## Concurrent Systems

Nebenläufige Systeme

X. Basics of Non-Blocking Synchronisation

Wolfgang Schröder-Preikschat

January 21, 2021



#### Outline

#### Preface

Constructional Axis

Exemplification

Transition

Transactional Axis





Preface

Constructional Axis

General

Exemplification

Transition

Transactional Axis

General

Case Study

Summary



CS (WS 2020/21, LEC 10)

# Subject Matter

- discussion on abstract concepts of synchronisation without lockout of critical action sequences of interacting processes (cf. [5])
  - attribute "non-blocking" here means abdication of mutual exclusion as the conventional approach to protect critical sections
  - note that even a "lock-free" solution may "block" a process from making progress, very well!
- develop an intuition for the dependency on process interleaving and contention rate when arguing on performance issues
  - what in case of high and what else in case of low contention?
  - what is the exception that proves the rule?
- follow suit, an explanation of the two-dimensional characteristic of non-blocking synchronisation is given
  - on the one hand, constructional, on the other hand, transactional
  - with different weighting, depending on the use case and problem size
- not least, engage in sort of tolerance to races of interacting processes while preventing faults caused by race conditions...



CS (WS 2020/21, LEC 10)

# Tolerance is the suspicion that the other person just might be right. 1



Source: Commemorative plaque, Berlin, Bundesallee 79



<sup>1</sup>(Ger.) Toleranz ist der Verdacht, dass der andere Recht hat.

CS (WS 2020/21, LEC 10)

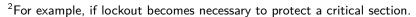
# Reentrancy

(Ger.) Eintrittsinvarianz

#### Definition

A program is re-entrant (Ger. ablaufinvariant) if, at execution time, its sequence of actions tolerates self-overlapping operation.

- those programs can be re-entered at any time by a new process, and they can also be executed by simultaneous processes
  - the latter is a logical consequence of the former: full re-entrant
  - but the former does not automatically imply the latter<sup>2</sup>
- originally, this property was typical for an **interrupt handler**, merely, that allows for nested execution—recursion not unresembling
  - each interrupt-driven invocation goes along with a new process
  - whereby the simultaneous processes develop **vertically** (i.e., stacked)
- generally, this property is typical for a large class of non-sequential programs whose executions may overlap each other
  - each invocation goes along with a new process, it must be "thread-safe"
  - whereby the simultaneous processes develop **horizontally**, in addition





#### Outline

Preface

Constructional Axis General Exemplification Transition

Transactional Axis



CS (WS 2020/21, LEC 10)

cf. [10, p. 22]

## Semaphore Revisited

devoid of an explicit protective shield all-embracing the semaphore implementation, i.e., the elementary operations P and V:

```
typedef struct semaphore {
                        /* value: binary or general */
    int gate;
    event t wait;
                        /* list of sleeping processes */
} semaphore_t;
```

- other than the original definition [1, p. 29], semaphore primitives are considered divisible operations in the following
  - merely single steps that are to be performed inside of these primitives are considered indivisible
  - these are operations changing the semaphore value (gate) and, as the case may be, the waitlist (wait)
  - but not any of these operations are secured by means of mutual exclusion at operating-system machine level
  - rather, they are safeguarded by falling back on ISA-level mutual exclusion in terms of atomic load/store or read-modify-write instructions



CS (WS 2020/21, LEC 10)

Constructional Axis - Exemplification

## Building Blocks for Barrier-Free Operation

- use of atomic (ISA-level) machine instructions for changing the semaphore value consistently (p. 11)
  - a TAS or CAS, resp., for a binary and a FAA for a general semaphore
  - instruction cycle time is bounded above, solely hardware-defined
  - wait-free [2, p. 124], irrespective of the number of simultaneous processes
- abolish abstraction in places, i.e., perform wait-action unfolding to prevent the lost-wakeup problem (p. 10)
  - make a process "pending blocked" before trying to acquire the semaphore
  - cancel that "state of uncertainty" after semaphore acquirement succeeded
  - wait- or lock-free [2, p. 142], depending on the waitlist interpretation
- accept dualism as to the incidence of processing states, i.e., tolerate a "running" process being seemingly "ready to run" (p. 12)
  - delay resolving until some process is in its individual idle state
  - have also other processes in charge of clearing up multiple personality
  - wait-free, resolution produces background noise but is bounded above
- forgo dynamic data structures for any type of waitlist or synchronise them using **optimistic concurrency control** (p. 16ff.)



