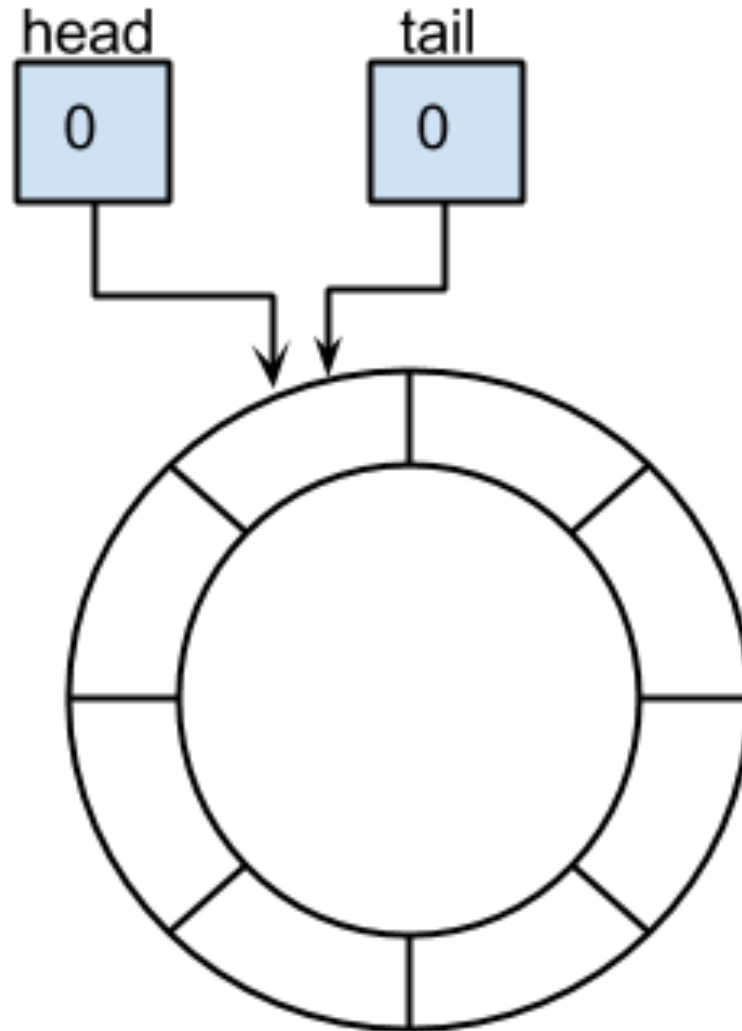
Disruptor

Why not used CLQ which is lock(wait)-free?

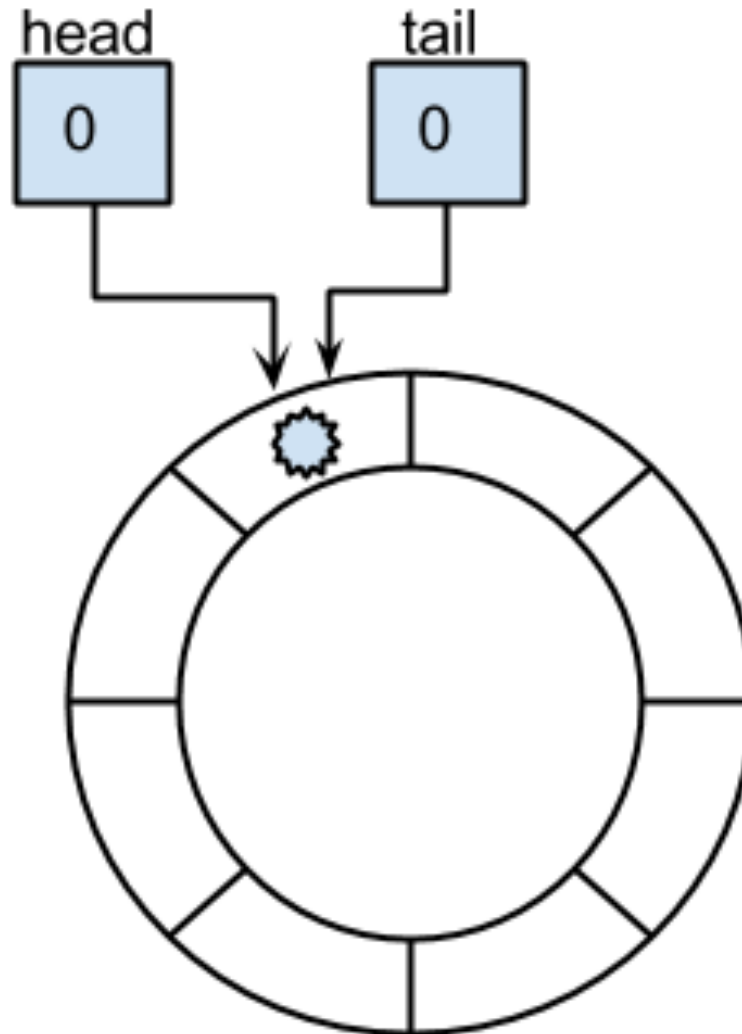
- Queue unbounded et non blocking
- Allocate a node at each insertion
- Not CPU cache friendly
- MultiProducer and MultiConsumer

