

# Concurrent Systems

## Exercise 05 – Non-Blocking Synchronisation

Stefan Reif

2019-01-14



# Agenda

---

Non-Blocking Data Structures

Lock-free Events

Worker Thread Notification

Assignment 5



# Agenda

---

## Non-Blocking Data Structures

Lock-free Events

Worker Thread Notification

Assignment 5



# Non-blocking queues

- Many queue algorithms have been published
  - ... through scientific channels, e.g. [2, 3, 6]
  - ... through non-scientific channels, e.g. [4]
- Minimal interface
  - At least `enqueue()` and `dequeue()`
  - Sometimes `is_empty()`, `make_empty()`
- Differences between lock-free queue algorithms [5]
  - Single-producer ↔ Multi-producer
  - Single-consumer ↔ Multi-consumer
  - Lock-free ↔ Wait-free
  - Nonintrusive ↔ intrusive
  - Node-based (“Unbounded”) ↔ Array-based (“bounded”) ↔ Hybrid
  - Consistency guarantees
  - ...



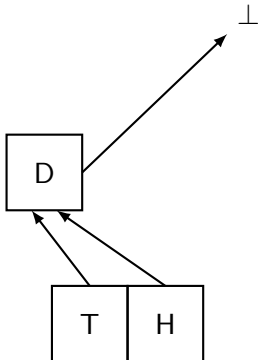
## Queue example

---

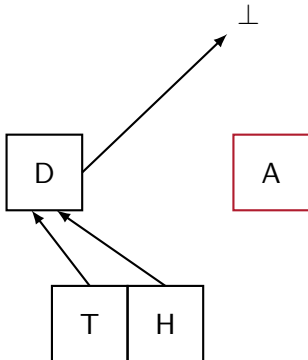
- Source code is online:  
<http://www.1024cores.net/home/lock-free-algorithms/queues/intrusive-mpsc-node-based-queue>
- Multi-producer Single-consumer (MPSC)
- Wait-free
- Intrusive
- Node-based, unbounded
- Inconsistent: not always connected



