

From the roots to the sky: end-to-end concurrency verification at DresdenRC



Jonas Oberhauser
Dresden Research Center

Selected achievements of DresdenRC in the recent years



non-safety-critical products



Fix bugs

concurrency bugs

- seL4**: 1 missing barrier and 1 misplaced barrier breaking integrity of critical sections. *3 months Arm-internal discussion whether the bug can happen on Arm*
- 10+** barrier and algorithmic bugs in various internal components. *Grandfather paradox*
- MariaDB**: Missing barrier leading to data corruption.
- HMCS lock**: Missing barrier in the HMCS lock paper leading to hang.

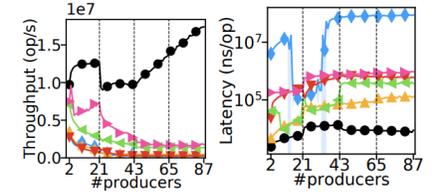
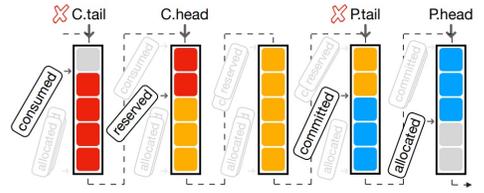


leave bug that increases throughput by 5% & causes 5 minutes downtime each month

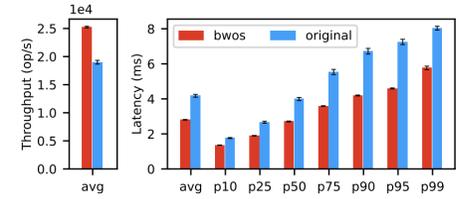
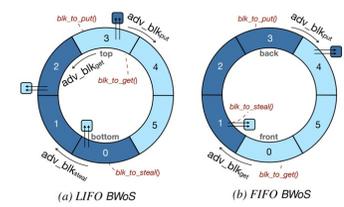


HMCS lock paper leading to hang

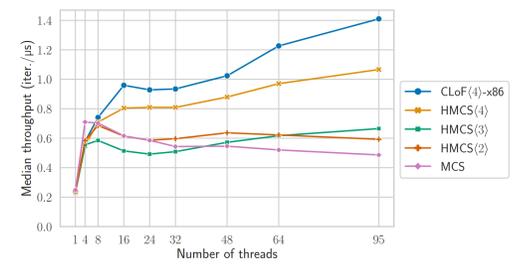
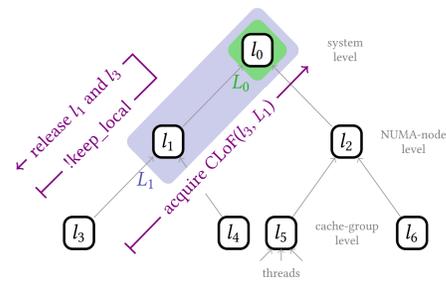
New concurrency algorithms