Array/LinkedBlockingQueue: Not lock-free

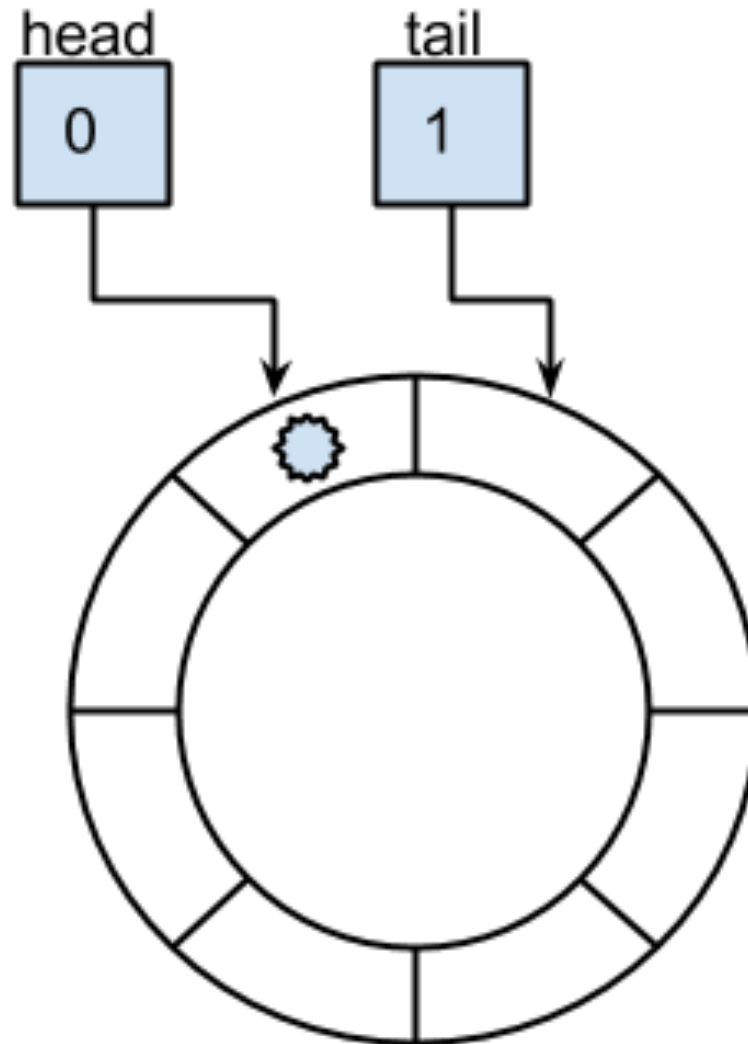
Ring Buffer: 1P 1C



Ring Buffer: 1P 1C



Ring Buffer: 1P 1C



Ring Buffer: 1P 1C

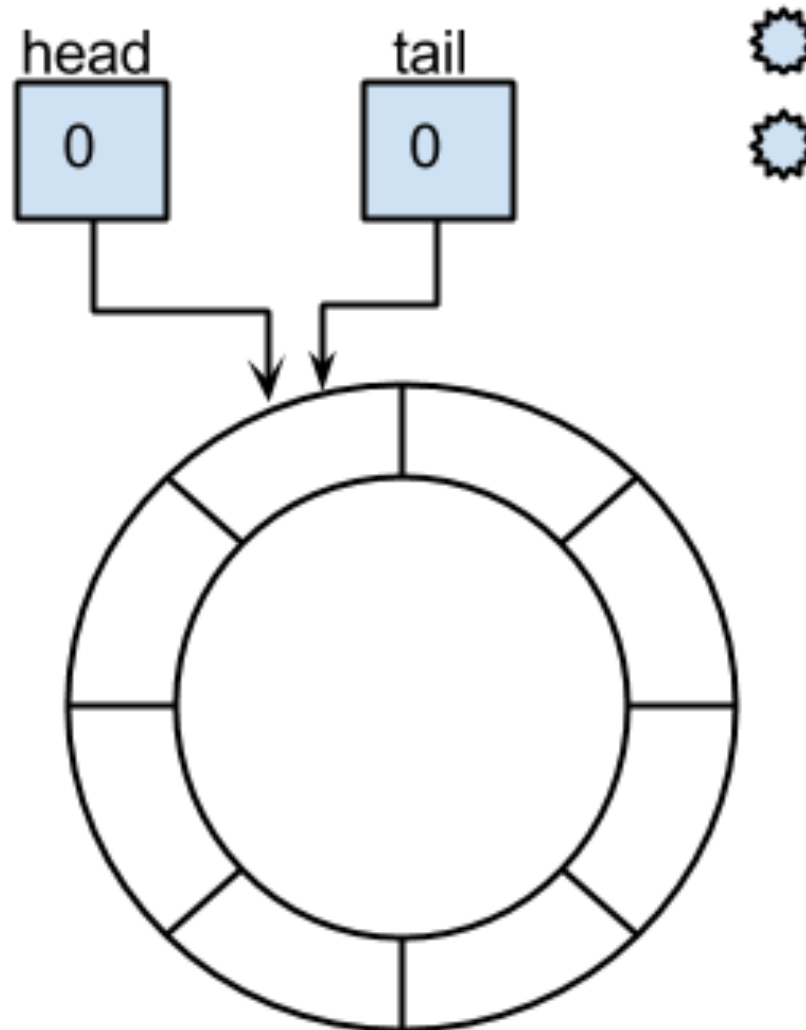
```
Object[] ringBuffer;
volatile int head;
volatile int tail;

public boolean offer(E e) {
    if (tail - head == ringBuffer.length)
        return false;
    ringBuffer[tail % ringBuffer.length] = e;
    tail++; // volatile write
    return true;
}
```

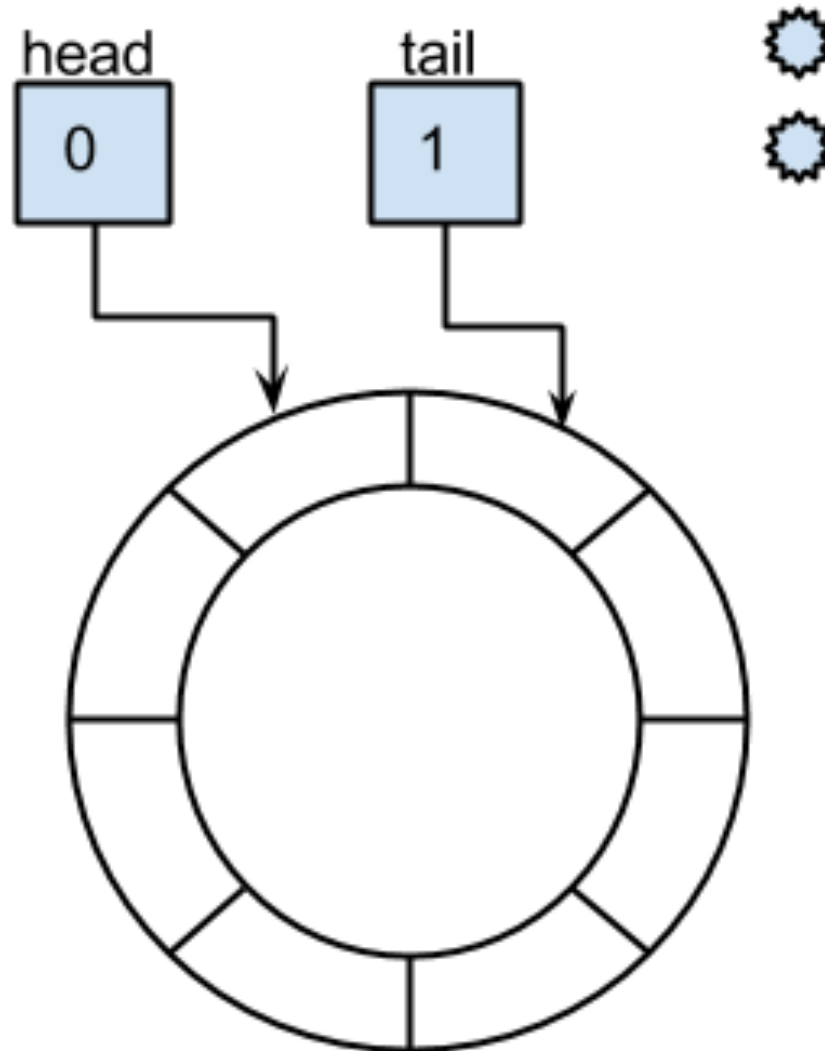

Ring Buffer: 1P 1C

```
public E poll() {  
    if (tail == head)  
        return null;  
    int idx = head % ringBuffer.length  
    E element = ringBuffer[idx];  
    ringBuffer[idx] = null;  
    head++; // volatile write  
    return element;  
}
```

