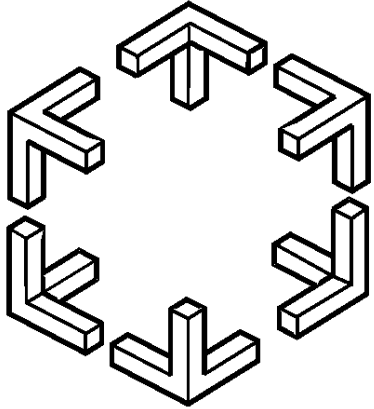




HÖGSKOLAN I BORÅS
VETENSKAP FÖR PROFESSION



Distributed Computing and Systems

Chalmers university of technology

Brushing the Locks out of the Fur:
A Lock-Free Work Stealing
Library Based on Wool

Håkan Sundell

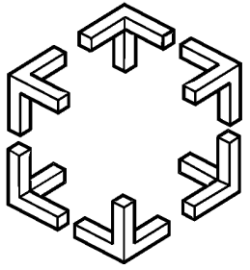
University College of Borås

Parallel Scalable Solutions AB

Philippas Tsigas

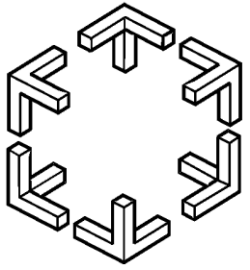
Chalmers University of Technology





Outline

- Synchronization of Shared Data
- Task Parallelism Library
 - Light-Weight
 - Previous Work
- The Wool Library
 - Architecture
 - Synchronization
- Wool with Lock-Free Synchronization
- Experiments
- Conclusions

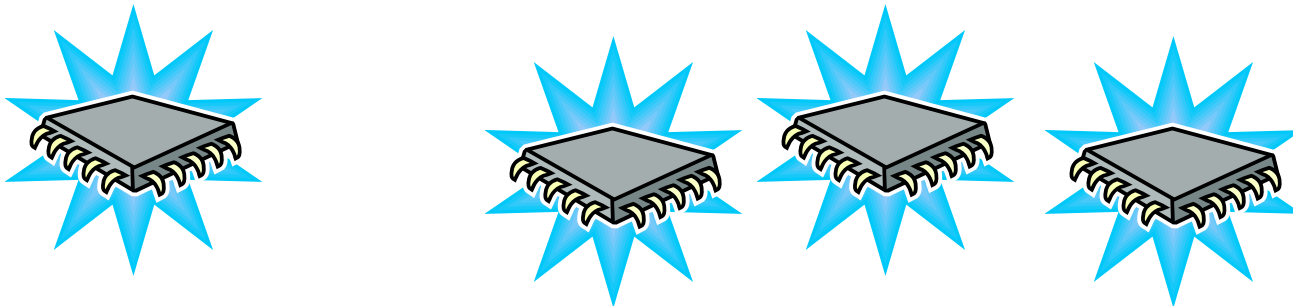


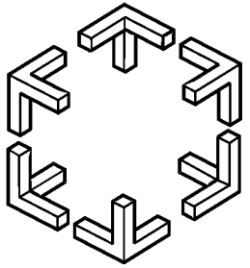
Parallel (e.g. Multi-threaded) Software

- Programs consist of many tasks (threads)



- That execute on one or more (logical) processors

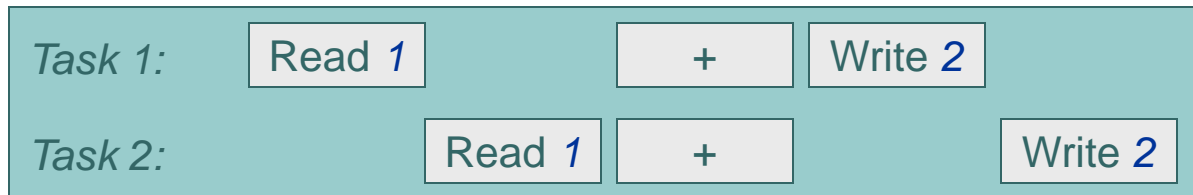




Critical Sections

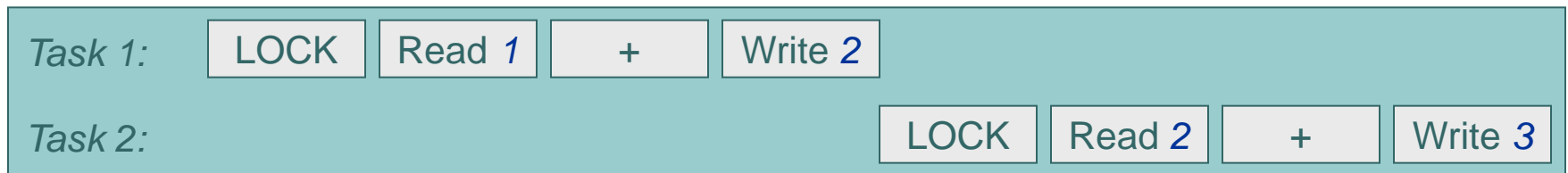
- Problem: operations on shared variables in programming languages are not atomic.

