

Atomics in Arm: Are they ruining your performance?

Ricardo Jesus and Michèle Weiland

EPCC, The University of Edinburgh

ISC 2022 AHUG Workshop

- Introduction
- The problem
- Benchmarking compare-and-swap operations
- Experiments and observations
- Conclusions

- **Atomics** are indivisible operations guaranteed to run to completion without interference from other threads
- **Compare-and-swap (CAS)** is one of the most common and general atomic operations
 - Allows a thread to atomically compare a memory location with a given value and, if the values match, to write another value in that memory location
 - Used to implement locks, semaphores, concurrent data structures, etc
- **In AArch64** CAS operations can be implemented via LL-SC pairs of instructions or explicit CAS instructions (LSE, Armv8.1+)

The performance of LL-SC/CAS approaches can be fundamentally different on different micro-architectures

- The RAJA Performance Suite (RAJAPerf)¹ is a benchmark suite consisting of loop-based kernels extracted from HPC applications, other benchmark suites, etc
- PI_ATOMIC is one of RAJAPerf's kernels. It computes π via

$$\pi = 4 \arctan(1) = 4 \int_0^1 \frac{1}{1+x^2} dx \quad (1)$$

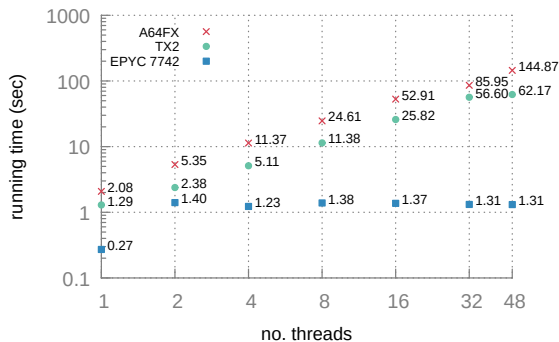
- Multiple implementations such as sequential, OpenMP, and RAJA
- Useful to assess the performance of OpenMP atomics in high contention

```
1  double *pi = ...;
2  #pragma omp parallel for
3  for(int i = 0; i < N; ++i) {
4      double x = ((double)i + 0.5) * dx;
5      #pragma omp atomic
6      *pi += dx / (1.0 + x * x);
7  }
8  pi *= 4.0;
```

¹<https://github.com/LLNL/RAJAPerf>

- We were using RAJAPerf to compare A64FX, ThunderX2, and AMD EPYC 7742 processors
- The PI_ATOMIC kernel stood out as the single worst-performing kernel on Arm
- **Fundamental difference** in the behaviours of the EPYC and the Arm-based CPUs
 - On the EPYC 7742, the running time **stays constant** for multi-threaded runs
 - On the A64FX/TX2, the running time **grows linearly** with the no. threads

- 50 repetitions of 1,000,000 iterations each (default parameters)
- Optimisation flags shown later



- Assembly¹ is identical for x86 and AArch64 targets
- Move value from GP to FP register, operate, move result back and attempt CAS
- Unnecessary jump present in both targets

```

1  /*      C/C++      */
2  for(...) {
3      #pragma omp atomic
4      *x += 1.0;
5  }

```

```

1  /*      ASM (x86)      */
2  .L2:
3      vmovq    %rdx, %xmm2    // move x to XMM
4      mov     %rdx, %rax      // create copy
5      vaddsd   %xmm0, %xmm2, %xmm1    // x+1.0
6      vmovq    %xmm1, %rcx    // move result back
7      lock cmpxchg %rcx, (%rdi) // attempt CAS
8
9      jne     .L4              // retry if fail
10     ...
11 .L4:
12     mov     %rax, %rdx
13     jmp     .L2

```

```

1  /*      ASM (A64)      */
2  .L2:
3      fmov     d1, x2    // move to FP register
4      mov     x3, x2      // create copy
5      fadd     d1, d0, d1    // compute x+1.0
6      fmov     x4, d1      // move result back
7      casal    x3, x4, [x0]    // attempt CAS
8      cmp     x2, x3
9      bne     .L4          // retry if fail
10     ...
11 .L4:
12     mov     x2, x3
13     b       .L2

```

¹Compiled with GCC 11.2 and -O3 -fopenmp -march=znver2/-mcpu=a64fx.