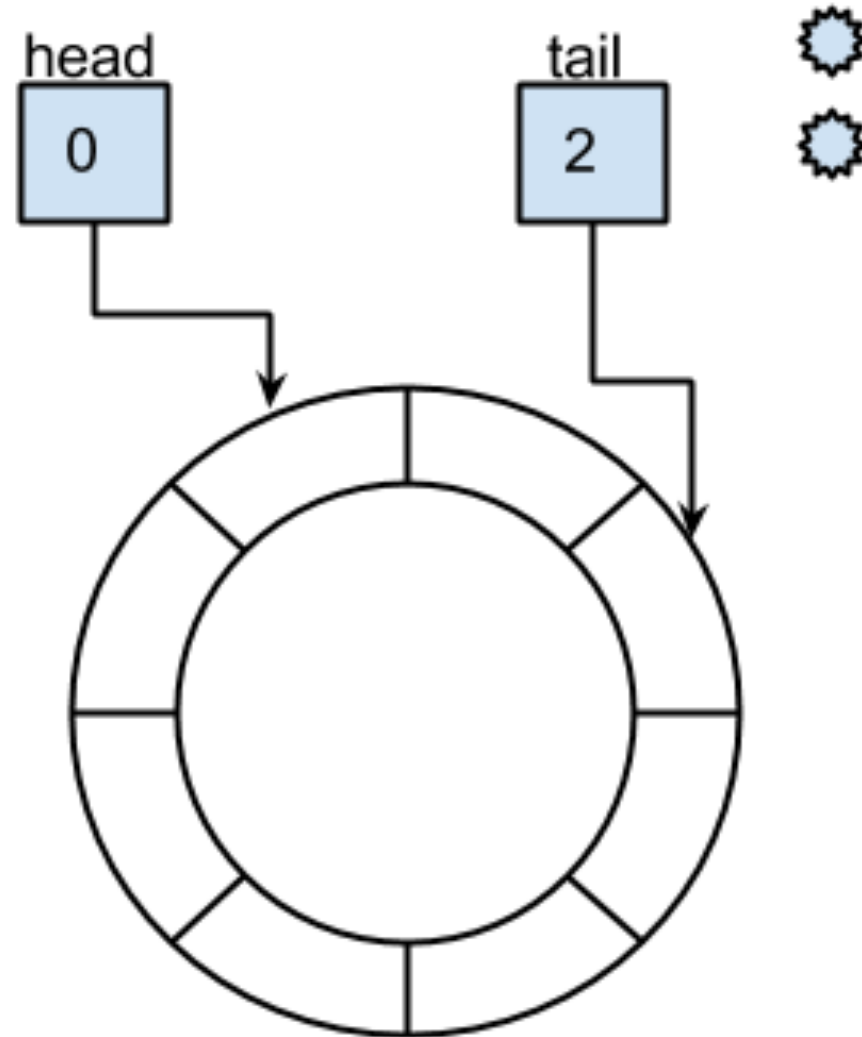
Ring Buffer: nP 1C



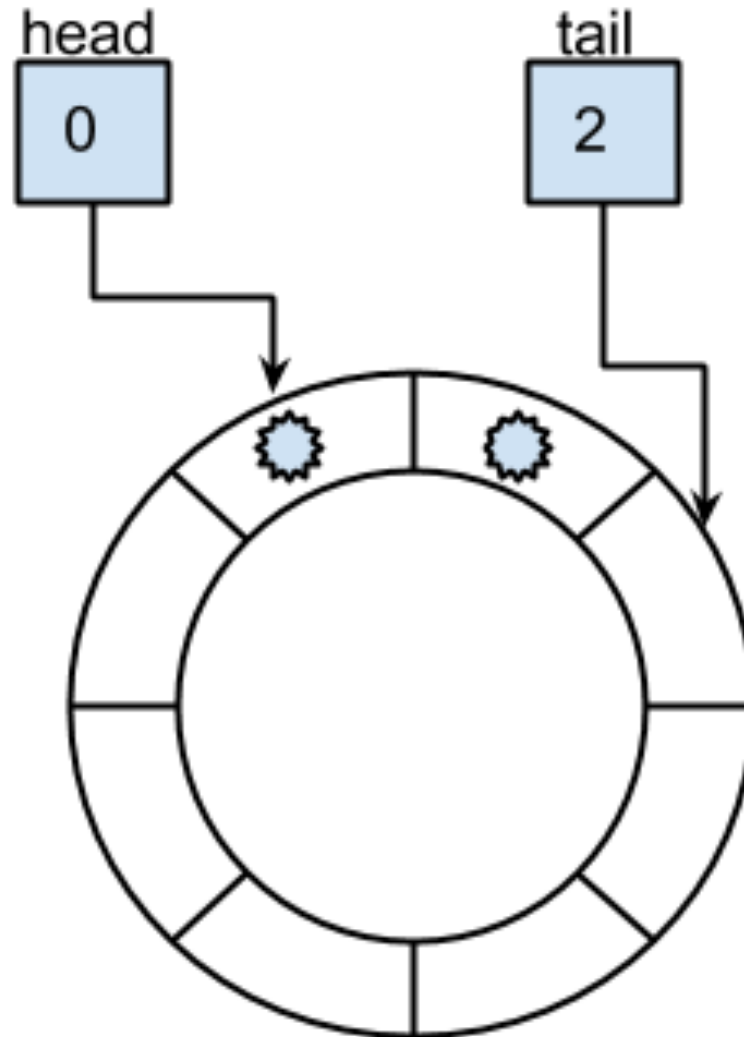
Ring Buffer: nP 1C



Ring Buffer: nP 1C



Ring Buffer: nP 1C



Ring Buffer: nP 1C

```
AtomicReferenceArray ringBuffer;  
volatile long head;  
AtomicLong tail;
```

Ring Buffer: nP 1C

```

public boolean offer(E e) {
    long curTail;
    do {
        curTail = tail.get();
        if (curTail - head == ringBuffer.length())
            return false;
    } while (!tail.compareAndSet(curTail, curTail+1));
    int idx = curTail % ringBuffer.length();
    ringBuffer.set(idx, e); // volatile write
    return true;
}

```

Ring Buffer: nP 1C

```

public E poll() {
    int index = head % ringBuffer.length();
    E element = ringBuffer.get(index);
    if (element == null)
        return null;
    ringBuffer.set(index, null);
    head++; // volatile write
    return element;
}

```


Disruptor

- Very flexible for different usages (strategies)
- Very good performance
- Data transfer from one thread to another (Queue)

Spinning

spinning

- Active wait
- very good for consumer reactivity
- Burns a cpu permanently

spinning

- Some locks are implemented with spinning (spinLock)
- Synchronized blocks spin a little bit on contention
- use of the pause instruction (x86)

spinning

How to avoid burning a core ?