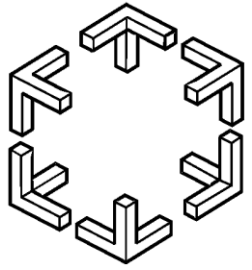
counter=counter+1; = Read + Write



! counter=2, but should be 3!

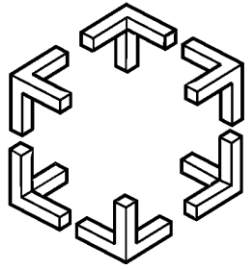
- Straightforward solution: Apply mutual exclusion





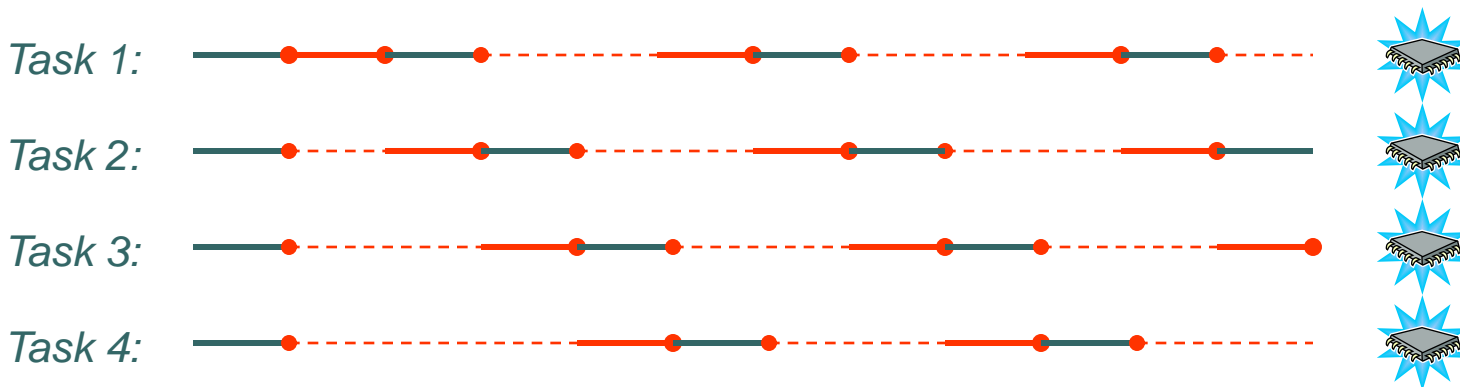
Critical Sections + Scheduling

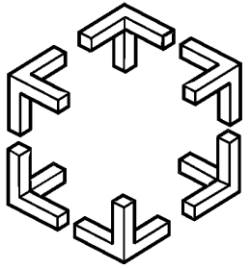
- **Blocking.** More advanced and pessimistic schedulability analysis.
- **Deadlocks.** Reduced fault-tolerance, if one task fails, other (even all) might also fail.
- **Priority Inversion.** Tasks might not execute with the proper priority even though it was set. Deadlines might be missed.



Critical Sections + Multiprocessors

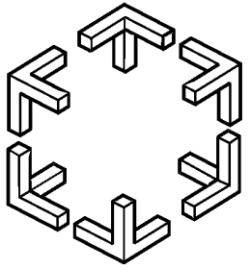
- **Reduced Parallelism.** Several tasks with overlapping critical sections will cause waiting processors to go idle.





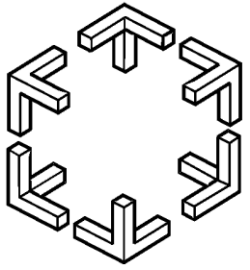
Avoid Critical Sections!

- **Avoid Blocking.** Easier and more optimistic analysis, i.e. less hardware needed.
- **Avoid Deadlocks.** Increased fault-tolerance as failed tasks can not affect others to fail.
- **Avoid Priority Inversion.** Easier and more reliable analysis, and avoids complex and high-overhead solutions.
- **Increased Parallelism.** Increased overall performance, more optimistic analysis, i.e. less hardware needed.



