

ETH zürich @spcl_eth DINFK

T. HOEFLER, M. PUESCHEL

Lecture 5: Fast practical locks, lock-free, consensus, and scalable locks

Teaching assistant: Salvatore Di Girolamo Motivational video: <https://www.youtube.com/watch?v=qx2dRIQXnbs>

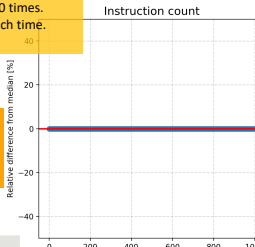


SPCL @spcl_eth ETH zürich

Nondeterminism in [most] performance measurements!

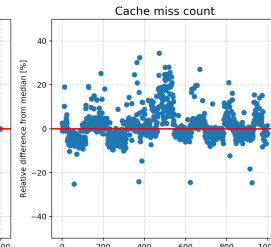
Same code executed 1000 times. Two metrics measured each time.

How do we report measurements showing high variation?



Instruction count

One is amazingly stable.



Cache miss count

The other—not at all!

```
const int n=1000;
volatile int a=0;
for (int i=0; i<n; ++i)
  a++;
```

TH, R. Belli: Scientific Benchmarking of Parallel Computing Systems, IEEE/ACM SC15 (full talk at <https://www.youtube.com/watch?v=6x01WAWU1U>)

SPCL @spcl_eth ETH zürich

Nondeterminism in [most] performance measurements!

Same code executed 1000 times. Two metrics measured each time.

How do we report measurements showing high variation?

Rule 5: Report if the measurement values are deterministic. For nondeterministic data, report confidence intervals of the measurement

```
const int n=1000;
volatile int a=0;
for (int i=0; i<n; ++i)
  a++;
```

One is amazingly stable. The other—not at all!

TH, R. Belli: Scientific Benchmarking of Parallel Computing Systems, IEEE/ACM SC15 (full talk at <https://www.youtube.com/watch?v=6x01WAWU1U>)

SPCL @spcl_eth ETH zürich

Administrivia

- Intermediate project presentation: next Monday 10/29 during lecture
 - Report will be due in January!
 - Starting to write early is very helpful --- write - rewrite - rewrite (no joke!)
- Coordinate your talk! You have 10 minutes (8 talk + 2 Q&A)
 - What will you be speaking about?
 - Focus on the key aspects (time is tight!)
 - Who will be speaking (up to you).
 - Engage the audience ☺
- Send slides by Sunday night (11:59pm CH time) to Salvatore!
 - We will use a single (windows) laptop to avoid delays when switching
 - Expect only Windows (powerpoint) or a PDF viewer
 - The order of talks will be randomized for fairness

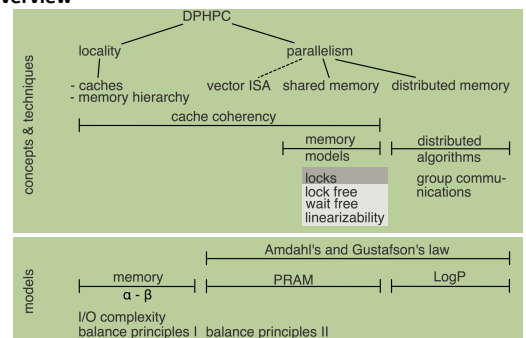
SPCL @spcl_eth ETH zürich

Review of last lecture

- Memory models in practical parallel programming
 - Synchronized programming
 - How locks synchronize processes and memory!
- Proving program correctness
 - Pre-/postconditions - sequential
 - Lifting to parallel
 - How to prove locked programs correct (nearly trivial)
- Lock implementation
 - Peterson lock - proof of correctness (using read/write histories, program and visibility orders)
 - With x86 memory model!
 - Lock performance
 - Simple x86 - how much does memory model correctness cost?

SPCL @spcl_eth ETH zürich

DPHPC Overview



concepts & techniques

- locality
 - caches
 - memory hierarchy
- parallelism
 - vector ISA
 - shared memory
 - distributed memory
- cache coherency
- memory models
 - locks
 - lock free
 - wait free
 - linearizability
- distributed algorithms
 - group communications

models

- memory $\alpha - \beta$
- PRAM
- LogP
- I/O complexity
- balance principles I
- balance principles II
- Little's Law
- scheduling

Amdahl's and Gustafson's law

SPCL ETH zürich

Goals of this lecture

- Fast and scalable practical locks!
 - Based on atomic operations
 - Why do we need atomic operations?
- Recap lock-free and wait-free programming
 - Proof that wait-free consensus is impossible without atomics
 - Valence argument: a proof technique similar to showing that atomics are needed for locks
- Locks in practical setting
 - How to block?
 - When to block?
 - How long to block?
 - Simple proof of competitiveness
- Case study: large-scale distributed memory locking
 - Problems and outline to next class