Backoff strategies:

- `Thread.yield()`
- `LockSupport.parkNanos(1)`
- `Object.wait()/Condition.await()/LockSupport.park()`

Ticketing: OrderedScheduler

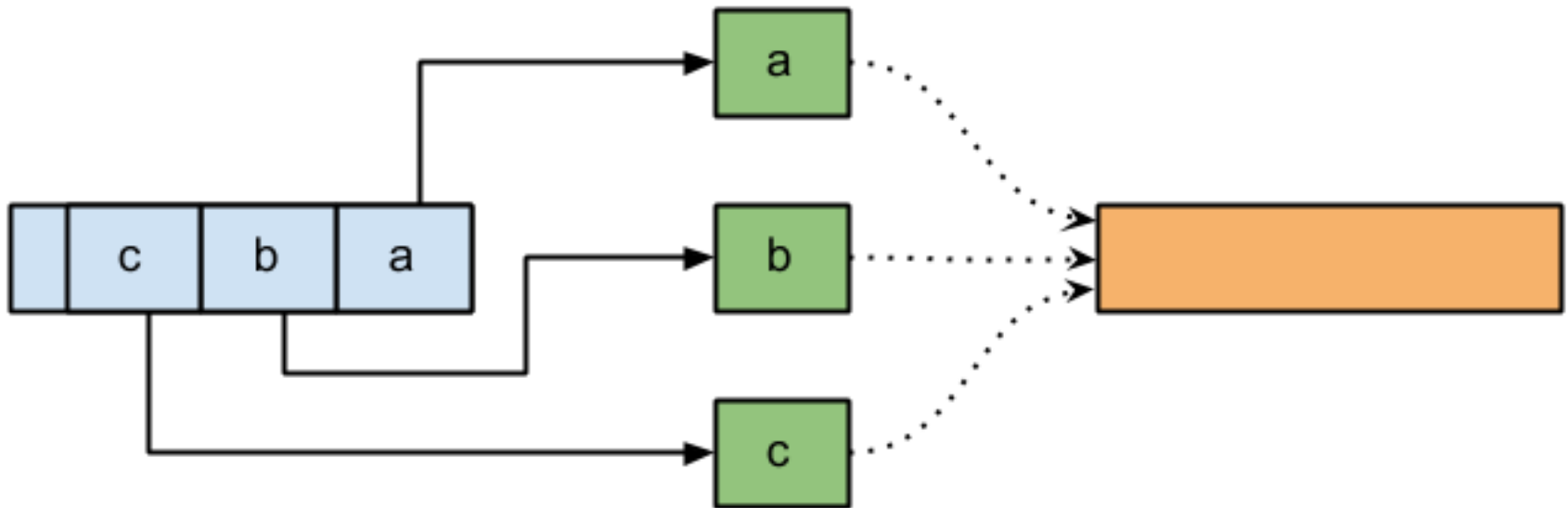
Ticketing

How to parallelize tasks while keeping ordering ?

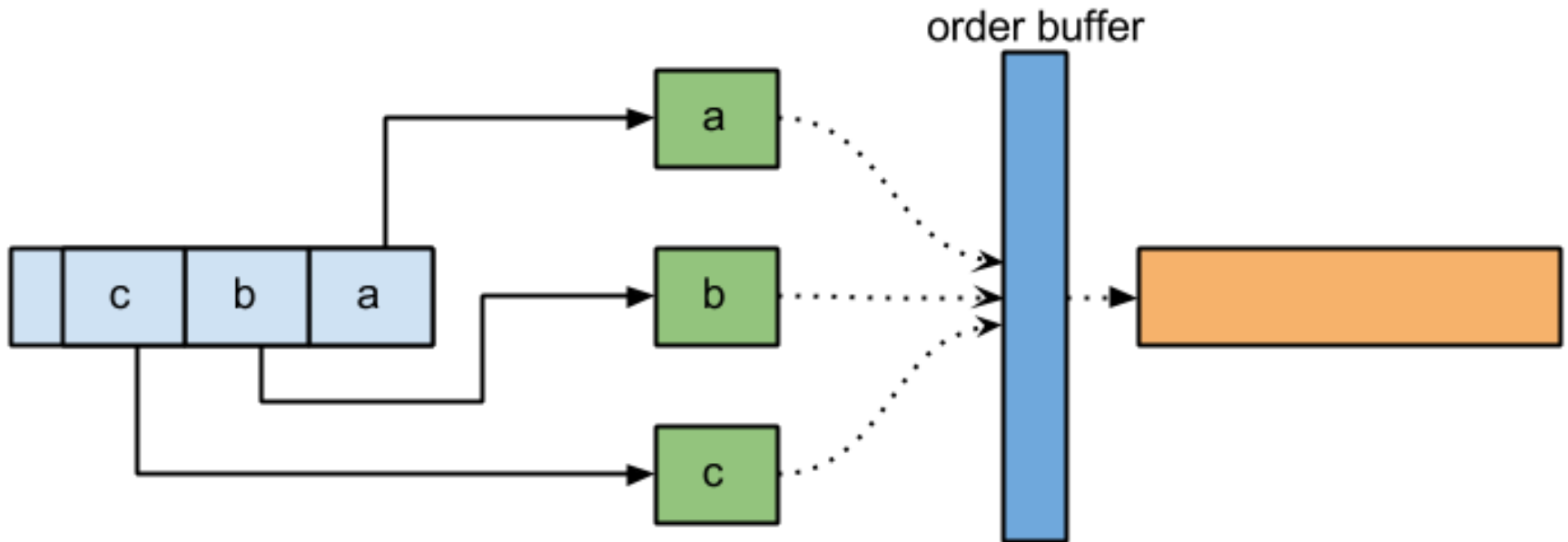
Example: Video stream processing

- read frame from the stream
- processing of the frame (parallelisable)
- writing into the output (in order)

