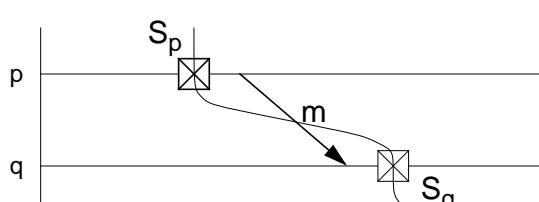
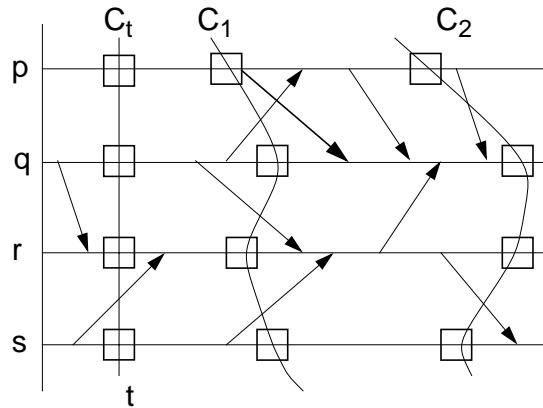


BP 1	Zustandssicherung: Inhalt
11	Zustandssicherung, Rücksetzen <i>Koo R.; Toueg, S.: Checkpointing and Rollback-Recovery for Distributed Systems. IEEE Trans. on Software Engineering, vol. SE-13, No. 1, Jan. 1987, pp. 23-31.</i>
11.1	Begriffsbildungen
11.2	Problem der Erzeugung von Sicherungspunkten
11.3	Algorithmus von Koo und Toueg
17.01.02	Universität Erlangen-Nürnberg, Lehrstuhl für Informatik 4 (Verteilte Systeme und Betriebssysteme), F. Hofmann Reproduktion jeder Art oder Verwendung dieser Unterlage zu Lehrzwecken außerhalb der Universität Erlangen-Nürnberg ist ohne Genehmigung des Autors unzulässig
	11.1-1

BP 1	Zustandssicherung: Begriffsbildungen
11.1	Begriffsbildungen Konsistenter globaler Zustand: Globaler Zustand, den das System nach bisheriger Beobachtung eingenommen haben kann.
	Beispiel eines inkonsistenten Zustands
	
	S_p: Zustand von p (das Versenden von m ist nicht aufgezeichnet)
	S_q: Zustand von q (der Empfang von m ist aufgezeichnet)
	Der globale Zustand (S_p, S_q) ist nicht konsistent.
17.01.02	Universität Erlangen-Nürnberg, Lehrstuhl für Informatik 4 (Verteilte Systeme und Betriebssysteme), F. Hofmann Reproduktion jeder Art oder Verwendung dieser Unterlage zu Lehrzwecken außerhalb der Universität Erlangen-Nürnberg ist ohne Genehmigung des Autors unzulässig
	11.1-2

Kanalzustände



C_t : Realzeit-Schnitt (Zustand des Systems zum Zeitpunkt t)

C_1 : Schnitt, der zu einem global konsistenten Zustand gehört

C_2 : Inkonsistenter Schnitt

Zustand des Kommunikationssystems bezüglich C_1

Je eine Nachricht in den Kanälen $\langle q, p \rangle$, $\langle q, r \rangle$ und $\langle s, r \rangle$

17.01.02

Universität Erlangen-Nürnberg, Lehrstuhl für Informatik 4 (Verteilte Systeme und Betriebssysteme), F. Hofmann
Reproduktion jeder Art oder Verwendung dieser Unterlage zu Lehrzwecken außerhalb der Universität Erlangen-Nürnberg
ist ohne Genehmigung des Autors unzulässig

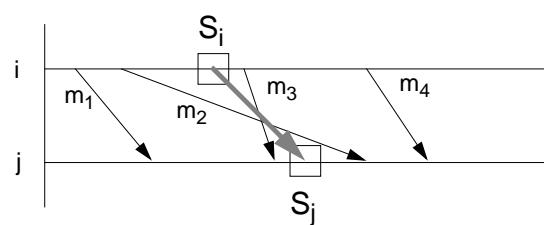
11.1-3

Informale Definition

Ein globaler Zustand $S = (S_1, S_2, \dots, S_i, \dots, S_j, \dots, S_n)$ heißt konsistent, wenn für alle i und j gilt, daß eine von i an j gesendete Nachricht, die "vor S_j " empfangen wurde, nicht "nach S_i " gesendet worden sein kann.

Zustand eines Kanals $\langle i, j \rangle$ bezüglich eines globalen Zustandes S:

Alle Nachrichten die "vor S_i " gesendet wurden, aber erst "nach S_j " empfangen werden.



17.01.02

Universität Erlangen-Nürnberg, Lehrstuhl für Informatik 4 (Verteilte Systeme und Betriebssysteme), F. Hofmann
Reproduktion jeder Art oder Verwendung dieser Unterlage zu Lehrzwecken außerhalb der Universität Erlangen-Nürnberg
ist ohne Genehmigung des Autors unzulässig

11.1-4

BP 1**Zustandssicherung: Erzeugung von Sicherungspunkten****11.2****Erzeugung von Sicherungspunkten****Methode 1:**

Prozesse sichern unabhängig voneinander lokale Zustände in einem stabilen Speicher. Nach einem Fehler wird eine konsistente Menge lokaler Zustandssicherungen bestimmt.