CS (WS 2020/21, LEC 10)

Constructional Axis - Exemplification

#### **Atomic Machine Instructions** differences to [10, p. 24/25]

load/store-based implementation for a binary semaphore:

```
inline bool avail(semaphore t *sema) {
    return CAS(&sema->gate, 1, 0);
}
```

■ both *lodge* and *unban* remain unchanged

enumerator-based implementation for a general semaphore:

```
inline int lodge(semaphore t *sema) {
    return FAA(&sema->gate, -1);
}
inline bool unban(semaphore_t *sema) {
    return FAA(&sema->gate, +1) < 0;
```

- avail remains unchanged
- note that both variants are insensitive to simultaneous processes
  - due to **indivisible operations** for manipulation of the semaphore value



## Wait-Action Unfolding

cf. [10, p. 23]

```
void prolaag(semaphore t *sema) {
       catch(&sema->wait);
                                /* expect notification */
       lodge(sema);
                                /* raise claim to proceed */
                                /* check for process delay */
       when (!avail(sema))
           coast():
                                /* accept wakeup signal */
       clean(&sema->wait);
                                /* forget notification */
   }
   void verhoog(semaphore t *sema) {
       if (unban(sema))
                               /* release semaphore */
10
           cause(&sema->wait); /* notify wakeup signal */
11
  }
12
```

implementation in the shape of a non-sequential program:

- 2 show interest in the receive of a notification to continue processing
- 3/4 draw on walkover, bethink and, if applicable, watch for notification
  - 5 either suspend or continue execution, depending on notification state
  - 6 drop interest in receiving notifications, occupy resource
  - 10 deregulate "wait-and-see" position above (I. 4), check for a sleeper
- 11 send notification to interested and, maybe, suspended processes



CS (WS 2020/21, LEC 10)

Constructional Axis - Exemplification

10

## **Dualism**

a process being in "running" state and, as the case may be, at the same time recorded on the waitlist of "ready to run" peers inline void catch(event t \*this) {

```
process t *self = being(ONESELF);
    self -> state |= PENDING;
                                     /* watch for event */
    apply(self, this);
                                     /* enter waitlist */
}
inline void clean(event t *this) {
    elide(being(ONESELF), this);
                                     /* leave waitlist */
}
```

- 3 prepares the "multiple personality" process to be treated in time
- 4 makes the process amenable to "go ahead" notification (p. 10, l. 11)
- 8 excludes the process from potential receive of "go ahead" notifications
- treatment of "multiple personality" processes is based on division of labour as to the different types of waitlist (cf. p. 34)
  - "ready" waitlist, the respective idle process of a processor (p. 33)
  - "blocked" waitlist, the semaphore increasing or decreasing process

11

catch of a "go ahead" event is by means of a per-process latch • i.e., a "sticky bit" holding member of the process control block (PCB) inline int coast() { /\* latch event \*/ stand(); return being(ONESELF)->merit; /\* signaller pid \*/ } int cause(event\_t \*this) { process\_t \*next; int done = 0; for (next = being(0); next < being(NPROC); next++)</pre> 10 if (CAS(&next->event, this, 0)) 11 done += hoist(next, being(ONESELF)->name); 12 13 14 return done; 15

11 • recognise willingness to catch a signal and continue execution

12 • notify "go ahead", pass own identification, and ready signallee

CS (WS 2020/21, LEC 10)

