

<b>BP 1</b>	<b>Einigungsprotokolle: Inhalt</b>
<b>12</b>	<b>Einigungsprotokolle</b>
<b>12.1</b>	<b>Erläuterung der Fragestellung</b>
<b>12.2</b>	<b>Unlösbarkeitsaussagen</b> <i>Fischer, M. J.; Lynch, N. A.; Merritt, M.: Easy impossibility proofs for distributed consensus problems. Distributed Computing, Vol. 1, 1986, pp. 26-39.</i> <i>Fischer, M. J.; Lynch, N. A.; Paterson, M. S.: Impossibility of Distributed Consensus with One Faulty Process. Journal of the ACM, Vol. 32, No. 2, 1985, pp. 374-382.</i>
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 12.1-1

<b>BP 1</b>	<b>Einigungsprotokolle: Inhalt</b>
<b>12.4</b>	<b>Paxos</b>
	<i>Lamport, L.: The Part-Time Parliament. Research Report 49, Digital Equipment Corporation Systems Research Center, Palo Alto, CA, September 1989.11</i>  <i>De Prisco, R.: Revisiting the Paxos Algorithm. M. S. Thesis, Massachusetts Institute of Technology, Laboratory for Computer Science, Cambridge, MA, June 1997. Technical Report MIT-LCS-TR-717, Lab. for Computer Science, MIT Cambridge, MA, USA, June 1997.</i>  <a href="http://theory.lcs.mit.edu/tds/paxos.html">http://theory.lcs.mit.edu/tds/paxos.html</a>
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 12.4-3

<b>BP 1</b>	<b>Einigungsprotokolle: Inhalt</b>
<b>12.3</b>	<b>Byzantinische Verständigung</b>
	<i>Lamport, L.; Shostak, R.; Pease, Marshall: The Byzantine Generals Problem. ACM Transactions on Programming Languages and Systems, Vol. 4, No. 3, July 1982, pp. 382-401.</i>  <i>Pease, M.; Shostak, R.; Lamport, L.: Reaching Agreement in the Presence of Faults. Journal of the ACM, Vol. 27, No. 2, April 1980, pp. 228-234.</i>  <i>Dolev, D.; Fischer, M. J.; Fowler, R.; Lynch, N.; Strong, H. R.: An Efficient Algorithm for Byzantine Agreement without Authentication. Information and Control, vol. 52 (1982), pp. 257-274.</i>  <i>Burns, J. E.; Neiger, G.: Fast and simple distributed consensus. Distributed Computing, Vol. 8 (1994), pp. 59-64.</i>
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 12.3-2

<b>BP 1</b>	<b>Einigungsprotokolle: Erläuterung der Fragestellung</b>
<b>12.1</b>	<b>Erläuterung der Fragestellung</b>
	<b>Systemmodell</b> <ul style="list-style-type: none"> <li>• n Prozessoren, davon maximal m fehlerhaft.</li> <li>• Voll vermaschtes Nachrichtensystem, nicht beeinträchtigt durch fehlerhafte Prozessoren.</li> <li>• Zu jeder empfangenen Nachricht ist dem Empfänger der Absender bekannt.</li> <li>• Das Kommunikationssystem ist zuverlässig (d. h. lediglich Prozessoren können fehlerhaft sein).</li> <li>• Einigung soll zwischen zwei Möglichkeiten (0 und 1) erzielt werden.</li> <li>• Abarbeitung in synchronisierten Runden oder asynchron.</li> </ul> <b>Prozessorfehler</b> <ul style="list-style-type: none"> <li>• Fail Stop</li> <li>• Unterlassung von Nachrichtenversand</li> <li>• Boshafte Fehlverhalten (Byzantinisches Fehlverhalten)</li> </ul>
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 12.1-4