Methode 2:

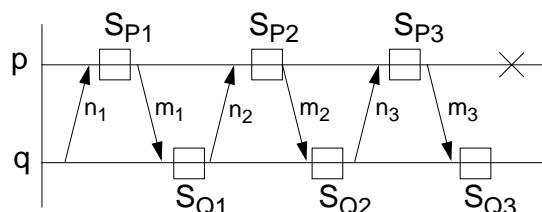
Die Prozesse koordinieren die Zeitpunkte für Zustandssicherungen derart, daß jeder Prozeß nur jeweils die jüngste Sicherung aufbewahrt. Durch die Koordinierung wird dafür gesorgt, daß die Gesamtheit der jüngsten Sicherungen einen konsistenten globalen Zustand darstellt.

17.01.02

Universität Erlangen-Nürnberg, Lehrstuhl für Informatik 4 (Verteilte Systeme und Betriebssysteme), F. Hofmann
Reproduktion jeder Art oder Verwendung dieser Unterlage zu Lehrzwecken außerhalb der Universität Erlangen-Nürnberg
ist ohne Genehmigung des Autors unzulässig

11.2-5**BP 1****Zustandssicherung: Erzeugung von Sicherungspunkten****Probleme der Methode 1:**

- Großer Speicheraufwand
- Dominoeffekt



Deshalb weitere Untersuchung konzentriert auf Methode 2

17.01.02

Universität Erlangen-Nürnberg, Lehrstuhl für Informatik 4 (Verteilte Systeme und Betriebssysteme), F. Hofmann
Reproduktion jeder Art oder Verwendung dieser Unterlage zu Lehrzwecken außerhalb der Universität Erlangen-Nürnberg
ist ohne Genehmigung des Autors unzulässig

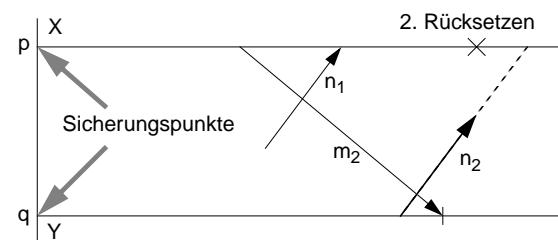
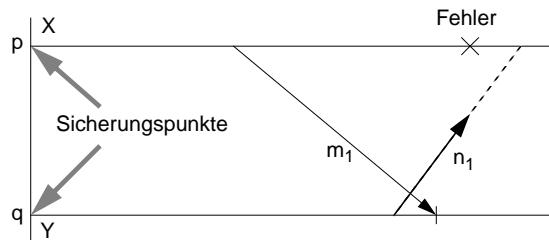
11.2-6

Rücksetzverfahren

Naive Vorgehensweise:

Jeder Prozeß setzt zu seinem zuletzt gesicherten Zustand zurück und nimmt dann die Arbeit wieder auf.

Problem: Livelock als Folge des nicht vollkommen zeitsynchronen Wiederanlaufs der Prozesse



Annahme:

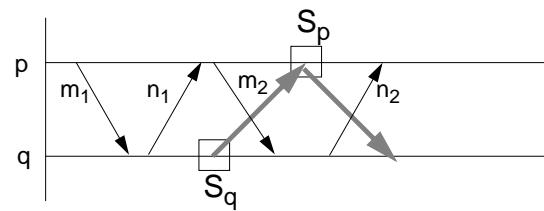
- Nur zuverlässige FIFO-Kanäle
- Keine Ausfälle während des Sicherns, Rücksetzens und Wiederanlaufs

Bemerkung:

- Der Algorithmus kann nebenläufig zur eigentlichen Anwendung abgewickelt werden!

Ablauf

- (a) Beim ersten Empfang einer Kontrollnachricht von q durch p und vor dem Senden weiterer Nachrichten:
- p sichert seinen Zustand S_p
 - p notiert den Kanal $\langle q, p \rangle$ als leer
 - p schickt eine Kontrollnachricht an alle seine Nachbarn (auch an den, von dem die betrachtete Kontrollnachricht kam)
- (b) Bei nachfolgendem Erhalt einer Kontrollnachricht von r :
- p notiert als Kanalzustand $\langle r, p \rangle$ die Sequenz der Nachrichten, die von r zwischen dem Festhalten von S_p und der eben erhaltenen Kontrollnachricht ankamen.

**Im weiteren:**

- Nur FIFO-Kanäle
- Ausfälle während des Rücksetzens und Wiederanlaufs sind möglich

Ein Algorithmus zur Bildung konsistenter Sicherungspunkte heißt **unverwüstlich** (resilient), wenn er auch bei Prozessorausfällen funktioniert.

