

Betriebssystemtechnik

Ablaufsteuerung

14. Juni 2010

Motiv: Synchronisation

Koordination der Kooperation und Konkurrenz gleichzeitiger Programmabläufe

ko-or-di-nie-ren *beiordnen*; in ein Gefüge einbauen; aufeinander abstimmen; nebeneinanderstellen; Termine ~.

- ▶ sich überlappen könnende Aktivitäten *der Reihe nach* ausführen
 - ▶ sicherstellen, **kritische Abschnitte konsistent** zu **durchlaufen**
- ▶ „der Reihe nach“ ~ die Verzögerung von Prozessen erzwingen
 - ▶ die überlappende oder die überlappte Aktivität, je nach Verfahren

Lernziel

- ▶ Techniken der Befehlssatzebene zur Synchronisation begreifen
- ▶ unteilbare (wettlaufintolerante) kritische Abschnitte eingrenzen
- ▶ teilbare (wettlauftolerante) kritische Abschnitte entwickeln

Überblick

Ablaufsteuerung

- Einleitung
- Wettlaufintoleranz
 - Unterbrechungssperre
 - Fortsetzungssperre
- Wettlauftoleranz
 - Unterbrechungstransparente Synchronisation
 - Nichtblockierende Synchronisation
- Zusammenfassung
- Bibliographie

Einordnung

Schicht	Funktion	Konzepte
12	Programmverwaltung	Text, Daten, Überlagerung
11	Dateiverwaltung	Dateisystem; Verzeichnis, Verknüpfung
10	Prozessverwaltung	Aktivitätsträger, Kontext, Stapel
9	Adressraumverwaltung	Arbeitsspeicher, Segment, Seite
8	Informationsaustausch	Paket, Nachricht, Kanal, Portal
7	Geräteprogrammierung	Kern; Signal, Zeichen, Block, Datenstrom
6	Platzanweisung	Hauptspeicher, Fragment, Seitenrahmen
5	Zugriffskontrolle	Subjekt, Objekt, Domäne, Befähigung
4	Betriebsmittelzugriff	Verdrängungs-/Vorgangssperre
3	Auftragseinplanung	Ereignis, Priorität, Zeitscheibe, Energie
2	Ablaufsteuerung	Unterbrechungs-/Fortsetzungssperre, Wettlauftoleranz
1	Kontrollflusswechsel	Koroutine, Unterbrechung, Fortsetzung
0	Stammprozessorabstraktion	Stammsystem
-1	Peripherie	MMU, (A)PIC, DMA, UART, ATA, SCSI, USB, ...
-2	Zentraleinheit	ARM, AVR, PowerPC, SPARC, x86, ...

Rekapitulation: SOS 1 [2] bzw. SP [3], BS [4]

Kritischer Abschnitt (KA) [1, S. 137]

- ▶ sich gegenseitig ausschließende Aktivitäten werden nie parallel ausgeführt \models **Elementaroperation** (ELOP)
 - ▶ sie verhalten sich zueinander, als seien sie unteilbar, weil keine Aktivität die andere unterbricht
- ▶ Anweisungen, deren Ausführung einen gegenseitigen Ausschluss erfordern (engl. *critical sections*, *critical regions*)

Beachte \leftrightarrow Abstraktionsniveau der ELOP

- ▶ setzen von „Synchronisationsklammern“ ist nicht zwingend, um einen kritischen Abschnitt zu schützen
- ▶ vielmehr gilt es sicherzustellen, dass die Ausführung eines solchen Abschnitts jederzeit ein **konsistentes Ergebnis** liefert

Integritätswahrung kritischer Abschnitte

Sicherstellung, dass die Ausführung eines kritischen Abschnitts jederzeit ein konsistentes Ergebnis liefert, ist auf zwei Wegen möglich:

- durch zeitweiligen **Ausschluss** der Möglichkeit einer überlappenden Ausführung der relevanten Programmanweisungen
 - ▶ Eintrittsanforderungen (problem- bzw. fallspezifisch) abwehren
- durch **Tolerierung** eben dieser Möglichkeit und Vorsehung gewisser Reparatur- bzw. Erholungsmaßnahmen
 - ▶ überlappend durchführbare, transaktionsartige Verfahren
 - ▶ mit/ohne spezielle Elementaroperationen der Befehlssatzebene

Beachte

- ▶ Überlappung meint je nach Prozessortyp zwei verschiedene Dinge:
 - Uniprozessor** \rightsquigarrow Unterbrechung und Wiedereintritt
 - Multiprozessor** \rightsquigarrow Parallelverarbeitung
- ▶ je nachdem fällt der Umgang mit Überlappung unterschiedlich aus

Integritätswahrung kritischer Abschnitte (Forts.)

zu (a) Eintrittsanforderungen **abwehren** ist *kein* gegenseitiger, sondern ein **einseitiger Ausschluss**

- ▶ von (hard- oder softwarebedingten) Unterbrechungen
- ▶ von (hardwarebedingten) Buszugriffen¹

zu (b) Prinzip: **nichtblockierende Synchronisation**

- lokale Kopie der zu aktualisierenden Variablen anlegen
- lokale Kopie aktualisieren
- Variable mit aktualisierter Kopie abgleichen \rightsquigarrow **Transaktion**
- scheitert 3., die Programmsequenz ab 1. erneut versuchen

Beachte

zu (a) privilegierte Befehle \rightsquigarrow privilegierter Arbeitsmodus der CPU

zu (b) nichtprivilegierte Befehle \rightsquigarrow beliebiger Arbeitsmodus

¹Der Schiedsrichter (engl. *arbiter*) wird angewiesen, Zugriffe anderer Prozessoren auf den gemeinsamen Daten-/Adressbus nicht durchzulassen.

Abwehr von Eintrittsanforderungen

Elementaroperationen schließen überlappte Ausführungen eines kritischen Abschnitts **ggf.** durch denselben oder einen anderen Prozessor aus

- selber Prozessor** \rightsquigarrow Unterbrechungssperre, Fortsetzungssperre
- anderer Prozessor** \rightsquigarrow Zugriffssperre (auf den gemeinsamen Bus)

Beachte \leftrightarrow Holzhammermethode

- ▶ setzen einer Sperre, obwohl die Eintrittsanforderung einen anderen kritischen Bereich betreffen kann
 - ▶ wenn KA und Anforderungsquelle (AQ) uneindeutig zugeordnet sind
 - ▶ genauer: der *Typ* des KA bzw. die *Klasse gemeinsamer Variablen*
 - ▶ d.h. jene Variablen, deren gemeinsame Aktualisierung kritisch ist
 - ▶ der Regelfall ist $N : M$ zwischen KA und AQ, mit $N \gg M$
- ▶ mögliche Nebenläufigkeit wird unnötig eingeschränkt und so letztlich auch ein wahrscheinlicher Verlust an Leistung hervorgerufen

Abwehr von Eintrittsanforderungen (Forts.)

```

orq_avert();
... /* critical section */
orq_admit();

```

ORQ ~ IRQ: Abk. für (engl.) *overlap request*

- pros**
- ▶ die einen KA bildenden Anweisungsfolgen bleiben unverändert
 - ▶ ist der Forderung nach Wiederverwendbarkeit zudenlich
 - ▶ sequentielles Programmierparadigma bleibt bestehen/erhalten
- cons**
- ▶ die Zulassung von Eintrittsanforderungen verzögert sich
 - ▶ um die ggf. nur (schwer) abschätzbare WCET² des KA
 - ▶ das Potential an Nebenläufigkeit wird nur suboptimal genutzt

Beachte

- ▶ einen KA durch „ORQ-Abwehr“ zu schützen, ist trivial
- ▶ dazu ist der KA aber überhaupt erst zu finden, was nicht trivial ist

²Ausführungszeit des schlimmsten Falls (engl. *worst-case execution time*).

Abwehr von Eintrittsanforderungen (Forts.)

Programmabläufe koordinieren sich nach dem Grundsatz: „wer zuerst kommt, mahlt zuerst“ (engl. *first come, first served*; FCFS)

- ▶ eine mit der Strategie zur Prozesseinplanung zumeist in Konflikt stehende „Anspruchshaltung“ zur Nutzung des Prozessors
 - ▶ Prozesseinplanung bestimmt die Reihenfolge der **Prozessorvergabe**
 - ▶ an welchen Prozess der Prozessor vergeben wird,
 - ▶ ob und ggf. wann der Prozess den Prozessor abgibt oder
 - ▶ ob dem Prozess der Prozessor entzogen werden kann, **ist Strategie**
 - ▶ den Prozessor nicht abzugeben, sollte dieser Strategie entsprechen
- ▶ wurden Eintrittsanforderungen abgewehrt, d.h., ist ein KA aktiv, beansprucht der laufende Prozess, den Prozessor nicht zu vergeben

Beachte

- ▶ die Entscheidungen zur Prozessorvergabe sind nicht mehr verbindlich
- ▶ Prioritätsverletzung bzw. Prioritätsumkehr kann die Folge sein

Abwehr von Unterbrechungen

Möglichkeiten der totalen oder partiellen **Unterbrechungssteuerung** —
Beispiel Intel Hardware:

- total**
- ▶ in der Senke der Unterbrechungsanforderung (UA): CPU
 - ▶ Unterbrechungsschalter im FLAGS Register, Bit 9 (IF)
 - ▶ 0 ↔ UA **nicht annehmen**: einstellen mit `cli`, `sti` oder `popf`
- partiell**
- ▶ im Mittler der UA: Gerät/PIC, 8259A
 - ▶ Unterbrechungsschalter im IMR, Bits 0–6
 - ▶ 1 ↔ UA **nicht weiterleiten**: einstellen mit `outb` (OCW1)
 - ▶ in der Quelle der UA: Gerät/UART, 8250 bzw. 16550
 - ▶ Unterbrechungsschalter im IER, Bits 0–3
 - ▶ 0 ↔ UA **nicht aussenden**: einstellen mit `outb`

Beachte ↔ verschachtelte kritische Abschnitte

- ▶ (1) abwehren, (2) KA, (3) wiederzulassen einer UA reicht nicht aus
- ▶ stattdessen: (1) sichern, (2) abwehren, (3) KA, (4) wiederherstellen

Abwehr von Unterbrechungen (Forts.)