8

SPCL ETH zürich

Back to Peterson in Practice ... on x86

- Implement and run our little counter on x86
- 100000 iterations
 - 1.6 · 10⁻⁶% errors
 - What is the problem?

```
volatile int flag[2];
volatile int victim;

void lock() {
  int j = 1 - tid;
  flag[tid] = 1; // I'm interested
  victim = tid; // other goes first
  while (flag[j] && victim == tid) {}; // wait
}

void unlock() {
  flag[tid] = 0; // I'm not interested
}
```

9

SPCL ETH zürich

Peterson in Practice ... on x86

- Implement and run our little counter on x86
- 100000 iterations
 - 1.6 · 10⁻⁶% errors
 - What is the problem?
 - No sequential consistency for W(v) and R(flag[j])

```
volatile int flag[2];
volatile int victim;

void lock() {
  int j = 1 - tid;
  flag[tid] = 1; // I'm interested
  victim = tid; // other goes first
  asm("mfence");
  while (flag[j] && victim == tid) {}; // wait
}

void unlock() {
  flag[tid] = 0; // I'm not interested
}
```

10

SPCL ETH zürich

Peterson in Practice ... on x86

- Implement and run our little counter on x86
- Many iterations
 - 1.6 · 10⁻⁶% errors
 - What is the problem?
 - No sequential consistency for W(v) and R(flag[j])
 - Still 1.3 · 10⁻⁶% Why?

```
volatile int flag[2];
volatile int victim;

void lock() {
  int j = 1 - tid;
  flag[tid] = 1; // I'm interested
  victim = tid; // other goes first
  asm("mfence");
  while (flag[j] && victim == tid) {}; // wait
}

void unlock() {
  flag[tid] = 0; // I'm not interested
}
```

11

SPCL ETH zürich

Peterson in Practice ... on x86

- Implement and run our little counter on x86
- Many iterations
 - 1.6 · 10⁻⁶% errors
 - What is the problem?
 - No sequential consistency for W(v) and R(flag[j])
 - Still 1.3 · 10⁻⁶% Why?
 - Reads may slip into CR!

```
volatile int flag[2];
volatile int victim;

void lock() {
  int j = 1 - tid;
  flag[tid] = 1; // I'm interested
  victim = tid; // other goes first
  asm("mfence");
  while (flag[j] && victim == tid) {}; // wait
}

void unlock() {
  asm("mfence");
  flag[tid] = 0; // I'm not interested
}
```

The compiler may inline this function ☺

12

SPCL ETH zürich

Correct Peterson Lock on x86

- Unoptimized (naïve sprinkling of mfences)
- Performance:
 - No mfence: 375ns
 - mfence in lock: 379ns
 - mfence in unlock: 404ns
 - Two mfence: 427ns (+14%)

```
volatile int flag[2];
volatile int victim;

void lock() {
  int j = 1 - tid;
  flag[tid] = 1; // I'm interested
  victim = tid; // other goes first
  asm("mfence");
  while (flag[j] && victim == tid) {}; // wait
}

void unlock() {
  asm("mfence");
  flag[tid] = 0; // I'm not interested
}
```

13

SPCL ETH zürich

Hardware Support?

- Hardware atomic operations:
 - Test&Set
Write const to memory while returning the old value
 - Atomic swap
Atomically exchange memory and register
 - Fetch&Op
Get value and apply operation to memory location
 - Compare&Swap
Compare two values and swap memory with register if equal
 - Load-linked/Store-Conditional LL/SC (or load-acquire (LDA) store-release (STL) on ARM)
Loads value from memory, allows operations, commits only if no other updates committed → mini-TM
 - Intel TSX (transactional synchronization extensions)
Hardware-TM (roll your own atomic operations)

14

SPCL ETH zürich

Relative Power of Synchronization

- Design-Problem I: Multi-core Processor
 - Which atomic operations are useful?
- Design-Problem II: Complex Application
 - What atomic should I use?
- Generally hard to answer ☹️
 - Depends on too many systems details (access patterns, CC implementation, contention, algorithm ...)
- Concept of “consensus number” C if a primitive can be used to solve the “consensus problem” in a finite number of steps (even if threads stop)
 - atomic registers have C=1 (thus locks have C=1!)
 - TAS, Swap, Fetch&Op have C=2
 - CAS, LL/SC, TM have C=∞

15

SPCL ETH zürich

Test-and-Set Locks

- Test-and-Set semantics
 - Memoize old value
 - Set fixed value TASval (true)
 - Return old value
- After execution:
 - Post-condition is a fixed (constant) value!

```
bool TestAndSet (bool *flag) {
    bool old = *flag;
    *flag = true;
    return old;
} // all atomic!
```

16

SPCL ETH zürich

Test-and-Set Locks

- Assume TASval indicates “locked”
- Write something else to indicate “unlocked”
- TAS until return value is != TASval (1 in this example)
- When will the lock be granted?
- Does this work well in practice?

```
bool TestAndSet (bool *flag) {
    bool old = *flag;
    *flag = true;
    return old;
} // all atomic!
```

```
volatile int lck = 0;

void lock() {
    while (TestAndSet(&lck) == 1);
}

void unlock() {
    lck = 0;
}
```