- Threads concurrently increment a counter iters times each using CAS operations
- Measure: **average time** and **average number of attempts** per increment
- A few parameters:
 - counter type: `uint64_t`, `double`
 - strategy: CAS, LL-SC
 - memory order: relaxed, acquire, release, acquire-release
- Kernels are implemented in assembly (next slide)

```
1  #include <stdint>
2
3  template <typename counter_t>
4  uint64_t kernel(uint64_t iters, counter_t* mem,
5                  bool weak, int memorder)
6  {
7      uint64_t attempts = 0;
8      do {
9          counter_t expected, desired;
10         do {
11             attempts++;
12             expected = *mem;
13             desired = expected + 1;
14         } while(!__atomic_compare_exchange(
15                 mem, &expected, &desired,
16                 weak, memorder, memorder));
17     } while(--iters);
18     return attempts;
19 }
```


- From left to right: double kernel for x86, uint64_t kernel for AArch64 with CAS strategy, and double kernel for AArch64 with LL-SC strategy

```

1  kernel_dbl_x86:
2      xor %r8d, %r8d
3      vmovsd one, %xmm0
4  loop_dbl_x86:
5      mov (%rsi), %rax
6  try_dbl_x86:
7      inc %r8
8      vmovq %rax, %xmm1
9      vaddsd %xmm0, %xmm1, %xmm2
10     vmovq %xmm2, %rcx
11     lock cmpxchg %rcx, (%rsi)
12     jne try_dbl_x86
13     dec %rdi
14     jne loop_dbl_x86
15     mov %r8, %rax
16     ret
17  one:
18     .double 1.0

```

lea 0x1(%rax), %rcx
for uint64_t kernel



```

1  kernel_u64_cas:
2      mov x9, 0
3
4  loop_u64_cas:
5      ldr x2, [x1]
6  try_u64_cas:
7      add x9, x9, 1
8
9      add x3, x2, 1
10
11     mov x4, x2
12     cas x2, x3, [x1]
13     cmp x2, x4
14     bne try_u64_cas
15     subs x0, x0, 1
16     bne loop_u64_cas
17     mov x0, x9
18     ret

```

```

1  kernel_dbl_ldxr_stxr:
2      mov x9, 0
3      fmov d0, 1.0
4  loop_dbl_ldxr_stxr:
5      ldxr x2, [x1]
6
7      add x9, x9, 1
8      fmov d1, x2
9      fadd d1, d0, d1
10     fmov x3, d1
11
12     stxr w4, x3, [x1]
13
14     cbnz w4, loop_dbl_ldxr_stxr
15     subs x0, x0, 1
16     bne loop_dbl_ldxr_stxr
17     mov x0, x9
18     ret

```

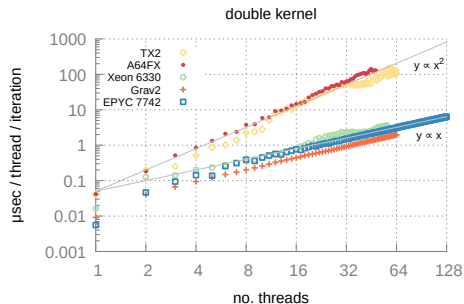
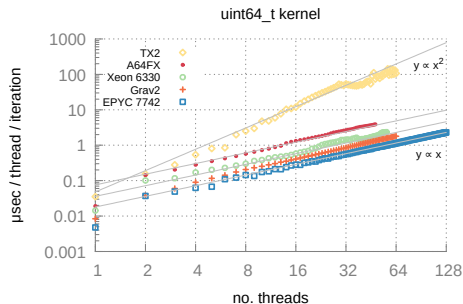