<b>BP 1</b>	<b>Einigungsprotokolle: Erläuterung der Fragestellung</b>
	<p><b>Beglaubigte versus unbeglaubigte Nachrichten (signed versus oral)</b></p> <ul style="list-style-type: none"> <li>• Beglaubigt meint, daß Nachrichten unterschrieben werden können und <ul style="list-style-type: none"> <li>• die Unterschriften fehlerfreier Prozessoren nicht fälschbar sind,</li> <li>• der Inhalt einer Nachricht, die von einem fehlerfreien Prozessor unterzeichnet ist, nicht verfälscht werden kann und</li> <li>• jeder die Authentizität bestimmter Unterschriften feststellen kann.</li> </ul> </li> </ul>
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>

12.1-5

<b>BP 1</b>	<b>Einigungsprotokolle: Erläuterung der Fragestellung</b>
	<p><b>Problemklassifikation</b></p> <p><b>Einigung aller fehlerfreien Prozessoren auf einen Wert, der von einem Prozessor vorgegeben wird (Verständigungsproblem, 'agreement problem').</b></p> <ul style="list-style-type: none"> <li>• <b>Einigungsziel:</b> Alle fehlerfreien Prozessoren einigen sich auf den gleichen Wert.</li> <li>• <b>Richtigkeit:</b> Falls der vorgegebende Prozessor fehlerfrei ist, erfolgt die Einigung auf seinen Wert.</li> </ul> <p><b>Einigung auf einen Wert, der eine Funktion der Anfangswerte aller Prozessoren ist (Einigungsproblem, 'consensus problem').</b></p> <ul style="list-style-type: none"> <li>• <b>Einigungsziel:</b> Alle fehlerfreien Prozessoren einigen sich auf den gleichen Wert.</li> <li>• <b>Richtigkeit:</b> Wenn alle fehlerfreien Prozessoren den gleichen Wert vorgeben, einigen sich alle auf diesen Wert.</li> </ul> <p><b>Einigung auf die Menge der Anfangswerte der Prozessoren (Konsistenzproblem, 'consistency problem').</b></p> <ul style="list-style-type: none"> <li>• <b>Einigungsziel:</b> Alle fehlerfreien Prozessoren einigen sich auf den gleichen Wertevektor.</li> <li>• <b>Richtigkeit:</b> Die Vorgabewerte fehlerfreier Prozessoren werden in dem Wertevektor richtig wiedergegeben.</li> </ul>
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>

12.1-6

<b>BP 1</b>	<b>Einigungsprotokolle: Erläuterung der Fragestellung</b>
	<p><b>Verwandschaft der Probleme</b></p> <p><b>Lösung des Verständigungsproblems --&gt; Lösung des Konsistenzproblems</b></p> <p>Jeder Knoten startet als Anführer eine Durchführung des Verständigungsalgorithmus, wodurch jeder Knoten einen Wertevektor mit den verlangten Eigenschaften ermittelt.</p> <p><b>Lösung des Konsistenzproblems --&gt; Lösung des Einigungsproblems</b></p> <p>Das Ergebnis wird als mehrheitlicher Wert der erlangten Wertevektoren bestimmt. Enthalten die Wertevektoren beide Werte gleich oft, wird als Ergebnis der Wert 1 genommen.</p> <p><b>Lösung des Einigungsproblems --&gt; Lösung des Verständigungsproblems</b></p> <ol style="list-style-type: none"> <li>1. Runde: Der Anführer sendet seinen Wert an alle einschließlich sich selbst.</li> <li>2. Runde: Alle Knoten einschließlich Anführer führen den Einigungsalgorithmus durch, wobei sie als Anfangswert den in der 1. Runde erhaltenen Wert nehmen.</li> </ol>
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>