17

SPCL ETH zürich

Cacheline contention (or: why I told you about MESI and friends)

- On x86, the XCHG instruction is used to implement TAS
 - x86 lock is implicit in xchg!
- Cacheline is read and written
 - Ends up in exclusive state, invalidates other copies
 - Cacheline is “thrown” around uselessly
 - High load on memory subsystem
x86 lock is essentially a full memory barrier ☹️

```
movl $1, %eax
xchg %eax, (%ebx)
```

18

SPCL ETH zürich

Test-and-Test-and-Set (TATAS) Locks

- Spinning in TAS is not a good idea
- Spin on cache line in shared state
 - All threads at the same time, no cache coherency/memory traffic
- Danger!
 - Efficient but use with great care!
 - Generalizations are very dangerous

```
volatile int lck = 0;

void lock() {
    do {
        while (lck == 1);
    } while (TestAndSet(&lck) == 1);
}

void unlock() {
    lck = 0;
}
```

19

Warning: Even Experts get it wrong!

- Example: Double-Checked Locking

1997

Double-Checked Locking
An Optimization Pattern for Efficiently Initializing and Accessing Thread-Safe Objects

Douglas C. Schmidt
Subhadra K. Sivasubramanian
Dirk Grunwaldt
Wan, C. C. Lo

Abstract

The paper presents an analysis of the lock-free initialization and accessing of thread-safe objects. It shows that the double-checked locking pattern is not as efficient as it is often claimed to be. In fact, it is shown that the double-checked locking pattern can be replaced by a more efficient pattern that avoids the overhead of acquiring a lock by first testing for locking contention.

Google double-checked locking

2018

About 9,500,000 results (0.41 seconds)

Double-checked locking - Wikipedia
https://en.wikipedia.org/wiki/Double-checked_locking

The "Double-Checked Locking is Broken" Declaration
https://www.inet4you.com/~jgriener/About/DoubleCheckedLocking.html

Double-Checked Locking with Singleton | Baeldung
https://www.baeldung.com/java-singleton-double-checked-locking

LOCK10-J Use a correct form of the double-checked locking idiom - SEI
https://wiki.sei.cmu.edu/view/fulltext.php?filename=lock10-j

Double-checked locking: Clever, but broken | JavaWorld
https://www.javaworld.com/article/2136744/core-java/DoubleCheckedLocking.html

Problem: Memory ordering leads to race-conditions!

Contention?

- Do TATAS locks still have contention?
- When lock is released, k threads fight for cache line ownership
 - One gets the lock, all get the CL exclusively (serially!)
 - What would be a good solution? (think "collision avoidance")

```
volatile int lck = 0;

void lock() {
  do {
    while (lck == 1);
  } while (TestAndSet(&lck) == 1);
}

void unlock() {
  lck = 0;
}
```

TAS Lock with Exponential Backoff

- Exponential backoff eliminates contention statistically
 - Locks granted in unpredictable order
 - Starvation possible but unlikely
 - How can we make it even less likely?

```
volatile int lck = 0;

void lock() {
  while (TestAndSet(&lck) == 1) {
    wait(time);
    time *= 2; // double waiting time
  }
}

void unlock() {
  lck = 0;
}
```

Similar to: T. Anderson: "The performance of spin lock alternatives for shared-memory multiprocessors", TPDS, Vol. 1 Issue 1, Jan 1990

TAS Lock with Exponential Backoff

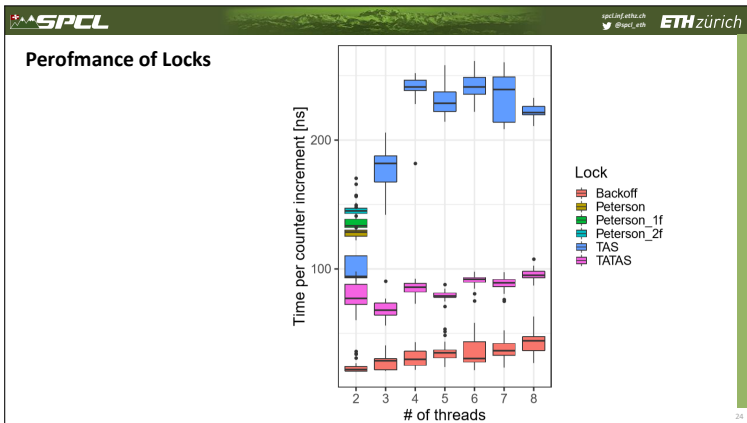
- Exponential backoff eliminates contention statistically
 - Locks granted in unpredictable order
 - Starvation possible but unlikely
 - Maximum waiting time makes it less likely

```
volatile int lck = 0;
const int maxtime=1000;

void lock() {
  while (TestAndSet(&lck) == 1) {
    wait(time);
    time = min(time * 2, maxtime);
  }
}

void unlock() {
  lck = 0;
}
```

Similar to: T. Anderson: "The performance of spin lock alternatives for shared-memory multiprocessors", TPDS, Vol. 1 Issue 1, Jan 1990



Improvements?

