

Wechselseitiger Ausschluß (mutual exclusion)

- Aufgabenstellung

„Mehrere Prozesse möchten etwas tun, wobei garantiert sein muß, daß es zu einem Zeitpunkt nur einer darf.“

- Voraussetzungen

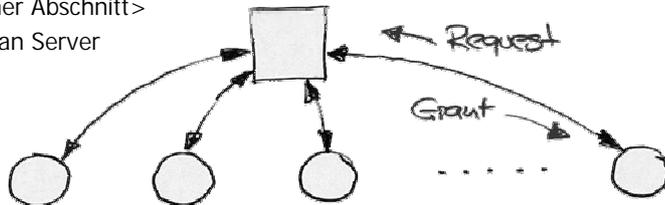
- exklusive Betriebsmittel
- „abstrakte“ exklusive Betriebsmittel (z.B. Termin in einem verteilten TVS)

- Lösung bei gemeinsamen Speicher

- Semaphore ⇒ Interessiert uns hier nicht!

Zentralisierter Ansatz

- Zentraler Server verwaltet Zustand
- Client
 - Request an Server
 - Warten auf Grant
 - <Kritischer Abschnitt>
 - Release an Server



- Bemerkungen
 - 3 Nachrichten
 - SPOF

Verteilte Systeme, Winter 2003

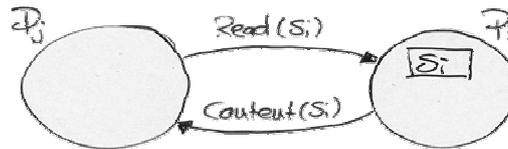
Folie 7.3

7.1 Lösungen mittels State-Variablen

State-Variable

- Jeder Prozeß P_i hat zwei Arten von Variablen
 - Lokale Variablen
 - Nur P_i kann auf sie zugreifen (lesend und schreibend)
 - State-Variablen
 - P_i darf lesend und schreibend zugreifen
 - $\forall j : j \neq i : P_j$ darf lesend zugreifen

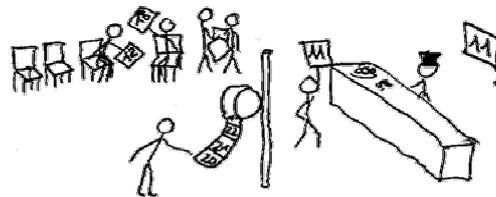
- Einfache Umsetzbarkeit in nachrichten-basierten Systemen



Verteilte Systeme, Winter 2003

Folie 7.5

„Bäckerei-Algorithmus“ (Lamport '74)



- Algorithmus für P_i

```

var number : array [0...n-1] of integer;
number[i] ← 1 + max(number[0], ..., number[n-1]);
for j=0 to n-1 (i ≠ j) {
    wait (number[j] + 0) D
    (number[i],i) < (number[j],j);
}

<kritischer Abschnitt>
number[i] ← 0;
                    
```

$$(a,b) < (c,d) \hat{=} (a < c) \hat{\vee} ((a = c) \hat{\vee} (b < d))$$

Verteilte Systeme, Winter 2003

Folie 7.6

Bemerkungen

- Berechnung der Reihenfolge durch
 - `number[i] = 1 + max(number[0], ..., number[n-1]);`
- Algorithmus garantiert wechselseitigen Ausschluß
 - `number[i]` besitzt eindeutiges Minimum
 - Sei P_i Prozeß mit diesem Minimum
 - \Rightarrow Für $P_i : \forall j : j \neq i : \text{wait}(\text{true})$
 - \Rightarrow Für $P_j : j \neq i : \exists \text{wait}(\text{false});$
- Algorithmus ist fair
 - Prozeß P_i muß maximal $(n-1)$ mal warten
- Algorithmus ist korrekt?
 - Fast
 - Problem: Berechnung von `number[i]` ist nicht atomar!