Intel Familie, Auswahl

Funktion	totale Abwehr		partielle Abwehr	
	Senke (CPU)		Mittler (PIC)	Quelle (UART)
IRQ_AVERT	<code>cli</code>	<code>pushf</code> <code>cli</code>	<code>inb IMR, %al</code> <code>pushl %eax</code> <code>movl mask, %eax</code> <code>outb %al, IMR</code>	<code>inb IER, %al</code> <code>pushl %eax</code> <code>movl mask, %eax</code> <code>outb %al, IER</code>
IRQ_ADMIT	<code>sti</code>	<code>popf</code>	<code>popl %eax</code> <code>outb %al, IMR</code>	<code>popl %eax</code> <code>outb %al, IER</code>

verschachtelte kritische Abschnitte

- sichern** ▶ FLAGS Register, IMR oder IER
- abwehren** ▶ nicht annehmen, weiterleiten oder aussenden
- zulassen** ▶ wieder annehmen, weiterleiten oder aussenden

Abwehr von Fortsetzungen

Fortsetzung (engl. *continuation*), auch: Einfädelung (engl. *merge*)

- ▶ Nachspann*anhang* einer Unterbrechungsbehandlung
 - ▶ ist grundsätzlich unterbrechbar: $IPL = 0$
- ▶ läuft *synchron* mit anderen Betriebssystemaktivitäten
 - ▶ im Gegensatz zum Nachspann: asynchron zum BS
- ▶ Auslösung (durch Nachspann) erfordert *Koordinierung*

Beachte \leftrightarrow Fortsetzung \subset Epilog [4]

- ▶ Epilog = Nachspann + Fortsetzung
 - ▶ der Nachspann ist „nach unten“ orientiert, fokussiert auf (s)ein Gerät
 - ▶ Betriebssystemfunktionen dürfen nicht aufgerufen werden
 - ▶ die Fortsetzung ist „nach oben“ orientiert, fokussiert auf das BS
 - ▶ dient (einzig) dem Zweck, Betriebssystemfunktionen aufzurufen
 - ▶ sie wird durch eine Art *Hochruf* (engl. *upcall*, [5]) aktiviert
- ▶ die Fortsetzung (eines Nachspanns) ist ein Epilog ohne AST

Abwehr, Zulassung und Weitergabe von Fortsetzungen

ICE (Abk. engl. *interrupt continuation executive*)

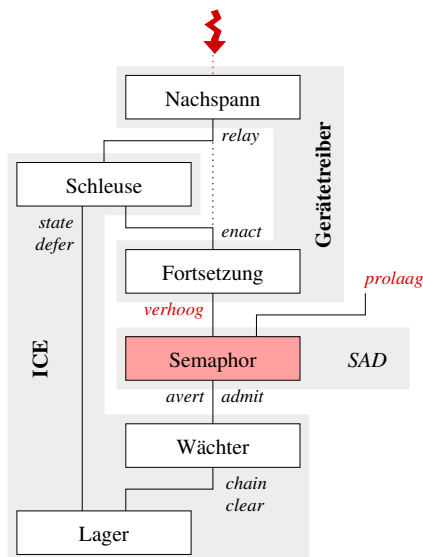
```
extern void ice_avert (ice_t *);           /* defer continuations */
extern void ice_treva (volatile ice_t *); /* allow continuations */
extern void ice_admit (ice_t *);         /* unleash continuations */
extern void ice_clear (ice_t *);        /* process continuations */
extern void ice_defer (ice_t *, job_t *); /* defer continuation */
extern void ice_relay (ice_t *, job_t *); /* enact/defer cont. */
```

```
extern char ice_state (ice_t *);         /* guard state */
extern job_t *ice_chain (volatile ice_t *); /* guard load */
extern ice_t *ice_guard ();              /* guard singleton */
```

```
typedef struct ice {
    char busy;           /* state */
    queue_t load;       /* jobs */
} ice_t;
```

```
typedef reinit_t job_t;
extern void job_enact (job_t *);
```

Fortsetzung der Ereigniszustellung: Entwurfsskizze



- Nachspann** ▶ SLIH
 - ▶ vgl. Kapitel 5
- Schleuse** ▶ serialisiert Fortsetzungen
- Wächter** ▶ kontrolliert krit. Abschnitt
- Lager** ▶ speichert Fortsetzungen

Beachte \leftrightarrow Semaphore

- ▶ blockadefrei synchronisiert
- ▶ ein überlappendes V kommt verzögert zur Ausführung

Serialisierte kritische Abschnitte: Semaphore

```
void ewd_prolaag (semaphore_t *) {
    ice_avert(ice_guard());
    /* ... */
    ice_admit(ice_guard());
}
```

```
void ewd_verhoog (semaphore_t *) {
    ice_avert(ice_guard());
    /* ... */
    ice_admit(ice_guard());
}
```

- guard** ▶ ist die Verwaltungsstruktur zur Steuerung von Fortsetzungen
- avert** ▶ signalisiert, dass Fortsetzungen zurückzustellen sind
- admit** ▶ signalisiert, dass Fortsetzungen nicht zurückzustellen sind
 - ▶ verarbeitet zurückgestellte Fortsetzungen, sofern erforderlich

Beachte

- ▶ zwischen *avert* und *admit* werden Fortsetzungen zurückgestellt
- ▶ die Funktionen arbeiten blockade- aber nicht verzögerungsfrei

Fortsetzung eines Nachspans: Treiber \iff Faden

```
typedef struct driver {
    job_t task;          /* SLIH continuation */
    semaphore_t port;   /* driver thread signalling socket */
} driver_t;

void job_signal (driver_t *this) { ewd_verhoog(&this->port); }

driver_t job_driver = {{{0}, &job_signal}, {0}};

void a_slih () { ice_relay(ice_guard(), (job_t*)&job_driver); }
void a_hils () { ewd_prolaag(&job_driver.port); }
```

- a_hils* ▶ Routine eines Fadens: konsumiert ein Gerätesignal
- a_slih* ▶ Nachspann des Treibers: löst eine Fortsetzung aus
- job_driver* ▶ Deskriptor einer Fortsetzung
- job_signal* ▶ Fortsetzung des Nachspans: produziert ein Gerätesignal

Serialisierung kritischer Abschnitte

```
void ice_await (ice_t *this) {
    this->busy = 1;
}
```

```
void ice_admit (ice_t *this) {
    ice_treva(this);
    if (ice_chain(this))
        ice_clear(this);
}
```

```
void ice_treva (volatile ice_t *this) {
    this->busy = 0;
}
```

- await* ▶ serialisieren
- admit* ▶ aufarbeiten
- treva* ▶ \neg *await*
- chain* ▶ etwas zu tun?

```
job_t *ice_chain (volatile ice_t *this) {
    return (job_t*)this->load.head.link;
}
```

Beachte \iff *admit*

- ▶ Fortsetzungen können sich überholen: kein striktes FCFS

Serialisierung kritischer Abschnitte (Forts.)

Fortsetzungen weiterleiten, d.h., bedingt aktivieren

```
void ice_relay (ice_t *this, job_t *task) {
    if (!ice_state(this)) job_enact(task); /* idle, call it */
    else ice_defer(this, task);          /* busy, store it */
}

char ice_state (ice_t *this) {
    return this->busy; /* 0 = idle, 1 = busy */
}

void ice_defer (ice_t *this, job_t *task) {
    fad_aback(&this->load, &task->next); /* append continuation */
}
```

- relay* ▶ durchschleusen einer Fortsetzung: aktivieren/zurückstellen
- state* ▶ Aktivierungszustand des serialisierten kritischen Abschnitts
- defer* ▶ zurückstellen (speichern) einer Fortsetzung

Serialisierung kritischer Abschnitte (Forts.)

Fortsetzungen freistellen und aktivieren

```
void ice_clear (ice_t *this) {
    job_t *task;
    while ((task = (job_t*)fad_fetch(&this->load)))
        job_enact(task);
}

void job_enact (job_t *task) {
    (task->work)(task); /* activate continuation object */
}
```

- clear* ▶ Warteschlange zurückgestellter Fortsetzungen abarbeiten
- enact* ▶ Fortsetzungsmethode auf Fortsetzungsobjekt applizieren

Beachte \iff *clear*

- ▶ die Freistellung von Fortsetzungen erfolgt außerhalb des KA
- ▶ die **Aktivierungsreihenfolge** von Fortsetzungen ist unbestimmt

Serialisierung kritischer Abschnitte (Forts.)