- Are TAS locks perfect?
 - What are the two biggest issues?
 - Cache coherency traffic (contending on same location with expensive atomics)
 - Critical section underutilization (waiting for backoff times will delay entry to CR)
- What would be a fix for that?
 - How is this solved at airports and shops (often at least)?
- Queue locks -- Threads enqueue
 - Learn from predecessor if it's their turn
 - Each threads spins at a different location
 - FIFO fairness

SPCL ETH zürich @spcl_eth

Array Queue Lock

- **Array to implement queue**
 - Tail-pointer shows next free queue position
 - Each thread spins on own location
 - CL padding!
 - index[] array can be put in TLS
- **So are we done now?**
 - What's wrong?
 - Synchronizing M objects requires $\Theta(NM)$ storage
 - What do we do now?

```

volatile int array[n] = {1,0,...,0};
volatile int index[n] = {0,0,...,0};
volatile int tail = 0;

void lock() {
  index[tid] = GetAndInc(tail) % n;
  while (!array[index[tid]]); // wait to receive lock
}

void unlock() {
  array[index[tid]] = 0; // I release my lock
  array[(index[tid] + 1) % n] = 1; // next one
}
  
```

26

SPCL ETH zürich @spcl_eth

CLH Lock (1993)

- **List-based (same queue principle)**
- **Discovered twice by Craig, Landin, Hagersten 1993/94**
- **2N+3M words**
 - N threads, M locks
- **Requires thread-local qnode pointer**
 - Can be hidden!

```

typedef struct qnode {
  struct qnode *prev;
  int succ_blocked;
} qnode;

qnode *lck = new qnode; // node owned by lock

void lock(qnode *lck, qnode *qn) {
  qn->succ_blocked = 1;
  qn->prev = FetchAndSet(lck, qn);
  while (qn->prev->succ_blocked);
}

void unlock(qnode **qn) {
  qnode *pred = (*qn)->prev;
  (*qn)->succ_blocked = 0;
  *qn = pred;
}
  
```

27

SPCL ETH zürich @spcl_eth

CLH Lock (1993)

- **Qnode objects represent thread state!**
 - succ_blocked == 1 if waiting or acquired lock
 - succ_blocked == 0 if released lock
- **List is implicit!**
 - One node per thread
 - Spin location changes
 - NUMA issues (cacheless)
- **Can we do better?**

```

typedef struct qnode {
  struct qnode *prev;
  int succ_blocked;
} qnode;

qnode *lck = new qnode; // node owned by lock

void lock(qnode *lck, qnode *qn) {
  qn->succ_blocked = 1;
  qn->prev = FetchAndSet(lck, qn);
  while (qn->prev->succ_blocked);
}

void unlock(qnode **qn) {
  qnode *pred = (*qn)->prev;
  (*qn)->succ_blocked = 0;
  *qn = pred;
}
  
```

28

SPCL ETH zürich @spcl_eth

MCS Lock (1991)

- **Make queue explicit**
 - Acquire lock by appending to queue
 - Spin on own node until locked is reset
- **Similar advantages as CLH but**
 - Only 2N + M words
 - Spinning position is fixed!
 - Benefits cache-less NUMA
- **What are the issues?**
 - Releasing lock spins
 - More atomics!

```

typedef struct qnode {
  struct qnode *next;
  int succ_blocked;
} qnode;

qnode *lck = NULL;

void lock(qnode *lck, qnode *qn) {
  qn->next = NULL;
  qnode *pred = FetchAndSet(lck, qn);
  if (pred != NULL) {
    qn->locked = 1;
    pred->next = qn;
    while (qn->locked);
  }
}

void unlock(qnode *lck, qnode *qn) {
  if (qn->next == NULL) { // if we're the last waiter
    if (CAS(lck, qn, NULL)) return;
    while (qn->next == NULL); // wait for pred arrival
  }
  qn->next->locked = 0; // free next waiter
  qn->next = NULL;
}
  
```

29

SPCL ETH zürich @spcl_eth

Lessons Learned!

- **Key Lesson:**
 - Reducing memory (coherency) traffic is most important!
 - Not always straight-forward (need to reason about CL states)
- **MCS: 2006 Dijkstra Prize in distributed computing**
 - "an outstanding paper on the principles of distributed computing, whose significance and impact on the theory and/or practice of distributed computing has been evident for at least a decade"
 - "probably the most influential practical mutual exclusion algorithm ever"
 - "vastly superior to all previous mutual exclusion algorithms"
 - fast, fair, scalable → widely used, always compared against!

30

SPCL ETH zürich @spcl_eth

Time to Declare Victory?