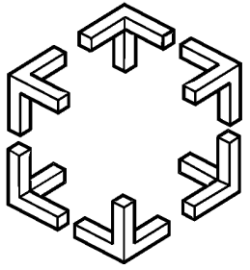
Non-Blocking Synchronization

- The key lies in how mutual exclusion (i.e. mutex, semaphore) is implemented in actual hardware (i.e. processors).
 - Atomic primitives in hardware can atomically update one memory word.
- Sophisticated solutions can exploit the same atomic primitives to support access to shared resources without locks, i.e. non-blocking.



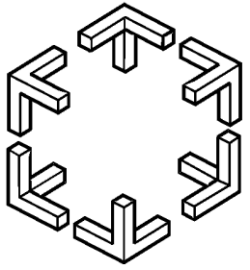
Non-Blocking Algorithms

- **Obstruction-Free.**
 - Guarantees progress in absence of contention.
 - Need extra module for contention management.
- **Lock-Free.**
 - Guarantees that always one operation is making progress.
 - Combined with scheduling information, schedulability analysis can be done.
- **Wait-Free.**
 - Guarantees that any operation will finish in a finite time.
 - Schedulability analysis can be done directly.



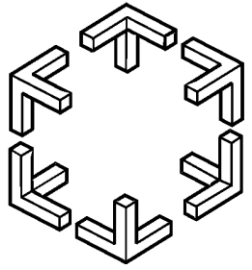
Task Parallelism Library

- Fine-grained parallelism is desired for achieving maximal speed-up.
- Spawning threads is expensive.
- Task-based approach:
 - Dynamically (recursively) spawn tasks.
 - Each *Task* contains a relatively small work-load.
 - Usually just a function call.
 - Side-effects are (usually) allowed.
- A *Task Parallelism Library* is usually a multi-threaded program (run-time system) together with a programming framework.



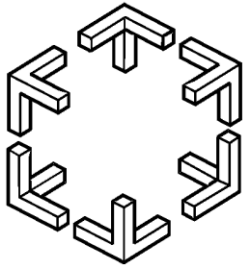
Fibonacci Example (Wool)

```
1 #include <stdio.h>
2 #include <stdlib.h>
3 #include "wool.h"
4
5 TASK_1 ( int , fib, int , n)
6 {
7     if(n<2) return n;
8     else {
9         int a,b;
10        SPAWN ( fib, n -2);
11        a= CALL ( fib, n -1);
12        b= SYNC ( fib);
13        return a+b;
14    }
15 }
16
17 TASK_2 ( int , main, int , argc,
18 char **, argv)
19 {
20     printf( "%d\n", CALL ( fib, atoi( argv[1])) );
21 }
```



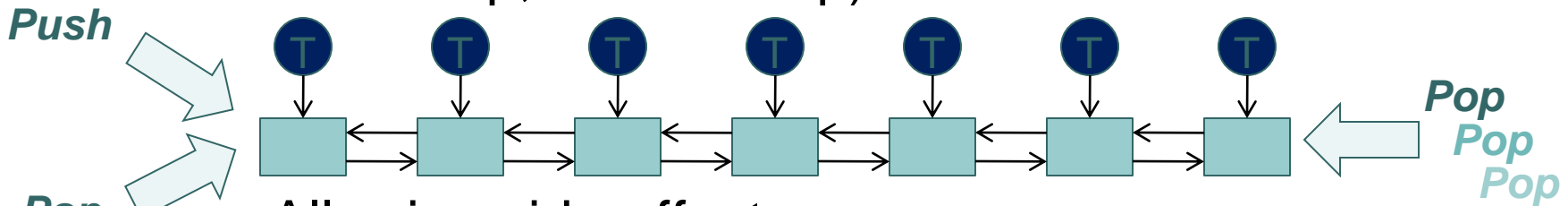
Light-Weight Task Management Libraries

- If considering a large number of tasks, handling costs becomes a bottleneck for efficiency (e.g. speed-up).
- Core issues:
 - Data structure in which the tasks are stored.
 - Strategy for load balancing between workers (i.e. threads)
 - Synchronization for moving tasks between workers and corresponding data structures in order to realize load balancing strategy.



Work-Stealing "Deque"

- Task objects stored in a "deque" (local Push/Pop, thieves Pop) data structure.

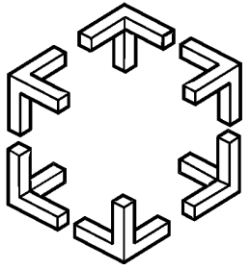


- Allowing side-effects:

- Arora et al. "*Thread scheduling for multiprogrammed multiprocessors*". 1998.
- ...
- Chase and Lev. "*Dynamic circular work-stealing deque*". 2005.