Plausibilitätskontrolle

relay stellt Fortsetzungen zurück \iff *state = BUSY*

- ▶ Programmausführung zwischen *avert* und *treva* unterbrochen

clear läuft mit *state = IDLE* \Rightarrow keine Fortsetzung geht verloren

enact aktiviert Fortsetzungen, die ggf. serialisierte KA ausführen

admit lässt Überlappungen von Fortsetzungen zu

- keine Überlappung \Rightarrow *clear* baut Fortsetzungsliste ab
- Überlappung \Rightarrow **indirekte Rekursion** ist möglich
 - ▶ sofern *enact* einen KA aktiviert, der mit *admit* endet
 - ▶ Fall (a) ist sodann Terminationsbedingung

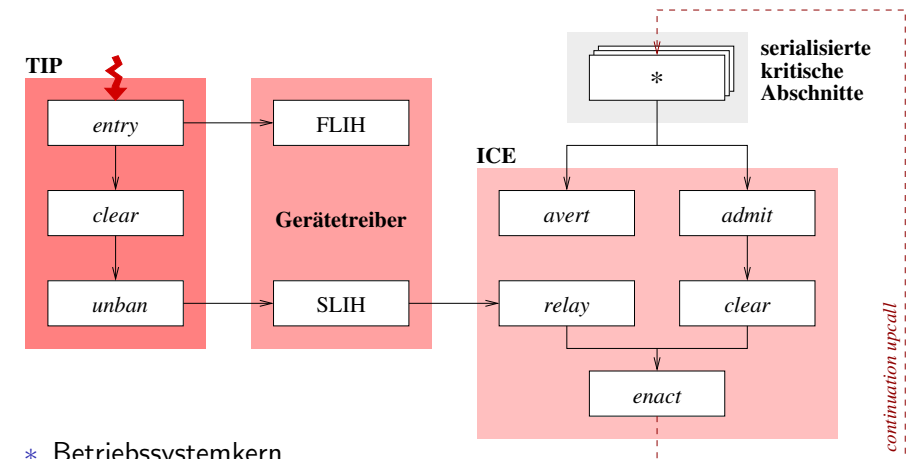
Beachte \iff Rekursionstiefe, hängt ab von:

(ein Fall für TAS, falls kritisch)

- ▶ Frequenz, Häufigkeit und Art von Unterbrechungen
- ▶ Funktion der Vor- und Nachspanne der Unterbrechungen
- ▶ Funktion der Fortsetzungen

Serialisierung kritischer Abschnitte (Forts.)

Gesamtzusammenhang: Aktivierungsfolge zentraler Funktionen



* Betriebssystemkern

- im normalen (synchronen) Fadenkontrollfluss aktivierte Funktionen
- asynchron ausgelöste, „einsynchronisierte“ Funktionen, z.B.:
 - ▶ *verhoog*, *preempt* – in normalen Ablauf einzufädelnde Aufrufe

Fortsetzungsliste — konstruktiv gesicherte Datenstruktur

Beachte \iff Operationsprinzip AST (vgl. Kap. 5)

- ▶ Nachspannfreistellung (TIP, *clear*) überlappt sich selbst nicht
 - ▶ Parallelverarbeitung (Zustellergruppe) einmal außen vor gelassen
- ▶ Folge: Fortsetzungweiterleitung (ICE, *clear*) überlappt sich nie
 - ▶ auch **Verdrängungsereignisse** werden sich hinten anstellen:

Verdrängung \iff Vorspann \mapsto Nachspann \mapsto Fortsetzung

- ▶ durch Konstruktion bedingte Wechselwirkung (engl. *interaction*)
 - ▶ beugt Wettlaufsituationen in Bezug auf die Fortsetzungsliste vor
 - ▶ macht explizite Synchronisation der Zugriffsoperationen überflüssig

FAD (Abk. für engl. *fundamental algorithms and data structures*)

```
extern void fad_reset (queue_t *);          /* clear queue */
extern void fad_aback (queue_t *, chain_t *); /* append chain item */
extern chain_t *fad_fetch (queue_t *);     /* remove chain item */
```

Fortsetzungsliste: FAD

```
typedef struct chain chain_t;

struct chain {
    chain_t *link;
};
```

```
void fad_reset (queue_t *this) {
    this->head.link = 0;
    this->tail = &this->head;
}
```

```
typedef struct queue queue_t;

struct queue {
    chain_t head;
    chain_t *tail;
};
```

```
void fad_aback (queue_t *this,
               chain_t *item) {
    item->link = 0;
    this->tail->link = item;
    this->tail = item;
}
```

```
chain_t *fad_fetch (queue_t *this) {
    chain_t *item;
    if ((item = this->head.link) && !(this->head.link = item->link))
        this->tail = &this->head;
    return item;
}
```

Unterbrechungstransparente Synchronisation [6]

Koordinierung unterbrechungsbedingter Aktivitäten, ohne **asynchrone Programmunterbrechungen** abzuwehren

- ▶ die Systemsoftware ist frei von Unterbrechungssperren³
- ▶ Synchronisation verwendet nur nichtprivilegierte Befehle der ISA
 - ▶ mit gewissen **Atomizitätseigenschaften** in Bezug auf Unterbrechungen
 - ▶ wie z.B. atomares Lesen/Schreiben von Speicherworten ($n > 1$ Bytes)
 - ▶ aber auch TAS, FAA oder CAS, d.h., **atomare Komplexbefehle**
- ▶ nur die Hardware selbst sperrt Unterbrechungen (zeitweilig) aus
 - ▶ gleichwohl bestimmt Software, wie lange diese Sperren aktiv sind !!!

Beachte ↔ Hilfestellung durch andere Konzepte

- ▶ um Unterbrechungssynchronisation zur *Ausnahme* werden zu lassen
 - ▶ z.B. Fortsetzungssperren, um KA konventionell schützen zu können
- ▶ die **Konzentration auf das Wesentliche** fördern: **Fortsetzungsliste**

³Erneute Abwehr von Unterbrechungen nach Unterbrechungsfreigabe, setzt keine Unterbrechungssperre im eigentlichen Sinn (vgl. Kapitel 5).

Fallstudie: Warteschlangensynchronisation

ITS (Abk. für engl. *interrupt transparent synchronization*)

```
extern void its_reset (queue_t *);           /* clear queue */
extern void its_aback (queue_t *, chain_t *); /* append chain item */
extern chain_t *its_fetch (queue_t *);      /* remove chain item */
```

```
void its_reset (queue_t *this) {
    fad_reset(this);
}
```

```
void its_aback (queue_t *, chain_t *) {
    /* ... */
}
```

```
chain_t *its_fetch (queue_t *this) {
    chain_t *item = fad_fetch(this);
    /* ... */
    return item;
}
```

ITS ↔ FAD

reset Wiederverwendung

- ▶ *chain_t*
- ▶ *queue_t*

aback Ersetzung

fetch Spezialisierung

Fallstudie: Warteschlangensynchronisation (Forts.)

Element wettlaufstolerant einfügen

```
void its_aback (queue_t *this, chain_t *item) {
    chain_t *last;

    item->link = 0;           /* make item last chain element */

    last = this->tail;        /* remember item insertion point */
    this->tail = item;        /* advance tail pointer, optimistically */

    while (last->link)        /* overlapping aback: find actual tail */
        last = last->link;

    last->link = item;        /* append item */
}
```

Atomizität ▶ Lesen/Schreiben von Zeigerwerten ist ELOP

Überlappungsmuster ▶ *aback* überlappt *aback* oder *fetch*

Fallstudie: Warteschlangensynchronisation (Forts.)

Element wettlaufstolerant entfernen

(*fad_fetch* expandiert)

```
chain_t *its_fetch (queue_t *this) {
    chain_t *item;

    if ((item = this->head.link) && !(this->head.link = item->link)) {
        this->tail = &this->head; /* point of problem! */
        if (item->link) {         /* race condition detected! */
            chain_t *help, *lost = item->link;
            do {
                help = lost->link; /* requeue lost elements */
                its_aback(this, lost);
            } while ((lost = help));
        }
    }
    return item;
}
```

Überlappungsmuster ▶ *fetch* wird nur von *aback* überlappt

Fallstudie: Warteschlangensynchronisation (Forts.)

Plausibilitätskontrolle

- aback**
- ▶ überlappt sich selbst immer nur **stapelweise**, wenn überhaupt
 - ▶ Wettlaufsituation $\iff last = tail$, jedoch $tail \neq item$
 - Normalfall $LINK(last) = 0 \Rightarrow$ kein *aback*-Wiedereintritt
 - Konfliktfall $LINK(last) \neq 0 \Rightarrow last$ falsch, korrigieren
 - ▶ Zuweisung an *link* (einfügen von *item*) $\iff last \neq tail$
- fetch**
- ▶ überlappt sich nie selbst, auch nicht durch Verdrängung
 - ▶ Wettlaufsituation $\iff head = 0$, jedoch $tail \neq \&head$
 - Normalfall $LINK(item) = 0 \Rightarrow$ *aback* überlappte nicht
 - Konfliktfall $LINK(item) \neq 0 \Rightarrow item \rightsquigarrow lost \ \& \ found$
 - ▶ umtragen der Einträge aus „*lost & found*“-Liste ist atomar

Beachte \leftrightarrow Eignung für TIP und ICE

- Nachspannliste** ▶ zutreffend, jedoch nicht bei Parallelverarbeitung
- Fortsetzungsliste** ▶ unzutreffend: AST serialisiert Nachspänne !!!