Constructional Axis – Exemplification

13

#### Outline

Preface

Constructional Axis

Exemplification

Transition

Transactional Axis

General

Case Study



#### A Means to an End...

- non-blocking synchronisation spans two dimensions of measures in the organisation of a non-sequential program
  - i a constructional axis, as was shown with the semaphore example, and
  - ii a transactional axis, which is coming up in the next section
  - in many cases, particularly given complex software structures such as operating systems, the former facilitates the latter
    - the building blocks addressed and drafted so far are not just dedicated to operating systems, but are suited for any kind of "threads package"
    - although quite simple, they still disclose handicaps as to legacy software
- reservation towards the exploitation of non-blocking synchronisation originates much more from the constructional axis
  - synchronisation is a typical cross-cutting concern of software and, thus, use case of aspect-oriented programming (AOP, [3])
  - but the semaphore example shows that even AOP is not the loophole here
- but note that the **transactional axis** does not suggest effortlessness and deliver a quick fix to the synchronisation problem
  - appropriate solutions, however, benefit from a much more localised view



CS (WS 2020/21, LEC 10)

Constructional Axis - Transition

# Optimistic Concurrency Control

cf. [7, p. 15]

## Definition (acc. [4])

Method of coordination for the purpose of updating shared data by mainly relying on transaction backup as control mechanisms.

```
do
    read phase:
         save a private copy of the shared data to be updated;
         compute a new private data value based on that copy;
     validation and, possibly, write phase:
         try to commit the computed value as new shared data;
while commit failed (i.e., transaction has not completed).
```

- during the read phase, all writes take place only on local copies of the shared data subject to modification
- a subsequent validation phase checks that the changes as to those local copies will not cause loss of integrity of the shared data
- if approved, the final write phase makes the local copies global, i.e., commits their values to the shared data

15

## Transactional Computation

CAS-oriented approach, value-based, typical for CISC:

```
word t any;
                                     /* shared data */
{
    word_t old, new;
                                     /* own data */
    do new = compute(old = any);
                                     /* read */
    while (!CAS(&any, old, new));
                                     /* validate/write */
}
```

LL/SC-oriented approach, reservation-based, typical for RISC:

```
/* shared data */
word_t any;
    word_t new;
                                     /* own data */
    do new = compute(LL(&any));
                                     /* read */
    while (!SC(&anv, new));
                                     /* validate/write */
}
```

#### CAS recreated using LL/SC (cf. [8, p. 16])

Reading phase carried out simultaneously remains undetected...



CS (WS 2020/21, LEC 10)

**Unsynchronised Operations** 

Transactional Axis - General

devoid of synchronisation

basic **precondition**: an item to be stacked is not yet stacked/queued

```
inline void push dos(stack t *this, chain t *item) {
    item->link = this->head.link;
    this->head.link = item;
}
2 • copy the contents of the stack pointer to the item to be stacked
```

3 • update the stack pointer with the address of that item

```
inline chain_t *pull_dos(stack_t *this) {
    chain t *node;
    if ((node = this->head.link))
        this->head.link = node->link;
    return node;
```

- 7 memorise the item located at the stack top, if any
- 8 update the stack pointer with the address of the next item

## Data Type I

let a very simple dynamic data structure be object of investigation • modelling a **stack** in terms of a single-linked list:

```
typedef struct stack {
     chain_t head;
                           /* top of stack: list head */
} stack t:
• whereby a single list element is of the following structure:
typedef struct chain {
     struct chain *link; /* next list element */
} chain t;
```

- stack manipulation by pushing or pulling an item involves the update of a single variable, only: the "stack pointer"
- when simultaneous processes are allowed to interact by sharing that stack structure, the update must be an indivisible operation



17

CS (WS 2020/21, LEC 10)

Transactional Axis - Case Study

18

# Lock-Free Synchronised Operations

benefit from the precondition: an item to be stacked is "own data"

```
inline void push_lfs(stack_t *this, chain_t *item) {
       do item->link = this->head.link;
       while (!CAS(&this->head.link, item->link, item));
   }
   2 • copy the contents of the stack pointer to the item to be stacked
   3 • attempt to update the stack pointer with the address of that item
   inline chain_t *pull_lfs(stack_t *this) {
       chain_t *node;
       do if ((node = this->head.link) == 0) break;
       while (!CAS(&this->head.link, node, node->link));
11
       return node;
  }
```