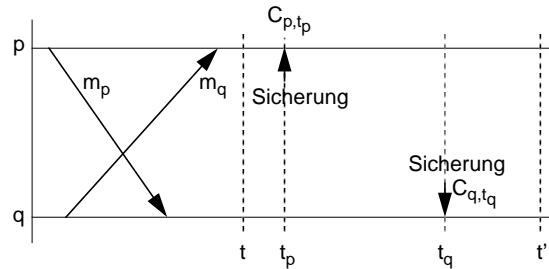
Eine (lokale) Zustandssicherung heißt **dauerhaft** (permanent), wenn sie nie rückgängig gemacht wird und (lokale) Berechnungen, die zur Erreichung des Sicherungspunktes notwendig waren, niemals wiederholt werden müssen.

Eine (lokale) Zustandssicherung heißt **versuchsweise** (tentative), wenn sie rückgängig oder dauerhaft gemacht werden kann.

Unverwüstliche Algorithmen zur Bildung konsistenter Sicherungspunkte können nicht nur dauerhafte lokale Zustandssicherungen vornehmen.

BP 1**Zustandssicherung: Erzeugung von Sicherungspunkten****Beweis:**

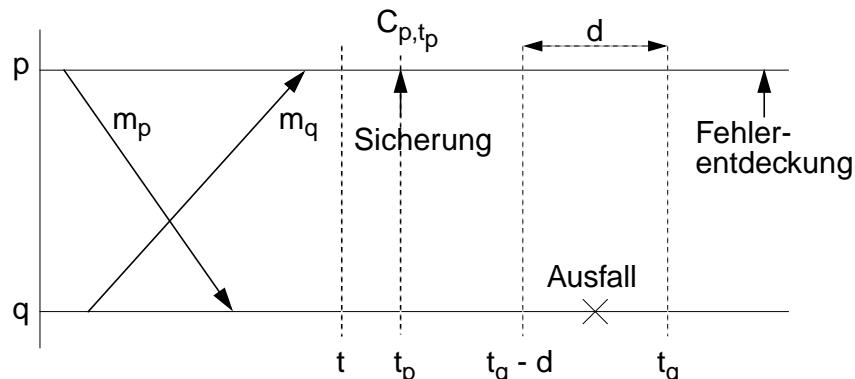
- Angenommen es existiere ein solcher Algorithmus A.
- Zum Zeitpunkt $t > 0$ habe p bereits die Nachricht m_q von q empfangen und q die Nachricht m_p von p.
- Prozeß p rufe zum Zeitpunkt t den Algorithmus A auf und er werde zum Zeitpunkt t' beendet.
- Der Prozeß p bilde zum Zeitpunkt t_p mit $t < t_p \leq t'$ die lokale Sicherung C_{p, t_p} .
- A unverwüstlich und C_{p, t_p} dauerhaft
 \Rightarrow Prozeß q besitzt dauerhafte Sicherung C_{q, t_q} mit $t < t_q \leq t'$.



17.01.02

Universität Erlangen-Nürnberg, Lehrstuhl für Informatik 4 (Verteilte Systeme und Betriebssysteme), F. Hofmann
 Reproduktion jeder Art oder Verwendung dieser Unterlage zu Lehrzwecken außerhalb der Universität Erlangen-Nürnberg ist ohne Genehmigung des Autors unzulässig

11.2-11

BP 1**Zustandssicherung: Erzeugung von Sicherungspunkten****Annahme: Einer der Prozesse fällt aus****Die minimale Zeit zur Entdeckung eines Fehlers sei d.**(a) $t_p \leq t_q$ **Man betrachte einen Ausfall von q in $(\max(t, t_q - d), t_q)$**  $\Rightarrow C_{p, t_p}$ angelegt, aber nicht C_{q, t_q} , da q vorher ausgefallen ist. \Rightarrow Dieser Fall kann nicht auftreten.

17.01.02

Universität Erlangen-Nürnberg, Lehrstuhl für Informatik 4 (Verteilte Systeme und Betriebssysteme), F. Hofmann
 Reproduktion jeder Art oder Verwendung dieser Unterlage zu Lehrzwecken außerhalb der Universität Erlangen-Nürnberg ist ohne Genehmigung des Autors unzulässig

11.2-12

BP 1	Zustandssicherung: Erzeugung von Sicherungspunkten <hr/> <p>(b) $t_p > t_q$</p> <p>Man betrachte einen Ausfall von p in $(\max(t, t_p - d), t_p)$. Eine zu a) analoge Überlegung führt auch in diesem Fall zu einem Widerspruch.</p>
17.01.02	Universität Erlangen-Nürnberg, Lehrstuhl für Informatik 4 (Verteilte Systeme und Betriebssysteme), F. Hofmann Reproduktion jeder Art oder Verwendung dieser Unterlage zu Lehrzwecken außerhalb der Universität Erlangen-Nürnberg ist ohne Genehmigung des Autors unzulässig

