

Better Mutual Exclusion on the Filesystem using Boost.AFIO

(asynchronous filesystem and file i/o)

Niall Douglas

Contents:

1. Quick overview of proposed Boost.AFIO, its motivation and its current status: v1 and v2
2. Overview of API design of Abstract Base Class
`afio::algorithm::shared_fs_mutex::shared_fs_mutex`
3. High-level overview of our memory map lock algorithm
4. Detail of the memory map lock algorithm implementation
5. Benchmarks comparing our memory map lock algorithm implementation to **lock files**, **byte ranges** and **atomic append** locks on various filing systems

History and status of proposed Boost.AFIO (2012-)

Motivation

1. Why is it still so hard in 2016 to **atomically** update more than file at once?
 - **Over usage of badly fitting SQLite**
2. Why is it still so hard in 2016 to write DMA-friendly **zero-copy** storage code?
 - **Ever more important with NVMe 4M IOPS SSDs ...**
3. Why is it still so hard in 2016 to write **bug-free** filesystem code?
 - **The infamous directory tree deletion anti-pattern on Windows**

Proposed Boost.AFIO v2 (> Oct 2015)

- “Bare metal” design
 - Exposes all the quirks of the host OS to user **unfiltered**
 - **No threads**, **no resource nor memory management**, **no exceptions** (we return lightweight **mongrel-monadic Outcomes** instead)
 - Performance **never measurably worse** than using host OS APIs directly - overhead often just **hundreds** of assembler opcodes
- Tight mapping between C++ type system and filing system object primitives
- Publicly exposes core file system template algorithms library into **afio::algorithm::*** (the AFIO “**FTL**”)
- Exemplar of simple, light, clean C++ 14/17 design

Proposed Boost.AFIO v2 - changes

Since ACCU 2016 (April):

- No AFIO v2 feature work apart from today's workshop
- BUT: Beginnings of a test suite based on **CTest**, **Boost.KernelTest** and threadsafe **CATCH**
 - KernelTest lets you write unit test suites *functionally* instead of imperatively
 - Parameter permutation lists, **preconditions**, **postconditions**
 - Maximum use of CPU cores, **bias towards CI driven soak testing**
 - In Progress: Automatic CI ASan, MSan, TSan, UBSan, valgrind, libfuzzer
 - Soon: Automatic **edge coverage** and **code bloat** calculation
 - Despite immaturity, lots of AFIO v2 bugs fixed already

Proposed Boost.AFIO v2 - changes

- Huge improvements to **Boost-lite**
 - AFIO now has a **cmake3** based build system yay!
 - As does **Outcome**, **KernelTest** and **Boost-lite** itself
 - Via Boost-lite build all Boost-lite libraries have:
 - Automatic CMake, CTest, CDash and (soon) CPack
 - Automatic binary and source uploads to **CDash**
 - Automatic **Travis** and **Appveyor** CI
 - Automatic use of **C++ Modularisation** & precompiled headers
 - Automatic submodularisation & **dependency tracking**
 - Automatic flat hierarchy or “**single download and drop in**” header-only installation
 - Automatic **ABI versioning** and hard version binds
 - Boost-lite remains very immature though!



Boost.AFIO

Dashboard Calendar Previous Current Project

No file changed as of **Wednesday, September 07 2016 - 01:00 EDT**



3 hours ago: 10 tests not run on Linux-3.13.0-40-generic-x86_64

3 hours ago: 12 errors introduced on Linux-3.13.0-40-generic-x86_64

17 hours ago: 10 tests not run on Linux-3.13.0-40-generic-x86_64

17 hours ago: 12 errors introduced on Linux-3.13.0-40-generic-x86_64

1 days ago: 1 test failed on Windows-AMD64

[See full feed](#)

Experimental

Site	Build Name	Update	Configure		Build		Test			Build Time
		Files	Error	Warn	Error	Warn	Not Run	Fail	Pass	
TravisLinuxWorker	Linux-3.13.0-40-generic-x86_64		0	0	12	0	10	0	0	2 hours ago
TravisDocumentation	Linux-3.13.0-40-generic-x86_64		0	0	0	0				2 hours ago
APPVYR-WIN	Windows-AMD64		0	6	0	0	0	0	30	3 hours ago