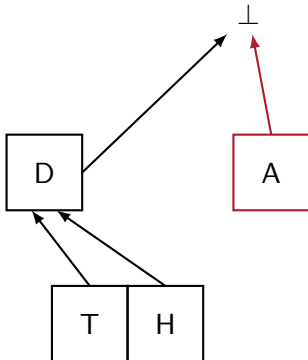
# Queue implementation



# Queue implementation

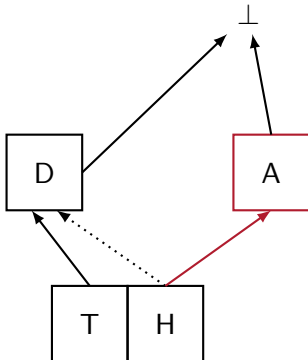


# Queue implementation

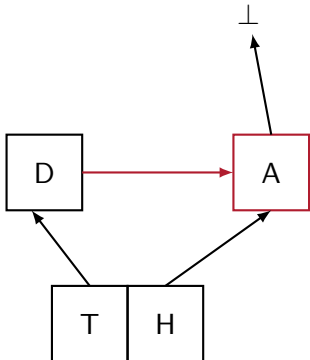




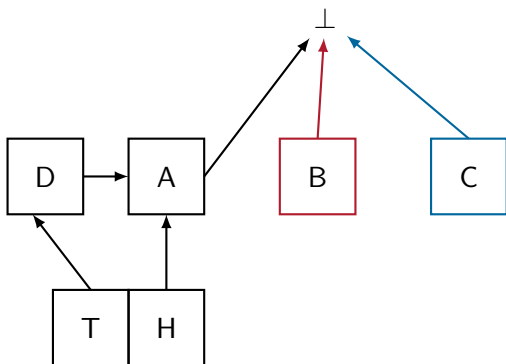
# Queue implementation



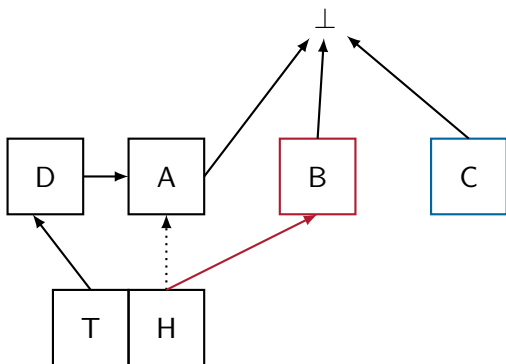
# Queue implementation



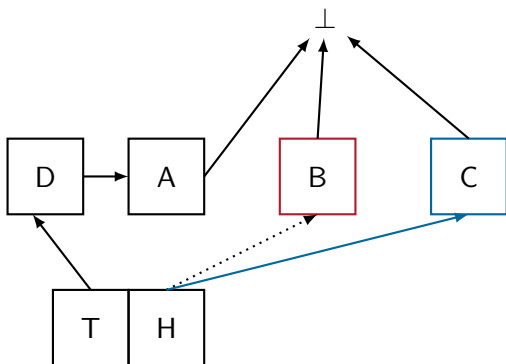
# Queue implementation



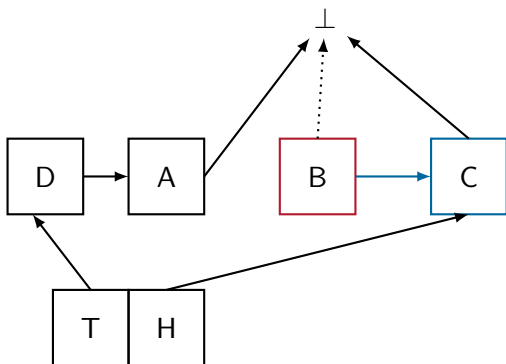
# Queue implementation



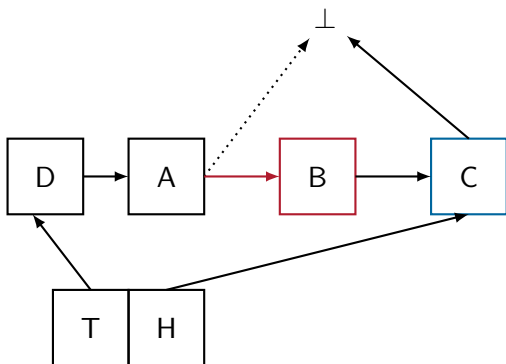
# Queue implementation



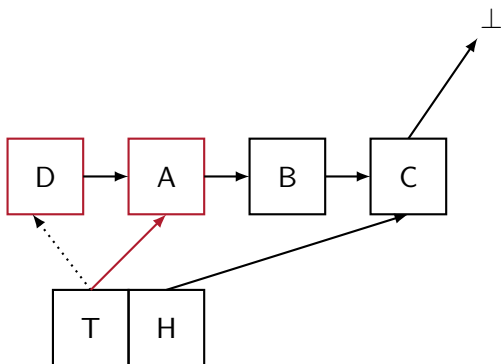
# Queue implementation



# Queue implementation

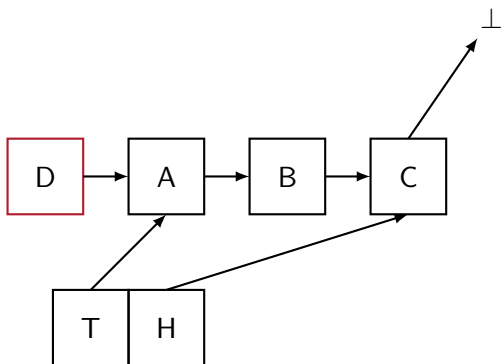


# Queue implementation

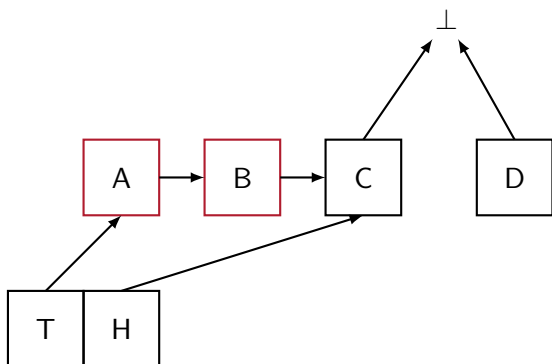




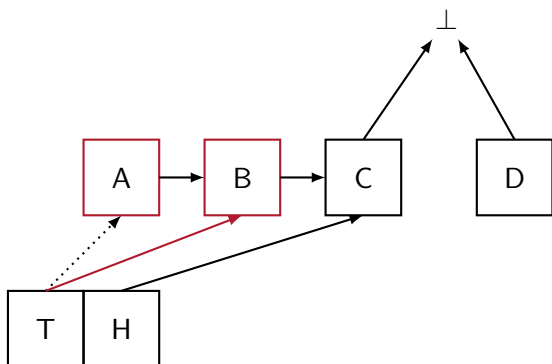
# Queue implementation



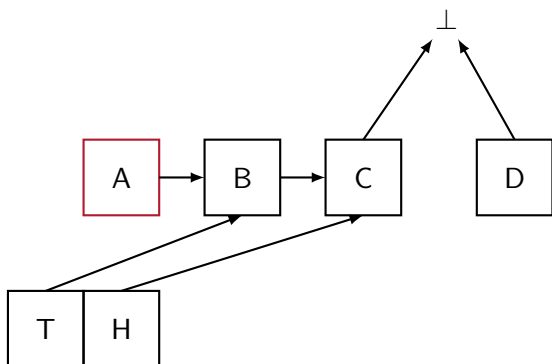
# Queue implementation



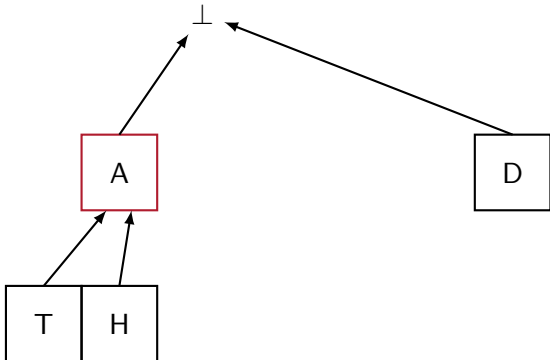
# Queue implementation



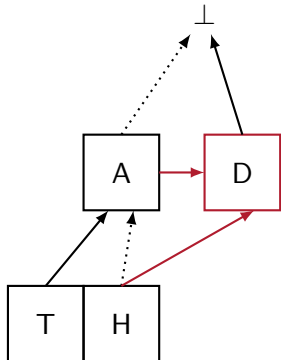
# Queue implementation



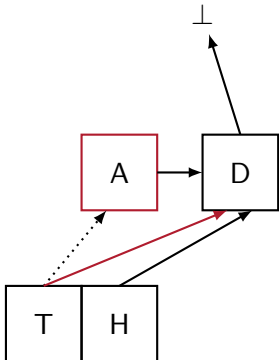
# Queue implementation



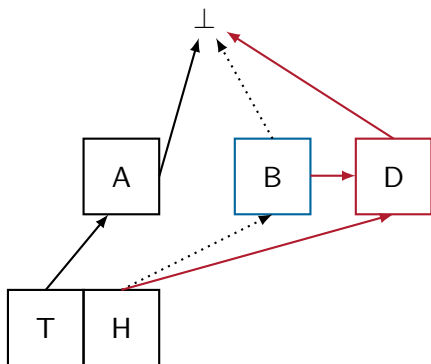
# Queue implementation



# Queue implementation



# Queue implementation





# Agenda

---

Non-Blocking Data Structures

Lock-free Events

Worker Thread Notification

Assignment 5



# Wait for a condition

---

- **Actively or Inactively (“passive”)**
  - Does the thread release the processor core?
  - Different performance, energy characteristics, ...
  
- **Active waiting ...**