- All measurements taken in exclusive mode
- Each thread runs for 100,000 iterations
- Threads pinned to cores and placed close to each other (OMP_PLACES=cores and OMP_PROC_BIND=close)

System	Processor	No. cores	Compiler	Flags
Fulham	ThunderX2	2×32	GCC 10.1.0	-O3 -fopenmp -mcpu=native
Isambard 2	A64FX	48	GCC 11.1.0	-O3 -fopenmp -mcpu=native
C6g	Graviton2	64	GCC 10.3.1	-O3 -fopenmp -mcpu=native
C7g	Graviton3	64	GCC 10.3.1	-O3 -fopenmp -mcpu=native
ARCHER2	EPYC 7742	2×64	GCC 11.2.0	-O3 -fopenmp -march=native
ICX	Xeon Gold 6330	2×28	GCC 11.2.0	-O3 -fopenmp -march=native

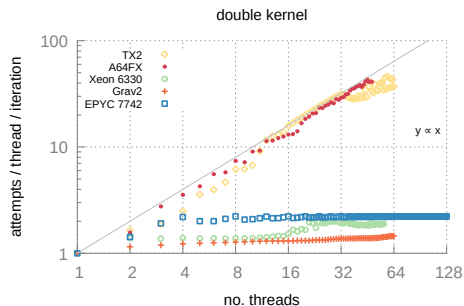
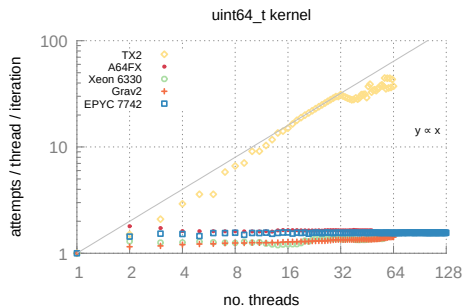
- Only show results for cas and ldxr+stxr are shown (i.e. relaxed memory order)
- Same trends with other memory orders (correlation coefficient below)

memory order	A64FX		Graviton2		ThunderX2	
	uint64_t	double	uint64_t	double	uint64_t	double
casa	0.9999	0.9993	1.0000	1.0000	0.9998	0.9998
casl	1.0000	0.9989	1.0000	1.0000	0.9999	0.9998
casal	0.9999	0.9989	1.0000	1.0000	0.9999	0.9997
ldaxr+stxr	1.0000	1.0000	0.9636	0.9559	0.9997	0.9994
ldxr+stlxx	0.9999	1.0000	0.9621	0.9380	0.9989	0.9987
ldaxr+stlxx	0.9999	1.0000	0.9638	0.9633	0.9983	0.9983

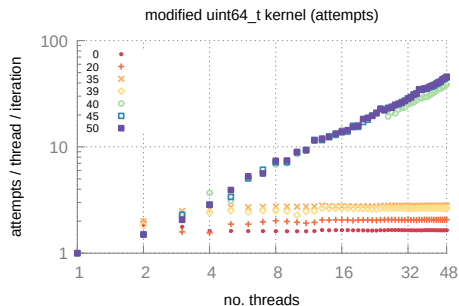
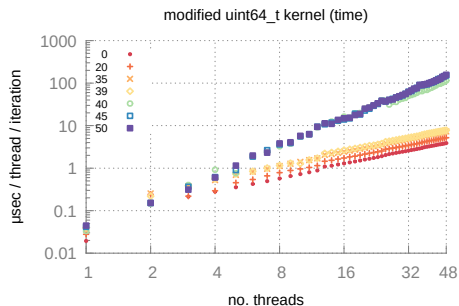
- Two types of behaviour in terms of average time per increment
 - **Linear scaling** on Graviton2, x86 CPUs, and A64FX (on the uint64_t kernel)
 - **Quadratic scaling** on ThunderX2 and A64FX (on the double kernel)
- Behaviour exhibited by the double kernel matches that seen in PI_ATOMIC