boost.afio-v2.0-binaries-win64-201608262229.zip	26-Aug-2016	22:29	5377132
boost.afio-v2.0-binaries-win64-201608270003.zip	27-Aug-2016	00:03	5377146
boost.afio-v2.0-binaries-win64-201608292307.zip	29-Aug-2016	23:07	5338070
boost.afio-v2.0-binaries-win64-201609030006.zip	03-Sep-2016	00:06	5338069
boost.afio-v2.0-binaries-win64-201609061654.zip	06-Sep-2016	16:54	5493508
boost.afio-v2.0-binaries-win64-201609070715.zip	07-Sep-2016	07:15	5494875
boost.afio-v2.0-source-201608131743.tar.xz	13-Aug-2016	17:43	7066248
boost.afio-v2.0-source-201608131810.tar.xz	13-Aug-2016	18:11	7066340
boost.afio-v2.0-source-201608131829.tar.xz	13-Aug-2016	18:30	7068964
boost.afio-v2.0-source-201608131935.tar.xz	13-Aug-2016	19:36	7069012
boost.afio-v2.0-source-201608141334.tar.xz	14-Aug-2016	13:35	7069032
boost.afio-v2.0-source-201608141813.tar.xz	14-Aug-2016	18:13	7081544
boost.afio-v2.0-source-201608171941.tar.xz	17-Aug-2016	19:41	7131376
boost.afio-v2.0-source-201608171943.tar.xz	17-Aug-2016	19:44	7131248
boost.afio-v2.0-source-201608171950.tar.xz	17-Aug-2016	19:50	7131292
boost.afio-v2.0-source-201608181945.tar.xz	18-Aug-2016	19:45	7130588
boost.afio-v2.0-source-201608192001.tar.xz	19-Aug-2016	20:02	7131556
boost.afio-v2.0-source-201608212029.tar.xz	21-Aug-2016	20:29	7153236
boost.afio-v2.0-source-201608222020.tar.xz	22-Aug-2016	20:20	7154696
boost.afio-v2.0-source-201608230917.tar.xz	23-Aug-2016	09:17	7154928
boost.afio-v2.0-source-201608232041.tar.xz	23-Aug-2016	20:41	7156384
boost.afio-v2.0-source-201608242239.tar.xz	24-Aug-2016	22:41	7159840
boost.afio-v2.0-source-201608262006.tar.xz	26-Aug-2016	20:06	7160264
boost.afio-v2.0-source-201608262054.tar.xz	26-Aug-2016	20:55	7160752
boost.afio-v2.0-source-201608262229.tar.xz	26-Aug-2016	22:29	7160276
boost.afio-v2.0-source-201609041922.tar.xz	04-Sep-2016	19:22	7171832
boost.afio-v2.0-source-201609041929.tar.xz	04-Sep-2016	19:30	7171916
boost.afio-v2.0-source-201609050824.tar.xz	05-Sep-2016	08:24	7171984
boost.afio-v2.0-source-201609060928.tar.xz	06-Sep-2016	09:28	7202656
boost.afio-v2.0-source-201609061652.tar.xz	06-Sep-2016	16:52	7203640
boost.afio-v2.0-source-201609070713.tar.xz	07-Sep-2016	07:13	7206252
boost.outcome-v1.0-source-201608131653.tar.xz	13-Aug-2016	16:54	614640
boost.outcome-v1.0-source-201608131738.tar.xz	13-Aug-2016	17:39	614600
boost.outcome-v1.0-source-201608131807.tar.xz	13-Aug-2016	18:08	614808
boost.outcome-v1.0-source-201608131827.tar.xz	13-Aug-2016	18:28	614820
boost.outcome-v1.0-source-201608141807.tar.xz	14-Aug-2016	18:08	615232
boost.outcome-v1.0-source-201608171939.tar.xz	17-Aug-2016	19:40	617892
boost.outcome-v1.0-source-201608181942.tar.xz	18-Aug-2016	19:43	617920
boost.outcome-v1.0-source-201608191957.tar.xz	19-Aug-2016	19:58	617876
boost.outcome-v1.0-source-201608212027.tar.xz	21-Aug-2016	20:28	623960

Top level AFIO v2

afio::algorithm::

shared_fs_mutex::*

API overview

Prereq: Boost.Outcome (mongel monads)

- Policy driven underlying implementation
`basic_monad<T, EC, E>` is aliased into convenience typedefs:
 - `outcome<T>` = empty OR `T` OR `std::error_code` OR `std::exception_ptr`
 - `result<T>` = empty OR `T` OR `std::error_code`
 - `option<T>` = empty OR `T`
- Has identical API to a `std::future<T>` or `std::optional<T>` (e.g. `get()`, `value()`, `has_value()` etc)

Prereq: AFIO v2 byte range lock API

```
class io_handle : public handle {  
    using extent_type = unsigned long long;  
    class extent_guard; // RAI guard  
    virtual result<extent_guard> lock(extent_type offset,  
    extent_type bytes, bool exclusive = true, deadline d =  
    deadline()) noexcept;  
};
```

Passes straight through to syscall, insane semantics and all if your POSIX OS is mad



afio::algorithm::shared_fs_mutex

There are these implementations of `shared_fs_mutex` in AFIO v2:

1. `shared_fs_mutex::lock_files`
2. `shared_fs_mutex::byte_ranges`
3. `shared_fs_mutex::atomic_append` (ACCU talk)
4. `shared_fs_mutex::memory_map` (today)