12.1-7

<b>BP 1</b>	<b>Einigungsprotokolle: Unlösbarkeitsaussagen</b>
12.2	<p><b>Unlösbarkeitsaussagen</b></p> <p><b>Unlösbarkeit bei asynchronen Nachrichten</b></p> <p><b>Modell</b></p> <ul style="list-style-type: none"> <li>• Das System besteht aus n Prozessoren.</li> <li>• Jeder Prozessor besitzt einen (allen anderen unbekannten) Anfangswert 0 oder 1.</li> <li>• Jeder Prozessor besitzt unbeschränkten lokalen Speicher.</li> <li>• Kommunikation erfolgt über ein zuverlässiges Nachrichtensystem.</li> </ul> <p>Die Operation <b>send(p, m)</b> übergibt die Nachricht (p, m) an das Nachrichtensystem. Dabei bezeichnet p den Empfänger und m den eigentlichen Nachrichteninhalt.</p> <p>Die Operation <b>receive(p)</b> entnimmt ein Paar (p, m) aus dem Nachrichtensystem und liefert m zurück oder sie liefert als Ergebnis <math>\Phi</math> und läßt das Nachrichtensystem unverändert (asynchrones Empfangen).</p> <p><b>S12.1 Satz</b></p> <p>Es existiert kein Protokoll, das unter Verwendung asynchroner Nachrichten das Einigungsproblem lösen kann, wenn auch nur ein Prozessor fehlerhaft ist.</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>

12.2-8

## BP 1 Einigungsprotokolle: Unlösbarkeitsaussagen

Beweis:

- **Bezeichnungen**
  - **Konfiguration:** Umfaßt die Zustände sämtlicher Prozessoren und des Nachrichtensystems.
  - **Schritt:** Überführt eine Konfiguration in die nächste durch eine Zustandsänderung, die ein einzelner Prozessor bewirkt durch Ausführung von receive oder durch lokalen Übergang mit eventuellem Aussenden von Nachrichten.
  - **Ereignis:** Die Prozessoren sind als deterministisch vorausgesetzt, so daß ein Übergang vollständig durch ein Ereignis  $e = (p, m)$  bestimmt ist. Eine Konfiguration  $C$  geht durch ein Ereignis  $e$  in die Konfiguration  $e(C)$  über.
  - **Lokales Ergebnis:** Wert für den sich der Prozessor unwiderruflich entschieden hat.
  - **Ablaufplan:** (Endliche oder unendliche) Folge  $\sigma$  von Ereignissen, die auf  $C$  angewandt werden kann.
  - **Lauf:** Schrittfolge, die bei Abwicklung eines Ablaufplans entsteht.

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

12.2-9

## BP 1 Einigungsprotokolle: Unlösbarkeitsaussagen

• **Bezeichnungen (Forts.)**

- **Erreichbar:** Wenn  $\sigma$  endlich ist, heißt die am Ende des Ablaufs erreichte Konfiguration erreichbar.
- **Zulässiger Lauf:** Lauf mit höchstens einem fehlerhaften Prozessor, in dem alle von fehlerfreien Prozessoren ausgesandten Nachrichten schließlich ankommen.

17.01.02

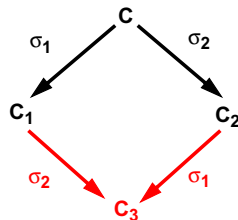
Universität Erlangen-Nürnberg, Lehrstuhl für Informatik 4 (Verteilte Systeme und Betriebssysteme), F. Hofmann  
Reproduktion jeder Art oder Verwendung dieser Unterlage zu Lehrzwecken außerhalb der Universität Erlangen-Nürnberg ist ohne Genehmigung des Autors unzulässig

12.2-10

## BP 1 Einigungsprotokolle: Unlösbarkeitsaussagen

### L12.1 Lemma

Ausgehend von der Konfiguration  $C$  führe  $\sigma_1$  bzw.  $\sigma_2$  zu  $C_1$  bzw.  $C_2$ . Falls die Prozesse, die in  $\sigma_1$  einen Schritt ausführen, verschieden sind von all denen, die in  $\sigma_2$  einen Schritt ausführen, dann kann  $\sigma_2$  auf  $C_1$  angewandt werden und  $\sigma_1$  auf  $C_2$  und beides führt zum gleichen Zustand  $C_3$ .



Beweis: Unmittelbare Folgerung der Definitionen.

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

12.2-11

## BP 1 Einigungsprotokolle: Unlösbarkeitsaussagen

Beweis des Satzes

- **Indirekt:** Angenommen es gäbe ein Protokoll, das auch bei einem fehlerhaften Prozessor das Einigungsproblem löst.