- 8 memorise the item located at the stack top, if any
- 9 attempt to update the stack pointer with the address of the next item

19

#### Shallowness: ABA Problem

cf. [8, p. 14 & 36-37]

given a LIFO list (i.e., stack) of following structure:  $head \Leftrightarrow A \Leftrightarrow B \Leftrightarrow C$ 

- with head stored at location  $L_i$  shared by processes  $P_1$  and  $P_2$
- furthermore assume actual parameter this is pointing to  $L_i$

```
inline chain_t *pull_lfs(stack_t *this) {
      chain_t *node;
      do if ((node = this->head.link) == 0) break;
      while (!CAS(&this->head.link, node, node->link));
      return node;
6 }
```

assuming that the following sequence of actions will take place:

P<sub>1</sub> reads head item A followed by B on the list, gets delayed at line 4

• remembers node = A, but has not yet done CAS:  $head \diamondsuit A \diamondsuit B \diamondsuit C$ 

 $P_2$  pulls head item A from the list:

head \$ B \$ C

pulls head item B from the list:

head \$ C

■ pushes item A back to the list, now followed by C: head  $\Rightarrow$  A  $\Rightarrow$  C

 $P_1$  resumes, CAS realises head = A (followed by B):  $head \Rightarrow B \Rightarrow \bigcirc$ 

■ list state  $head \diamondsuit A \diamondsuit C$  as left behind by  $P_2$  is lost...



11

12 13

15 16 CS (WS 2020/21, LEC 10)

Transactional Axis - Case Study

21

## ABA Problem Tackled I

... as ugly as sin

```
typedef chain_t* chain_l;
                                    /* labelled pointer! */
#define BOX (sizeof(chain_t) - 1) /* tag-field mask */
inline void push_lfs(stack_t *this, chain_l item) {
   do ((chain_t *)raw(item, BOX))->link = this->head.link;
   while (!CAS(&this->head.link, ((chain_t *)raw(item, BOX))->link, tag(item, BOX)));
chain_l pull_lfs(stack_t *this) {
   chain_l node;
   do if (raw((node = this->head.link), BOX) == 0) break;
   while (!CAS(&this->head.link. node. ((chain t *)raw(node. BOX))->link)):
   return node;
```

aggravating side-effect of the solution is the loss of transparency

- the pointer in question originates from the environment of the critical operation (i.e., push and pull in the example here)
- tampered pointers must not be used as normal  $\sim$  derived type
- language embedding and compiler support would be of great help...

Hint (CAS vs. LL/SC)

The ABA problem does not exist with LL/SC!



## Approach to Solving the ABA Problem

```
workaround using a change-number tag as pointer label:
inline void *raw(void *item, long mask) {
    return (void *)((long)item & ~mask);
}
inline void *tag(void *item, long mask) {
    return (void *)
        ((long)raw(item, mask) | ((long)item + 1) & mask);
```

- alignment of the data structure referenced by the pointer is assumed
  - an **integer factor** in accord with the data-structure size (in bytes)
  - rounded up to the next **power of two**:  $2^N > sizeof(datastructure)$
- zeros the N low-order bits of the pointer—and discloses the tag field
- rather a kludge (Ger. Behelfslösung) than a clearcut solution<sup>3</sup>
  - makes ambiguities merely unlikely, but cannot prevent them
  - "operation frequency" must be in line with the **finite values margin**
- if applicable, attempt striving for problem-specific frequency control

<sup>3</sup>This also holds for DCAS when using a "whole word" change-number tag.