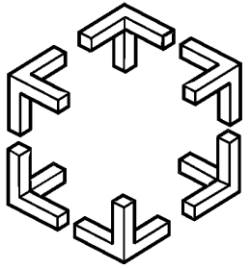
- Disallowing side-effects:

- Michael et al. "*Idempotent work stealing*". 2009.



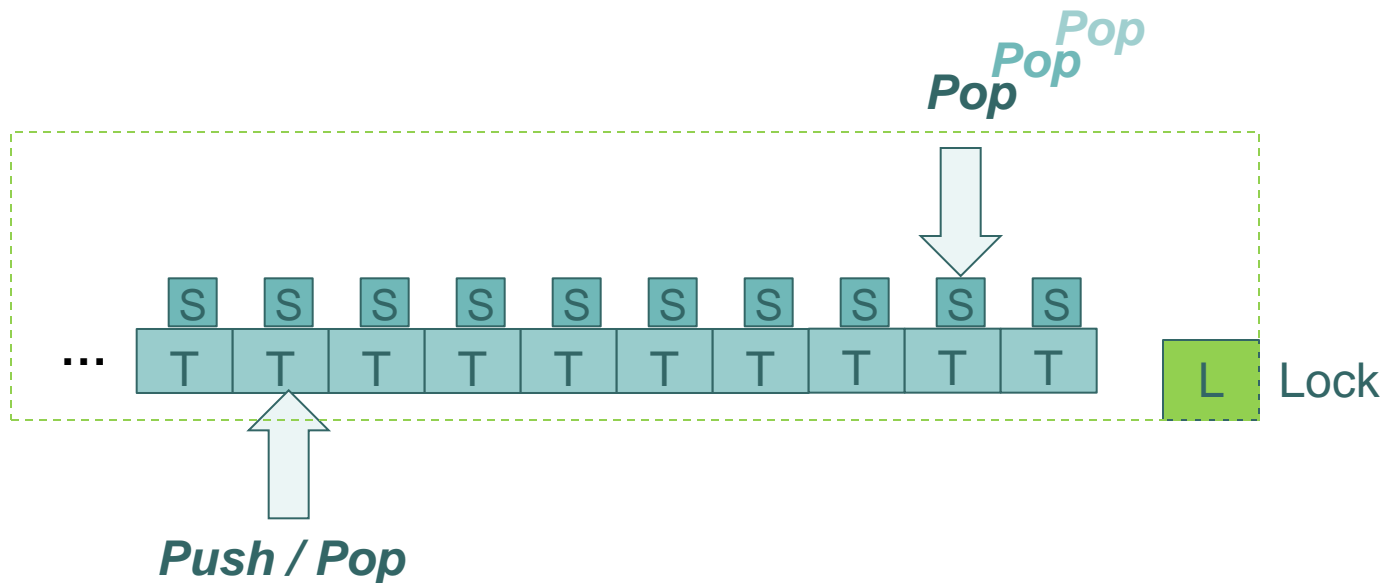
The Wool Library (v.0.1.1)

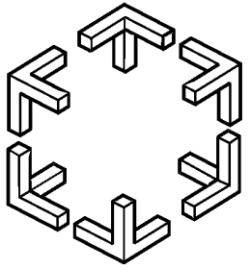
- Karl-Filip Faxén, "*Wool – A work stealing library*", MCC 2008.
- Really light-weight.
 - Simplified framework.
- Efficient synchronization
 - Tasks and "deque" data structure is the same ("collapsed layers").
 - Un-even synchronization
 - Optimizes for the average case.



Wool: Architecture

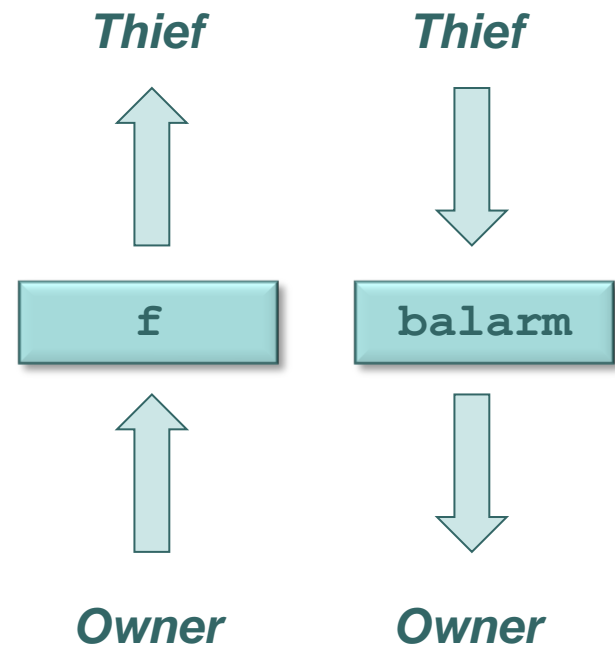
- Each worker has a large array of Tasks.
- Each Task includes stealing/availability status.