BP 1	Zustandssicherung: Algorithmus von Koo und Toueg <hr/> <p>11.3</p> <p>Algorithmus von Toueg</p> <p>Eigenschaften des Algorithmus:</p> <p>Unverwüstlich, d.h. Ausfälle bei Zustandssicherung oder Rücksetzen werden toleriert</p> <p>Minimale Zahl von Prozessen ist involviert</p> <p>Grundmuster des Algorithmus:</p> <p>Analog 2-Phasen-Commit</p> <ol style="list-style-type: none"> 1. Phase: Veranlassung versuchsweiser Sicherungspunkte 2. Phase: Verbreitung und Durchsetzung der Entscheidung <p>Spezielle Zustandsgröße:</p> <p>$willing_to_ckpt_p$ ($\in \{"yes", "no"\}$)</p> <p>Beschränkung in der Versendung von Basisnachrichten:</p> <p>Zwischen dem Anlegen einer lokalen Sicherungskopie und dem Empfang der Entscheidung des Initiators werden keine Basisnachrichten versandt.</p>
17.01.02	Universität Erlangen-Nürnberg, Lehrstuhl für Informatik 4 (Verteilte Systeme und Betriebssysteme), F. Hofmann Reproduktion jeder Art oder Verwendung dieser Unterlage zu Lehrzwecken außerhalb der Universität Erlangen-Nürnberg ist ohne Genehmigung des Autors unzulässig

Zustandssicherung

Benutzernachrichten m enthalten eine vom Sender aufsteigend vergebene Sequenznummer $m.l$ mit $\text{MIN} < m.l < \text{MAX}$

$R_{q,p} :=$ Menge aller Nachrichten, die q seit der letzten dauerhaften Sicherung seines lokalen Zustandes von p erhalten hat.

$$\text{last_rmsg}_q(p) := \begin{cases} \max(\{m.l \mid m \in R_{q,p}\}) & \text{falls } \{m.l \mid m \in R_{q,p}\} \neq \emptyset \\ \text{MIN} & \text{falls } \{m.l \mid m \in R_{q,p}\} = \emptyset \end{cases}$$

$S_{q,p} :=$ Menge aller Nachrichten, die q seit der letzten versuchsweisen oder dauerhaften Sicherung seines lokalen Zustandes an p gesendet hat.

$$\text{first_smsg}_q(p) := \begin{cases} \min\{m.l \mid m \in S_{q,p}\} & \text{falls } \{m.l \mid m \in S_{q,p}\} \neq \emptyset \\ \text{MIN} & \text{falls } \{m.l \mid m \in S_{q,p}\} = \emptyset \end{cases}$$

$p \in \text{ckpt_cohort}_q := q$ hat eine versuchsweise lokale Sicherungskopie angelegt und $\text{last_rmsg}_q(p) > \text{MIN}$

Kontrollnachrichten:

- "take a tentative checkpoint and $\text{last_rmsg}_q(p)$ "
- "make tentative checkpoint permanent"
- "undo tentative checkpoint"

Daemon process

```

send(initiator, "take a tentative checkpoint and MAX");
await(initiator, willing_to_ckptinitiator);
if (willing_to_ckptinitiator == "yes")
    send(initiator, "make tentative checkpoint permanent");
else
    send(initiator, "undo tentative checkpoint");

```

INITIAL STATE

```

first_smsgp(daemon) = MAX;
willing_to_ckptp = "yes", falls p gewillt ist eine Sicherungskopie anzulegen, und
                    "no" sonst;

```

All processes p

```

UPON RECEIPT OF "take a tentative checkpoint and last_rmsgq(p)" from q DO {
    if (willing_to_ckptp and last_rmsgq(p) ≥ first_smsgp(q) > MIN) {
        take_a_tentative_checkpoint;
        for (all r ∈ ckpt_cohortp)
            send(r, "take a tentative checkpoint and last_rmsgp(r)");
        for (all r ∈ ckpt_cohortp) await(r, willing_to_ckptr);
        if (exists r ∈ ckpt_cohortp (willing_to_ckptr == "no"))
            willing_to_ckptp = "no";
    }
    send(q, willing_to_ckptp);
}

```

UPON RECEIPT OF m == "make tentative checkpoint permanent"

or m == "undo tentative checkpoint" DO {

if (m == "make tentative checkpoint permanent")

make_tentative_checkpoint_permanent;

else

undo_tentative_checkpoint;

for (all r ∈ ckpt_cohort_p) **send**(r, m);

}

BP 1	Zustandssicherung: Koo; Toueg - Zustandssicherung
<p>o erbt eine Sicherungsaufforderung von q, wenn er von q eine Nachricht "take a tentative checkpoint and last_rmsg_q(p)" erhält und das Prädikat $willing_to_ckpt_p$ und $last_rmsg_q(p) \geq first_smsg_p(q) > MIN$ erfüllt ist.</p>	
17.01.02	<p>L11.1 Jeder Prozeß akzeptiert höchstens eine Aufforderung zum Anlegen einer versuchsweisen Sicherungskopie.</p> <p>Beweis:</p> <p>Prozesse versenden vom Zeitpunkt des Anlegens einer Sicherungskopie bis zum Empfang der vom Initiator ausgehenden Entscheidung keine Benutzernachrichten, so daß in dieser Zeit $first_smsg_p(q) = MIN$ ist. Also werden weitere Aufforderungen zum Anlegen einer Sicherungskopie nicht akzeptiert.</p>