Ticketing



Ticketing



Ticketing

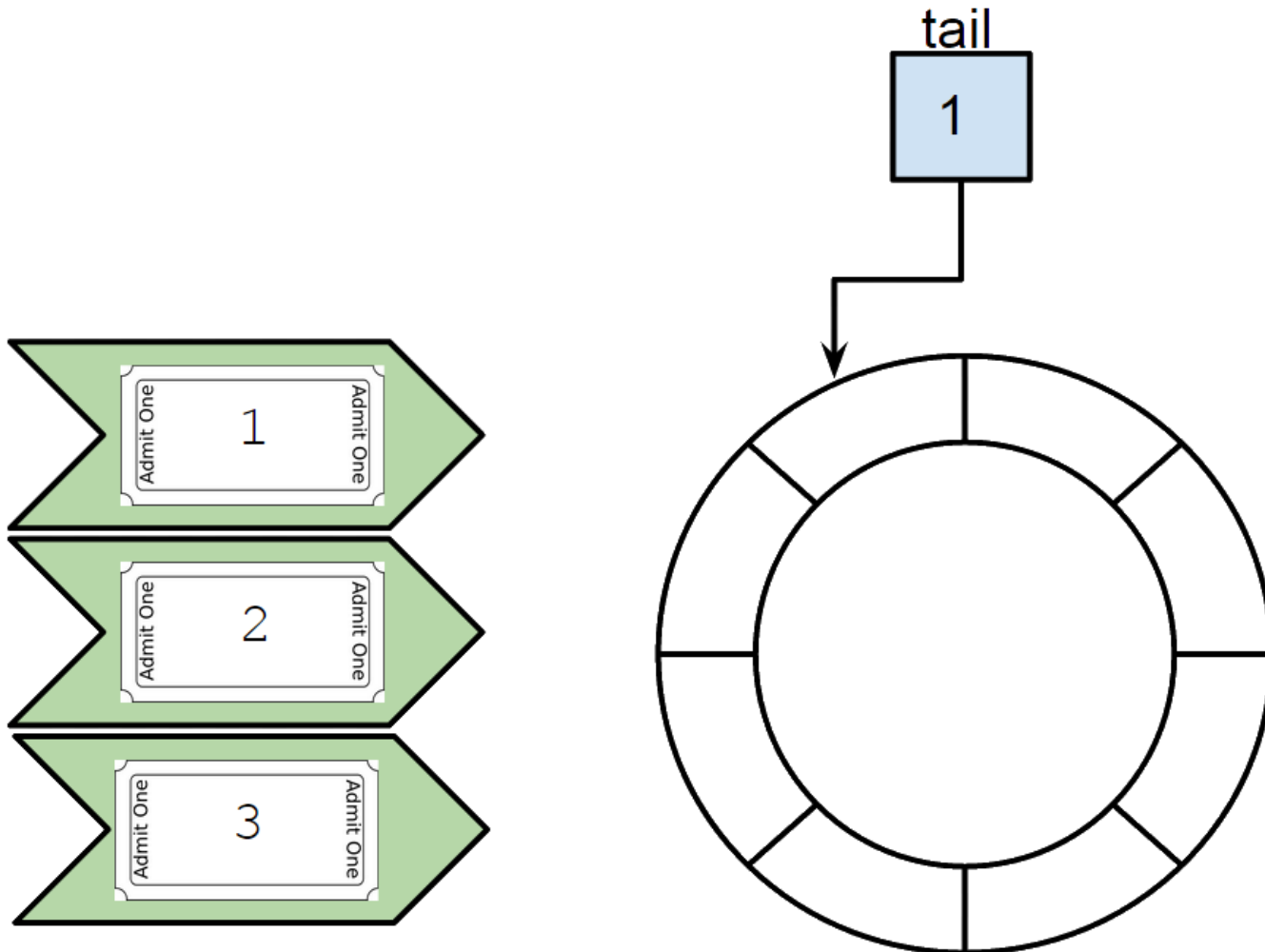
Can do this with Disruptor, but with a consumer thread

OrderedScheduler can do the same but:

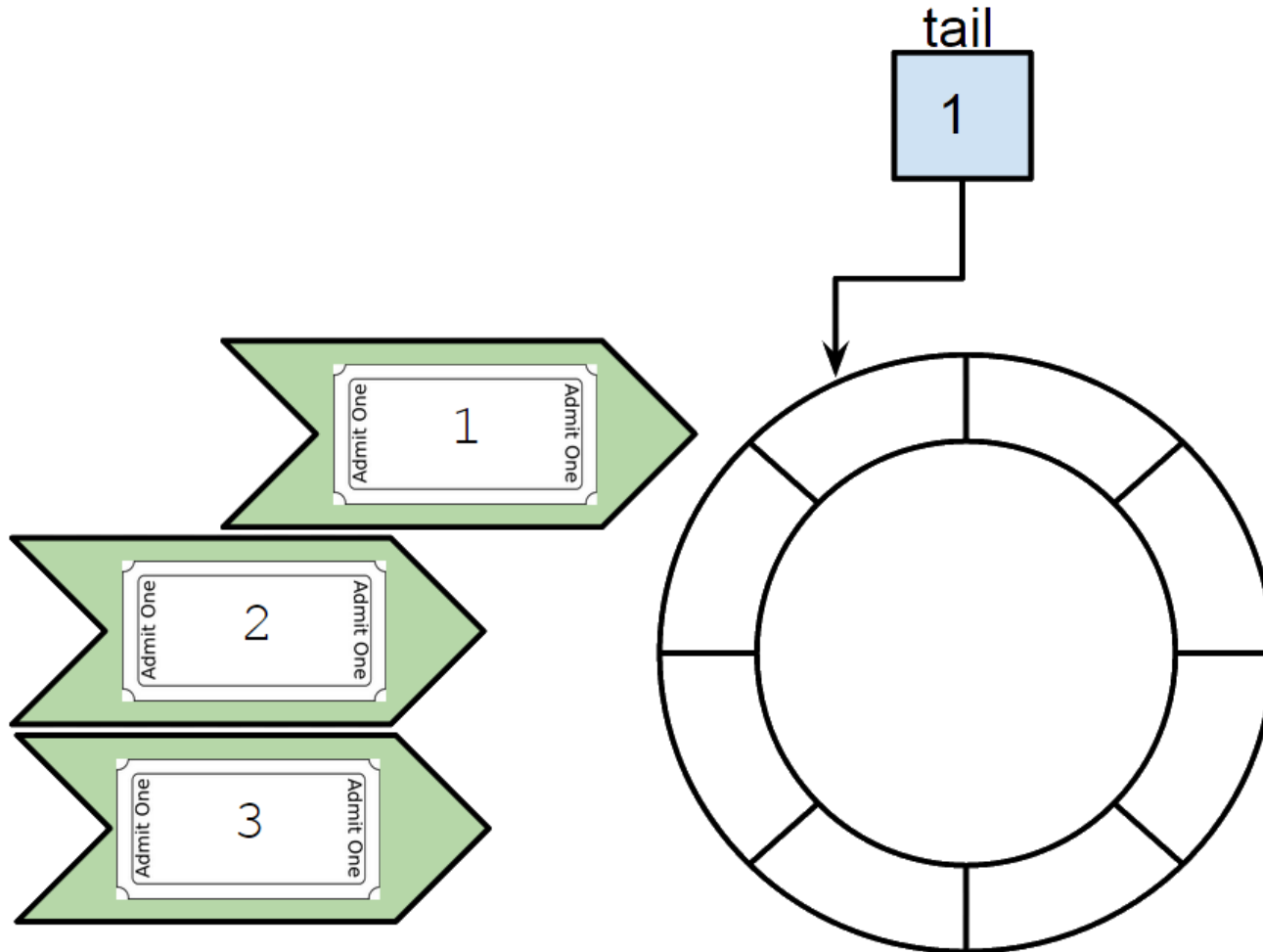
- no inter-thread communication overhead
- no additional thread
- no wait strategy

Take a ticket...