- **Down to memory complexity of 2N+M**
 - Probably close to optimal
- **Only local spinning**
 - Several variants with low expected contention
- **But: we assumed sequential consistency**
 - Reality causes trouble sometimes
 - Sprinkling memory fences may harm performance
 - Open research on minimally-synching algorithms!
 - Come and talk to me if you're interested

31

SPCL ETH zürich

Fighting CPU waste: Condition Variables

- **Allow threads to yield CPU and leave the OS run queue**
 - Other threads can get them back on the queue!
- **cond_wait(cond, lock) – yield and go to sleep**
- **cond_signal(cond) – wake up sleeping threads**
- **Wait and signal are OS calls**
 - Often expensive, which one is more expensive?
Wait, because it has to perform a full context switch

32

SPCL ETH zürich

When to Spin and When to Block?

- **Spinning consumes CPU cycles but is cheap**
 - “Steals” CPU from other threads
- **Blocking has high one-time cost and is then free**
 - Often hundreds of cycles (trap, save TCB ...)
 - Wakeup is also expensive (latency)
Also cache-pollution
- **Strategy:**
 - Poll for a while and then block
But what is a “while”??

33

SPCL ETH zürich

When to Spin and When to Block?

- **Optimal time depends on the future**
 - When will the active thread leave the CR?
 - Can compute optimal offline schedule
Q: What is the optimal offline schedule (assuming we know the future, i.e., when the lock will become available)?
 - Actual problem is an online problem
- **Competitive algorithms**
 - An algorithm is *c*-competitive if for a sequence of actions *x* and a constant *a* holds:
 $C(x) \leq c \cdot C_{opt}(x) + a$
 - What would a good spinning algorithm look like and what is the competitiveness?

34

SPCL ETH zürich

Competitive Spinning

- **If *T* is the overhead to process a wait, then a locking algorithm that spins for time *T* before it blocks is 2-competitive!**
 - Karlin, Manasse, McGeoch, Owicki: “Competitive Randomized Algorithms for Non-Uniform Problems”, SODA 1989
- **If randomized algorithms are used, then $e/(e-1)$ -competitiveness (~1.58) can be achieved**
 - See paper above!

35

SPCL ETH zürich

Remember: lock-free vs. wait-free


- **A lock-free method**
 - guarantees that infinitely often **some** method call finishes in a finite number of steps
- **A wait-free method**
 - guarantees that **each** method call finishes in a finite number of steps (implies lock-free)
- **Synchronization instructions are not equally powerful!**
 - Indeed, they form an infinite hierarchy; no instruction (primitive) in level *x* can be used for lock-/wait-free implementations of primitives in level $z > x$.

36

SPCL ETH zürich

Concept: Consensus Number

- **Each level of the hierarchy has a “consensus number” assigned.**
 - Is the maximum number of threads for which primitives in level *x* can solve the consensus problem
- **The consensus problem:**
 - Has single function: `decide(v)`
 - Each thread calls it at most once, the function returns a value that meets two conditions:
consistency: all threads get the same value
validity: the value is some thread's input
 - Simplification: binary consensus (inputs in {0,1})



37

SPCL ETH zürich

Understanding Consensus

- Can a particular class solve n-thread consensus wait-free?
 - A class C solves n-thread consensus if there exists a consensus protocol using **any number** of objects of class C and **any number** of atomic registers
 - The protocol has to be wait-free (bounded number of steps per thread)
 - The consensus number of a class C is the largest n for which that class solves n-thread consensus (may be infinite)
 - Assume we have a class D whose objects can be constructed from objects out of class C. If class C has consensus number n, what does class D have?

38

SPCL ETH zürich

Starting simple ...

- Binary consensus with two threads (A, B)!
 - Each thread moves until it decides on a value
 - May update shared objects
 - Protocol state = state of threads + state of shared objects
 - Initial state = state before any thread moved
 - Final state = state after all threads finished
 - States form a tree, wait-free property guarantees a finite tree
 - Example with two threads and two moves each!*

39

SPCL ETH zürich

Atomic Registers

- Theorem [Herlihy'91]: Atomic registers have consensus number one
 - I.e., they cannot be used to solve even two-thread consensus! Really?
- Proof outline:
 - Assume arbitrary consensus protocol, thread A, B
 - Run until it reaches critical state where next action determines outcome (show that it must have a critical state first)
 - Show all options using atomic registers and show that they cannot be used to determine one outcome for all possible executions!
 - Any thread reads (other thread runs solo until end)
 - Threads write to different registers (order doesn't matter)
 - Threads write to same register (solo thread can start after each write)

40

SPCL ETH zürich

Atomic Registers

- Theorem [Herlihy'91]: Atomic registers have consensus number one
- Corollary: It is impossible to construct a wait-free implementation of any object with consensus number of >1 using atomic registers
 - "perhaps one of the most striking impossibility results in Computer Science" (Herlihy, Shavit)
 - We need hardware atomics or Transactional Memory!
- Proof technique borrowed from:
 - Impossibility of distributed consensus with one ... - ACM Digital Library
 - by M.J. Fischer - 1985 - Cited by 4189 - Related articles
 - Apr 1, 1985 - The consensus problem involves an asynchronous system of processes, some of which may be unreliable. The problem is for the reliable ...
- Very influential paper, always worth a read!
 - Nicely shows proof techniques that are central to parallel and distributed computing!