queue – **BBQ**: up to 10x improved throughput, latency



work stealing – **BWoS**: better latency, more throughput



lock – **CloF**: automatically configurable, improved throughput

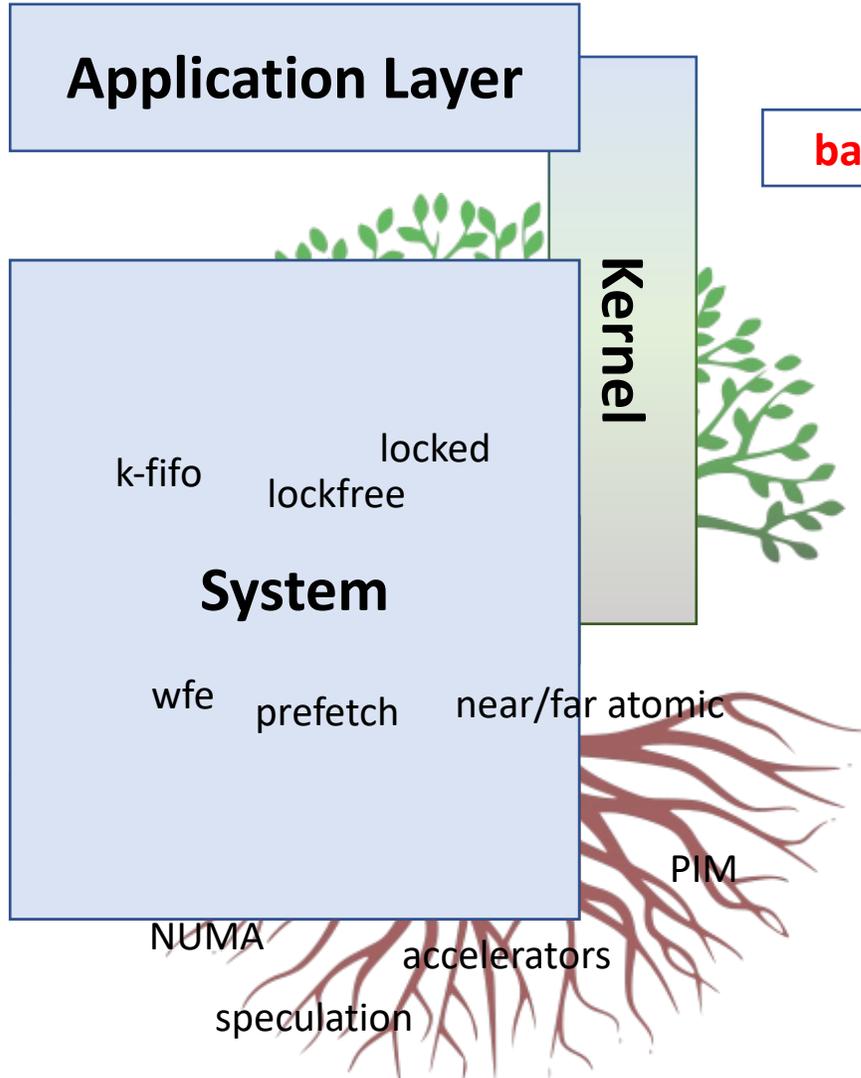


Concurrency

Thread-Safe Object

channels, maps, lists

Shared memory



barrier

```
void T1() {
  n = 20;
  release(flag = true);
}
```

```
void T2() {
  while (! acquire(flag));
  assert(n==20);
}
```

→ **Weak Memory Model (WMM)**

Highly optimized –

Challenge

- program now crashes after *only 100µsec!*
design application layer & system to leverage evolving hardware & compiler optimizations
- manage resulting **complexity**
- **disable optimizations selectively** when they make the application incorrect (e.g. with **barriers**)