Verteilte Systeme, Winter 2003

Folie 7.7

Szenario: Ein wechselseitiger Nicht-Ausschluß

Initialisierung: `number[i] = 0`

P_1

```
number[0] == 0
number[1] == 0
```

⋮

```
number[0] = 1 + max(0,0);
```

```
for ...
wait (number[1] != 0) true
P (number[0],0) < number[1],1) true
```

⋮

Im kritischen Abschnitt

P_2

```
number[0] == 0
number[1] == 0
number[1] = 1 + max(0,0);
```

⋮

```
for ...
wait (number[0] != 0) false
```

⋮

Im kritischen Abschnitt



Verteilte Systeme, Winter 2003

Folie 7.8

„Bäckerei-Algorithmus“, iterierte Fassung

- Algorithmus für Prozeß i

```

var number : array [0...n-1] of integer;
    choice : array [0...n-1] of boolean;

choice[i] ← true;
number[i] ← 1 + max (number[0],...,number[n-1]);
choice[i] ← false;

for j=0 to n-1, i≠j {
    wait ¬choice[j];
    wait (number[j] ≠ 0)
        ▷ (number[i],i) < (number[j],j);
}

<kritischer Abschnitt>

number[i] ← 0;
    
```

Szenario: Ein wechselseitiger Ausschluß

Initialisierung:
 number[i] = 0
 choice[i] = false



```

choice[0] = true
number[0] == 0
number[1] == 0
    
```



```

number[0] = 1 + max(0,0);
choice[0] = false
    
```

```

for ...
    wait ¬choice[1] true
    wait (number[1] != 0) true
    ▷ (number[0],0) < (number[1],1) true
    ...
    Im kritischen Abschnitt
    
```



```

choice[1] = true
number[0] == 0
number[1] == 0
number[1] = 1 + max(0,0);
choice[1] = false
    
```



```

for ...
    wait ¬choice[0] false
    wait (number[0] != 0) false
    
```

P2 blockiert

Bemerkungen

- Algorithmus ist symmetrisch
- Fehlertoleranz
 - Es muß „bloß“ garantiert sein, daß bei Absturz von P_i jeder andere Prozeß P_j davon ausgeht: P_i ist nicht im kritischen Abschnitt
 - Deadlock-Gefahr
 - Prozeß P_i stürzt immer bei der Maximumbildung ab!
 - ⇒ $choice_i = true$
 - ⇒ Alle Prozesse warten an $wait \neg choice_i$;
- Nachrichtenaufwand
 - Hängt davon ab, wie man die State-Variablen realisiert
 - $\sim 2n + 4 \cdot (n - 1)$

Verteilte Systeme, Winter 2003

Folie 7.11

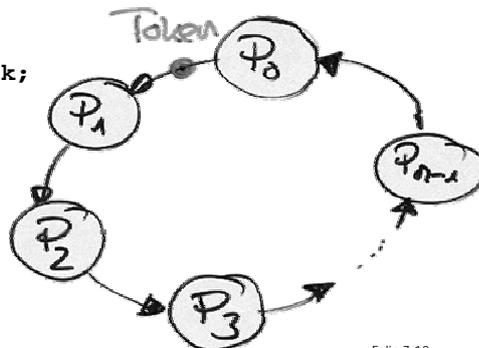
Dijkstra's selbst-stabilisierender Algorithmus (1974)

- Sei $k > 1$

```
var flag : array [0...n-1] of 0...k;
```
- P_0

```
wait flag[0] = flag[n-1];
<kritischer Abschnitt>
flag[0] ← (flag[0]+1) mod k;
```
- P_i ($i \neq 0$)

```
wait flag[i] ≠ flag[i-1];
<kritischer Abschnitt>
flag[i] ← flag[i-1];
```