  - Is usually fast, except for cache contention
  - Causes relatively high energy consumption
  - Probably wastes CPU time
  
- **Passive waiting ...**
  - Tends to be more energy-efficient
  - Supports many-threaded system
    - More threads than cores



# Eager event

- No thread management data

```
1 typedef struct event {  
2     atomic_int state;  
3 } event;
```

- Eager check for condition

```
1 void event_await(event *e) {  
2     while (PENDING == atomic_load(&e->state)); ← acquire  
3 }
```

- Notification by assignment


```
1 void event_signal(event *e) {  
2     atomic_store(&e->state, PAST); ← release  
3 }
```



# Passive waiting done wrong

## ■ How not to wait for a condition

```
1 void event_await(event *e) {  
2     while (PENDING == atomic_load(&e->state))  
3         sleep();  
4 }
```



## ■ Lost-wakeup problem

- Critical section between `load()` and `sleep()`
- Concurrent notification is lost
- Need for atomic `test_and_sleep()` operation



# Passive waiting

- How to wait for a condition passively?
  - Data race between condition check and sleep instruction
  - Atomic operation required
    - x86: `monitor`, `mwait` (privileged [1])
    - Linux: `futex()` system call

## Linux futex interface

```
1 syscall(SYS_futex, addr, FUTEX_WAIT, exp, &timeout);
2 syscall(SYS_futex, addr, FUTEX_WAIT, exp, NULL);
3 syscall(SYS_futex, addr, FUTEX_WAKE, num_threads);
4 ...
```

- Problem of spurious wakeup
  - Threads can wake up despite not being signalled
  - Need for re-evaluation of waiting conditions