11.3-19

BP 1	Zustandssicherung: Koo; Toueg - Zustandssicherung
<p>L11.2 Jede lokale Ausführung des Algorithmus terminiert.</p> <p>Beweis:</p> <p>(a) Falls die Aufforderung zum Anlegen einer Sicherungskopie nicht akzeptabel ist, terminiert der Algorithmus</p> <p>(b) Falls p den Auftrag akzeptiert, legt es genau eine versuchsweise Sicherungskopie an. Es ist also zu zeigen, daß p nach Anlegen der Sicherungskopie terminiert.</p> <ul style="list-style-type: none"> - $q \in ckpt_cohort_p \Rightarrow$ auf jeden Fall $send(p, willing_to_ckpt)$, bevor q auf die Entscheidung des Initiators wartet. Es kann daher keinen Wartezyklus geben, in dem p auf die Antwort von q wartet und q auf die Entscheidung des Initiators. - Wenn q von p den Auftrag zum Anlegen einer Sicherungskopie akzeptiert, so hat vorher p akzeptiert. Da wegen L11.1 die Akzeptiert-von-Relation nicht zirkulär sein kann, ist eine zyklische Wartesituation nicht möglich. <p>(c) Durchsetzen der Entscheidung ist trivial, da die Kanäle und Prozesse als zuverlässig vorausgesetzt sind.</p>	
17.01.02	<p>Universität Erlangen-Nürnberg, Lehrstuhl für Informatik 4 (Verteilte Systeme und Betriebssysteme), F. Hofmann Reproduktion jeder Art oder Verwendung dieser Unterlage zu Lehrzwecken außerhalb der Universität Erlangen-Nürnberg ist ohne Genehmigung des Autors unzulässig</p>

11.3-20

BP 1	Zustandssicherung: Koo; Toueg - Zustandssicherung
<p>L11.3 Wenn die dauerhaften Sicherungskopien vor Beginn des Algorithmus konsistent waren, so gilt dies auch nach Beendigung des Algorithmus.</p> <p>Beweis:</p> <p>(a) Falls keine neuen dauerhaften Sicherungskopien entstanden, ist die Aussage trivial.</p> <p>(b) Annahme: Es wurden neue dauerhafte Sicherungskopien gebildet und der Zustand ist inkonsistent.</p> <p>Dann existieren p und q derart, daß p nach Bildung seiner letzten dauerhaften Sicherungskopie eine Benutzernachricht m an q gesendet hat, die q vor Bildung seiner letzten dauerhaften Sicherungskopie erhalten hat. Da die alten Sicherungskopien konsistent waren, muß q seine neueste Sicherungskopie während des Ablaufs des Algorithmus gebildet haben. Vor der versuchsweisen Bildung muß demnach $\text{last_rmsg}_q(p) \geq \text{m.l}$ gewesen sein.</p> <p>Also galt $p \in \text{ckpt_cohort}_q$ und p wurde von q aufgefordert eine versuchsweise Kopie anzulegen. Als p diese Aufforderung erhielt, muß $\text{willing_to_ckpt}_p = \text{"yes"}$ gewesen sein, da sonst q seinerseits seine versuchsweise Sicherungskopie nicht dauerhaft gemacht hätte.</p>	

17.01.02

Universität Erlangen-Nürnberg, Lehrstuhl für Informatik 4 (Verteilte Systeme und Betriebssysteme), F. Hofmann
 Reproduktion jeder Art oder Verwendung dieser Unterlage zu Lehrzwecken außerhalb der Universität Erlangen-Nürnberg ist ohne Genehmigung des Autors unzulässig

11.3-21

BP 1	Zustandssicherung: Koo; Toueg - Zustandssicherung
<p>Weiter muß p eine versuchsweise Kopie nach dem Versenden von m angelegt haben oder es muß $\text{last_rmsg}_q(p) \geq \text{first_smsg}_p(q) > \text{MIN}$ gewesen sein. In beiden Fällen muß p seine Sicherungskopie nach dem Versenden von m gebildet haben. Diese Sicherungskopie wird dauerhaft gemacht, wenn die von q dauerhaft gemacht wird. Also machte p nach dem Versenden von m seine versuchsweise Kopie dauerhaft im Widerspruch zur Annahme.</p>	

17.01.02

Universität Erlangen-Nürnberg, Lehrstuhl für Informatik 4 (Verteilte Systeme und Betriebssysteme), F. Hofmann
 Reproduktion jeder Art oder Verwendung dieser Unterlage zu Lehrzwecken außerhalb der Universität Erlangen-Nürnberg ist ohne Genehmigung des Autors unzulässig

11.3-22

BP 1**Zustandssicherung: Koo; Toueg - Zustandssicherung**

Es sei $P = \{p_0, p_1, \dots, p_k\}$ die Menge der Prozesse, die während der Durchführung des Algorithmus ihren lokalen Zustand sichern, wobei p_0 Initiator ist.