Rückblick: Nachspann-/Fortsetzungsliste

Querschneidende Belange durch Verdrängung und Parallelverarbeitung

- Nachspannliste**
- ▶ Aufbau auf Anforderung durch einen Vorspann
 - ▶ Vorspänne überlappen sich nur stapelweise
 - ▶ dito: Fließbandmodell der Parallelverarbeitung
 - ▶ Abbau im Rahmen der Behandlung eines AST
 - ▶ Freistellung der Nachspänne ist unteilbar und überlappt deren Zurückstellung nicht
 - ▶ gilt bei Parallelverarbeitung nicht
- Fortsetzungsliste**
- ▶ Aufbau auf Anforderung durch einen Nachspann
 - ▶ Nachspänne überlappen sich nie: **AST serialisiert**
 - ▶ gilt bei Parallelverarbeitung nur bedingt
 - ▶ Abbau beim Verlassen eines serialisierten KA
 - ▶ Freistellung der Fortsetzungen ist **bedingt teilbar** und überlappt deren Zurückstellung nicht
 - ▶ gilt bei Parallelverarbeitung nicht

Beachte \leftrightarrow Parallelverarbeitung (Multiprozessoren)

- ▶ lässt Annahmen zum Überlappungsmuster nicht aufrechterhalten

Nichtblockierende Synchronisation

Koordinierung sich einander ggf. überlappender Aktivitäten, ohne dabei **gleichzeitige Prozesse**⁴ auszuschließen

- ▶ toleriert (pseudo-) parallele Programmausführungen
 - parallel** ▶ Multiprozessor, wirkliche Parallelität
 - pseudoparallel** ▶ Uniprozessor, Parallelität durch Unterbrechungen
- ▶ die Verfahren greifen auf nichtprivilegierte Befehle der ISA zurück
 - CISC** ▶ TAS, FAA, CAS bzw. CMPXCHG
 - RISC** ▶ LL/SC
- ▶ d.h., sie funktionieren im Benutzer- wie auch im Systemmodus

Beachte

- ▶ kein gegenseitiger Ausschluss \Rightarrow **Verklemmungsvorbeugung**
- ▶ die *benutzten* Befehle sind „echte“ Elementaroperationen der ISA

⁴Prozesse, deren Ausführung sich zeitlich überschneidet.

CAS: Prinzip

- compare and swap** ▶ Elementaroperation der Befehlssatzebene⁵
- ▶ unteilbar: Unterbrechungs- und Zugriffssperre (Datenbus)
 - ▶ implementiert eine „Transaktion“ für Uni- und Multiprozessoren

```
int cas (word_t *ref, word_t exp, word_t val) {
    unsigned char aux;

    orq_avert();
    if (aux = (*ref == exp)) *ref = val;
    orq_admit();

    return aux;
}
```

Operationsergebnis

- true**
- ▶ gleich
 - ▶ geschrieben
 - ▶ gelungen
- false**
- ▶ ungleich
 - ▶ gelesen
 - ▶ gescheitert

Generalisierte Schnittstelle

```
#define CAS(r,e,v)    cas((word_t*)r, (word_t)e, (word_t)v)
```

⁵Lässt sich durch passende Auslegung von *avert* und *admit* ggf. nachbilden.

CAS: Nachbildung für x86

```
int cas (word_t *ref, word_t exp, word_t val) {
    unsigned char aux;

    __asm__ __volatile__(
        "lock\n\t"                /* prefix next instruction */
        "cmpxchgl %2,%1\n\t"      /* (ref) == exp ? (ref) = val */
        "sete %0"                 /* extend ZF into aux */
        : "=q" (aux), "=m" (*ref)
        : "r" (val), "m" (*ref), "a" (exp) /* %eax loaded with exp */
        : "memory");

    return aux;
}
```

Beachte \leftrightarrow *lock*

optional für Uniprocessorsysteme: Befehle der ISA sind ununterbrechbar
zwingend für Multiprocessorsysteme: setzt Zugriffssperre (Datenbus)

Nichtblockierende Synchronisation mit CAS

erledige NBS mit CAS;

wiederhole

ziehe *lokale Kopie* des Inhalts der *Adresse* einer globalen Variablen;
 verwende die Kopie, um einen neuen *lokalen Wert* zu berechnen;
 versuche CAS: an *Adresse*, die *lokale Kopie* mit dem *lokalen Wert*;

solange CAS scheitert;

basta.

- pros
- ▶ Tolerierung beliebiger Überlappungsmuster
 - ▶ transparent für die Einplanung: keine Prioritätsumkehr
 - ▶ Vorbeugung von Verklemmungen: Bedingung 1 entkräftet [3]
 - ▶ Robustheit: keine hängenden Sperren bei Programmabbrüchen
 - ▶ in funktionaler Hinsicht wiederverwendbar und komponierbar
- cons
- ▶ Gefahr von **Verhungering** (engl. *starvation*)
 - ▶ Wiederverwendung sequentieller Altsoftware unmöglich
 - ▶ **Entwicklung** nebenläufiger Varianten im Regelfall **nicht trivial**

Fallstudie: Zählermanipulation

- fetch and add* ▶ atomare Manipulation eines Speicherzelleninhalts
- ▶ Nachbildung von XADD mittels CAS: **unteilbares Zählen**
 - ▶ einer mit 80486 eingeführten ELOP für die x86-Familie
 - ▶ geeignet zur (blockadefreien) Implementierung von Semaphore

NBS (Abk. für engl. *non-blocking synchronization*)

```
extern int nbs_count (int *, int);          /* fetch and add number */
```

```
int nbs_count (int *this, int rate) {
    int copy;

    do copy = *this;                /* fetch contents */
    while (!CAS(this, copy, copy + rate)); /* add value if unvaried */

    return copy;                    /* return (old) contents */
}
```

Fallstudie: Zählermanipulation (Forts.)

Plausibilitätskontrolle

- copy = *this* ▶ zieht eine lokale Kopie der zu manipulierenden Zahl
- copy + rate* ▶ berechnet den neuen Wert auf Basis dieser Kopie
- CAS ▶ versucht, den neuen Wert zu binden (engl. *commit*)
- true* \iff kein Zugriffskonflikt, neuer Wert gültig
- false* \iff **Zugriffskonflikt**, neuer Wert verworfen
- do ... while* ▶ terminiert nur, falls der neue Wert gebunden wurde

Beachte \leftrightarrow ABA-Fall (vgl. S. 6-40)

- ▶ die Feststellung des Zugriffskonflikts basiert auf eine Überprüfung des Inhalts einer Speicherstelle, nicht auf der Anwendung ihrer **Adresse**
- ▶ nicht jeder **überlappende Schreibzugriff** ist daher wirklich erkennbar
 - ▶ z.B., wenn zwischenzeitlich eine Änderung um n und $-n$ erfolgt

Fallstudie: Wartestapelmanipulation

LCFS (Abk. für engl. *last come, first served*)

- ▶ mögliche Grundlage für stapelorientierte Fadenverarbeitung [7]
 - Fadeneinplanung \mapsto Fadenkontrollblock *push*
 - Fadeneinlastung \mapsto Fadenkontrollblock *pull*
- ▶ atomare Manipulation einer einfach verketteten (LIFO) Liste

NBS — minimale Erweiterung

```
extern void nbs_flush (chain_t *);          /* clear chain */
extern void nbs_ahead (chain_t *, chain_t *); /* push chain item */
extern chain_t *nbs_strip (chain_t *);     /* pull chain item */
```

```
void nbs_flush (chain_t *this) {
    this->link = 0;
}
```

Fallstudie: Wartestapelmanipulation (Forts.)

Aufnahme in die Liste nach LCFS

```
void nbs_ahead (chain_t *this, chain_t *item) {
    do item->link = this->link;          /* is elected head */
    while (!CAS(&this->link, item->link, item)); /* try push item */
}
```

Entnahme aus der Liste nach LCFS

```
chain_t *nbs_strip (chain_t *this) {
    chain_t *node;

    do if ((node = this->link) == 0) break; /* access head */
    while (!CAS(&this->link, node, node->link)); /* try pull node */

    return node;
}
```

Fallstudie: Wartestapelmanipulation (Forts.)

Plausibilitätskontrolle

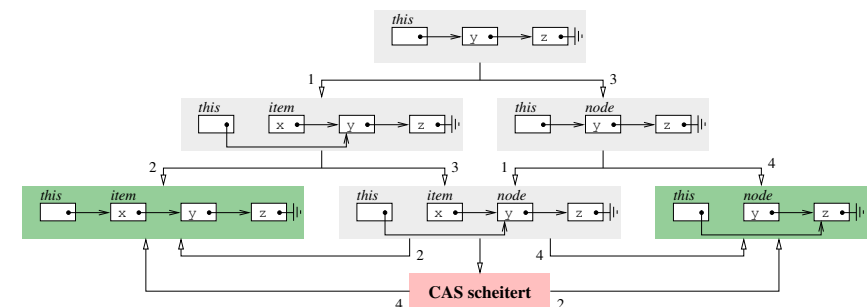
- ahead*
- ▶ das durch *item* adressierte Element ist noch nicht in der Liste und es wird auch nicht gleichzeitig in diese Liste eingetragen
 - ▶ kritischer Datenbestand ist $LINK(this)$, der **Listenkopf**
 - ▶ dieser wird in $LINK(item)$ als Kopie angelegt
 - ▶ neuer Listenkopf ist *item*, dessen Bindung CAS versucht
 - ▶ *do ... while* terminiert, falls *item* als Kopf gebunden wurde
- strip*
- ▶ kritischer Datenbestand ist $LINK(this)$, der **Listenkopf**
 - ▶ dieser wird in *node* als (lokale) Kopie angelegt
 - ▶ neuer Kopf ist $SUCC(node)$, dessen Bindung CAS versucht
 - ▶ *do ... while* terminiert, falls $SUCC(node)$ gebunden wurde

Beachte

- ▶ auch wenn mehrere Fäden auf denselben Kopfzeiger gleichzeitig zugreifen, wird CAS für nur einen Faden die Manipulation zulassen
- ▶ je nach Nutzung von *ahead* und *strip* droht das „ABA-Problem“

Fallstudie: Wartestapelmanipulation (Forts.)