# Linux futex semantics

- Wait until a variable changes

```
1 void futex_wait(int *addr, int exp) {  
2     if (*addr == exp)  
3         sleep_on(addr);  
4 }
```

- Address-dependent wakeup function

```
1 void futex_wake(int *addr) {  
2     for (Thread *t ← Threads)  
3         if (t->waiting_on == addr)  
4             wakeup(t);  
5 }
```



# Linux futex semantics

## ■ Wait until a variable changes

```
1 void futex_wait(int *addr, int exp) {  
2     if (*addr == exp) ] atomic  
3         sleep_on(addr);  
4 }
```

## ■ Address-dependent wakeup function

```
1 void futex_wake(int *addr) {  
2     for (Thread *t ← Threads)  
3         if (t->waiting_on == addr)  
4             wakeup(t);  
5 }
```



# Linux futex usage pattern

---

## ■ futex usage pattern

```
int i = atomic_load(&v);    atomic_store(&v, j);  
futex_wait(&v, i);        futex_wake(&v);
```

## ■ Critical Section

- Between `atomic_load()` and `futex_wait()`
- Concurrent signal must be caught





# Linux futex usage pattern

## ■ futex usage pattern

```
int i = atomic_load(&v);  
futex_wait(&v, i);
```

```
atomic_store(&v, j);  
futex_wake(&v);
```

## ■ Critical Section

- Between `atomic_load()` and `futex_wait()`
- Concurrent signal must be caught

## ■ futex overlapping patterns

```
int i = atomic_load(&v);  
futex_wait(&v, i);
```

```
atomic_store(&v, j);  
futex_wake(&v);
```



# Linux futex usage pattern

## ■ futex usage pattern

```
int i = atomic_load(&v);  
futex_wait(&v, i);
```

```
atomic_store(&v, j);  
futex_wake(&v);
```

## ■ Critical Section

- Between `atomic_load()` and `futex_wait()`
- Concurrent signal must be caught

## ■ futex overlapping patterns

```
int i = atomic_load(&v);  
  
futex_wait(&v, i);
```

```
atomic_store(&v, j);  
  
futex_wake(&v);
```



# Linux futex usage pattern

## ■ futex usage pattern

```
int i = atomic_load(&v);  
futex_wait(&v, i);
```

```
atomic_store(&v, j);  
futex_wake(&v);
```

## ■ Critical Section

- Between `atomic_load()` and `futex_wait()`
- Concurrent signal must be caught

## ■ futex overlapping patterns

```
int i = atomic_load(&v);  
  
futex_wait(&v, i);
```

```
atomic_store(&v, j);  
futex_wake(&v);
```



## Simple futex-based event

- Glibc has no direct `SYS_futex` wrapper

```
1 #define futex(...) syscall(SYS_futex, __VA_ARGS__)
```

- wait-operation using a futex

```
1 void event_await(event *e) {  
2     while (PENDING == atomic_load(&e->state))  
3         futex((int *) &e->state,  
4             FUTEX_WAIT, PENDING, NULL);  
5 }
```

- signal-operation on a futex

```
1 void event_signal(event *e) {  
2     atomic_store(&e->state, PAST);  
3     futex((int *) &e->state, FUTEX_WAKE, INT_MAX);  
4 }
```



## Simple futex-based event

- Glibc has no direct `SYS_futex` wrapper

```
1 #define futex(...) syscall(SYS_futex, __VA_ARGS__)
```

- wait-operation using a futex

```
1 void event_await(event *e) {  
2     while (PENDING == atomic_load(&e->state))  
3         futex((int *) &e->state,  
4             FUTEX_WAIT, PENDING, NULL);  
5 }
```

- signal-operation on a futex

```
1 void event_signal(event *e) {  
2     atomic_store(&e->state, PAST);  
3     futex((int *) &e->state, FUTEX_WAKE, INT_MAX);  
4 }
```



# ABA-safe futex-based event

- Use larger state space

```
1  #define is_pending(x) (!((x)&1))
```

- Conditionally wait on variable change

```
1  void event_await(event *e) {  
2      int val = atomic_load(&e->state);  
3      if (!is_pending(val))  
4          return;  
5      do {  
6          futex((int *) &e->state, FUTEX_WAIT, val, NULL);  
7      } while (atomic_load(&e->state) == val);  
8  }
```