Verteilte Systeme, Winter 2003

Folie 7.12

Bemerkungen

- Annahme: sei $\forall i : \mathbf{flag}[i] = 0$
- Algorithmus garantiert wechselseitigen Ausschluß
 - Schritt 1
 - P0: **wait** ($\mathbf{flag}[0] = \mathbf{flag}[n-1]$) true
 - Pi: **wait** ($\mathbf{flag}[i] \neq \mathbf{flag}[i-1]$) false $\forall i : i \neq 0$
 - P0 betritt kritischen Abschnitt
 - danach $\mathbf{flag}[0] = 1$
 - Schritt 2
 - P0 : **wait** ($\mathbf{flag}[0] = \mathbf{flag}[n-1]$) false
 - P1 : **wait** ($\mathbf{flag}[1] \neq \mathbf{flag}[0]$) true
 - Pi : **wait** ($\mathbf{flag}[i] \neq \mathbf{flag}[i-1]$) false $\forall i : i \neq 0, i \neq 1$
 - P1 betritt kritischen Abschnitt
 - ...
- Algorithmus verhindert Deadlocks ✓
- Algorithmus ist fair

Verteilte Systeme, Winter 2003

Folie 7.13

... aber

- Algorithmus ist nicht symmetrisch
- Prozesse bekommen das Privileg, obwohl sie gar nicht in den kritischen Abschnitt wollen!
 - ⇒ Prozesse müssen Token immer weitergeben
- Fehlertoleranz
 - Abstürzende Prozesse müssen aus dem Ring entfernt werden!
 - Wie man entfernt, hängt vom Prozeß P ab
 - $P_0 \leftarrow \rightarrow$ Anderer wird P0
 - $P_i \leftarrow \rightarrow$ Raus aus Ring
 - ⇒ Zusätzliche Informationen für den Ringaustag

Verteilte Systeme, Winter 2003

Folie 7.14

Selbst-Stabilisierung

	P ₀	P ₁	P ₂	P ₃	
■ Initiale Bedingung $\forall i: \text{flag}[i] = 0$ unnötig	0	1	0	1	
					3 ↯
	0	0	1	0	
					3 ↯
■ Korrektheit hängt von k ab	1	0	0	1	
– k ≤ n: Kein WA garantiert					3 ↯
■ Beispiel: n = 4, k = 2	0	1	0	0	
					3 ↯
– k > n: WA garantiert	1	0	1	0	
■ P ₀ inkrementiert (wenn er darf)					3 ↯
■ P _i schieben Token nach rechts	1	1	0	1	
■ „Müll“ wird nach rechts geschoben und durch sinnvolle Inhalte in P ₀ ergänzt					3 ↯
■ Initialisierungsdauer berechenbar	0	1	1	0	
					3 ↯
	0	1	0	1	

Verteilte Systeme, Winter 2003
Folie 7.15

Bemerkungen: State-Variablen

- Einfache, verhältnismäßig fehlertolerante Algorithmen
- 2 Ebenen werden durcheinander gebracht
 - Logik
 - Implementierung

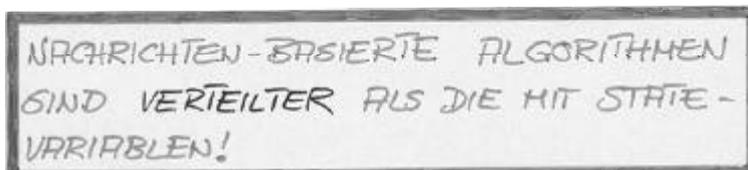
⇒ schwer beschreibbar & schwer verständlich
- Zugriff auf State-Variablen
 - Blindes Testen von Statevariablen
 - Hohe Anzahl an unnötigen Nachrichten
- Effiziente Realisierung in „shared-memory“-System
- Vorteil gegenüber z.B. Semaphore
 - keine globalen Variablen
 - Erhöhte Fehlertoleranz