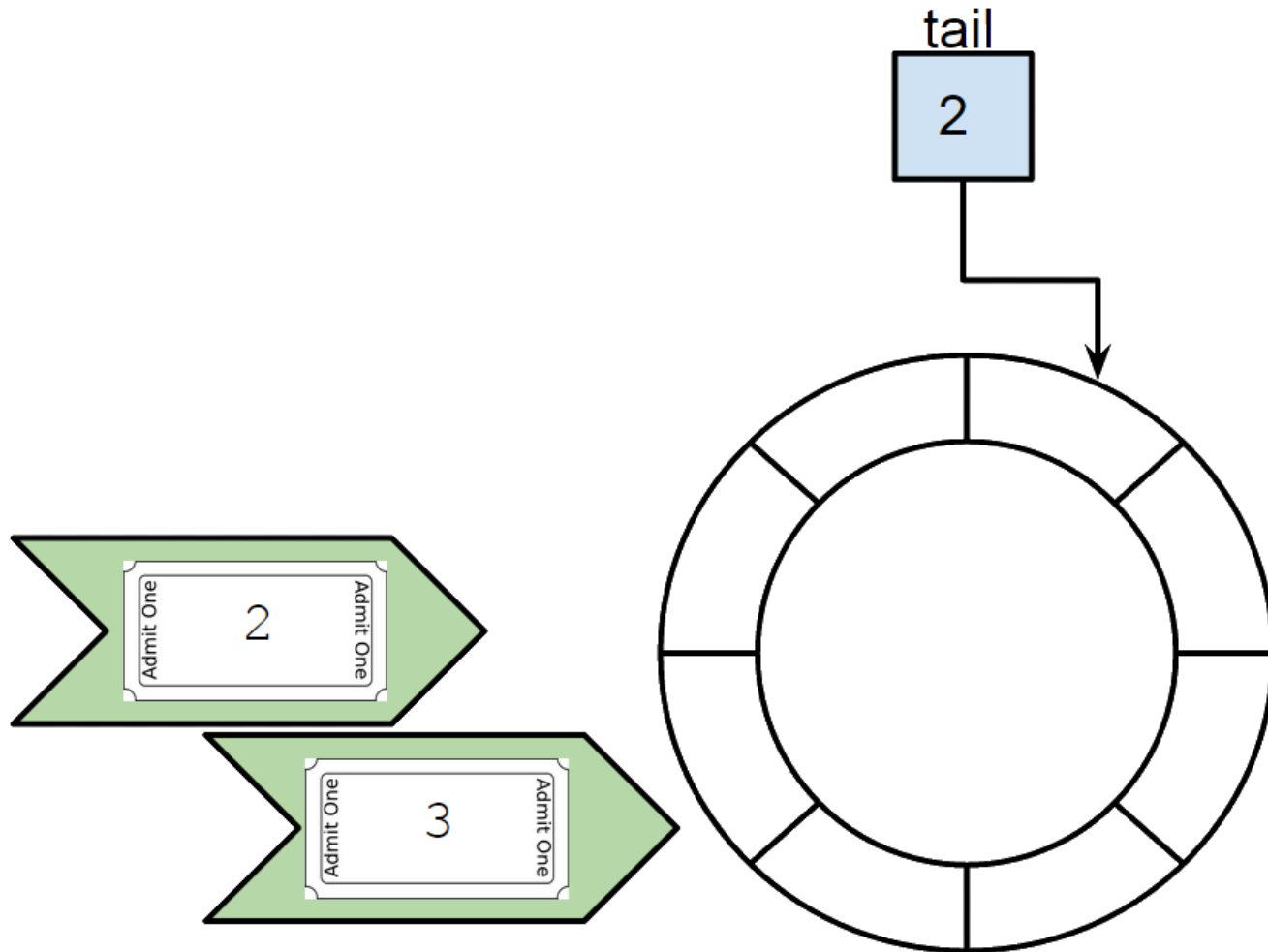
OrderedScheduler



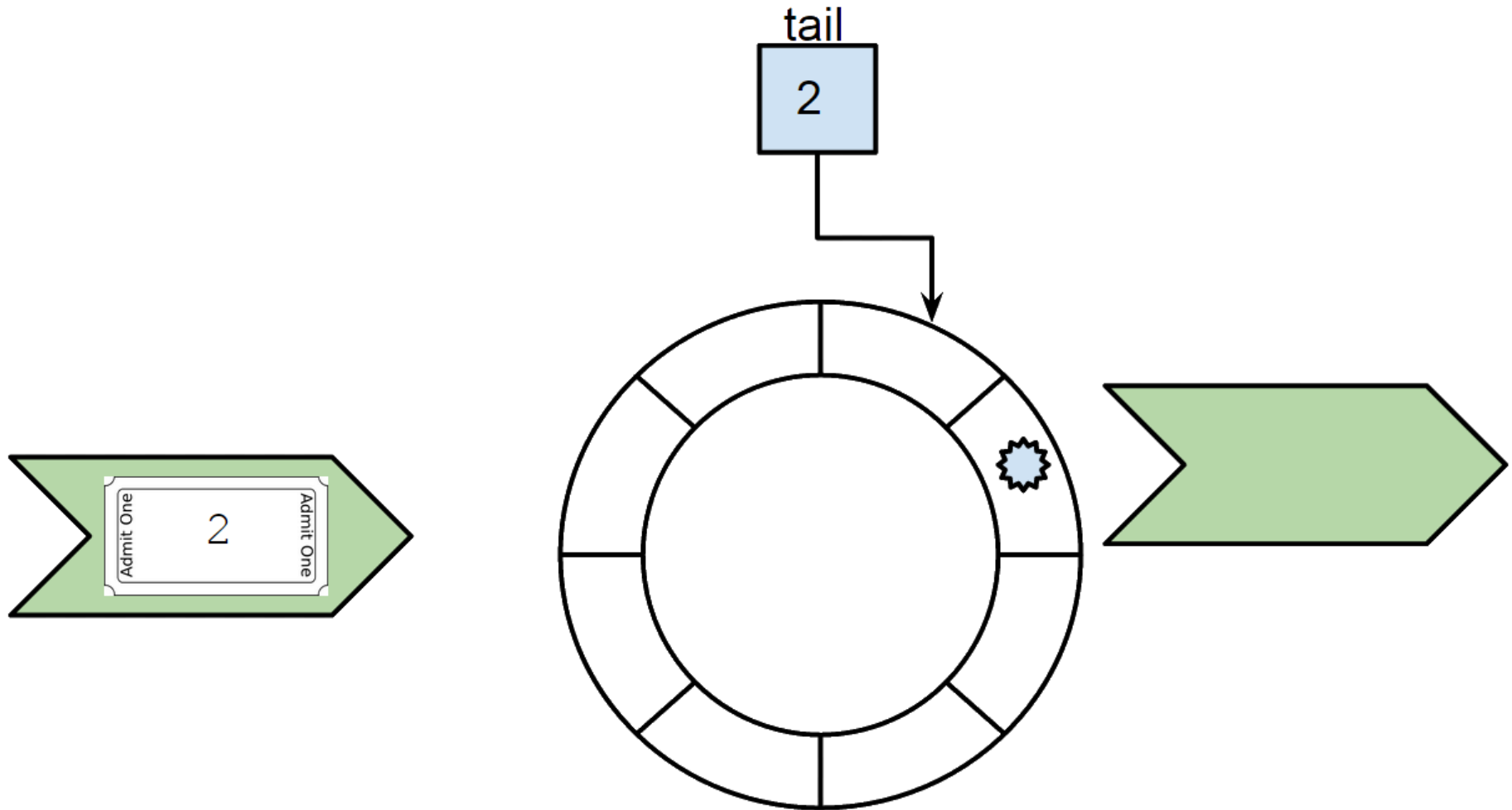
OrderedScheduler



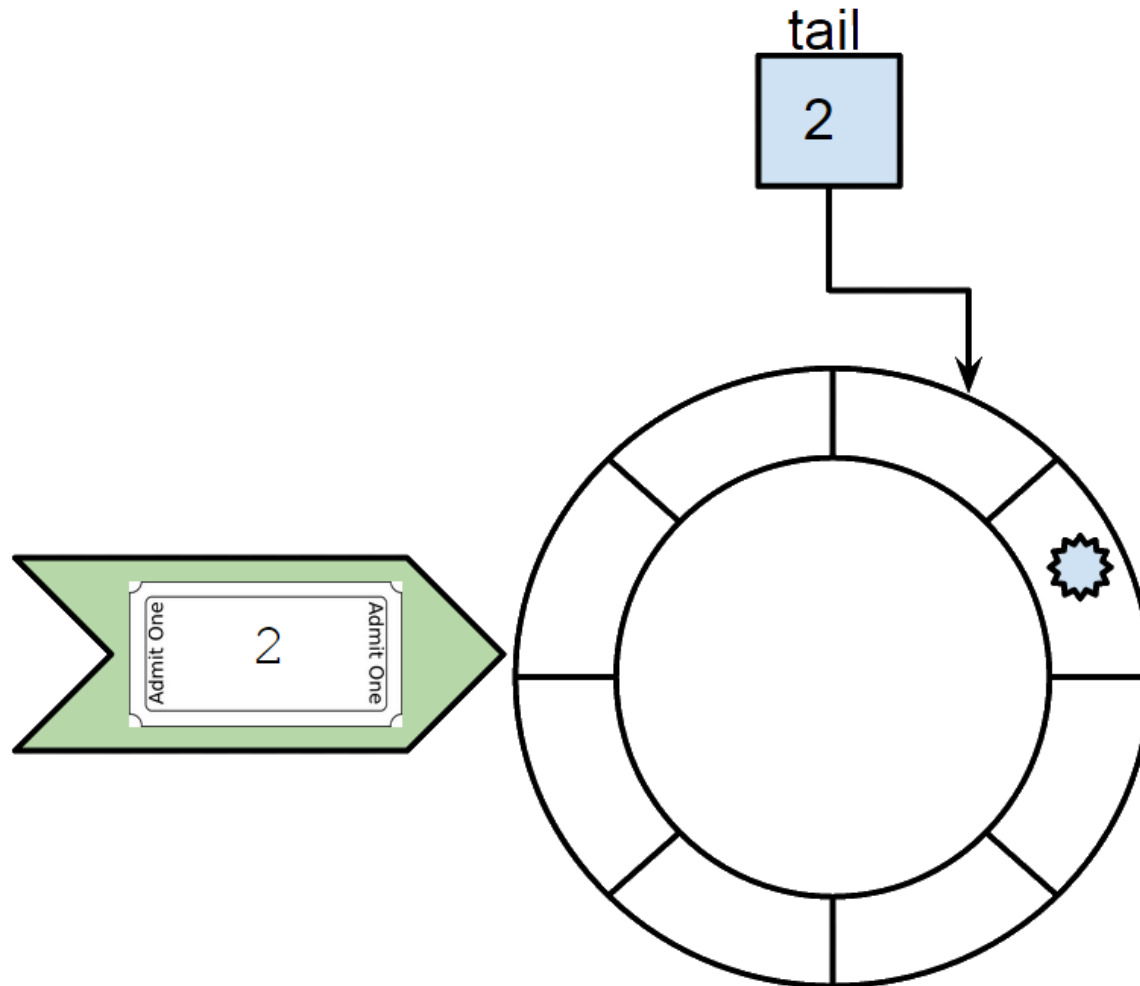
OrderedScheduler



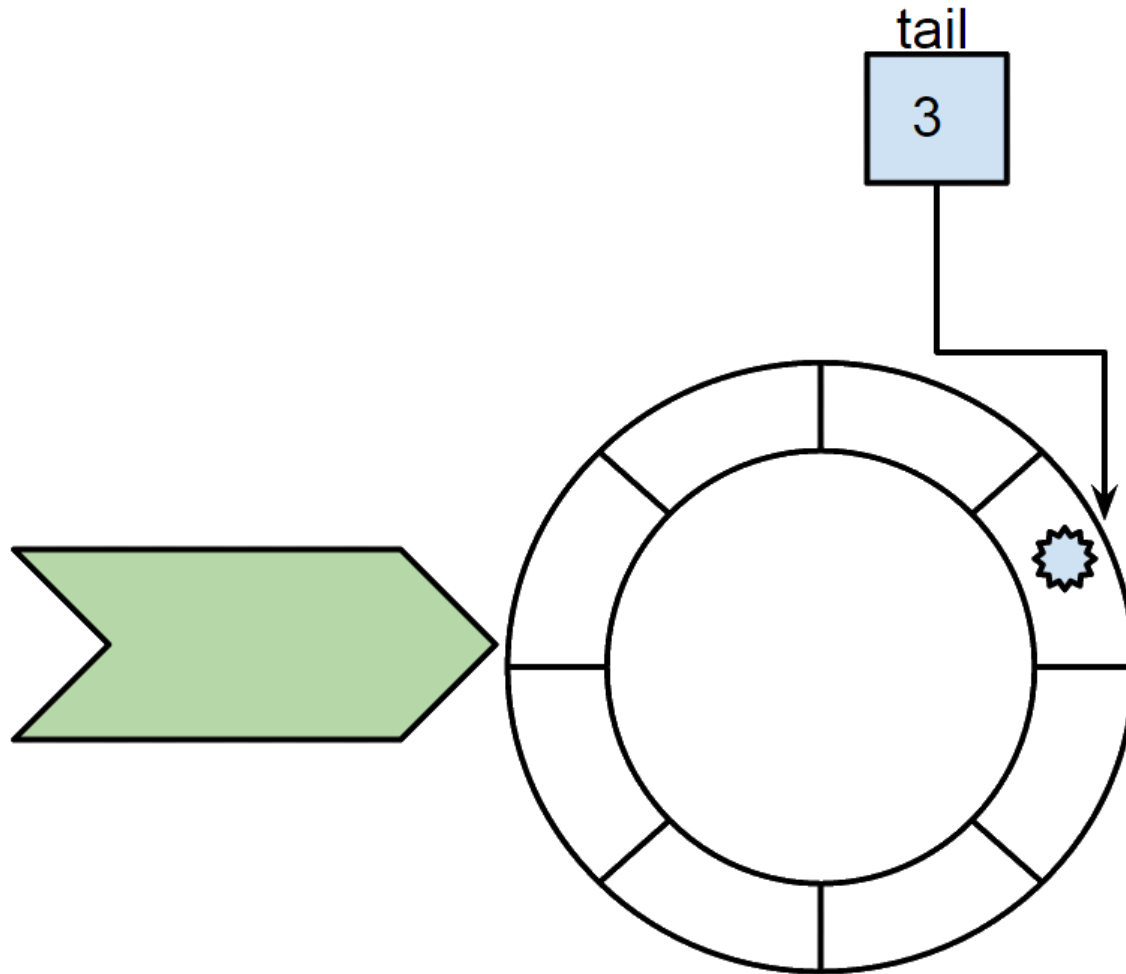
OrderedScheduler



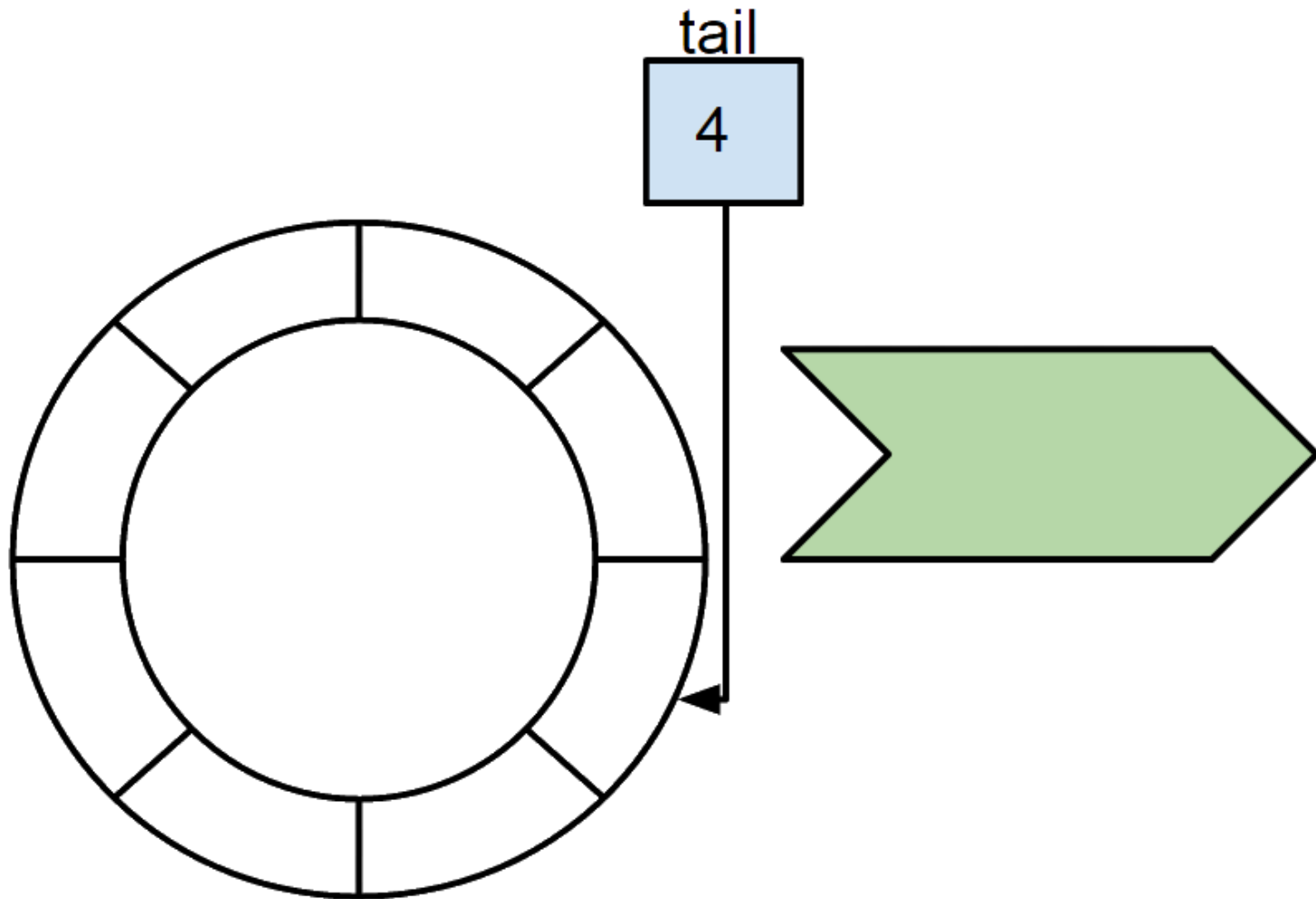
OrderedScheduler



OrderedScheduler



OrderedScheduler



OrderedScheduler

```
public void execute() {  
    synchronized (this) {  
        FooInput input = read();  
        BarOutput output = process(input);  
        write(output);  
    }  
}
```

OrderedScheduler

```
OrderedScheduler scheduler = new OrderedScheduler();
```

```
public void execute() {  
    FooInput input;  
    long ticket;  
    synchronized (this) {  
        input = read();  
        ticket = scheduler.getNextTicket();  
    }  
    [...]
```

OrderedScheduler

```
public void execute() {  
    [...]  
    BarOutput output;  
    try {  
        output = process(input);  
    }  
    catch (Exception ex) {  
        scheduler.trash(ticket);  
        throw new RuntimeException(ex);  
    }  
    scheduler.run(ticket, { () => write(output); });  
}
```

Ticketing

- Open sourced on GitHub