Plausibilitätskontrolle: Datenstrukturentwicklung je nach Überlappungsfall



```
void nbs_ahead (chain_t *this, chain_t *item) {
    1 do item->link = this->link;
    2 while (!CAS(&this->link, item->link, item));
}
```

```
chain_t *nbs_strip (chain_t *this) {
    chain_t *node;
    3 do if ((node = this->link) == 0) break;
    4 while (!CAS(&this->link, node, node->link));
    return node;
}
```

Problem ABA

Phänomen der nichtblockierenden Synchronisation auf Basis eines CAS, d.h., einer ELOP, die inhaltsbasiert arbeitet⁶

- ▶ angenommen zwei Fäden, F_1 und F_2 , stehen im Wettstreit um eine gemeinsame Variable V
 - F_1 ▶ liest den Wert A von V , speichert diesen als Kopie, wird dann allerdings vor dem CAS_V für unbestimmte Zeit verzögert
 - F_2 ▶ durchläuft dieselbe Sequenz, schafft jedoch mittels CAS_V den Wert B an V zuzuweisen
 - ▶ anschließend wird (in einem weiteren Durchlauf dieser Sequenz) wieder der ursprüngliche Wert A an V zugewiesen
 - F_1 ▶ setzt seine Ausführung mit CAS_V fort, erkennt, dass V den Wert A seiner Kopie speichert und überschreibt V
- ▶ im Ergebnis kann dieses Überlappungsmuster dazu führen, dass F_1 mittels CAS_V einen falschen Wert nach V transferiert

⁶Bei *Adressreservierung* wie z.B. mit LL/SC besteht dieses Problem nicht.

Problem ABA (Forts.)

Wartestapelmanipulation

Wartestapelzustand: $LINK(this) \rightarrow A \rightarrow B \rightarrow C$

	F_1	F_2	F_2	F_2	F_1
	<i>st. . .</i>	<i>strip</i>	<i>strip</i>	<i>ahead</i>	<i>. . .rip</i>
$LINK(this)$	A	A	B	C	A
$node_{strip}/item_{ahead}$	A	A	B	A	A
$LINK(node)$	B	B	C	–	B
$CAS_{LINK(this)}$. . .	$A \rightsquigarrow B$	$B \rightsquigarrow C$	$C \rightsquigarrow A$	$A \rightsquigarrow B$

Beachte die Zerstückelung der Liste: $LINK(this) \rightarrow B \rightarrow ? \wedge A \rightarrow C$

- ▶ eine Teilliste (A , Kopf) ist aus dem Wartestapel verschwunden
- ▶ eine andere Teilliste (B , Kopf) gelangte zurück in den Wartestapel
 - ▶ nachdem F_2 Eintrag B mittels *strip* regulär entnommen hatte
 - ▶ jedoch ohne dass B danach mittels *ahead* wieder aufgenommen wurde

Problem ABA (Forts.)

Abhilfe besteht darin, den umstrittenen Zeiger (nämlich *item* bzw. *node*) um einen problemspezifischen **Generationszähler** zu erweitern

- Etikettieren**
- ▶ Zeiger mit einem Anhänger (engl. *tag*) versehen
 - ▶ Ausrichtung (engl. *alignment*) ausnutzen, z.B.:

$$\begin{aligned} \text{sizeof}(\text{chain}_t) &\rightsquigarrow 4 = 2^2 \Rightarrow n = 2 \\ &\Rightarrow \text{chain}_t * \text{ ist Vielfaches von } 4 \\ &\Rightarrow \text{chain}_t * \text{Bits}_{[0:1]} \text{ immer } 0 \end{aligned}$$

- ▶ Platzhalter für n -Bit Marke/Zähler in jedem Zeiger

- DCAS**
- ▶ Abk. für (engl.) *double compare and swap*
 - ▶ Marke/Zähler als elementaren Datentyp auslegen
 - ▶ *unsigned int* hat Wertebereich von z.B. $[0, 2^{32} - 1]$
 - ▶ zwei Maschinenworte (Zeiger, Marke/Zähler) ändern

Problem ABA (Forts.)

Abhilfe (engl. *workaround*)

- ▶ bei Software, umgehen **unlösbarer Fehler** [8]

Beachte \leftrightarrow **Generationszähler**

- ▶ eine Lösung für CAS-artige Verfahren ist nicht wirklich gegeben
- ▶ die Wahrscheinlichkeit, dass das Problem auftritt, wird verringert
- ▶ Überlappungsmuster haben Einfluss auf den Wertebereich
 - ▶ bestimmt durch Zusammenspiel *und* Anzahl der wettstreitigen Fäden
 - ▶ ein Bit kann reichen, ebenso, wie ein *unsigned int* zu klein sein kann

Vorbeugung (engl. *prevention*) ist der richtige Ansatz — sofern möglich

- ▶ beliebige Überlappungsmuster konstruktiv (Entwurf) ausschließen
- ▶ auf **Adressreservierungsverfahren** der Hardware zurückgreifen
 - ▶ unterstützt nicht jede Hardware, ist nur typisch für RISC
 - ▶ z.B. CAS mit *load linked*, *store conditional* (LL/SC) implementieren

Fallstudie: Wartestapelmanipulation (Forts.)

Vorbeugung bzw. Abhilfe zum ABA-Problem für die Nachspannliste

Abläufe zur Propagierung von asynchronen Systemsprüngen aufgreifen:

- Zurückstellung**
- ▶ baut die Nachspannliste nur auf, niemals ab
 - ▶ ruft nur *ahead* auf, niemals *strip*
- Freistellung**
- ▶ baut die Nachspannliste nur ab, niemals auf
 - ▶ ruft nur *strip* auf, niemals *ahead*

Beachte \leftrightarrow Rechnerbetriebsart

- Uniprozessor**
- ▶ Freistellung schreitet sequentiell voran
 - ▶ *strip* überlappt nie \Rightarrow kein ABA-Problem
- Multiprozessor**
- ▶ Freistellung schreitet bedingt parallel voran
 - ▶ *strip* überlappt nur im Falle **Zustellergruppe: ahead**
 - ▶ **Mehrfachauslösung** eines Nachspanns ist kritisch
 - ▶ Freiliste \rightsquigarrow eindeutiger *tip_t** \Rightarrow kein ABA-Problem
 - ▶ nicht so jedoch bei **Ereigniszähler/Bitschalter**
 - ▶ ABA-Problem $\iff f_{\text{Zurückstellung}(x)} < f_{\text{Freistellung}(x)}$

Fallstudie: Wartestapelmanipulation (Forts.)

Abhilfe zum ABA-Problem: Etikettierung

Abstrakter Datentyp *chain_p**: Spezialisierung von *chain_t**

```
typedef chain_t chain_p;
```

```
extern chain_p *aba_wheel (chain_p *); /* rotate pointer tag bit(s) */
extern chain_t *aba_index (chain_p *); /* return pointer value */
```

Etikett anheften

```
chain_p *aba_wheel (chain_p *item) {
    return (chain_p *)((unsigned)item ^ 1);
}
```

Etikett entfernen

```
chain_t *aba_index (chain_p *item) {
    return (chain_t *)((unsigned)item & ~1);
}
```

Verwendung

- wheel**
- ▶ in *ahead*
 - ▶ markieren
 - ▶ Zeiger färben
- index**
- ▶ in beiden
 - ▶ bereinigen
 - ▶ Zeiger liefern

Fallstudie: Wartestapelmanipulation (Forts.)

Lebensdauer des Generationszählers bzw. der Etiketten

Abhilfe gegen des Problem der Mehrdeutigkeit (hier: ABA) von Zeigern ist **nicht transparent** für die Programme, die die Zeiger verwenden

- ▶ dies trifft insbesondere auch zu auf die Etikettierung von Zeigern
 - ▶ damit bleiben lediglich *Einzelwort-CAS* weiterhin möglich
 - ▶ Transparenz durch Beibehaltung der Zeigergröße ist nicht das Ziel
- ▶ voll ausgeprägte Generationszähler sind offensichtlich intransparent
 - ▶ Zeiger und Generationszähler müssen eine Einheit bilden
 - ▶ beide zusammen verdoppeln die Zeigergröße \rightsquigarrow *Doppelwort-CAS*

Zeiger samt Generationszähler/Etikett sind Exemplar eines Typs

- ▶ problemspezifische Auslegung, je nach Wertebereich des Zählers
 - Einzelwort* *chain_t**, falls einfache Etikettierung genügt
 - Doppelwort* *chain_t** und *unsigned int*, sonst
- ▶ Repräsentation als **abstrakter Datentyp** [9] \Rightarrow Anpassung von NBS

Fallstudie: Wartestapelmanipulation (Forts.)

Etikettierter Zeigertyp *chain_t**

```
void nbs_ahead (chain_t *this, chain_p *item) {
    chain_p *turn = aba_wheel(item); /* new pointer generation */

    do aba_index(item)->link = this->link;
    while (!CAS(&this->link, aba_index(item)->link, turn));
}

chain_p *nbs_strip (chain_t *this) {
    chain_p *node;

    do if (aba_index((node = this->link)) == 0) break;
    while (!CAS(&this->link, node, aba_index(node)->link));

    return node;
}
```

Beachte \leftrightarrow abstrakter Datentyp *chain_p** \mapsto *chain_t**

- ▶ Exemplare dieses Typs dürfen nur mittels *index* benutzt werden

Fallstudie: Warteschlangenmanipulation

NBS — minimale Erweiterung (Forts.)

```
extern void nbs_reset (queue_t *);          /* clear queue */
extern void nbs_aback (queue_t *, chain_t *); /* append chain item */
extern chain_t *nbs_fetch (queue_t *);     /* remove chain item */
```

```
void nbs_reset (queue_t *this) {
    fad_reset(this);
}
```

```
void nbs_aback (queue_t *, chain_t *) {
    /* ... */
}
```

```
chain_t *nbs_fetch (queue_t *this) {
    chain_t *node;
    /* ... */
    return node;
}
```

NBS \iff FAD/ITS*reset* Wiederverwendung*▶ *chain_t*▶ *queue_t**aback* Ersetzung*fetch* Ersetzung

Beachte *

- ▶ ggf. nicht, sollte das ABA-Problem bestehen

Fallstudie: Warteschlangenmanipulation (Forts.)