All implement the abstract base class

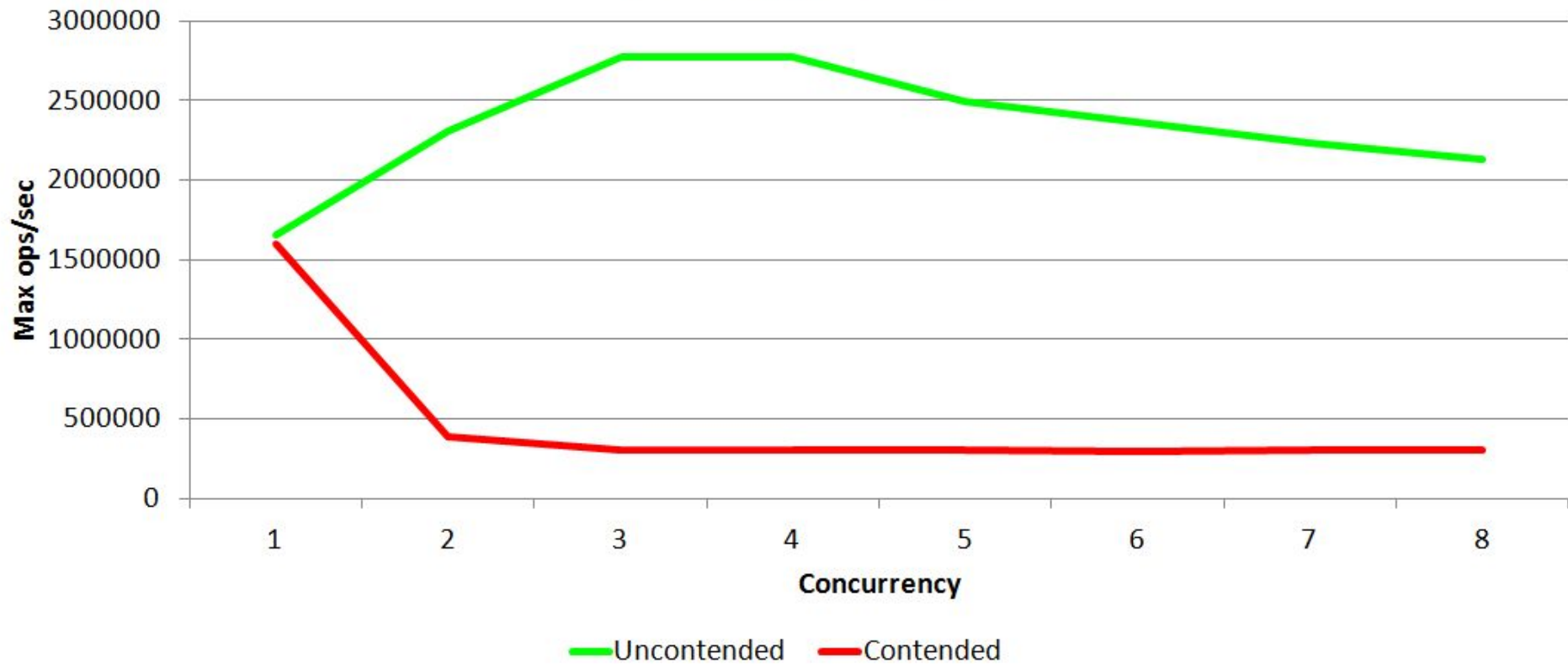
`shared_fs_mutex::shared_fs_mutex`

shared_fs_mutex abstract base class

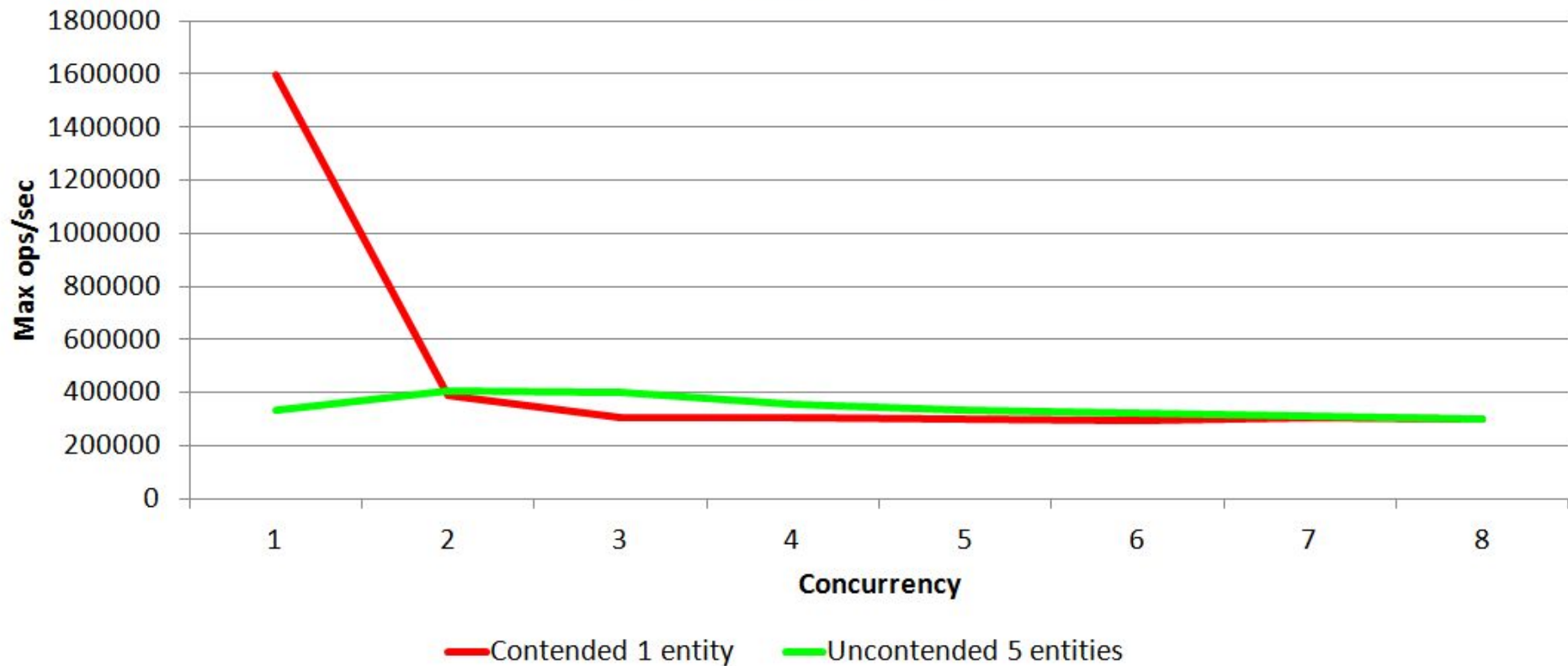
```
class shared_fs_mutex {  
    struct entity_type { value : 63; exclusive: 1; };  
    using entities_type = gsl::span<entity_type>;  
    class entities_guard; // RAII guard  
    (virtual) result<entities_guard> lock(entities_type  
entities, deadline d = deadline(), bool spin_not_sleep =  
false) noexcept;  
    (virtual) result<entities_guard> try_lock(entities_type  
entities) noexcept { return lock(std::move(entities),  
deadline(std::chrono::seconds(0))); }  
};
```

**Why does the shared FS
mutex API lock many
entities?**

Scaling of byte range locks to concurrency on NTFS (4 core 8 thread CPU)



How many uncontended entities equals single contended byte range lock (4 core 8 thread CPU)