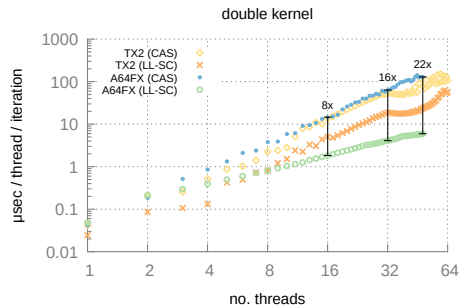
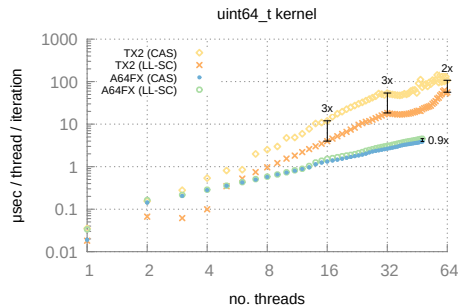
- Analogous trends in terms of average number of attempts per increment
 - **Constant** on Graviton2, x86 CPUs, and A64FX (on the uint64_t kernel)
 - **Linear scaling** on ThunderX2 and A64FX (on the double kernel)
- Explains the slowdown
- Suggests a lack of built-in backoff mechanisms in the A64FX and ThunderX2



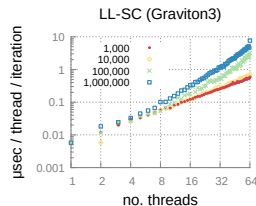
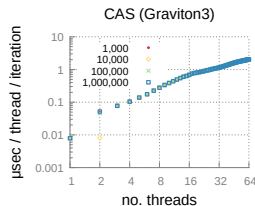
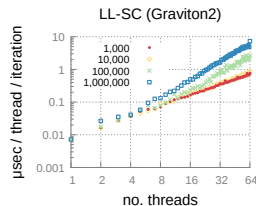
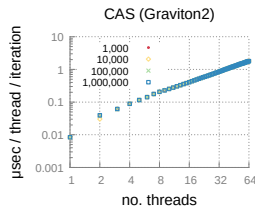
- Tested by inserting add instructions to the body of the integer kernel (“additional latency of loop body”)
- **Linear behaviour** if number of instructions inserted ≤ 39
- Sharp transition to **quadratic behaviour** if ≥ 40 instructions inserted
⇒ Matches double kernel



- Flat 2–3x speedup on ThunderX2 over CAS
- The **double LL-SC kernel exhibits linear scaling on the A64FX**
 - For high thread counts \Rightarrow speedup of more than 20x vs. CAS



- With LL-SC, average time per increment increases as the number of increments performed increases
- This contrasts with CAS, which is **invariant** to the number of increments per thread
- For a relatively small no. increments (<100,000), LL-SC approach can still bring significant speedups (2–3x) over CAS



- The newer LSE CAS instructions do not perform optimally on either the A64FX or ThunderX2
 - They can cause very significant slowdowns whereby the average time per successful CAS operation **quadruples** as the no. threads attempting the operation doubles
 - On the A64FX this behaviour arises when the time to prepare the CAS operation rises above a critical value around the 40 clock cycles mark
- Despite being older, the LL-SC instructions can provide significant speedups when compared to the newer CAS instructions on all Arm-based CPUs we tested
 - Flat 2–3x speedups on ThunderX2 and on Graviton2/3 (for low iteration counts)
 - On the A64FX, **LL-SC reduces** the quadratic scaling seen with CAS instructions **to the expected linear scaling**
⇒ For high thread counts, this results in speedups of over 20x
- We are currently working with vendors to understand the reasons leading to these behaviours and to develop optimised atomic primitives for their microarchitectures

Questions?