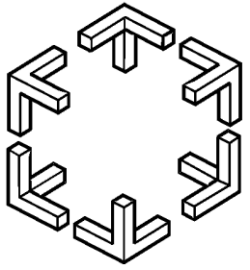




Modified Task data structure

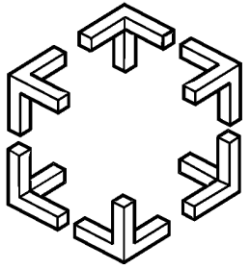
- Thieves synchronize through lock.
- Thief and owner synchronize through both `f` and `balarm`.





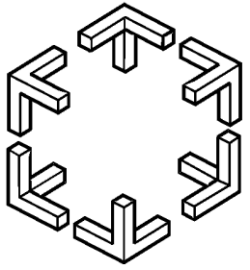
Wool: Stealers

```
1 bool steal( Worker *victim )
2 {
3     lock( victim->lck );
4     Task *t = victim->bot;
5     t->balarm = STOLEN;
6     memory_barrier();
7     if( t->f == INLINED ) {
8         unlock( victim->lck );
9         t->balarm = READY;
10        return false;
11    } else {
12        victim->bot++;
13        unlock( victim->lck );
14        ... // Run the task
15        memory_barrier();
16        t->balarm = DONE;
17        return true;
18    }
19 }
```



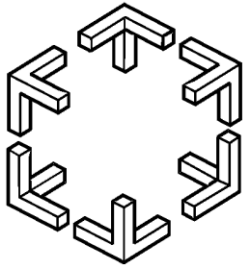
Wool: Task owners

```
21 void sync( Task *t )
22 {
23     t->f = INLINED;
24     memory_barrier();
25     if( t->balarm != READY ) {
26         // Wait for thief to fully decide
27         lock( self->lck );
28         if( t->balarm == READY ) {
29             unlock( self->lck );
30             ... // Run the task
31         } else {
32             unlock( self->lck );
33             ... // Wait for thief to finish
34             self->bot--;
35         }
36     }
37 }
```



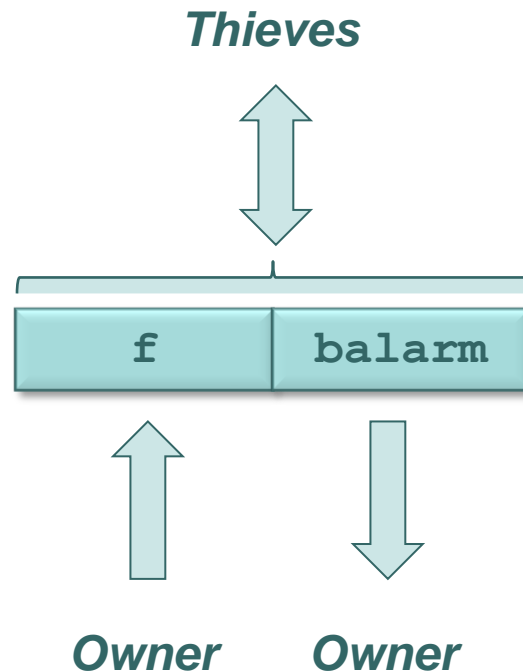
Lock-Free Approach: Atomic Primitives

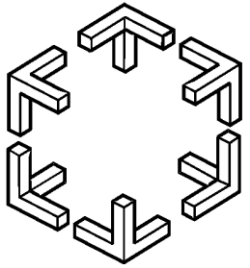
```
1 void FAA( int volatile *address, int number ) atomically do {
2     *address = *address + number;
3 }
4 //
5 bool CAS( int volatile *address, int oldvalue, int newvalue ) atomically do {
6     if( *address == oldvalue ) {
7         *address = newvalue;
8         return true;
9     }
10    else return false;
11 }
12 //
13 bool DWCAS( int volatile *address, int oldvalue1, int oldvalue2, int
newvalue1, int newvalue2) atomically do {
14     if( address[0] == oldvalue1 && address[1] == oldvalue2 ) {
15         address[0] = newvalue1;
16         address[1] = newvalue2;
17         return true;
18     }
19     else return false;
20 }
```



Modified Task data structure

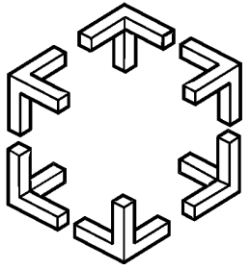
- Place both `f` and `balarm` into same double-word.