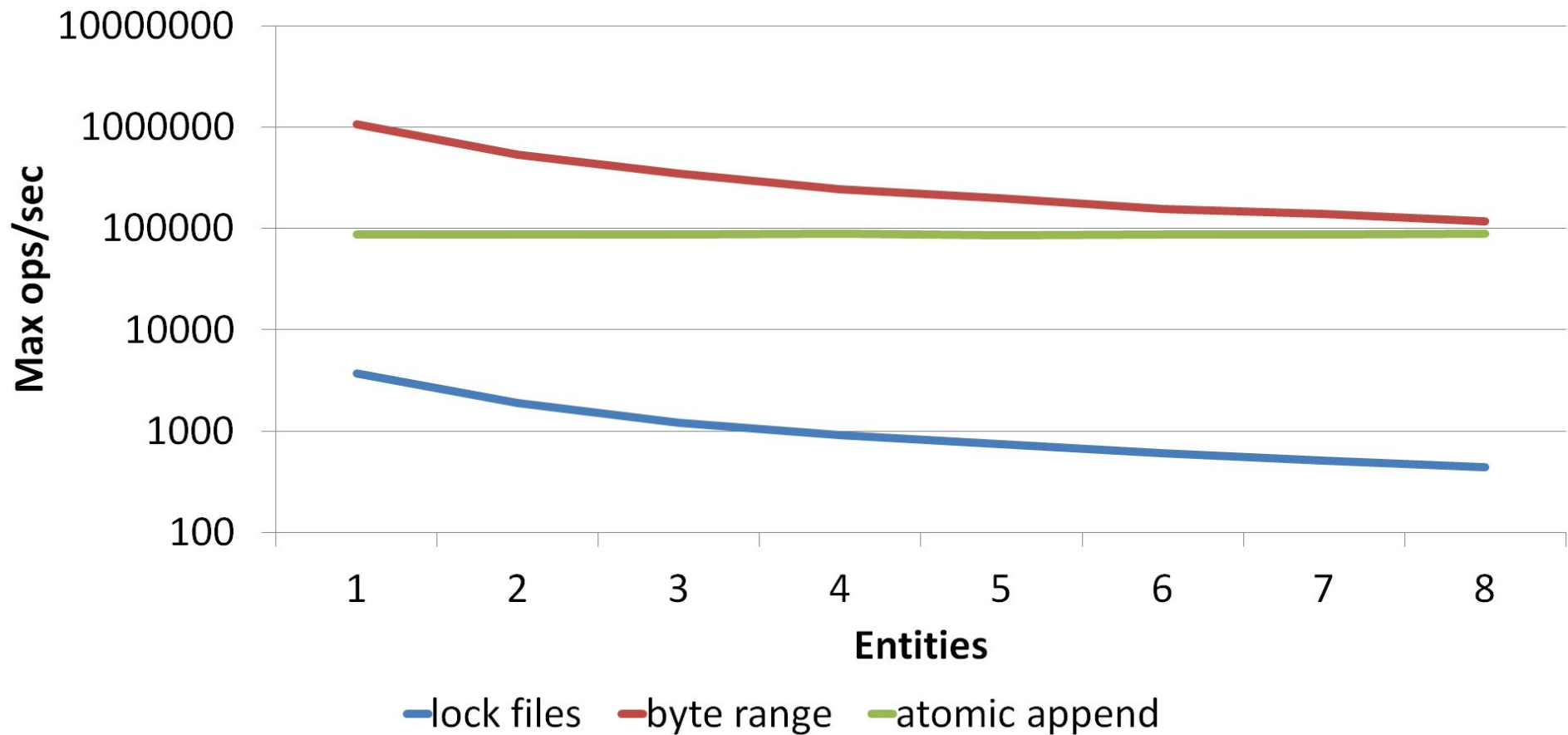
afio::algorithm::shared_fs_mutex

- How entities are mapped onto the file system depends on the implementation, so:
 - `shared_fs_mutex::lock_file` maps entities into 16 character hexadecimal files opened with `O_EXCL`
 - `shared_fs_mutex::byte_ranges` maps entities into single byte offsets into the shared lock file
 - Contention is handled by backing off all preceding locks and randomising the list before trying again, starting with the entity which was contended last try. This allows sleeping until contended lock becomes free

afio::algorithm::shared_fs_mutex

- `shared_fs_mutex::atomic_append` (see my ACCU 2016 talk) solves the problem of *inverse log scaling* to entity count in `lock_files` and `byte_ranges`
 - `atomic_append` pretty much always wins once you are locking ≥ 15 entities uncontended
 - Has easily best worst case performance in most heavy concurrency use cases, avoiding the “scalability hole” which most kernel filesystem lock implementations fall into after a certain load of entities and concurrency

Scaling of four concurrent users of the lock algorithms to entities on NTFS



afio::algorithm::shared_fs_mutex

Last remaining class of shared FS mutex is a very high performance one based entirely in userspace using **shared memory** atomics

- Comes with *lots* of caveats ... more on that shortly ...
- But it is probably likely most folk will default to it due to its sheer raw performance which is 10x-20x more than any other shared FS mutex in AFIO

(I have noticed people tend to default to **best average performance** rather than **best worst case performance**)

Further reading:

AFIO v2 API reference:

<https://ned14.github.io/boost.afio>

Abstract base class

`algorithm::shared_fs_mutex` API reference:

<http://goo.gl/s23ecD>

ACCU 2016 atomic append talk video:

<https://www.youtube.com/watch?v=elegewDwm64>

The memory map lock algorithm - Overview

Questions before we begin this section?

The shared memory map algorithm

- Only viable if there are no networked drive users
 - Therefore will also need a fallback lock which IS compatible with networked drive users
 - And the ability for all users of the lock to jointly degrade to that fallback lock in safe fashion if a networked drive user turns up
 - Which in turns means we need some way of portably detecting the arrival of a networked drive user

The memory_map algorithm

To construct the shared mutex instance:

1. If the lockfile doesn't exist, race free create a mapfile in `/tmp` and write its path into the lockfile. Otherwise open the `lockfile` and read the path of the `mapfile` somewhere in `/tmp`. Try opening that, if we can't then **degrade the lock**
2. Map the lockfile into memory read-only and the `/tmp` `mapfile` into memory read-write

(Big assumption here is that whatever is said to be `/tmp` by the OS is visible to all local processes, but not to any networked drive users)

The memory_map algorithm

To lock:

1. Has the shared memory map been **degraded** by a networked drive user? If so, pass through lock request to **fallback lock**
2. Hash each entity and modulus by total map entries
3. Try locking the **spin lock** at each entity's hash index. If fail to lock, unlock everything, randomise and retry

To unlock:

1. Unlock all the spin locks locked earlier (or **fallback lock**)

The memory_map algorithm

Characteristics:

- Does at least one atomic operation per entity
 - But cascades atomic ops if any contended, and SMP and especially NUMA have a very finite total system atomic operation bandwidth
- No ability to sleep the CPU
 - Spin locks do have a counter and do call sleep(1ms) eventually
- Quality of hash function very important
 - Collision probing left and right substantially increases cache line bouncing, whole system performance damage is very severe
- This algorithm is fundamentally anti-social to all other code on a machine ... but it's fast, very fast

The memory map lock algorithm - Detail

Questions before we begin this section?

```
memory_map<...>
```

```
template <
```

```
    template <class> class Hasher = fnv1a_hash,
```

```
    size_t HashIndexSize = 4096,
```

```
    class SpinlockType = shared_spinlock<>
```

```
> class memory_map : public shared_fs_mutex
```

- `fnv1a_hash` and `shared_spinlock<Policy>` come from Boost-lite <https://github.com/ned14/boost-lite>
- `1Mb + 0` = lock in use, `1Mb + 1` = map in use

memory_map<>::lock()

1. Is lock degraded (**lockfile**[0] is zero)? If so:
 - a. Have we only just degraded? If so:
 - i. Release shared lock map-in-use byte **1Mb + 1**
 - ii. Exclusive lock map-in-use byte **1Mb + 1** and immediately release
 - b. Redirect to **fallback lock**
2. Hash and modulus to **mapfile** entries every entity in the lock list in an auto-vectorising SIMD friendly fashion, eliding index collisions (with exclusive bit upgrade)

memory_map<>::lock()

3. For every entity index:

- a. If exclusively locking, try to exclusive lock the spinlock at that index, else try to shared lock it
- b. If lock fails:
 - i. Unlock everything locked so far
 - ii. Check if deadline has passed, if so return **ETIMEDOUT**
 - iii. Swap contended index with first index
 - iv. `std::random_shuffle()` rest of index list
 - v. Yield thread context and retry sequence

Questions?

memory_map<>::unlock()

1. If lock degraded, pass to **fallback lock**
2. Hash and modulus to **mapfile** entries every entity in the lock list in an auto-vectorising SIMD friendly fashion, eliding index collisions (with exclusive bit upgrade)
3. For every entity index:
 - a. If exclusively unlocking, try to exclusive unlock the spinlock at that index, else try to shared unlock it
4. Is lock newly degraded (**lockfile**[0] is zero)? If so, release shared lock map-in-use byte **1Mb +1**

Questions?

memory_map<>::fs_mutex_map()

Remember:

- $1\text{Mb} + 0$ = “lock in use” byte range lock offset
 - Used to detect other users of the lock
 - $1\text{Mb} + 1$ = “map in use” byte range lock offset
 - Used to detect other users of the map
1. Open the **lockfile**, creating if needed, **caching only reads (not writes)**
 2. Try to exclusive lock in-use byte at $1\text{Mb} + 0$

memory_map<>::fs_mutex_map()

3. If success (i.e. lock is NOT in use):
 - a. Truncate **lockfile** to zero
 - b. Create a randomly named **mapfile** in **/tmp**
 - c. Truncate it to **HashIndexSize**
 - d. Write path of **/tmp mapfile** into **lockfile**
 - e. Shared lock map-in-use byte at **1Mb + 1**
 - f. Convert exclusive lock in-use byte at **1Mb + 0** into a shared lock
 - g. Map **mapfile** (rw) and **lockfile** (ro) into memory

memory_map<>::fs_mutex_map()

3. If failure (i.e. lock IS in use):
 - a. Shared lock in-use byte at **1Mb + 0** (blocks)
 - b. Read path from **lockfile** and try to open
 - c. If not found, write zeroes all over **lockfile** if needed and return **EBUSY**
 - d. If found, shared lock map-in-use byte at **1Mb + 1**
 - e. Map **mapfile** (rw) and **lockfile** (ro) into memory

Questions?

memory_map<>::~~memory_map()

1. Unlock in-use byte at 1Mb + 0 and map-in-use byte at 1Mb + 1
2. Exclusive lock in-use byte at 1Mb + 0. If success:
 - a. Write 4Kb of zeroes to lockfile
 - b. Truncate lockfile to zero
 - c. Unlink /tmp mapfile
 - d. Release all locks, lock is reset

Questions?

Benchmarks!

The following benchmarks are for:

2 core 4 thread 4.125Gb/sec main memory bandwidth

GenuineIntel Intel(R) Core(TM) i5 CPU M 540 @

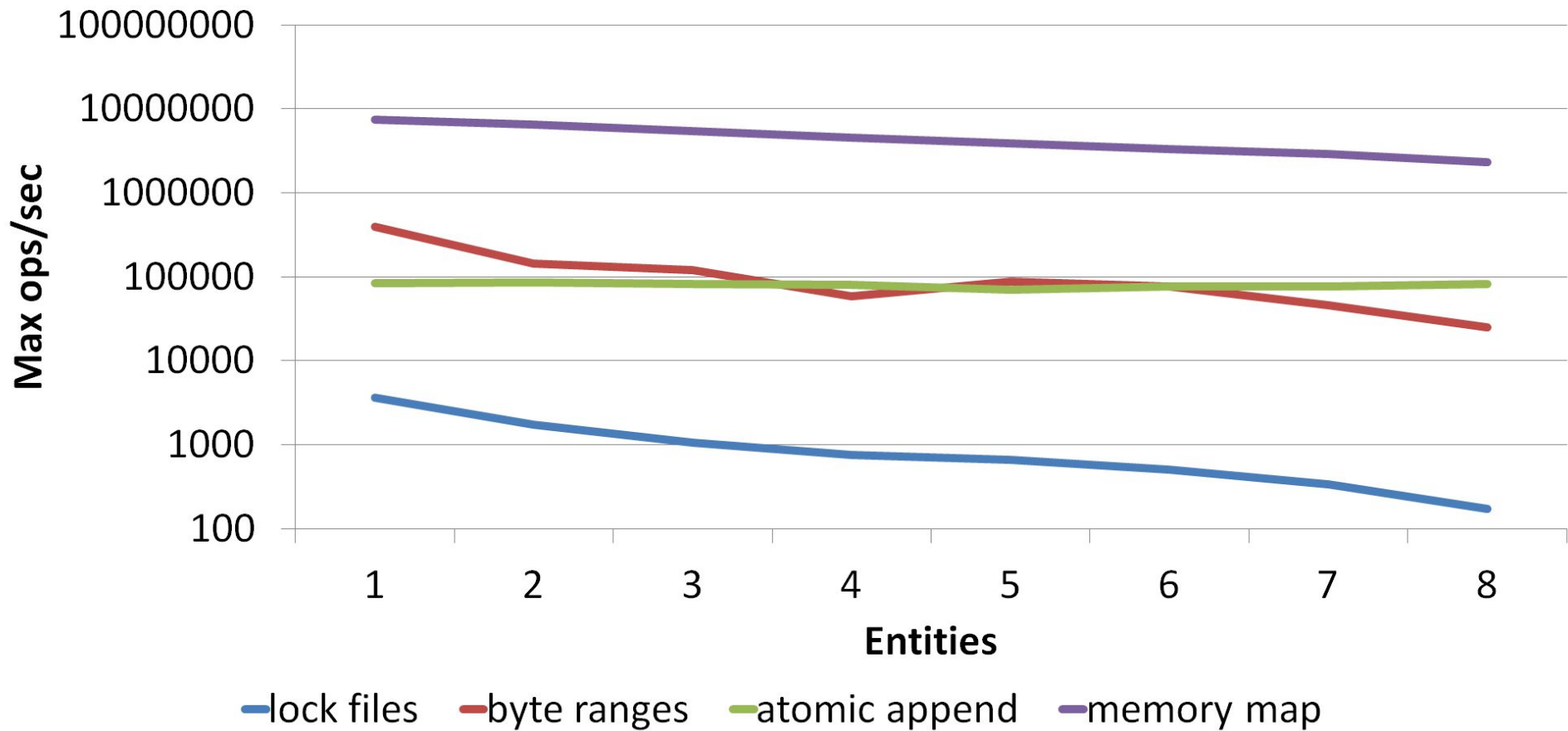
2.53GHz (2008 era laptop)

Microsoft Windows 10.0.14393.103

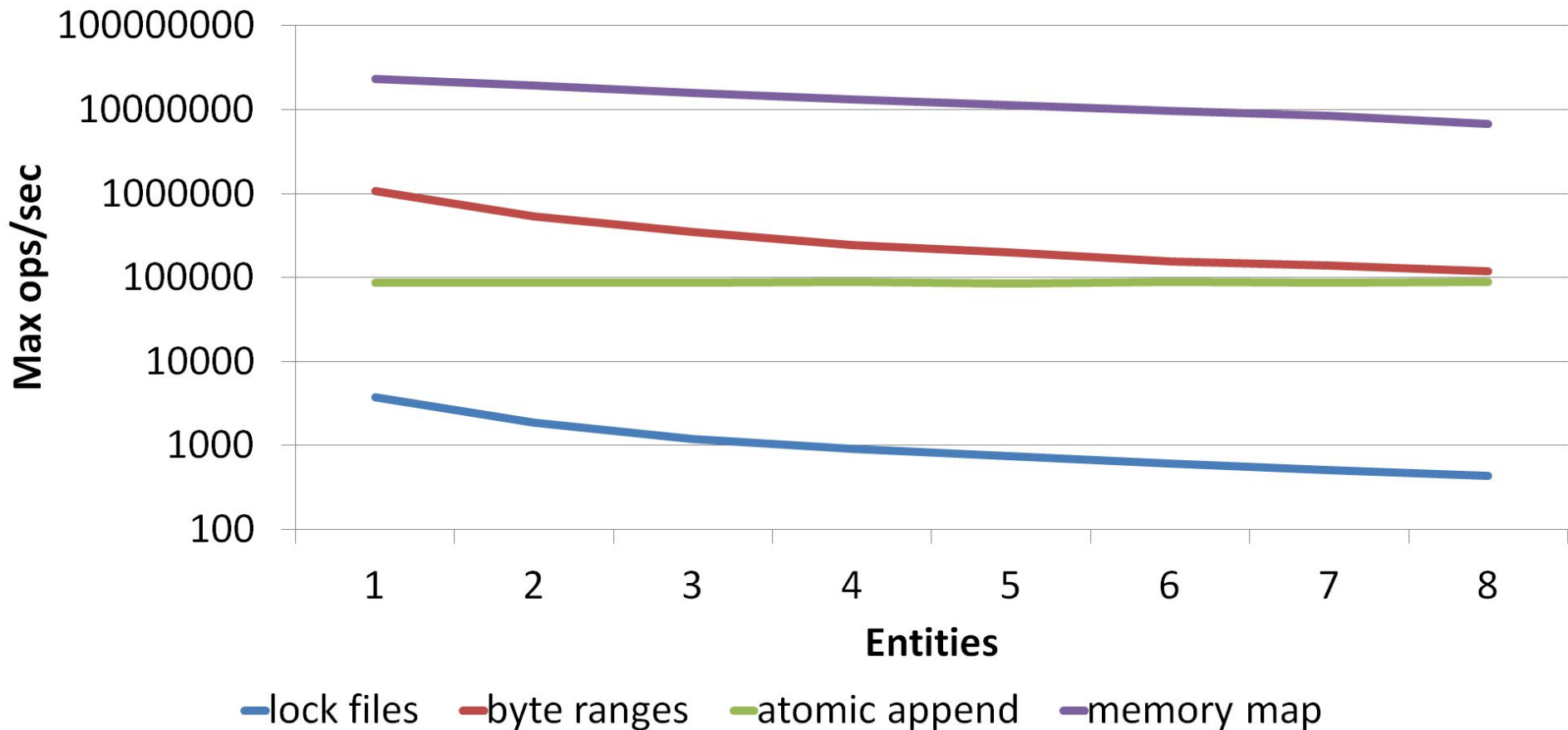
Known deficiencies in this memory_map implementation:

1. ~~Uses `std::random_shuffle` (deprecated)~~
2. Lock degrade currently racy under multiple thread users of same object instance
3. Networked user lock degrade blocks all lock user until all users realise degrade is happening
4. Multiple object instances racy on POSIX systems with insane byte locks