$C(P) = \{c(p_0), c(p_1), \dots, c(p_k)\}$ sei die Menge der neuen zugehörigen Sicherungen.

$C'(P) = \{c'(p_0), c'(p_1), \dots, c'(p_k)\}$ besitze die Eigenschaft, daß $c(p_0) = c'(p_0)$ ist und daß für $1 \leq i \leq k$ entweder $c'(p_i)$ gleich $c(p_i)$ oder $c'(p_i)$ gleich der letzten Sicherung vor Ausführung des Algorithmus ist.

17.01.02

Universität Erlangen-Nürnberg, Lehrstuhl für Informatik 4 (Verteilte Systeme und Betriebssysteme), F. Hofmann
Reproduktion jeder Art oder Verwendung dieser Unterlage zu Lehrzwecken außerhalb der Universität Erlangen-Nürnberg
ist ohne Genehmigung des Autors unzulässig

11.3-23

BP 1**Zustandssicherung: Koo; Toueg - Zustandssicherung****S11.2**

$C'(P)$ ist genau dann konsistent, wenn $C'(P) = C(P)$ ist.

Beweis:

\Leftarrow : **L11.3**

\Rightarrow : Die Relation **erbt eine Sicherungsaufforderung** ist nach L11.1 hierarchisch und, da es nur einen Initiator gibt, ein Baum T.

Auf Grund der Konstruktion enthält T alle Prozesse.

Annahme: $C'(P) \neq C(P)$ und $C'(P)$ konsistent. Sei $r \in P$ und $c'(r) \neq c(r)$.

Dann ist $(r \neq p_0)$ und es existiert in T ein Pfad von r nach p_0 . Auf diesem Pfad müssen p und q existieren derart, daß gilt:

p erbt eine Sicherungsaufforderung von q, $c'(p) \neq c(p)$ und $c'(q) = c(q)$.

Da p erbt, muß gegolten haben

$willing_to_ckpt_p$ and $last_rmsg_q(p) \geq first_smsg_p(q) > MIN$.

Sei m die Nachricht, die q von p erhalten hat und für die $last_rmsg_q(p) = m.l$ war.

Wegen $m.l \geq first_smsg_p(q)$ ist das Senden dieser Nachricht in $c'(p)$ nicht aufgezeichnet, aber ihr Empfang in $c'(q)$. Also ist $C'(P)$ nicht konsistent.

17.01.02

Universität Erlangen-Nürnberg, Lehrstuhl für Informatik 4 (Verteilte Systeme und Betriebssysteme), F. Hofmann
Reproduktion jeder Art oder Verwendung dieser Unterlage zu Lehrzwecken außerhalb der Universität Erlangen-Nürnberg
ist ohne Genehmigung des Autors unzulässig

11.3-24

BP 1**Zustandssicherung: Koo; Toueg - Zustandssicherung****Berücksichtigung von Prozessor-Ausfällen****Daemon process**

```

send(initiator, "take a tentative checkpoint and MAX");
await(initiator, willing_to_ckptinitiator);
if (willing_to_ckptinitiator == "yes")
    send(initiator, "make tentative checkpoint permanent");
else
    send(initiator, "undo tentative checkpoint");

```

INITIAL STATE

first_smsg_p(daemon) = MAX;
 willing_to_ckpt_p = "yes", falls p gewillt ist eine Sicherungskopie anzulegen, und
 "no" sonst;

17.01.02

Universität Erlangen-Nürnberg, Lehrstuhl für Informatik 4 (Verteilte Systeme und Betriebssysteme), F. Hofmann
 Reproduktion jeder Art oder Verwendung dieser Unterlage zu Lehrzwecken außerhalb der Universität Erlangen-Nürnberg
 ist ohne Genehmigung des Autors unzulässig

11.3-25

BP 1**Zustandssicherung: Koo; Toueg - Zustandssicherung****All processes p**

UPON RECEIPT OF "take a tentative checkpoint and last_rmsg_q(p)" from q DO {
if (willing_to_ckpt_p **and** last_rmsg_q(p) ≥ first_smsg_p(q) > MIN) {
 take_a_tentative_checkpoint;
for (all r ∈ ckpt_cohort_p)
send(r, "take a tentative checkpoint and last_rmsg_p(r)");
for (all r ∈ ckpt_cohort_p) **await**(r, willing_to_ckpt_r);
if (exists r ∈ ckpt_cohort_p ((willing_to_ckpt_r == "no") **or** has_failed(r)))
 willing_to_ckpt_p = "no";
 }
send(q, willing_to_ckpt_p);
}

UPON RECEIPT OF "make tentative checkpoint permanent"

or m == "undo tentative checkpoint" DO {

if (m == "make tentative checkpoint permanent")

make_tentative_checkpoint_permanent;

else

undo_tentative_checkpoint;

for (all r ∈ ckpt_cohort_p) **send**(r, m);

}

17.01.02

Universität Erlangen-Nürnberg, Lehrstuhl für Informatik 4 (Verteilte Systeme und Betriebssysteme), F. Hofmann
 Reproduktion jeder Art oder Verwendung dieser Unterlage zu Lehrzwecken außerhalb der Universität Erlangen-Nürnberg
 ist ohne Genehmigung des Autors unzulässig

11.3-26

Rücksetzen und Wiederanlauf

$S_{q,p}$:= Menge aller Nachrichten, die q vor der letzten dauerhaften Sicherung seines lokalen Zustandes an p gesendet hat.

$$\text{last_smsg}_q(p) := \begin{cases} \max(\{m.l \mid m \in S_{q,p}\}) & \text{falls } \{m.l \mid m \in S_{q,p}\} \neq \emptyset \\ \text{MAX} & \text{falls } \{m.l \mid m \in S_{q,p}\} = \emptyset \end{cases}$$

$p \in \text{roll_cohort}_q$:= q kann an p Nachrichten senden.