Szenario 1: *aback* toleriert Überlappungen nur von sich selbst

```
void nbs_aback (queue_t *this, chain_t *item) {
    chain_t *last;

    item->link = 0;

    do last = this->tail;
    while (!CAS(&this->tail, last, &item->link));

    last->link = item;
}
```

Plausibilitätskontrolle \iff kritisch ist *tail*, der Listenfuß

- ▶ die Kopie des aktuellen Wertes von *tail* wird in *last* vermerkt
- ▶ neuer Wert von *tail* ergibt sich aus $\&LINK(item)$
- ▶ CAS gelingt \iff *tail* hält den alten Wert $\Rightarrow tail = \&LINK(item)$
- ▶ *do ... while* terminiert \iff CAS gelingt

Fallstudie: Warteschlangenmanipulation (Forts.)

Szenario 2: *fetch* toleriert Überlappungen nur von sich selbst

```
chain_t *nbs_fetch (queue_t *this) {
    chain_t *node;

    do if ((node = this->head.link) == 0) return 0;
    while (!CAS(&this->head.link, node, node->link));

    if (node->link == 0)
        this->tail = &this->head;

    return node;
}
```

Plausibilitätskontrolle \iff kritisch ist $LINK(head)$, der Listenkopf

- ▶ analog zu *aback*, jedoch Sonderbehandlung wenn gilt:
 - ▶ $LINK(head) = node \wedge tail = \&LINK(node) \Rightarrow$ einziges Element
- ▶ im Falle des einzigen Elements gilt nach weiteren Verlauf:
 - ▶ $LINK(head) = LINK(node) = 0 \Rightarrow tail = \&LINK(head)$!!!

Fallstudie: Warteschlangenmanipulation (Forts.)

Szenario 3: *fetch* toleriert Überlappungen von sich selbst und (vermeintlich) von *aback*

```
chain_t *nbs_fetch (queue_t *this) {
    chain_t *node;

    do if ((node = this->head.link) == 0) return 0;
    while (!CAS(&this->head.link, node, node->link));

    if (node->link == 0)
        CAS(&this->tail, &node->link, &this->head);

    return node;
}
```

Beachte \iff hier gilt: $LINK(node) = 0 \Rightarrow LINK(head) = 0$

- ▶ überlappendes *aback* hat ggf. *item* angehängt: $tail = \&LINK(item)$
- ▶ dann müsste jedoch auch gelten: $LINK(head) = item$
- \Rightarrow Widerspruch: *fetch* toleriert überlappendes *aback* nicht

Fallstudie: Warteschlangenmanipulation (Forts.)

Rekapitulation der Szenarien 1–3

- aback*||*aback* ▶ Integrität von *tail* ist sichergestellt
- fetch*||*fetch* ▶ Integrität von *LINK(head)* ist sichergestellt
- fetch*||*aback* ▶ Integrität von *LINK(head)* bzw. *tail* ist sichergestellt
- ▶ Integrität von beiden zusammen ist **nicht sichergestellt**
- aback*||*fetch* ▶ Integrität \sim *fetch*||*aback* \Rightarrow **nicht sichergestellt**

Neuralgischer Punkt \leftrightarrow Kopfelement wird zum Verkettungsglied

- aback* ▶ hängt *item* ggf. an ein Kopfelement (*node = last*) an, das ggf. schon nicht mehr auf der Liste steht
- ▶ muss die Aktualisierung von *LINK(last)* in Abhängigkeit vom Ausführungsverlauf von *fetch* vornehmen
- fetch* ▶ setzt *tail* ggf. auf *&LINK(head)*, obwohl ggf. ein weiteres (zweites) Element an die Liste angehängt wurde
- ▶ muss das Zurücksetzen des Fußzeigers (*tail*) in Abhängigkeit vom Ausführungsverlauf von *aback* vornehmen

Fallstudie: Warteschlangenmanipulation (Forts.)

Lösungsansatz: Problemspezifisches Protokoll zwischen *aback* und *fetch*

Idee ist es, den Verkettungszeiger (*link*) eines Listenelements auch zur „Signalisierung“ zwischen *aback* und *fetch* zu nutzen

- ▶ sei *that* Zeiger (*chain_t**) auf ein Listenelement, dann soll gelten:

$$LINK(that) = \begin{cases} that & \text{that ist gültiges Verkettungsglied} \\ 0 & \text{Verkettungsglied } that \text{ wurde entfernt} \\ \dots & \text{Listenelement } that \text{ mit Nachfolger} \end{cases}$$

- ▶ für die beiden Funktionen lässt sich dies dann wie folgt ausnutzen:
- aback* ▶ hängt das neue Element nicht an $\Leftrightarrow LINK(that) \neq that$
- ▶ beachte: *that* \mapsto *last* \Rightarrow *last* ist Verkettungsglied
- fetch* ▶ nullt *LINK(that)* $\Leftrightarrow LINK(that) = that$, und
- ▶ versucht *tail* zurückzusetzen $\Leftrightarrow LINK(that) = 0$
- ▶ beachte: *that* \mapsto *node* \Rightarrow *node* ist Verkettungsglied
- ▶ beachte: *last = node*, beide Funkt. sehen dasselbe Verkettungsglied

Fallstudie: Warteschlangenmanipulation (Forts.)

Korrektur des Kopfzeigers

Aufmerksamkeit ist noch dem **Kopfzeiger** (*LINK(head)*) zu geben, dessen Integrität im Konfliktfall (*last = node*) wieder herzustellen ist

- fetch*||*aback* ▶ *fetch* stellt fest, dass *node* das einzige Element war
- ▶ dann gilt $LINK(head) = 0 \wedge LINK(node) = node$
- ▶ zuvor jedoch hängt *aback* noch ein neues Element an
- ▶ dann gilt $LINK(node) \neq node \Rightarrow$ wird nicht genullt
- ▶ *LINK(node)* zeigt auf das soeben angehängte Element
- \Rightarrow *fetch*: *LINK(head)* ist auf *LINK(node)* zu korrigieren
- aback*||*fetch* ▶ *last* zeigt auf das Verkettungsglied, *tail* ist umgesetzt
- ▶ *aback* stellt fest, dass *last* soeben entfernt wurde
- ▶ dann gilt $LINK(last) = 0 \Rightarrow$ *item* hängt nicht an
- ▶ *fetch* setzt *tail* in der Situation jedoch nicht zurück
- ▶ es gilt: $tail = \&LINK(item) \wedge LINK(head) = 0$
- \Rightarrow *aback*: *LINK(head)* ist auf *item* zu korrigieren

Fallstudie: Warteschlangenmanipulation (Forts.)

Zusammenfügen: Wettlaufstolerantes *aback*

```
void nbs_aback (queue_t *this, chain_t *item) {
    chain_t *last, *self;

    item->link = item;    /* item becomes new tail resp. next last */

    do self = (last = this->tail)->link; /* draw copies as needed */
    while (!CAS(&this->tail, last, &item->link));

    if (!CAS(&last->link, self, item)) /* last removed by fetch */
        this->head.link = item;    /* item becomes new head */
}
```

Beachte \leftrightarrow Problem ABA

(Selbststudium)

- ▶ ggf. ist *item* ein *chain_p** und *wheel* entsprechend zu applizieren
- ▶ ebenso wäre dann *index* auf allen relevanten Stellen anzuwenden

Fallstudie: Warteschlangenmanipulation (Forts.)

Zusammenfügen: Wettlaufstolerantes *fetch*

```

chain_t *nbs_fetch (queue_t *this) {
    chain_t *node, *next;

    do if ((node = this->head.link) == 0) return 0;
    while (!CAS(&this->head.link, node,
                ((next = node->link) == node ? 0 : next)));

    if (next == node) { /* last element just removed, be careful */
        if (!CAS(&node->link, next, 0)) this->head.link = node->link;
        else CAS(&this->tail, &node->link, &this->head);
    }

    return node;
}

```

Beachte ↔ Problem ABA

(Selbststudium)

- ▶ ggf. ist *node* ein *chain_p** und *index* entsprechend zu applizieren

Fallstudie: Warteschlangenmanipulation (Forts.)