Scaling of two contended users of the lock algorithms to entities on NTFS



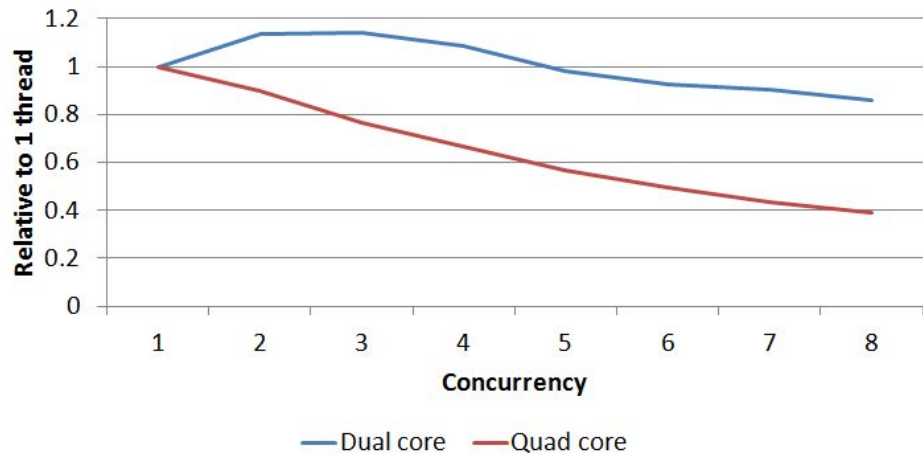
Scaling of four concurrent users of the lock algorithms to entities on NTFS



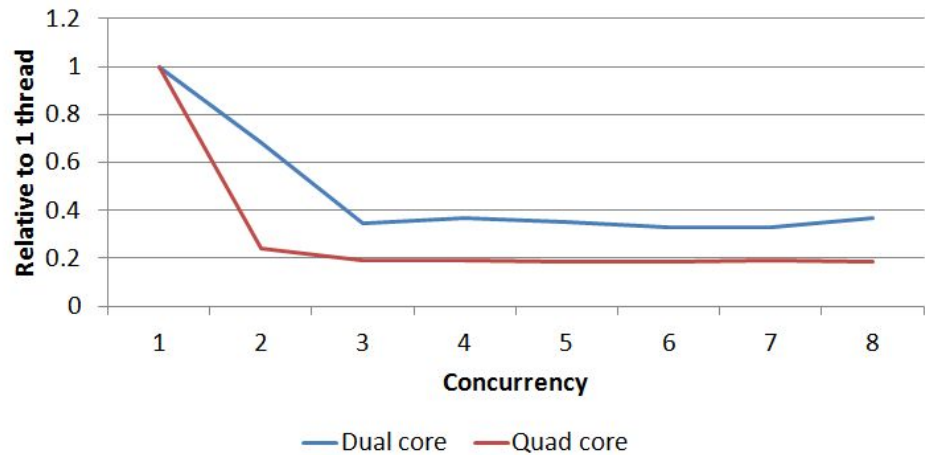
Concurrency scaling

Questions before we begin this section?