41

SPCL ETH zürich

Other Atomic Operations

- Simple RMW operations (Test&Set, Fetch&Op, Swap, basically all functions where the op commutes or overwrites) have consensus number 2!
 - Similar proof technique (bivalence argument)
- CAS and TM have consensus number ∞
 - Constructive proof!

42

SPCL ETH zürich


Compare and Set/Swap Consensus

```

const int first = -1;
volatile int thread = -1;
int proposed[n];

int decide(v) {
    proposed[tid] = v;
    if(CAS(thread, first, tid))
        return v; // I won!
    else
        return proposed[thread]; // thread won
}

```



- CAS provides an infinite consensus number
 - Machines providing CAS are asynchronous computation equivalents of the Turing Machine
 - I.e., any concurrent object can be implemented in a wait-free manner (not necessarily fast!)

43

SPCL ETH zürich

Now you know everything ☺

- **Not really ... :-)**
 - We'll argue more about **performance** now!
- **But you have all the tools for:**
 - Efficient locks
 - Efficient lock-based algorithms
 - Efficient lock-free algorithms (or even wait-free)
 - Reasoning about parallelism!
- **What now?**
 - A different class of problems
Impact on wait-free/lock-free on actual performance is not well understood
 - Relevant to HPC, applies to shared and distributed memory
→ Group communications

44

SPCL ETH zürich

Case study: Fast Large-scale Locking in Practice

Inuitive semantics

Various performance penalties

SPCL ETH zürich

Locks: Challenges

Calciu et al.: NUMA-aware reader-writer locks, PPOPP'13

SPCL ETH zürich

Locks: Challenges

We need intra- and inter-node topology-awareness

We need to cover arbitrary topologies

SPCL ETH zürich

Locks: Challenges

We need to distinguish between readers and writers

We need flexible performance for both types of processes

[1] V. Venkatesaramani et al. Tsai: How facebook serves the social graph, SIGMOD'12

SPCL ETH zürich

Locks: Challenges

What will we use in the design?

Ingredient 1 - MCS Locks

Mellor-Crummey and Scott: Algorithms for Scalable Synchronization on Shared-Memory Multiprocessors, ACM TOCS'91

Ingredient 2 - Reader-Writer Locks

How to manage the design complexity?

How to ensure tunable performance?

What mechanism to use for efficient implementation?

REMOTE MEMORY ACCESS (RMA) PROGRAMMING

TH, J. Dinan, R. Thakur, B. Barrett, P. Balaji, W. Gropp, K. Underwood: Remote Memory Access Programming in MPI-3, ACM TOPC'15

REMOTE MEMORY ACCESS PROGRAMMING

- Implemented in hardware in NICs in the majority of HPC networks (RDMA support).

RMA-RW - Required Operations

MPI RMA primer ([much] more in the recitation sessions)

- Windows expose memory**
 - Created explicitly
- Remote accesses**
 - Put, get
 - Atomics
 - Accumulate (also atomic Put)
 - Get_accumulate (also atomic Get)
 - Fetch and op (faster single-word get_accumulate)
 - Compare and swap
- Synchronization**
 - Two modes: passive and active target
 - We use passive target today, similar to shared memory!
 - Synchronization: flush, flush_local
- Memory model**
 - Unified (coherent) and separate (not coherent) view - it's complicated but versatile

How to manage the design complexity?

How to ensure tunable performance?

What mechanism to use for efficient implementation?

How to manage the design complexity?

- each element has its own distributed MCS queue (DQ) of writers
- Readers and writers synchronize with a distributed counter (DC)
- Modular design
- MCS queues form a distributed tree (DT)

P. Schmid, M. Besta, TH: High-Performance Distributed RMA Locks, ACM HPDC'16, best paper

How to ensure tunable performance?

- Each DQ: fairness vs throughput of writers
- DC: a parameter for the latency of readers vs writers
- A tradeoff parameter for every structure
- DT: a parameter for the throughput of readers vs writers

P. Schmid, M. Besta, TH: High-Performance Distributed RMA Locks, ACM HPDC'16, best paper

Distributed MCS Queues (DQs) - Throughput vs Fairness

- Larger $T_{L,i}$: more throughput at level i . Smaller $T_{L,i}$: more fairness at level i .
- Each DQ: The maximum number of lock passings within a DQ at level i , before it is passed to another $T_{L,i}$ DQ at i .

P. Schmid, M. Besta, TH: High-Performance Distributed RMA Locks, ACM HPDC'16, best paper

Distributed Tree of Queues (DT) - Throughput of readers vs writers

- DT: The maximum number of consecutive lock passings within readers (T_R).