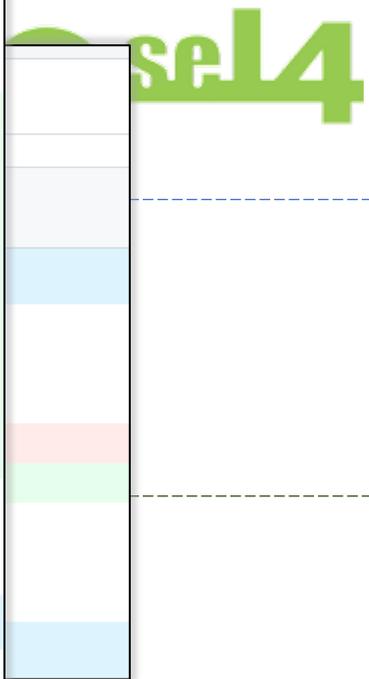


16 include/arch/arm/arch/model/smp.h

```

@@ -22,6 +22,16 @@ static inline bool_t try_arch_atomic_exchange(void *ptr, void *new_val, void **p
22 22     uint32_t atomic_status;
23 23     void *temp;
24 24
25 +     /* On failure there is no store.
26 +     Thus we can ignore release in the failure memory order,
27 +     and place a release barrier only if the success memory order implies
28 +     such an ordering. */
29 +     if (success_memorder == __ATOMIC_RELEASE || success_memorder == __ATOMIC_ACQ_REL) {
30 +         __atomic_thread_fence(__ATOMIC_RELEASE);
31 +     } else if (success_memorder == __ATOMIC_SEQ_CST) {
32 +         __atomic_thread_fence(__ATOMIC_SEQ_CST);
33 +     }
34 +
25 35     asm volatile(
26 36         LD_EX "[%prev_output], [%[ptr_val]]           \n\t" /* ret = *ptr */
27 37         ST_EX "%" OP_WIDTH "[atomic_var], %[new_val] , [%[ptr_val]] \n\t" /* *ptr = new */
@@ -34,7 +44,11 @@ static inline bool_t try_arch_atomic_exchange(void *ptr, void *new_val, void **p
34 44
35 45     /* Atomic operation success */
36 46     if (likely(!atomic_status)) {
37 -         __atomic_thread_fence(success_memorder);