Spezielle Zustandsgröße:

willing_to_roll_p ($\in \{\text{"yes"}, \text{"no"}\}$)

Kontrollnachrichten:

- "prepare to roll back and $\text{last_smsg}_q(p)$ "
- "roll back"
- "do not roll back"

17.01.02

Universität Erlangen-Nürnberg, Lehrstuhl für Informatik 4 (Verteilte Systeme und Betriebssysteme), F. Hofmann
Reproduktion jeder Art oder Verwendung dieser Unterlage zu Lehrzwecken außerhalb der Universität Erlangen-Nürnberg
ist ohne Genehmigung des Autors unzulässig

11.3-27

Daemon process

```
send(initiator, "prepare to roll back and MIN");
await(initiator, willing_to_rollinitiator);
if (willing_to_rollinitiator == "yes")
    send(initiator, "roll back");
else
    send(initiator, "do not roll back");
```

INITIAL STATE

```
ready_to_rollp = "yes";
last_rmsgp(daemon) = MAX;
willing_to_rollp = "yes", falls p gewillt ist zurückzusetzen, und
                    "no" sonst;
```

17.01.02

Universität Erlangen-Nürnberg, Lehrstuhl für Informatik 4 (Verteilte Systeme und Betriebssysteme), F. Hofmann
Reproduktion jeder Art oder Verwendung dieser Unterlage zu Lehrzwecken außerhalb der Universität Erlangen-Nürnberg
ist ohne Genehmigung des Autors unzulässig

11.3-28

BP 1**Zustandssicherung: Koo; Toueg - Rücksetzen und Wiederanlauf****All processes p**

```

UPON RECEIPT OF "prepare to roll back and last_smsgq(p)" from q DO {
    if (willing_to_rollp and last_rmsgp(q) > last_smsgq(p) and ready_to_rollp) {
        ready_to_rollp = "no";
        for (all r ∈ roll_cohortp)
            send(r, "prepare to roll back and last_smsgp(r)");
        for (all r ∈ roll_cohortp) await(r, willing_to_rollr);
        if (∃ r ∈ roll_cohortp (willing_to_rollr == "no" or has_failed(r)))
            willing_to_rollp = "no";
    }
    send(q, willing_to_rollp);
}

UPON RECEIPT OF m == "roll back" or m == "do not roll back" and ready_to_rollp == "no" DO {
    if (m == "roll back")
        restart_from_p's_permanent_checkpoint;
    else
        resume_execution;
    for (all r ∈ roll_cohortp) send(r, m);
}

```

17.01.02

Universität Erlangen-Nürnberg, Lehrstuhl für Informatik 4 (Verteilte Systeme und Betriebssysteme), F. Hofmann
Reproduktion jeder Art oder Verwendung dieser Unterlage zu Lehrzwecken außerhalb der Universität Erlangen-Nürnberg
ist ohne Genehmigung des Autors unzulässig

11.3-29

BP 1**Zustandssicherung: Koo; Toueg - Rücksetzen und Wiederanlauf****L11.4**

Jeder Prozeß akzeptiert höchstens eine Aufforderung zur Rücksetzvorbereitung ("prepare to roll back").

Beweis:

Prozesse setzen nach Empfang einer Aufforderung zur Rücksetzvorbereitung ready_to_roll = false und ändern es bis zum Abschluß des Algorithmus nicht mehr. Also werden weitere Aufforderungen zum Anlegen einer Sicherungskopie nicht akzeptiert.

17.01.02

Universität Erlangen-Nürnberg, Lehrstuhl für Informatik 4 (Verteilte Systeme und Betriebssysteme), F. Hofmann
Reproduktion jeder Art oder Verwendung dieser Unterlage zu Lehrzwecken außerhalb der Universität Erlangen-Nürnberg
ist ohne Genehmigung des Autors unzulässig