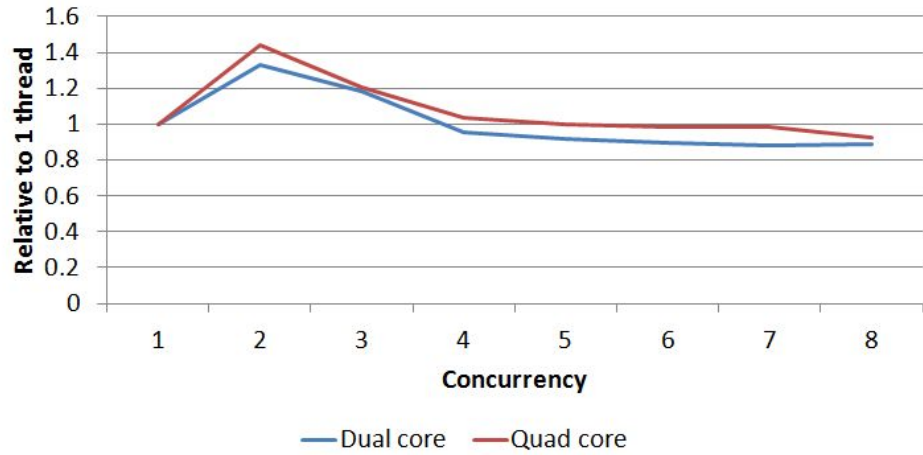
Contended lock files



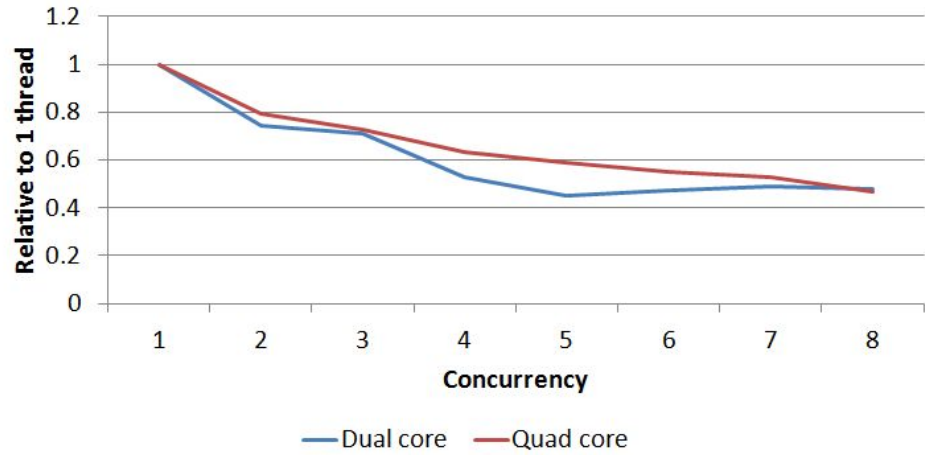
Contended byte ranges



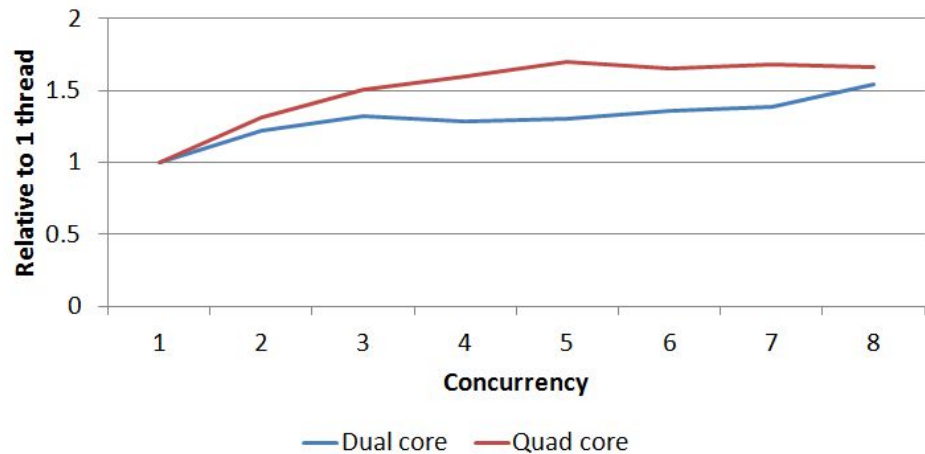
Contended atomic append



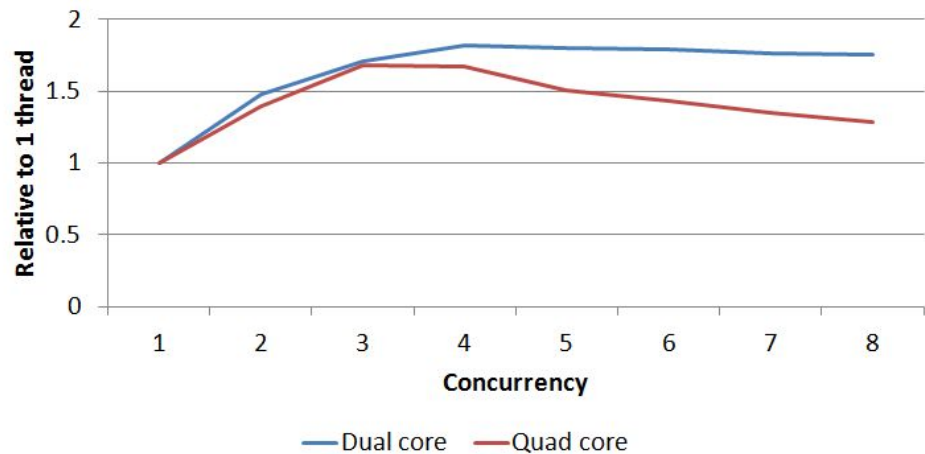
Contended memory maps



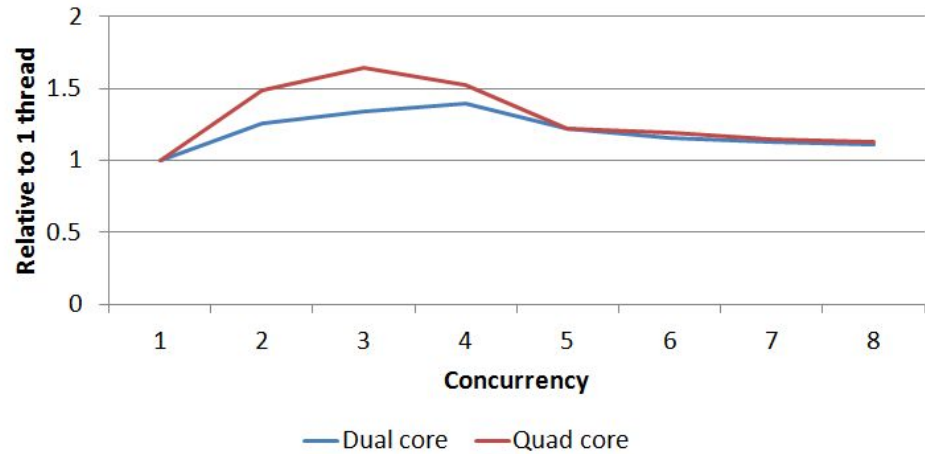
Uncontended lock files



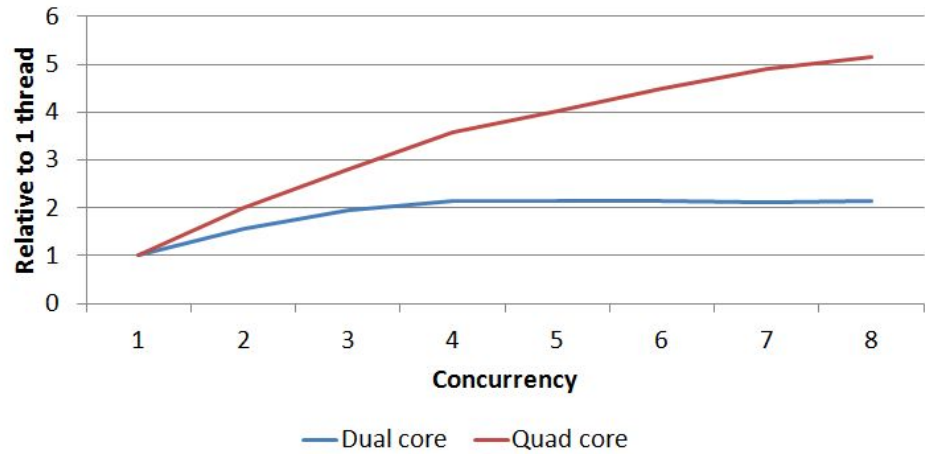
Uncontended byte ranges



Uncontended atomic append



Uncontended memory maps



Thank you

And let the questions begin!

Github: <https://github.com/ned14/boost.afio>

Ref docs: <https://ned14.github.io/boost.afio/>