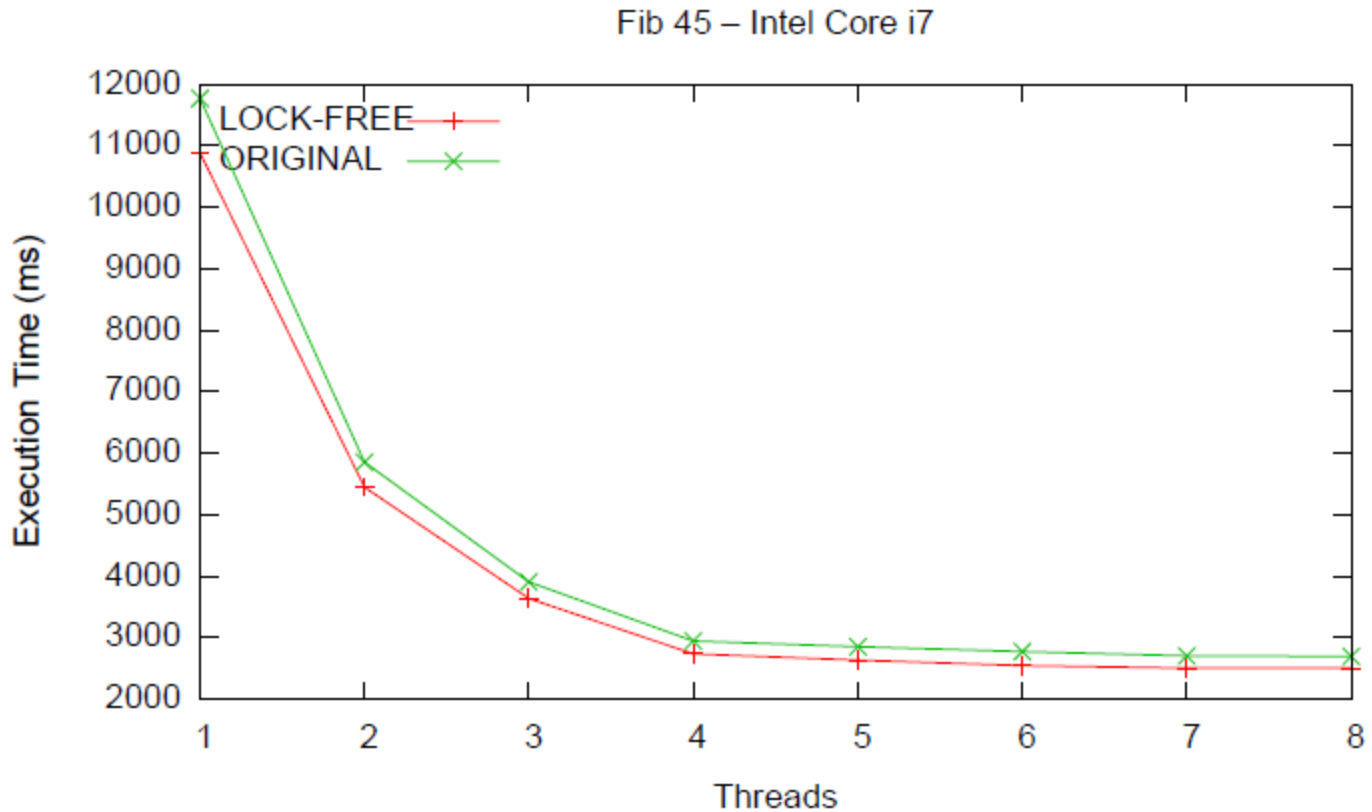


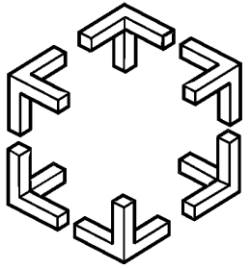
Lock-Free Wool

```
1 bool steal( Worker *victim )
2 {
3     Task *t = victim->bot;
4     f = t->f;
5     if( f != INLINED && DWCAS( &t->f, f, READY, f, STOLEN ) ) {
6         FAA( &victim->bot, 1 );
7         ... // Run the task
8         memory_barrier();
9         t->balarm = DONE;
10        return true;
11    }
12    else return false;
13 }
14
15 void sync( Task *t )
16 {
17     t->f = INLINED;
18     memory_barrier();
19     if( t->balarm == READY ) {
20         ... // Run the task
21     }
22     else {
23         ... // Wait for thief to finish
24         FAA( &self->bot, -1 );
25     }
26 }
```