P. Schmid, M. Besta, TH: High-Performance Distributed RMA Locks, ACM HPDC'16, best paper

Distributed Counter (DC) - Latency of readers vs writers

DC: every k th compute node hosts a partial counter, all of which constitute the DC.
 $k = T_{DC}$

A writer holds the lock $b|x|y$
 Readers that arrived at the CS
 Readers that left the CS

$T_{DC} = 1$
 $T_{DC} = 2$

P. Schmid, M. Besta, TH: High-Performance Distributed RMA Locks, ACM HPDC'16, best paper

Design space

Locality vs fairness (for writers) $T_{L,i}$
 Higher throughput of writers vs readers T_R
 Lower latency of writers vs readers T_{DC}

Design B
 Design A

P. Schmid, M. Besta, TH: High-Performance Distributed RMA Locks, ACM HPDC'16, best paper

Lock Acquire by Readers

A lightweight acquire protocol for readers: only one atomic fetch-and-add (FAA) operation

A writer holds the lock $b|x|y$
 Readers that arrived at the CS
 Readers that left the CS

FAA
 FAA
 FAA

P. Schmid, M. Besta, TH: High-Performance Distributed RMA Locks, ACM HPDC'16, best paper

Lock Acquire by Writers

Acquire the main MCS lock
 Acquire the main lock
 Acquire MCS
 Acquire MCS

P. Schmid, M. Besta, TH: High-Performance Distributed RMA Locks, ACM HPDC'16, best paper

EVALUATION

- CSCS Piz Daint (Cray XC30)
- 5272 compute nodes
- 8 cores per node
- 169TB memory
- Microbenchmarks: acquire/release: latency, throughput
- Distributed hashtable

P. Schmid, M. Besta, TH: High-Performance Distributed RMA Locks, ACM HPDC'16, best paper

Evaluation - Distributed Counter Analysis

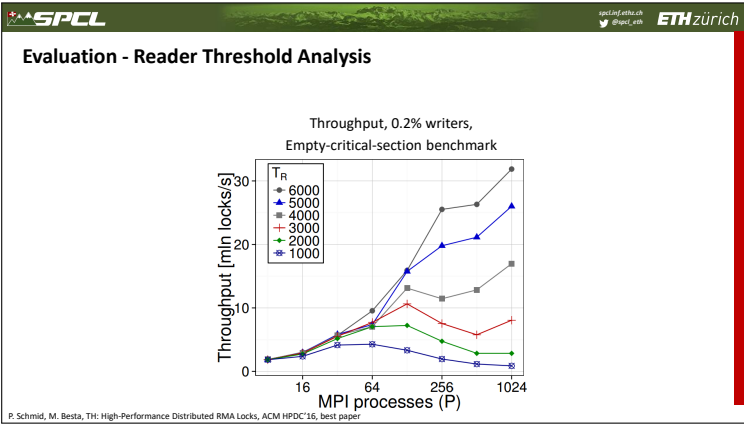
Throughput, 2% writers
 Single-operation benchmark

Throughput [min locks/s]

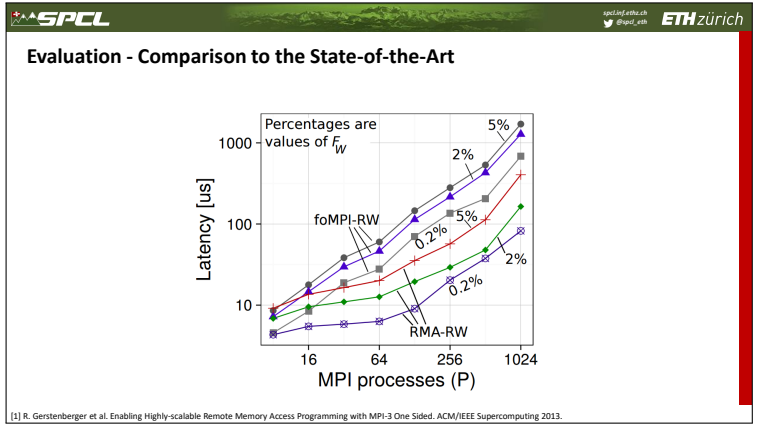
MPI processes (P)

T_{DC}
 64
 32
 16
 8
 4
 2

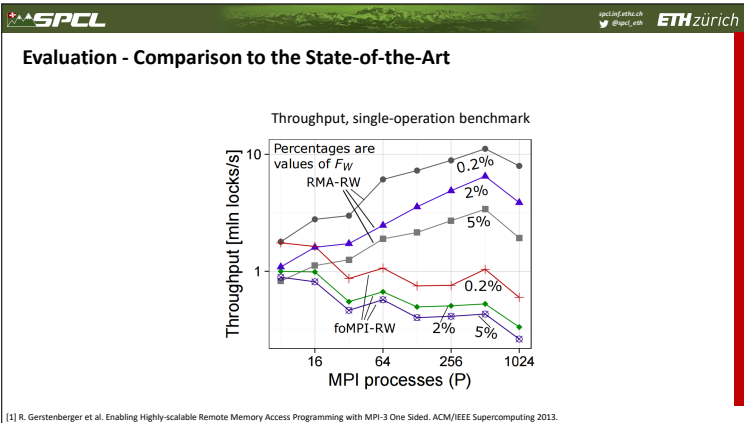
P. Schmid, M. Besta, TH: High-Performance Distributed RMA Locks, ACM HPDC'16, best paper