- Opened to PR & discussion on the design
- Used internally

Takeaways

Takeaways

- Measure
- Distribute
- Atomically update
- Order
- RingBuffer: easy lock-free
- Ordered without consumer thread

References

- jucProfiler: <http://www.infoq.com/articles/jucprofiler>
- Java Memory Model Pragmatics: <http://shipilev.net/blog/2014/jmm-pragmatics/>
- Memory Barriers and JVM Concurrency: http://www.infoq.com/articles/memory_barriers_jvm_concurrency
- JSR 133 (FAQ): <http://www.cs.umd.edu/~pugh/java/memoryModel/jsr-133-faq.html>
- CPU cache flushing fallacy: <http://mechanical-sympathy.blogspot.fr/2013/02/cpu-cache-flushing-fallacy.html>
- atomic<> Weapons: <http://channel9.msdn.com/Shows/Going+Deep/Cpp-and-Beyond-2012-Herb-Sutter-atomic-Weapons-1-of-2>

References

- circular buffers: <http://lwn.net/Articles/378262/>
- Mr T queues: https://github.com/mjpt777/examples/tree/master/src/java/uk/co/real_logic/queues
- Lock-free algorithms by Mr T: <http://www.infoq.com/presentations/Lock-free-Algorithms>
- Futex are tricky U. Drepper: <http://www.akkadia.org/drepper/futex.pdf>
- JCTools: <https://github.com/JCTools/JCTools>
- Nitsan Wakart blog: <http://psy-lob-saw.blogspot.com>

References

- Exploiting data parallelism in ordered data streams: <https://software.intel.com/en-us/articles/exploiting-data-parallelism-in-ordered-data-streams>
- OrderedScheduler: <https://github.com/Ullink/ordered-scheduler>
- Concurrency Freaks: <http://concurrencyfreaks.blogspot.com>
- Preshing on programming: <http://preshing.com/>
- Is Parallel Programming Hard, And If So, What Can You Do About It? <https://kernel.org/pub/linux/kernel/people/paulmck/perfbook/perfbook.2015.01.31a.pdf>

ABOUT ULLINK

ULLINK's electronic financial trading software solutions and services address the challenges of new regulations and increased fragmentation. Emerging markets require fast, flexible and compliant solutions for buy-side and sell-side market participants; providing a competitive advantage for both low touch and high touch trading.

www.ullink.com

FIND OUT MORE

Contact our Sales Team to learn more about the services listed herein as well as our full array of offerings:



+81 3 3664 4160 (Japan)

+852 2521 5400 (Asia Pacific)

+1 212 991 0816 (North America)

+55 11 3171 2409 (Latin America)

+44 20 7488 1655 (UK and EMEA)



connect@ullink.com