Verteilte Systeme, Winter 2003
Folie 7.16

7.2 Nachrichten-basierte Lösungen

Nachrichten-basierte Ansätze

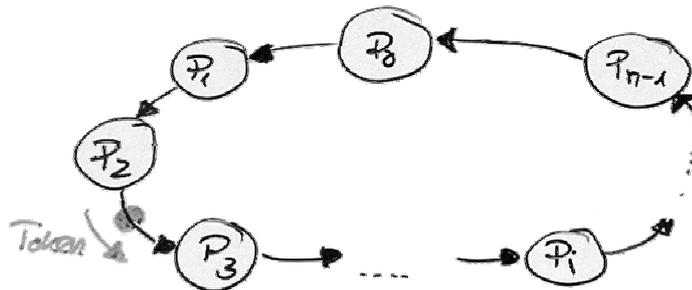
- State-Variablen können in einer nachrichten-basierten Umgebung realisiert werden
 - „Request Informationen“-Nachrichten (Pull)
 - Viele z.T. unnötige Nachrichten
- Push-Ansätze
 - „Send Informationen“-Nachrichten
 - Anzahl der Nachrichten wird minimiert
 - Keine unnützen Nachrichten
 - Alle Nachrichten signalisieren Zustandswechsel



NACHRICHTEN-BASIERTE ALGORITHMEN
SIND VERTEILT ER ALS DIE MIT STATE-
VARIABLEN!

Le Lann's Algorithmus (1977)

- Natürlich gibt es auch hier eine Token-Ring Lösung



- Prozeß i
 - Warte auf Token von P_{i-1}
 - <Kritischer Abschnitt>
 - Sende Token an P_{i+1}

Verteilte Systeme, Winter 2003

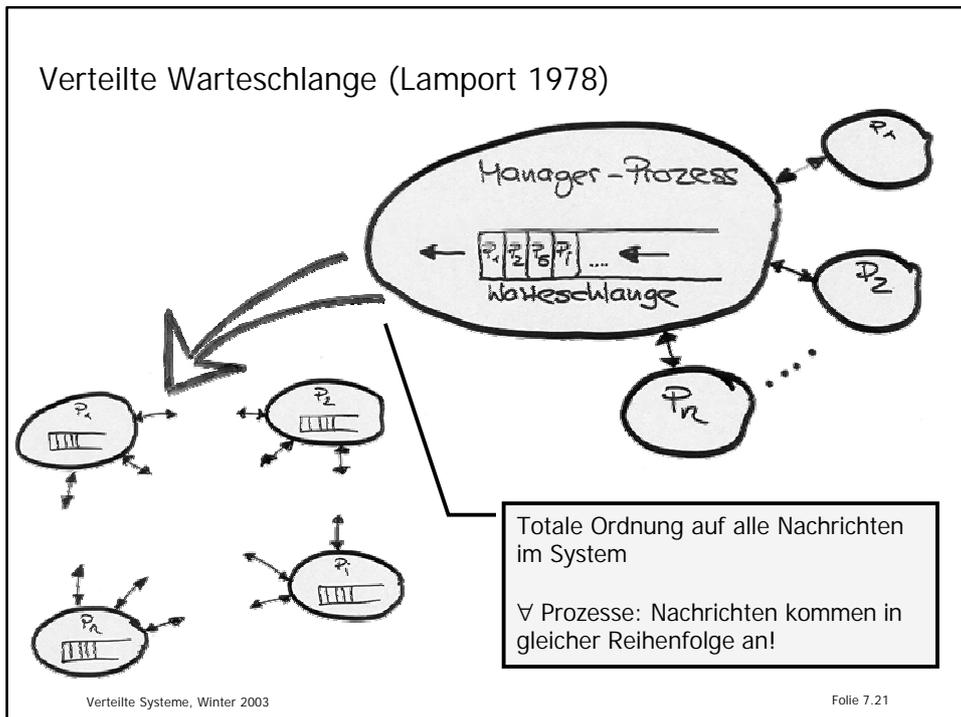
Folie 7.19

Bemerkungen

- Klar
 - Garantiert wechselseitigen Ausschluß
 - Verhindert Deadlocks
 - Fair
- Nachteile
 - Anzahl Nachrichten: $1 \dots + \infty$
 - Prozeß bekommt Privileg, obwohl er es nicht braucht!
 - Fehlertoleranz
 - Token-Verlust \Rightarrow Genau 1 neues Token
 - Prozeß-Absturz \Rightarrow Ring verkürzen