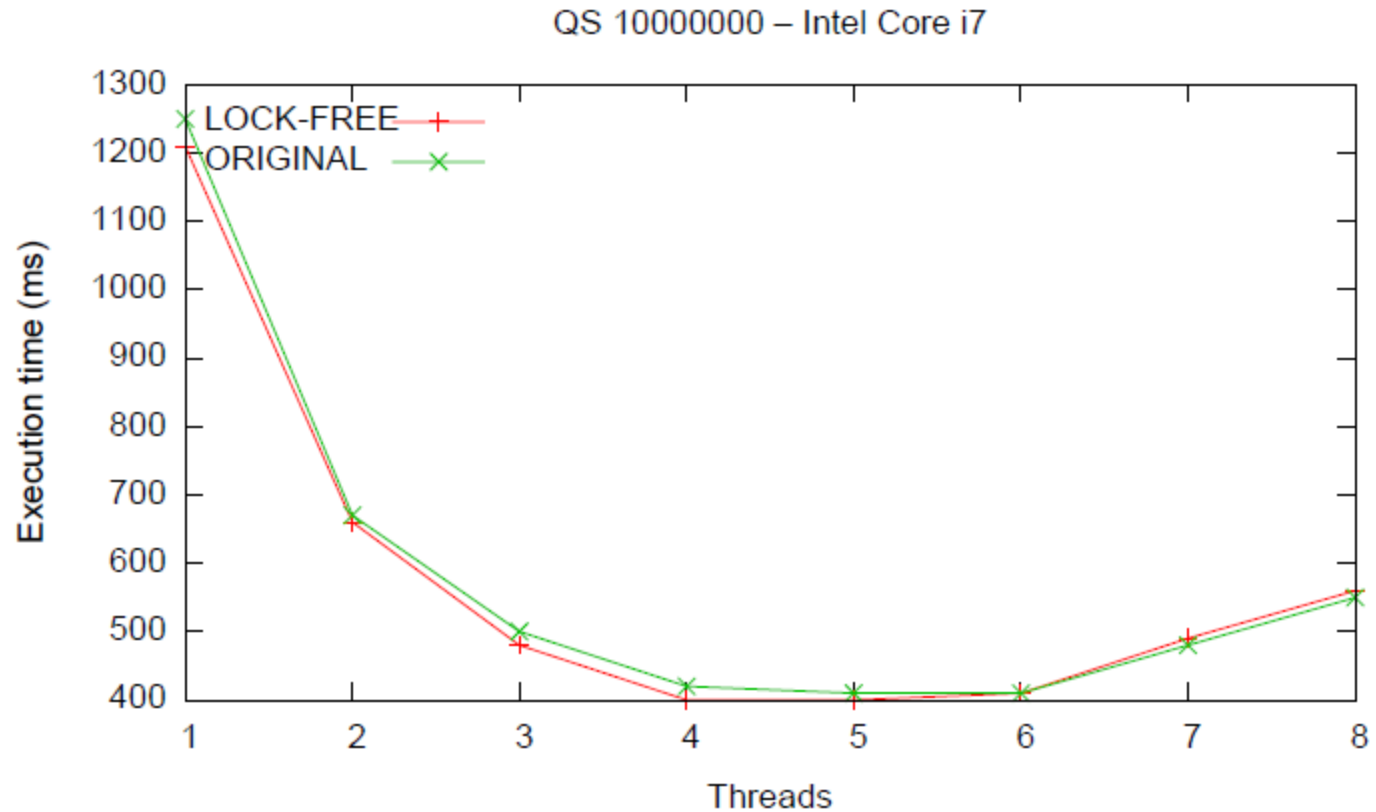


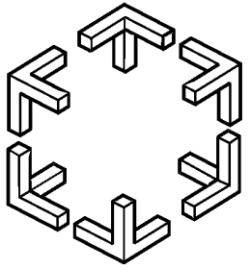
Experiments (Intel core i7): Fibonacci, fully expanded spawn-tree