11.3-30

BP 1	Zustandssicherung: Koo; Toueg - Rücksetzen und Wiederanlauf
L11.5	<p>Jede lokale Ausführung des Algorithmus R terminiert.</p> <p>Beweis:</p> <p>(a) Falls die Aufforderung zum Rücksetzen ("prepare to roll back") nicht akzeptabel ist, terminiert der Algorithmus.</p> <p>(b) Es ist also nur zu zeigen, daß p nach Annahme einer Rücksetzaufforderung terminiert (d. h. keine Verklemmungen auftreten).</p> <ul style="list-style-type: none"> - $q \in \text{roll_cohort}_p$ wartet auf die Entscheidung des Initiators ⇒ q hat <code>willing_to_roll</code> bereits gesendet. Es kann daher keinen Wartezyklus geben, in dem p auf die Antwort von q wartet und q auf die Entscheidung des Initiators. - Wenn q von p den Auftrag zum Rücksetzen akzeptiert, so hat vorher p akzeptiert. Da wegen L11.4 die Akzeptiert-von-Relation nicht zirkulär sein kann, ist eine zyklische Wartesituation nicht möglich. <p>(c) Durchsetzen der Entscheidung ist trivial, da die Kanäle und Prozesse als zuverlässig vorausgesetzt sind.</p>
17.01.02	<small>Universität Erlangen-Nürnberg, Lehrstuhl für Informatik 4 (Verteilte Systeme und Betriebssysteme), F. Hofmann Reproduktion jeder Art oder Verwendung dieser Unterlage zu Lehrzwecken außerhalb der Universität Erlangen-Nürnberg ist ohne Genehmigung des Autors unzulässig</small>

11.3-31

BP 1	Zustandssicherung: Koo; Toueg - Rücksetzen und Wiederanlauf
L11.6	<p>Wenn der Zustand des Systems vor Ausführung von R konsistent ist, so befindet es sich auch nach Ausführung von R in einem konsistenten Zustand.</p> <p>Beweis:</p> <p>Annahme: Entgegen der Behauptung sei nach Ablauf von R der Zustand inkonsistent.</p> <p>Dann müssen Prozesse p und q und eine Basisnachricht m von q an p existieren derart, daß während der Durchführung von R das Senden von m rückgängig gemacht wurde, aber nicht das Empfangen.</p> <p>q hat nach Empfang der Rücksetzaufforderung keine Basisnachrichten mehr verschickt. Also hat q an p einen Rücksetzauftrag erst nach dem Versenden von m verschickt.</p> <p>Da q das Senden von m rückgängig machte, muß $m.i > \text{last_msg}_q(p)$ sein.</p>
17.01.02	<small>Universität Erlangen-Nürnberg, Lehrstuhl für Informatik 4 (Verteilte Systeme und Betriebssysteme), F. Hofmann Reproduktion jeder Art oder Verwendung dieser Unterlage zu Lehrzwecken außerhalb der Universität Erlangen-Nürnberg ist ohne Genehmigung des Autors unzulässig</small>

11.3-32

BP 1	<p>Zustandssicherung: Koo; Toueg - Rücksetzen und Wiederanlauf</p> <p>Andererseits kann p keine dauerhafte Sicherungskopie angelegt haben nach Empfang von m und vor q's Rücksetzaufforderung: Die Erzeugung dieser Sicherungskopie zusammen mit der Tatsache, daß q keine neue dauerhafte Sicherungskopie angelegt hat, widersprüchen nämlich L10.3.</p> <p>Als q's Rücksetzaufforderung bei p ankam, mußte also p schon eine solche Aufforderung erhalten haben oder es akzeptierte die von q.</p> <p>Da p und q die gleiche Aufforderung bekamen, muß p zurückgesetzt haben.</p> <p>Fall 1: m erreichte p nach der als erste erhaltenen Rücksetzaufforderung .</p> <p>Da es sich um FIFO-Kanäle handelt, kam m bei p vor der Rücksetzaufforderung von q an. Der Initiator traf seine Entscheidung nachdem p auf q's Anfrage geantwortet hatte. Also setzte p nach Erhalt von m zurück. Widerspruch!</p> <p>Fall 2: m erreichte p vor der als erste erhaltenen Rücksetzaufforderung.</p> <p>Wie schon gezeigt, legte p keine dauerhafte Sicherungskopie nach Erhalt von m an. Also wird beim Rücksetzen von p der Empfang von m rückgängig gemacht.</p>
17.01.02	Universität Erlangen-Nürnberg, Lehrstuhl für Informatik 4 (Verteilte Systeme und Betriebssysteme), F. Hofmann Reproduktion jeder Art oder Verwendung dieser Unterlage zu Lehrzwecken außerhalb der Universität Erlangen-Nürnberg ist ohne Genehmigung des Autors unzulässig

BP 1	<p>Zustandssicherung: Koo; Toueg - Rücksetzen und Wiederanlauf</p> <p>S11.3 Es sei E eine Ausführung von R, in der Initiator p_0 und eine Menge P weiterer Prozesse zurücksetzen. Sei E' eine Ausführung, die mit E übereinstimmt bis auf die Tatsache, daß nur eine echte Teilmenge von P zurücksetzt. Dann hinterläßt E' das System in einem inkonsistenten Zustand.</p>
17.01.02	Universität Erlangen-Nürnberg, Lehrstuhl für Informatik 4 (Verteilte Systeme und Betriebssysteme), F. Hofmann Reproduktion jeder Art oder Verwendung dieser Unterlage zu Lehrzwecken außerhalb der Universität Erlangen-Nürnberg ist ohne Genehmigung des Autors unzulässig