Verteilte Systeme, Winter 2003

Folie 7.20



- ### Implementierungsansatz
- Verteilung der Nachrichten
 - Request: Multicast an alle Teilnehmer
 - Grant: Unicast Jedes Teilnehmers zurück an Sender
 - Release: Multicast an alle Teilnehmer
 - Nachrichten enthalten Zeitstempel
 - Totale Ordnung \Rightarrow Erweiterte Lamport-Zeit
 - Einsortieren gemäß Zeitstempel
 - Annahmen
 - FIFO-Eigenschaft jeweils zwischen P_i und P_k
 - Lamport-Algorithmus
 - Implementierung?
- Verteilte Systeme, Winter 2003 Folie 7.22

Bemerkungen

- Algorithmus ist
 - Fair
 - Deadlock frei (Totale Ordnung enthält keine Zyklen)
 - Korrekt (Garantiert wechselseitigen Ausschluß)
- Nachrichten-Komplexität
 - 2 mal Multicast + (n-1) mal Unicast
 - Multicast = n Unicast ~ 3(n-1) Nachrichten
 - Multicast = 1 Unicast ~ n+1 Nachrichten
- Allgemein $O(n)$

Ein optimaler Algorithmus

- Ricart & Agrawala's Algorithmus (1981)
 - Optimierung von Lamport's Algorithmus
- Prozess P_i möchte in den kritischen Abschnitt
 - **Sende (Request, clock_i, i) an alle**
 - **Warte auf n-1 Grants**
 - **<kritischer Abschnitt>**
- Prozess P_k empfängt Request von P_i
 - Sendet Reply sofort, wenn
 - P_k nicht in den kritischen Abschnitt will!
 - P_k zwar in den kritischen Abschnitt will, aber P_i höhere Rechte hat
 - Verzögert das Reply, wenn
 - P_k in den kritischen Abschnitt will und höhere Rechte als P_i hat

Ein optimaler Algorithmus

- Ricart & Agrawala's Algorithmus (1981)
 - Optimierung von Lamport's Algorithmus
- Prozess P_i möchte in den kritischen Abschnitt
 - **Sende (Request, clock $_i$, i) an alle**
 - **Warte auf n-1 Grants**
 - **<kritischer Abschnitt>**
 - **Gib alle verzögerten Replies raus**
- Nachrichtenkomplexität: $2(n-1)$

Verteilte Systeme, Winter 2003

Folie 7.25

Ein optimaler Algorithmus

- In dieser Richtung sind dann einige Verbesserungen vorgenommen worden

– Le Lann (1977)	1 ... + ∞
– Lamport (1978)	3 (n-1)
– Ricart & Agrawala (1981)	2 (n-1)
– Carvalho & Roucairol	0 ... (2(n-1))
– Suzuki & Kasami (1982)	n
– Ricart & Agrawala (1983)	n