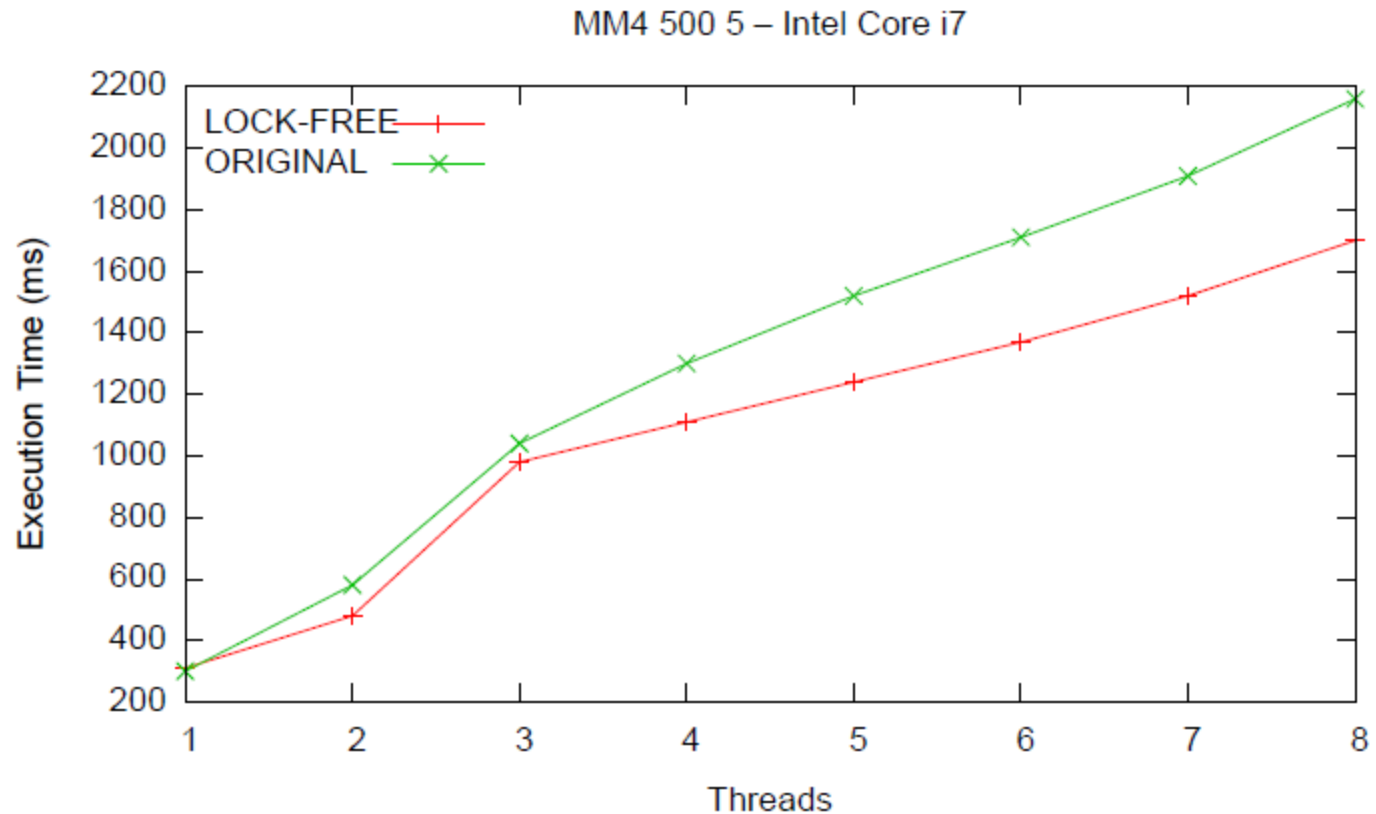
Experiments: Quicksort using shared memory

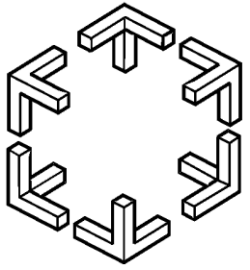




Experiments:

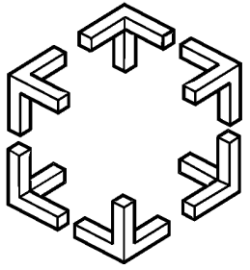
Matrix multiplication using "parallel for"





Conclusions

- Although Wool was highly optimized, adding Lock-Free synchronization could improve (absolute) performance.
- "Un-even" synchronization is an interesting technique for optimizing the average case.
- Task "size" is significant for performance.
 - "parallel for" is especially sensitive for task size, due to relatively high overhead.



Questions?

Thank You for listening!

www.pss-ab.com

www.adm.hb.se/~hsu

www.cse.chalmers.se/~tsigas