Plausibilitätskontrolle: Neuralgische Punkte

```

void nbs_aback (queue_t *this, chain_t *item) {
    chain_t *last, *self;

    item->link = item;

    t do self = (last = this->tail)->link;
    t while (!CAS(&this->tail, last, &item->link));

    1 if (!CAS(&last->link, self, item))
    2   this->head.link = item;
    }

chain_t *nbs_fetch (queue_t *this) {
    chain_t *node, *next;

    h do if ((node = this->head.link) == 0) return 0;
    h while (!CAS(&this->head.link, node,
                ((next = node->link) == node ? 0 : next)));

    3 if (next == node) {
    4   if (!CAS(&node->link, next, 0))
    5     this->head.link = node->link;
    6   else CAS(&this->tail, &node->link, &this->head);
    }

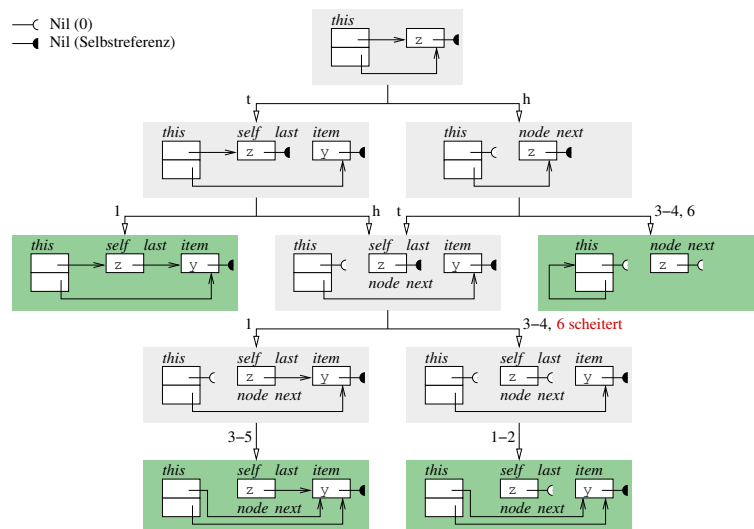
    return node;
}

```

- t ▶ *tail* weitersetzen ✓
- h ▶ *head* weitersetzen ✓
- 1 ▶ Verkettungsglied gültig?
 - ▶ ja, *item* angehängt
 - ▶ *fetch* wurde signalisiert
- 2 ▶ nein, *fetch* kam vorbei
- ▶ *head* korrigieren
- 3 ▶ letztes Element raus?
- 4 ▶ wirklich keins mehr da?
 - ▶ ja, *aback* signalisiert: 6
- 5 ▶ nein, *aback* kam vorbei
- ▶ *head* korrigieren
- 6 ▶ *tail* ggf. zurücksetzen

Fallstudie: Warteschlangenmanipulation (Forts.)

Plausibilitätskontrolle: Datenstrukturentwicklung je nach Überlappungsfall



Nichtsequentielle Programme sind nicht einfach...

Rekapitulation am Beispiel wichtiger, noch nicht behandelter Aspekte:

Ausführungspfadanalyse

- ▶ welcher Aufwand trifft den *Normalfall*?
- ▶ welcher Zusatzaufwand den *Ausnahmefall*?
- ▶ wann „lohnt“ sich welches Verfahren?

Fortschrittsgarantie

- ▶ ist das Vorankommen nichtsequentieller Programmausführung sichergestellt?

Anmerkung

- ▶ die Untersuchungen sind nicht exakt: keine Taktzyklen/Laufzeiten
- ▶ um ein Gefühl zu bekommen reicht es, Maschinenbefehle zu zählen
- ▶ auch soll es genügen, sich Umgebungseinflüsse vor Augen zu halten

Ausführungspfadanalyse

Referenz zum Vergleich mit NBS: FAD mit „annotierten“ kritischen Abschnitten

```
void fad_aback (queue_t *this, chain_t *item) {
    item->link = 0;

    ENTER(fad);
    this->tail->link = item;
    this->tail = item;
    LEAVE(fad);
}
```

Schutzoptionen

nil unsynchronisiert
ice Fortsetzungssperre
irq Unterbrechungssperre

```
chain_t *fad_fetch (queue_t *this) {
    chain_t *item;

    ENTER(fad);
    if ((item = this->head.link) && !(this->head.link = item->link))
        this->tail = &this->head;
    LEAVE(fad);

    return item;
}
```

Ausführungspfadanalyse (Forts.)

Anzahl der Maschinenbefehle (x86): gesamt/davon im kritischen Abschnitt

Konfiguration	<i>aback</i>		<i>fetch</i>	
	-	—	-	—
unsynchronisiert	7/0	7/0	10/0	11/0
Unterbrechungssperre	9/4	9/4	12/8	14/10
Fortsetzungssperre	12/4 + δ_1	12/4 + δ_1	19/8 + δ_1	21/10 + δ_1
NBS	24/0 + δ_2	25/0 + δ_2	25/0 + δ_2	35/0 + δ_2

- kürzester Pfad, — längster Pfad

δ_1 Aufwand für die Abarbeitung zurückgestellter Fortsetzungen: *clear*

2 Befehle für die Aufrufsequenz

12 Befehle Verwaltungsgemeinkosten

6 Befehle für Entnahme und Aktivierung *einer* Fortsetzung

▶ Entnahme: plus 10/11 Befehle des unsynchronisierten *fetch*

δ_2 Aufwand für die Wiederholungsversuche bei gescheitertem CAS

8/14 Befehle pro Versuch bei *aback*/*fetch*

Fortschrittsgarantie

Aussagen zur **Lebendigkeit** (engl. *liveliness*)⁷ von Programmen zu treffen, die nichtsequentiell ablaufen, braucht **Umgebungswissen**

- ▶ alle hier betrachteten Fälle zeigten **nichtsequentielle Programme**
- ▶ in ihnen spiegeln sich Annahmen zu Überlappungsmustern wider
- ▶ die durch das Operationsprinzip des Rechensystems begründet sind
 - ▶ Unterbrecher- und Uni- bzw. Multiprozessorbetrieb

Ansätze zum Nachweis von Lebendigkeitseigenschaften nichtsequentieller Programme gibt es bereits seit geraumer Zeit

There is a rather large body of sad experience to indicate that a concurrent program can withstand very careful scrutiny without revealing its errors. [10, S. 456]

- ▶ formale Methoden sind eigentlich der einzig richtige Weg zum Ziel
 - ▶ wären Betriebssysteme (universal/spezial) nur nicht so komplex

⁷In der Informatik (fälschlicherweise) auch: *liveness*.

Unterbrechungs-/Fortsetzungssperren

Fortschrittsgarantie für die Weiterleitung zurückgestellter Arbeitsaufträge (d.h., Nachspänne oder Fortsetzungen/Einfädelungen)

- ▶ leitet sich ab aus der WCET des gesperrten kritischen Abschnitts, für den eine Eintrittsanforderung vorliegt⁸
- ▶ ergibt sich aus der Annahme, dass die diesen Abschnitt ausmachende Anweisungsfolge des nichtsequentiellen Programms terminiert

Beachte ↔ δ_1 verlängert die Laufzeit serialisierter KA nicht

- ▶ die betrachteten Sperrverfahren garantieren Fortschritt, da ihre Korrektheit nicht von der der kritischen Abschnitte abhängt
- ▶ die Sperrverfahren „benutzen“ die kritischen Abschnitte nicht

⁸Die in Kap. 5 vorgestellte Unterbrechungsweitergabe sowie die damit verbundene Nachspannfreistellung (*clear*) zeigt Merkmal eines kritischen Abschnitts, der nämlich zu einem Zeitpunkt nur einmal aktiv sein darf, um Stapelüberlauf vorzubeugen.

Unterbrechungstransparente Synchronisation

Fortschrittsgarantie für Warteschlangenoperationen in einer Umgebung, die ein rollenspezifisches bzw. -verteiltes Überlappungsmuster vorgibt:

$\left. \begin{array}{l} \text{aback} \\ \text{fetch} \end{array} \right\}$ ist nur von Überlappung durch *aback* betroffen, ggf.

- ▶ die Ausführung von *aback* wird ereignisbedingt, mehrstufig ausgelöst
- ▶ und: nichtsequentielle Programmausführung im Uniprozessorbetrieb

Beachte \leftrightarrow **Programmschleifen** (S. 6-26/6-27)

- ▶ aus den Algorithmen allein heraus lässt sich keine WCET bestimmen
- ▶ sie ist erst durch Analyse des Unterbrecherbetriebs abschätzbar

Unterbrechungstransparente Synchronisation (Forts.)

Analyse des Unterbrecher- und Treiberbetriebs

- IPL_{max}
- ▶ höchst mögliche Unterbrechungsebene der Hardware
 - ▶ $0 < max < N$
 - ▶ $0 \rightsquigarrow 1 \rightsquigarrow 2 \rightsquigarrow \dots \rightsquigarrow N$ ist schlimmster Fall
 - ▶ d.h., vertikal kompletter Durchlauf, von unten nach oben
- IPL_{tip}
- ▶ Anzahl der Gerätetreiber einer Ebene tip
 - ▶ $0 \leq tip \leq IPL_{max}$
 - ▶ $0 \rightsquigarrow 1 \rightsquigarrow 2 \rightsquigarrow \dots \rightsquigarrow IPL_{tip}$ ist schlimmster Fall in Ebene tip
 - ▶ d.h., horizontal kompletter Durchlauf (Abfragebetrieb)
- TIP_{dev}
- ▶ Nachspannauslösungen von Treiber dev pro Unterbrechung
 - ▶ $0 \leq dev \leq M$: $M = 1$ bei Einfachauslösung, $\delta > 1$ sonst
- TIP_{max}
- ▶ maximale Länge der Warteschlange von Nachspännen
 - ▶ $\sum TIP_{dev}$ aller Treiber aller Unterbrechungsebenen

Beachte $\leftrightarrow \delta_1$

- ▶ TIP_{max} ist auch untere Grenze für die Anzahl von Fortsetzungen

Unterbrechungstransparente Synchronisation (Forts.)