## ■ Conditionally increment state variable

```
1 void event_signal(event *e) {
2     int val = atomic_load(&e->state);
3     if (!is_pending(val)) return;
4     if (atomic_compare_exchange_strong(&e->state,
5                                       &val, val + 1))
6         futex((int *) &e->state, FUTEX_WAKE, INT_MAX);
7 }
8 void event_reset(event *e) {
9     int val = atomic_load(&e->state);
10    if (is_pending(val)) return;
11    atomic_compare_exchange_strong(&e->state,
12                                  &val, val + 1);
13 }
```



# Agenda

---

Non-Blocking Data Structures

Lock-free Events

Worker Thread Notification

Assignment 5





# Worker thread waiting condition

---

- Similar to event implementation
  - Unilateral synchronisation
  - Correctness depends on job queue consistency
- Implementation variants
  - sensitivity flag
    - Wait while sensitive  $\neq 0$
  - load counter
    - Wait for load  $\neq 0$
  - generation counter
    - Wait for generation increment
  - ...



## Worker thread waiting condition – sensitivity flag

### ■ Get next Job

```
1 job *next() {
2     atomic_store(&sensitive, 1); // single consumer!
3     job *n = dequeue();
4     if (n)
5         return atomic_store(&sensitive, 0), n;
6     futex((int *) &sensitive, FUTEX_WAIT, 1, NULL);
7     return next();
8 }
```

### ■ Submit a job

```
1 void submit(job *j) {
2     enqueue(j); int old = 1;
3     if (atomic_compare_exchange_strong(&sensitive,
4                                         &old, 0))
5         futex((int *) &sensitive, FUTEX_WAKE, 1);
```



## Worker thread waiting condition – load counter

### ■ Get next Job

```
1 job *next(){
2     while (!atomic_load(&load))
3         futex((int *) &load, FUTEX_WAIT, 0, NULL);
4     atomic_fetch_sub(&load, 1); // single consumer!
5     return dequeue();
6 }
```

### ■ Submit a job

```
1 void submit(job *j) {
2     enqueue(j);
3     if (!atomic_fetch_add(&load, 1))
4         futex((int *) &sensitive, FUTEX_WAKE, 1);
5 }
```



## Worker thread waiting condition – generation counter

### ■ Get next Job

```
1 job *next() {
2     int now = atomic_load(&generation);
3     job *n = dequeue();
4     if (n)
5         return n;
6     futex((int *) &generation, FUTEX_WAIT, now, NULL);
7     return next();
8 }
```

### ■ Submit a job

```
1 void submit(job *j) {
2     enqueue(j);
3     atomic_fetch_add(&generation, 1);
4     futex((int *) &generation, FUTEX_WAKE, INT_MAX);
5 }
```



# Agenda

---

Non-Blocking Data Structures

Lock-free Events

Worker Thread Notification

Assignment 5



# Assignment 5

---

- Rewrite LWT without locks
  - Remove all scheduler locks ...
  - ... but still provide high-level synchronisation for LWT threads
- Tackle the ABA problem
  - Use version counters where needed
- All LWT synchronisation mechanisms use futexes internally
  - Internal `lwt_futex` data structure
  - Used by mutexes, condition variables, semaphore, `join()`, ...
  - Only one mechanism for passive waiting / thread blocking
- Logical cores use the futex system call for passive waiting
  - Sleep when no `lwt_thread` is ready



# Reference List I

---

- [1] INTEL CORPORATION (Hrsg.):  
*Intel Architecture Software Developer's Manual.*  
Santa Clara, California, USA: Intel Corporation, 2016
- [2] MICHAEL, M. M. ; SCOTT, M. L.:  
Simple, Fast, and Practical Non-blocking and Blocking Concurrent Queue Algorithms.  
In: *Proceedings of the 15th Annual ACM Symposium on Principles of Distributed Computing (PODC)*, ACM Press, 1996, S. 267–275
- [3] MORRISON, A. ; AFEK, Y. :  
Fast Concurrent Queues for x86 Processors.  
In: *Proceedings of the 18th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP)*, ACM Press, 2013, S. 103–112
- [4] VYUKOV, D. :  
*Intrusive MPSC node-based queue.*  
<http://www.1024cores.net/home/lock-free-algorithms/queues/intrusive-mpsc-node-based-queue>,
- [5] VYUKOV, D. :  
*Producer-Consumer Queues.*  
<http://www.1024cores.net/home/lock-free-algorithms/queues>,



## Reference List II

---

- [6] YANG, C. ; MELLOR-CRUMMEY, J. :  
A Wait-free Queue As Fast As Fetch-and-add.  
In: *Proceedings of the 21st ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP)*, ACM Press, 2016, S. 16:1–16:13