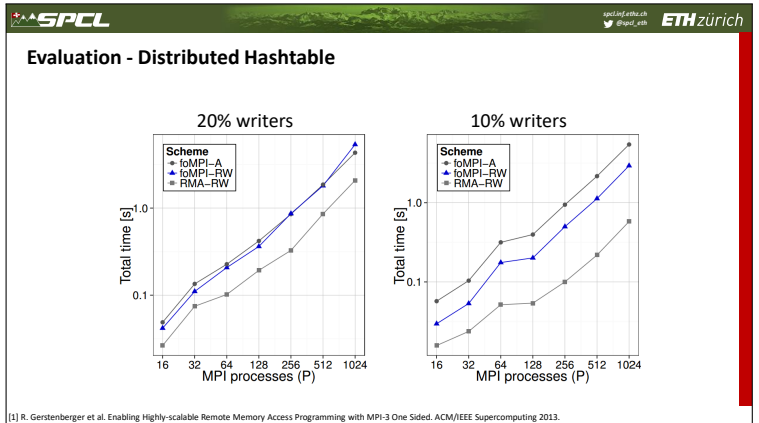
[1] P. Schmid, M. Besta, TH: High-Performance Distributed RMA Locks, ACM HPDC'16, best paper



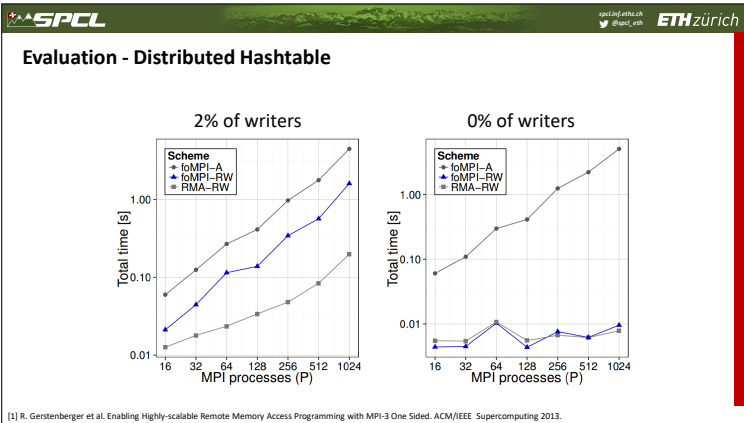
[1] R. Gerstenberger et al. Enabling Highly-scalable Remote Memory Access Programming with MPI-3 One Sided. ACM/IEEE Supercomputing 2013.



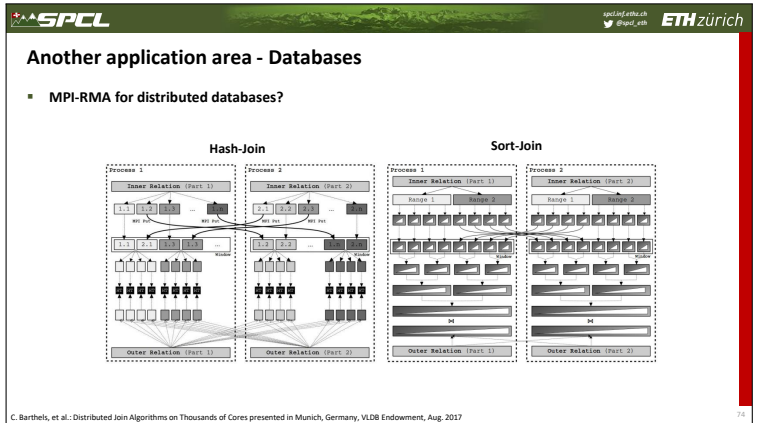
[1] R. Gerstenberger et al. Enabling Highly-scalable Remote Memory Access Programming with MPI-3 One Sided. ACM/IEEE Supercomputing 2013.



[1] R. Gerstenberger et al. Enabling Highly-scalable Remote Memory Access Programming with MPI-3 One Sided. ACM/IEEE Supercomputing 2013.



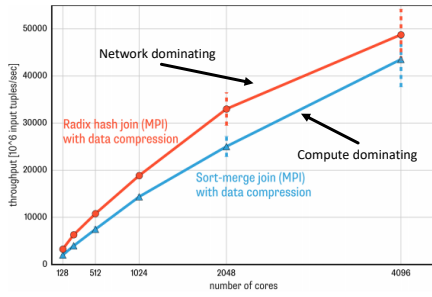
[1] R. Gerstenberger et al. Enabling Highly-scalable Remote Memory Access Programming with MPI-3 One Sided. ACM/IEEE Supercomputing 2013.



C. Barthels, et al.: Distributed Join Algorithms on Thousands of Cores presented in Munich, Germany, VLDB Endowment, Aug. 2017

Another application area - Databases

- MPI-RMA for distributed databases on Piz Daint



Another application area - Databases

- MPI-RMA for distributed databases on Piz Daint