CS (WS 2020/21, LEC 10) Transactional Axis – Case Study

22

# ABA Problem Tackled II

... provided the processor plays along

```
same precondition (cf. p. 20): an item to be stacked is "own data"
inline void push lfs(stack t *this, chain t *item) {
    do item->link = LL(&this->head.link);
```

- while (!SC(&this->head.link, item)); }
- 2 copy the head pointer and make a reservation to his address
- update the head pointer if the reservation still exists

```
inline chain_t *pull_lfs(stack_t *this) {
       chain_t *node;
       do if ((node = LL(&this->head.link)) == 0) break;
       while (!SC(&this->head.link, node->link)):
10
11
       return node;
  }
12
```

- memorise the head pointer and make a reservation to his address
- 9 update the head pointer if the reservation still exists

#### Outline

Preface

Constructional Axis

General

Exemplification

Transition

Transactional Axis

Genera

Case Study

Summary



wosch CS (WS 2020/21, LEC 10)

Summary

mmary

25

#### Reference List I

[1] DIJKSTRA, E. W.:

Cooperating Sequential Processes / Technische Universiteit Eindhoven.

Eindhoven, The Netherlands, 1965 (EWD-123). -

Forschungsbericht. -

(Reprinted in *Great Papers in Computer Science*, P. Laplante, ed., IEEE Press, New York, NY, 1996)

[2] Herlihy, M.:

Wait-Free Synchronization.

In: ACM Transactions on Programming Languages and Systems 11 (1991), Jan., Nr. 1, S. 124–149

[3] KICZALES, G.; LAMPING, J.; MENDHEKAR, A.; MAEDA, C.; LOPES, C. V.; LOINGTIER, J.-M.; IRWIN, J.:

Aspect-Oriented Programming.

In: AKSIT, M. (Hrsg.); MATSUOKA, S. (Hrsg.): Proceedings of the 11th European Conference on Object-Oriented Programming (ECOOP'97) Bd. 1241, Springer-Verlag, 1997 (Lecture Notes in Computer Science). – ISBN 3-540-63089-9, S. 220-242

[4] Kung, H.-T.; Robinson, J. T.:

On Optimistic Methods for Concurrency Control.

In: ACM Transactions on Database Systems 6 (1981), Jun., Nr. 2, S. 213–226



#### Résumé

- non-blocking synchronisation  $\mapsto$  abdication of mutual exclusion
- systems engineering makes a two-dimensional approach advisable
  - the constructional track brings manageable "complications" into being
  - these "complications" are then subject to a *transactional track*

The latter copes with *non-blocking synchronisation* "in the small", while the former is a *state-machine outgrowth* using atomic instructions, sporadically, and enables barrier-free operation "in the large".

- no bed of roses, no picnic, no walk in the park—so is non-blocking synchronisation of reasonably complex simultaneous processes
  - but it constrains sequential operation to the absolute minimum and,
  - thus, paves the way for parallel operation to the maximum possible

#### Hint (Manyfold Update)

Solutions for twofold updates already are no "no-brainer", without or with special instructions such as CDS or DCAS. Major updates are even harder and motivate techniques such as **transactional memory**.



wosch CS (WS 2020/21, LEC 10)

Summ

2

#### Reference List II

[5] Moir, M.; Shavit, N.:

"Concurrent Data Structures".

In: Mehta, D. P. (Hrsg.); Sahni, S. (Hrsg.): Handbook of Data Structure and Applications.

CRC Press, Okt. 2004, Kapitel 47, S. 1–32

[6] Schröder-Preikschat, W.; Lehrstuhl Informatik 4 (Hrsg.):

Concurrent Systems.

FAU Erlangen-Nürnberg, 2014 (Lecture Slides)

[7] Schröder-Preikschat, W.:

Critical Sections.

In: [6], Kapitel 4

[8] Schröder-Preikschat, W.:

Elementary Operations.

In: [6], Kapitel 5

[9] Schröder-Preikschat, W.:

Monitor.