47 +         if (success_memorder == __ATOMIC_ACQUIRE || success_memorder == __ATOMIC_ACQ_REL) {
48 +             __atomic_thread_fence(__ATOMIC_ACQUIRE);
49 +         } else if (success_memorder == __ATOMIC_SEQ_CST) {
50 +             __atomic_thread_fence(__ATOMIC_SEQ_CST);
51 +         }
38 52     } else {
39 53         /* Atomic operation failure */
40 54         __atomic_thread_fence(failure_memorder);

```



missing release-barrier – two threads can enter critical section!

release-barrier is implemented wrong – two threads can still enter critical section!

ARM_ERRATA_742230
 bool "ARM errata: DMB operation may be faulty"
 depends on CPU_V7 && SMP
 depends on !ARCH_MULTIPLATFORM
 help
 This option enables the workaround for the 742230 Cortex-A9 (r1p0..r2p2) erratum. Under rare circumstances, a DMB instruction between two write operations may not ensure the correct visibility ordering of the two writes. This workaround sets a specific bit in the diagnostic register of the Cortex-A9 which causes the DMB instruction to behave as a DSB, ensuring the correct behaviour of the two writes.



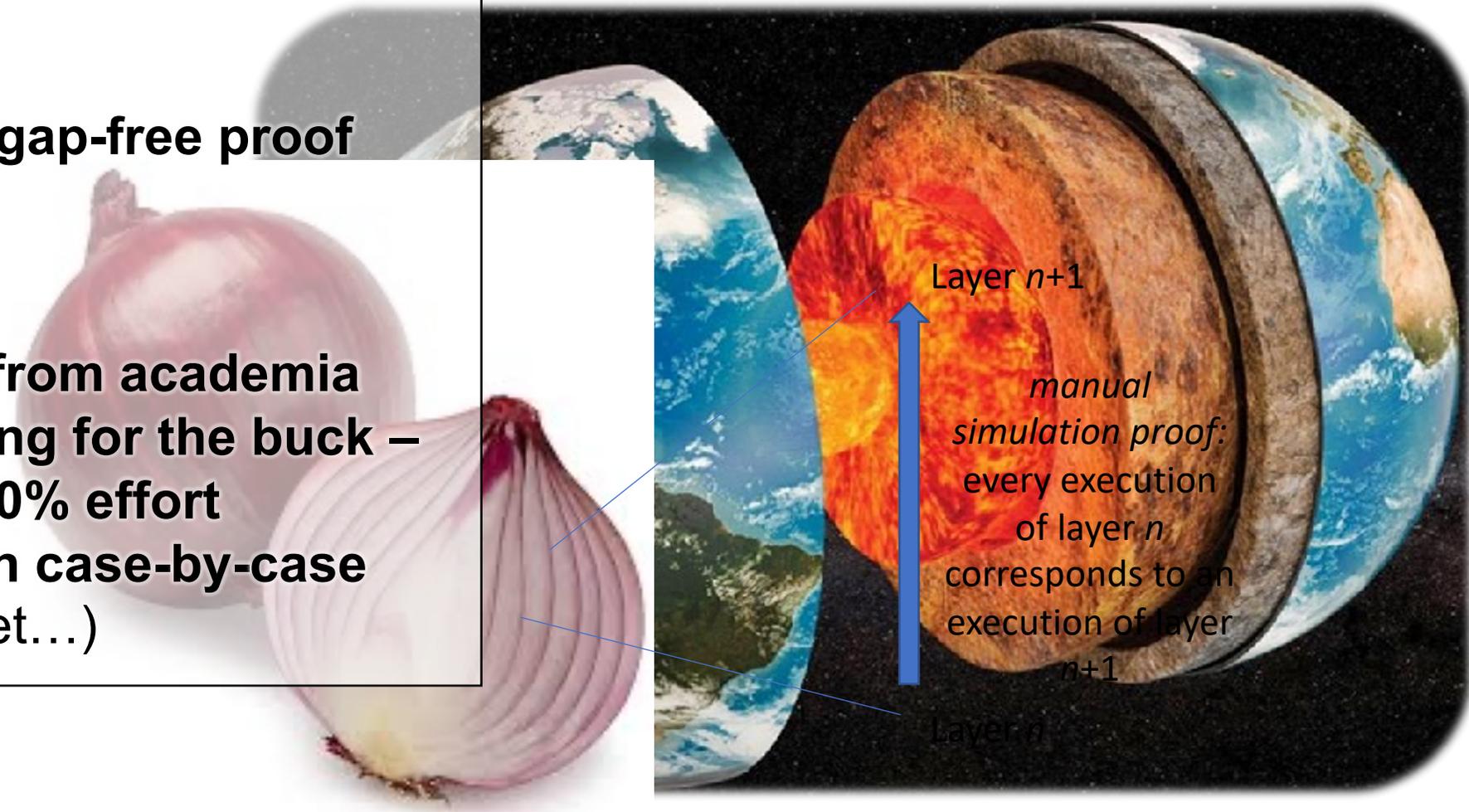
How to do end-to-end Verification?