Obergrenze für die Anzahl der Durchläufe der kritischen Programmschleifen

Konfiguration	IPL_{max}	$\sum IPL_{tip}$	TIP_{dev}	TIP_{max}
\emptyset PC _{Linux}	15	17	1	17
\approx MacBook	15	20	1	20
AX*	7	7	1	7
PEACE* _{SUPRENUM}	8	8	1	8
PEACE _{MANNA}	1	3	1	3
PURE _{LKW}	\sim	10	1	10
PURE _{Wetterstation}	2	2	1	2
CiAO-AS	2	2	1	2
OOSTuBS	15	2	1	2
MPStuBS	15	3	1	3
EZStubs	15	1	1	1

* Nicht jede Unterbrechungsebene hat Treiber, mindestens aber einen Abfänger

Nichtblockierende Synchronisation

Fortschrittsgarantie für Warteschlangenoperationen in einer Umgebung, die ein rollenspezifisches Überlappungsmuster vorgibt:

- ▶ Operationsausführungen können sich beliebig überlappen
 - Wartestapel \mapsto *ahead* und *strip*
 - Warteschlange \mapsto *aback* und *fetch*
- ▶ nichtsequentielle Programmausführung: Uni-/Multiprozessorbetrieb

Beachte $\leftrightarrow \delta_2$ und **Programmschleifen** (S. 6-37 bzw. 6-55/6-56)

- ▶ aus den Algorithmen allein heraus lässt sich keine WCET bestimmen
 - ▶ sie ist erst durch Analyse der Programmabläufe abschätzbar
- ▶ festgehalten werden kann, dass sich die Verfahren **sperrfrei** verhalten
 - ▶ die Ausführung des nichtsequentiellen Programms schreitet voran
 - ▶ allerdings lässt die Aushungerung von Fäden nicht ausschließen
- ▶ bei bekannter WCET aller Fadenabläufe wären sie sogar **wartefrei**

Nichtblockierende Synchronisation (Forts.)

Differenzierungen von Fortschrittsgarantien

behinderungsfrei (engl. *obstruction-free*)

- ▶ ein einzelner, in Isolation ablaufender Faden wird seine Operation in begrenzter Anzahl von Schritten beenden
- ▶ ein Faden läuft in Isolation ab, wenn alle ihn behindern könnenden anderen Fäden in der Ausführung zurückgestellt sind

sperrfrei (engl. *lock-free*), umfasst Behinderungsfreiheit

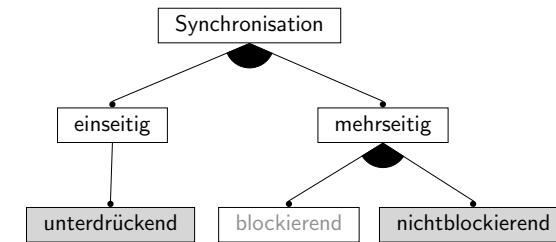
- ▶ jeder bei Ablauf eines Fadens auszuführende Schritt trägt dazu bei, dass das nichtsequentielle Programm insgesamt voranschreitet
- ▶ garantiert systemweiten Durchsatz, erlaubt jedoch Aushungerung

wartefrei (engl. *wait-free*), umfasst Sperrfreiheit

- ▶ die Anzahl der zur Beendigung einer Operation bei Fadenabläufen auszuführenden Schritte ist begrenzt
- ▶ garantiert systemweiten Durchsatz und ist frei von Aushungerung

Merkmaldiagramm

Synchronisation mit Elementaroperationen der Ebene₂

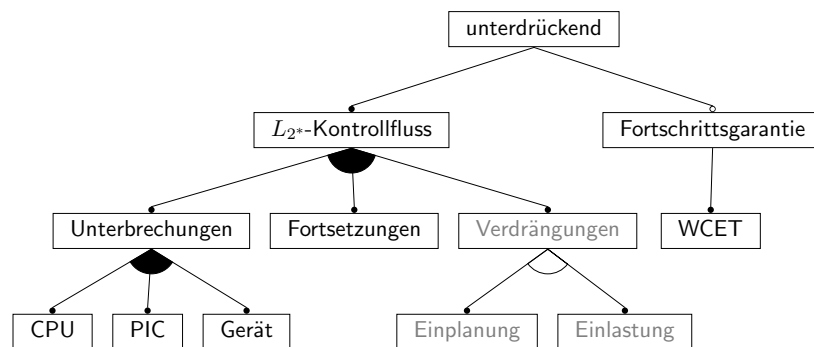


Beachte ↔ **blockierende Synchronisation**

- ▶ bezogen auf das Abstraktionsniveau der Befehlssatzebene heißt das:
 - Schlossvariable (engl. *lock variable*) und
 - aktives Warten (engl. *busy waiting, spin locking*)
- ▶ Blockadefreiheit ist damit grundsätzlich ausgeschlossen
 - ▶ weshalb Verfahren dazu auch nicht weiter untersucht wurden...

Merkmaldiagramm (Forts.)

Einseitiger Ausschluss der Ausführung eines Ebene₂-Programms

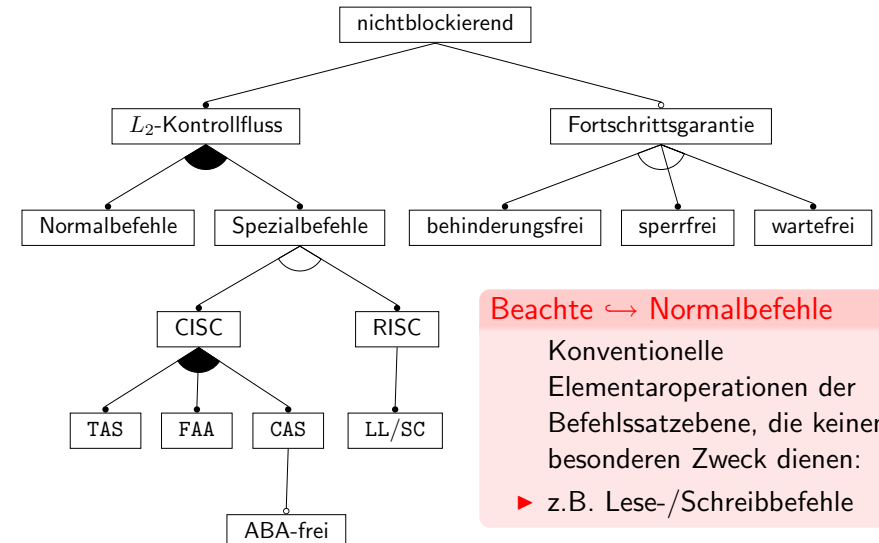


Beachte ↔ **Verdrängungen**

- ▶ diese „benutzen“ Abstraktionen zur Prozessverwaltung (z.B. Fäden)
- ▶ die auf Ebene der Ablaufsteuerung jedoch nicht bekannt sind
 - ▶ es sei denn, die Befehlssatzebene würde ein Prozesskonzept bieten

Merkmaldiagramm (Forts.)

Einseitiger Ausschluss der Ausführung eines Ebene₂-Befehls



Beachte ↔ **Normalbefehle**

- Konventionelle Elementaroperationen der Befehlssatzebene, die keinem besonderen Zweck dienen:
 - ▶ z.B. Lese-/Schreibbefehle

Resümee

Wettlaufintoleranz ↔ Abwehr von Ereignisanforderungen

- ▶ einseitiger Ausschluss der Ausführung eines Ebene₂-Programms
- ▶ Unterbrechungssperre, total/partiell
- ▶ Fortsetzungssperre, Fortsetzungen/Einfädelungen

Wettlauftoleranz ↔ nichtblockierende Synchronisation

- ▶ Spezialfall: Unterbrechungstransparente Synchronisation
- ▶ einseitiger Ausschluss der Ausführung eines Ebene₂-Befehls
- ▶ Einzelwort „*compare and swap*“
- ▶ Problem „ABA“, Etikettierung von Zeiger, Generationszähler

Fortschrittsgarantien ↔ WCET

- ▶ Weiterleitung zurückgestellter Arbeitsaufträge
- ▶ Reparatur inkonsistent gewordener Datenstrukturen
- ▶ Behinderungs-, Sperr- und Wartefreiheit

Literaturverzeichnis

- [1] Ralf Guido Herrtwich and Günter Hommel.
Kooperation und Konkurrenz — Nebenläufige, verteilte und echtzeitabhängige Programmsysteme.
Springer-Verlag, 1989.
- [2] Wolfgang Schröder-Preikschat.
Softwaresysteme 1.
<http://www4.informatik.uni-erlangen.de/Lehre/SOS1>, 2004.
- [3] Wolfgang Schröder-Preikschat.
Systemprogrammierung.
<http://www4.informatik.uni-erlangen.de/Lehre/SP>, 2008.
- [4] Daniel Lohmann.
Betriebssysteme.
<http://www4.informatik.uni-erlangen.de/Lehre/BS>, 2007.

Literaturverzeichnis (Forts.)

- [5] David D. Clark.
The structuring of systems using upcalls.
ACM SIGOPS Operating Systems Review, 19(5):171–180, 1985.
- [6] Friedrich Schön, Wolfgang Schröder-Preikschat, Olaf Spinczyk, and Ute Spinczyk.
On interrupt-transparent synchronization in an embedded object-oriented operating system.
In *Proceedings of the 3rd IEEE International Symposium on Object-Oriented Real-Time Distributed Computing (ISORC '00)*, pages 270–277. IEEE Computer Society Press, March 2000.
- [7] Theodore P. Baker.
Stack-based scheduling of realtime processes.
Real-Time Systems, 3(1):67–99, 1991.

Literaturverzeichnis (Forts.)

- [8] LEO GmbH.
LEO Deutsch-Englisches Wörterbuch.
<http://dict.leo.org>.
- [9] Barbara H. Liskov and Stephen N. Zilles.
Programming with abstract data types.
ACM SIGPLAN Notices, 9(4):50–59, April 1974.
- [10] Susan Owicki and Leslie Lamport.
Proving liveness properties of concurrent programs.
ACM Transactions on Programming Languages and Systems (TOPLAS), 4(3):455–495, July 1982.