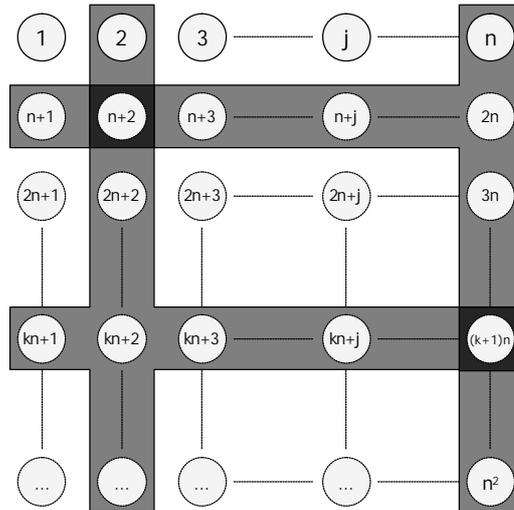


Verteilte Systeme, Winter 2003

Folie 7.26

Ein optimaler Algorithmus: Maekawa (1985)

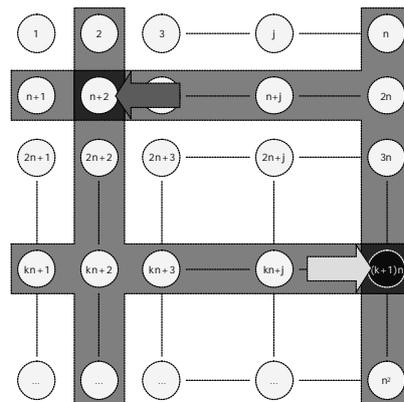
- Idee
 - Anordnung in Gitter
 - Prozeß i fragt Gruppe S_i von anderen Prozessen um Erlaubnis
 - $S_i \cap S_j \neq \emptyset \iff \forall i \neq j$



Verteilte Systeme, Winter 2003

Folie 7.27

Ein optimaler Algorithmus



- Deadlock-Gefahr
- Nachrichten für wechselseitigen Ausschluß
 - REQUEST, LOCKED, RELEASE
- Nachrichten zur Deadlock-Vermeidung
 - INQUIRE, RELINQUISH, FAILED

Verteilte Systeme, Winter 2003

Folie 7.28

Algorithmus

- P_i möchte kritischen Abschnitt betreten
 - **Sende REQUEST an alle Elemente von S_i ;**
 - **Warte auf LOCKED von allen Prozessen aus S_i ;**
 - **<kritischer Abschnitt>**
 - **Sende RELEASE an alle Prozesse aus S_i ;**
- $P_j \in S_i$ empfängt REQUEST-Nachricht
 - **if (frei) {**
 - Zustand \rightarrow locked;**
 - Sende LOCKED-Nachricht an P_i ;**
 - }**
 - else { /* Abschnitt ist von $P_k \hat{=} S_j$ belegt */**
 - if (P_j kennt aussichtsreicheren REQUEST-Kandidaten)**
 - Sende FAILED-Nachricht an P_i ;**
 - else**
 - Sende INQUIRE-Nachricht an P_k ;**
 - /* Warte auf RELINQUISH oder RELEASE von P_k */**
 - }**

Verteilte Systeme, Winter 2003

Folie 7.29

Algorithmus (contd.)

- $P_k \in S_i$ empfängt INQUIRE-Nachricht von P_j
 - **if (Prozess $\hat{=} S_k$ hat FAILED an P_k geschickt) {**
 - Hebe bereits angelaufene LOCK's auf;**
 - Sende RELINQUISH an P_j ;**
 - }**
 - else**
 - Sende RELEASE nach Verlassen des K_A ;**
- weitere Funktionen?
- Eigenschaften
 - Algorithmus garantiert wechselseitigen Ausschluß
 - Auflösung von Deadlocks
 - Fairness
 - Fehlertoleranz
 - Übernahme der RG-Menge eines Prozesses P_i
 - Das Wissen von P_i ist in $O(\sqrt{n})$ Nachrichten wieder herleitbar

Verteilte Systeme, Winter 2003

Folie 7.30

Nachrichtenkomplexität

- Grundversion

$$c \cdot \sqrt{n} \quad c \in [6 \dots 10]$$

- Verbesserungen
 - „Request-Grant“-Mengen im Gitter sind nicht minimal
 - Problem aus der projektiven Geometrie
 - Für bestimmte n sind minimale Si möglich
($|S_i|$ = Potenz einer Primzahl)
- Im minimalen Fall ist die Nachrichten-Komplexität