- Much larger problem scale
- High speed of change

→ **manual, unbounded, gap-free proof can't keep up!**

What to do?

- 1. learn how to think from academia**
- 2. look for biggest bang for the buck – 80% verified with 20% effort**
- 3. choose methods on case-by-case basis (no silver bullet...)**



Level of detail at each level



Application Layer

Concurrency Library

HW Mapping

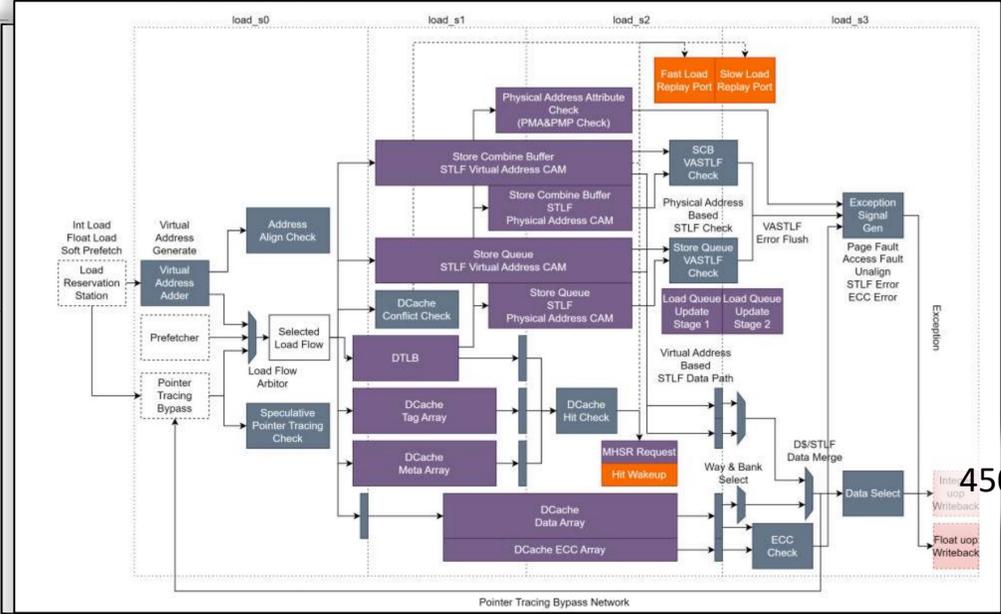
Hardware

Kernel

LEVEL OF DETAIL



50-200 LoC



```

red-before *)
[A | Q]

red-before *)
dmb.full]; po
e ([A];amo;[L]); po
[A]
; po; [dmb.ld]; po
po
[dmb.st]; po; [W]

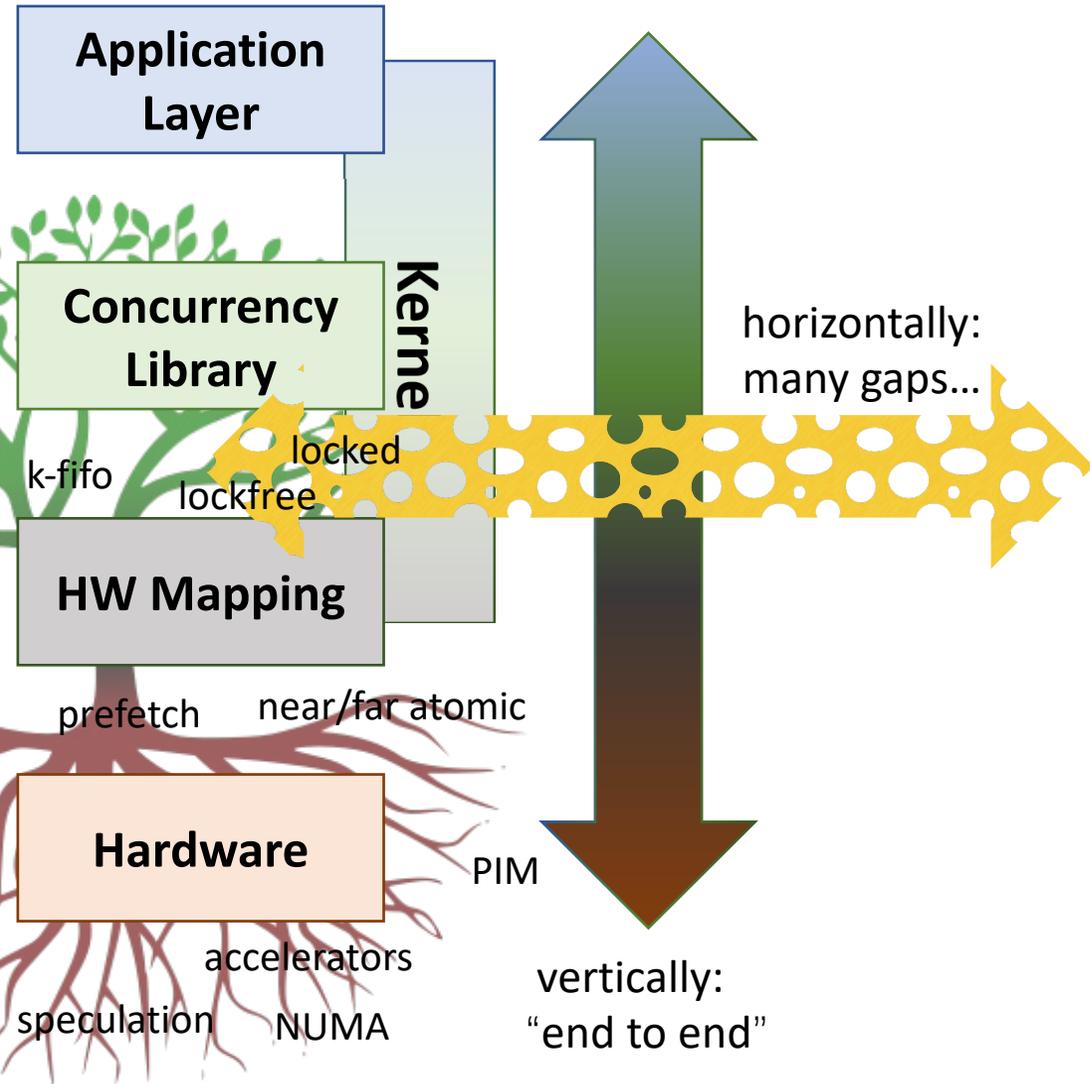
red-before *)
ws; si
lob
lob
lob; lob

sibility requirement *)
s; si | lob)+ as external
    
```





Summary & Lessons Learned



- Each layer has a different point in “detail” vs “degrees of freedom” making different methodologies and different amount of gaps effective
- Learned a lot about the architecture at each layer by opening “system box” and doing verification
- Could build more efficient concurrency algorithms due to that knowledge
- More efficient concurrency **needs** verification to manage complexity