In: [6], Kapitel 8



Summary – Bibliography

#### Reference List III

```
[10] Schröder-Preikschat, W.:
    Semaphore.
    In: [6], Kapitel 7
```



CS (WS 2020/21, LEC 10)

Summary – Bibliography

29

# Receive-Side "Sticky Bit" Operations

cf. p. 13

a simple mechanism that allows a process to "latch onto" an event:

```
inline void shade(process_t *this) {
                                        /* clear latch */
       this->latch.flag = false;
  }
   inline void stand() {
       process_t *self = being(ONESELF);
       if (!self->latch.flag)
                                        /* inactive latch */
           block();
                                        /* relinguish... */
       shade(self);
                                        /* reset latch */
10
11
   inline void latch() {
       being(ONESELF)->state |= PENDING; /* watch for */
13
                                            /* & latch */
       stand();
14
15 }
```

- 8 either suspend or continue the current process (cf. p. 33)
  - was marked "pending" to catch a "go ahead" notification (cf. p.12)



#### **Propagate Notifications**

```
int cause(event_t *this) {
       chain t *item;
       int done = 0;
       if ((item = detach(&this->wait)))
           do done += hoist((process_t *)
                coerce(item, (int)&((process_t *)0)->event),
                    being(ONESELF)->name);
           while ((item = item->link));
       return done;
11
12
```

- variant relying on a **dynamic data structure** for the waitlist
  - 5 adopt the waitlist on the whole, indivisible, and wait-free
  - 6-8 notify "go ahead", pass own identification, and ready signallee
    - 7 pattern a dynamic type-cast from the chain\_t\* member event to the process t\* of the enclosing process structure (i.e., PCB)
    - 9 notify one process at a time, bounded above, N-1 times at worst



CS (WS 2020/21, LEC 10)

Addendum - Re-Entrant Operations

cf. p. 13

# Send-Side "Sticky Bit" Operations

non-blocking measure to signal a single process, one-time, and keep signalling effective, i.e., "sticky" (Ger. klebrig) until perceived<sup>4</sup>

```
inline void punch(process t *this) {
      if (!this->latch.flag) {
                                       /* inactive latch */
           this->latch.flag = true;
                                       /* activate it */
           if (this->state & PENDING) /* is latching */
                                       /* set readu */
               yield(this);
      }
  }
   inline int hoist(process t *next, int code) {
       next->merit = code:
                                     /* pass result */
       punch(next);
                                       /* send signal */
       return 1;
13 }
```

- 2–3 assuming that the PCB is not shared by simultaneous processes
  - otherwise, replace by TAS(&this->latch.flag) or similar
  - 5 makes the process become a "multiple personality", possibly queued

<sup>4</sup>In contrast to the signalling semantics of monitors (cf. [9, p. 8]).



31

CS (WS 2020/21, LEC 10) Addendum – Re-Entrant Operations

## Resolving Multiple Personality

```
cf. [10, p. 37]
```

```
void block() {
       process_t *next, *self = being(ONESELF);
3
       do {
                             /* ... become the idle process */
           while (!(next = elect(hoard(READY))))
                             /* enter processor sleep mode */
               relax():
       } while ((next->state & PENDING)
                                              /* clean-up? */
            && (next->scope != self->scope));
8
9
       if (next != self) { /* it's me who was set ready? */
10
           self->state = (BLOCKED | (self->state & PENDING));
11
           seize(next);
                              /* keep pending until switch */
12
13
       self -> state = RUNNING;
                                    /* continue cleaned... */
14
15 }
```

- a "pending blocked" process is still "running" but may also be "ready to run" as to its queueing state regarding the ready list
  - such a process must never be received by another processor (1. 7–8)



CS (WS 2020/21, LEC 10)

Addendum - Re-Entrant Operations

#### Waitlist Association

depending on the waitlist interpretation, operations to a greater or lesser extent in terms of non-functional properties:

```
inline void apply(process_t *this, event_t *list) {
   #ifdef __FAME_EVENT_WAITLIST__
       insert(&list->wait, &this->event);
   #else
       this->event = list;
   #endif
   }
   inline void elide(process_t *this, event_t *list) {
   #ifdef __FAME_EVENT_WAITLIST__
       winnow(&list->wait, &this->event);
11
   #else
12
13
       this->event = 0;
   #endif
   }
15
```

3/11 ■ dynamic data structure, bounded above, lock-free, lesser list walk

5/13 ■ elementary data type, constant overhead, atomic, larger table walk



33

CS (WS 2020/21, LEC 10) Addendum – Re-Entrant Operations