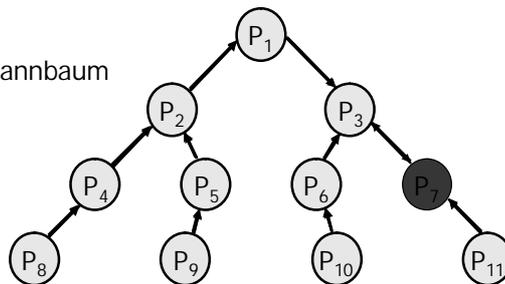
$$c \cdot \sqrt{n} \quad c \in [3 \dots 5]$$

Verteilte Systeme, Winter 2003

Folie 7.31

Ein optimaler Algorithmus: Raymond (1989)

- Idee: Topologie ist ein Spannbaum
- Ein Prozess P_i besitzt das Privileg („Token“)
 - kann den kritischen Abschnitt betreten
- Andere Prozesse kennen Richtung zum Token
- Beim Tokenwechsel werden Richtungen aktualisiert
 - Einfache Pointeroperationen beim Wechseln
- Requests werden in Tokenrichtung resendet
 - Letztendlich kommt Token an
 - Token kann dem Request nicht ewig „davonrennen“



Verteilte Systeme, Winter 2003

Folie 7.32

Nachrichten und Datenstrukturen

- Nachrichten
 - REQUEST : Ein Prozess will das Token
 - PRIVILEG : Hier ist das Privileg!

- Lokale Variablen für Prozeß Pi
 - HOLDER: „self“ oder Name eines direkten Nachbarn in dessen Richtung das Privileg liegt
 - USING : bool = true wenn Pi im kritischen Abschnitt ist (USING = true \Rightarrow Holder = self)
 - REQUEST-QUEUE = FIFO-Schlange angekommener Request-Nachrichten von Prozessen, die das Privileg noch nicht hatten
 - Enthält nur „self“ und Namen direkter Nachbarn
 - ASKED : bool = true \Leftrightarrow Request nach Privileg in Richtung HOLDER ist bereits verschickt worden.

Algorithmus

- Prozess Pi will Privileg
 - Füge „self“ an REQUEST-QUEUE an;
 - if (\emptyset ASKED) {
 - sende REQUEST an HOLDER;
 - ASKED \leftarrow true;
 - }

- Prozess Pk empfängt Request von Nachbar Pi
 - Füge „Pi“ an REQUEST-QUEUE an;
 - if (\emptyset ASKED) {
 - sende REQUEST an HOLDER;
 - ASKED \leftarrow true;
 - }

Algorithmus (contd.)

- Prozess P_k hat Privileg und will es loswerden
 - `if (|REQUEST-QUEUE| > 0) {`
 - `X ← Head(REQUEST-QUEUE);`
 - `Holder ← X;`
 - `sende PRIVILEGE an HOLDER;`
 - `asked ← using ← false;`
 - `}`
- Prozess P_k empfängt Privileg von P_m
 - `X ← Head(REQUEST-QUEUE);`
 - `if (X = self) {`
 - `Using ← true;`
 - `<kritischer Abschnitt>`
 - `}`
 - `else {`
 - `HOLDER ← X;`
 - `sende Privilege an HOLDER;`
 - `}`

Verteilte Systeme, Winter 2003

Folie 7.35

Bemerkungen

- Allgemein
 - Garantiert wechselseitigen Ausschluß
 - Verhindert Deadlock's: Spannbaum ist azyklisch
⇒ Privileg trifft nach endlicher Zeit auf einen Request.
 - Fair! (FIFO-Queue's)
- Nachrichten-Komplexität

$$O(\log_k n)$$
 - k = Grad des Baumes

Verteilte Systeme, Winter 2003

Folie 7.36

Bemerkungen (contd.)

- Im Artikel werden weitere Verbesserungen vorgeschlagen
- Initialisierung
 - Ein Prozeß wird privilegiert
 - Sendet INITIALIZE an alle Nachbarn;
 - Alle Prozesse, die INITIALIZE empfangen
 - Setzen $HOLDER \leftarrow \text{INIT-Sender}$;
 - Senden INITIALIZE an alle Nachbarn;
- Fehlertoleranz

Verteilte Systeme, Winter 2003

Folie 7.37

Vergleich der vorgestellten Algorithmen

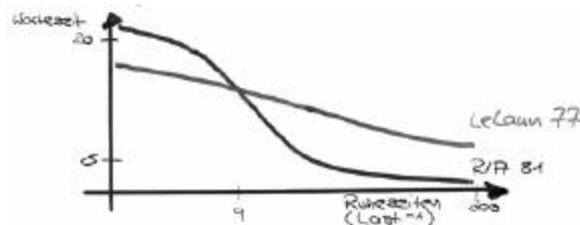
- | | |
|----------------------------|--------------------|
| ■ Le Lann (1977) | $1 \dots + \infty$ |
| ■ Lamport (1978) | $3 (n-1)$ |
| ■ Ricart & Agrawala (1981) | $2 (n-1)$ |
| ■ Carvalho & Roucairol | $0 \dots (2(n-1))$ |
| ■ Ricart & Agrawala (1983) | n |
| ■ Maekawa (1985) | $c \cdot \sqrt{n}$ |
| ■ Raymond (1989) | $\log_k n$ |
- Unterschiede
 - Symmetrische und asymmetrische Algorithmen

Verteilte Systeme, Winter 2003

Folie 7.38

Welchen Algorithmus soll man nehmen

- DEN optimalen Algorithmus gibt es nicht!
- Güte hängt nicht allein vom „Worst-Case“-Verhalten ab
 - Empirische Untersuchungen, da „mittleres“ Verhalten fast nie theoretisch bestimmt werden kann
 - z.B. Ausführliche Simulation aller Algorithmen bis 1983
- Simulationsparameter
 - Anzahl der Knoten, Ruhezeiten, Belegzeiten



Verteilte Systeme, Winter 2003

Folie 7.39

Nochmal wechselseitigen Ausschluß

- Idee der „Request-Grant“-Mengen kann erweitert werden!
 - B. A. Sanders (1987)
- 2 Mengen/Prozeß
 - Inform-Menge
 - Alle Prozesse, die informiert werden müssen, wenn ein bestimmter Prozeß den kritischen Abschnitt verläßt
 - Request-Menge
 - Alle Prozesse, die gefragt werden müssen, wenn ein bestimmter Prozeß in den kritischen Abschnitt will.
- Anpassbar an bestimmte Anwendungen
- Bestimmte Anforderungen an die Inform- und Request-Mengen, um wechselseitigen Ausschluß zu garantieren

Verteilte Systeme, Winter 2003

Folie 7.40

Literatur

- Michel Raynal
Algorithmus for mutual Exclusion
MIT Press 1986
(State-Variablen und nachrichten-
basierte Algorithmen)
- Mamoru Maekawa
*A Θ n Algorithm for Mutual Exclusion
in Decentralized Systems*
ACM Transactions on Computer
Systems, Vol. 3, No. 2, May 1985,
pp. 145 – 159
- Kerry Raymond
*A Tree-Based Algorithm for
Distributed Mutual Exclusion*
ACM Transactions on Computer
Systems, Vol. 7, No. 1, February
1989, pp. 61 – 77
- Beverly A. Sanders
*The Information Structure of Distributed
Mutual Exclusion Algorithms*, ACM
Transactions on Computer Systems,
Vol. 5, No. 3, August 1987,
pp. 284 – 299