**Definition:**

Sei  $C$  eine Konfiguration und  $V$  die Menge lokaler Ergebnisse von  $C$  aus erreichbarer Konfigurationen.

- $C$  heißt **bivalent**, wenn  $|V| = 2$  ist.
- $C$  heißt **i-valent**, wenn  $|V| = 1$  ist und die lokalen Ergebnisse, den Wert  $i$  besitzen.

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

12.2-12

## BP 1 Einigungsprotokolle: Unlösbarkeitsaussagen

### L12.2 Lemma

Es gibt eine bivalente Anfangskonfiguration.

Beweis:

- Angenommen es gibt nur 0- und 1-valente Anfangskonfigurationen.
- Anfangskonfigurationen differieren nur in den Eingaberegistern. Also müssen zwei Anfangskonfigurationen  $C_0$  und  $C_1$  existieren, die sich nur im Wert des Eingaberegisters eines Prozessors  $P_p$  unterscheiden und von denen eine 0-valent und die andere 1-valent ist. (Man betrachte die Anordnung der Anfangskonfigurationen, in der die Vektoren aus den Eingaberegisterinhalten nach Gray-Code angeordnet sind.)
- $\sigma$  sei ein zulässiger Lauf, der ohne Mitwirken von  $P_p$  von  $C_0$  ausgehend zur Entscheidung führt (das muß gehen, da  $P_p$  der ausgefallene Prozessor sein könnte).
- $\sigma$  kann auch auf  $C_1$  angewandt werden.
- Die resultierenden Konfigurationen unterscheiden sich nur im lokalen Zustand von  $P_p$ .
- Gemäß Voraussetzung ermitteln  $\sigma(C_0)$  und  $\sigma(C_1)$  die gleichen Entscheidungswerte im Widerspruch dazu, daß  $C_0$  0-valent ist und  $C_1$  1-valent.

## BP 1 Einigungsprotokolle: Unlösbarkeitsaussagen

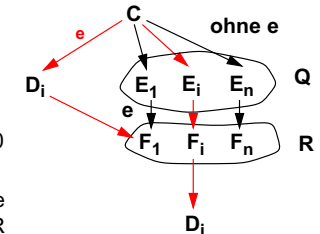
### L12.3 Lemma

Sei  $C$  eine bivalente Konfiguration und  $e = (p, m)$  ein in dieser Konfiguration anwendbares Ereignis. Weiter sei  $Q$  die Menge der Konfigurationen, die von  $C$  aus ohne Mitwirkung von  $P_p$  erreicht werden können, und es sei  $R = e(Q)$ . Dann enthält  $R$  eine bivalente Konfiguration.

Beweis:

Angenommen  $R$  enthalte **keine** bivalente Konfiguration.

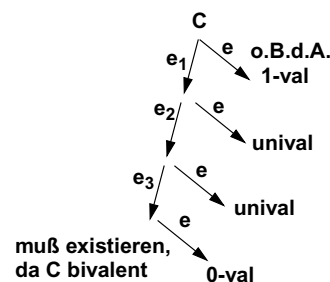
- $D_i$  sei eine von  $C$  aus erreichbare  $i$ -valente Konfiguration (die nach Voraussetzung für  $i = 0$  und  $i = 1$  existieren muß).
- Dann gilt entweder  $D_i \in Q \wedge F_i = e(D_i) \in R$  oder  $e$  wurde angewandt, um  $D_i$  zu erreichen, und  $D_i \in R$  oder es existiert  $F_i \in R$ , das von  $E_i$  aus erreichbar ist.
- $R$  muß also auf jeden Fall 0- und 1-valente Konfigurationen enthalten.



## BP 1 Einigungsprotokolle: Unlösbarkeitsaussagen

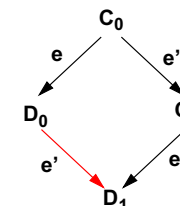
Beweis (Forts.)

- Zwei Konfigurationen werden als Nachbarn bezeichnet, wenn die eine in einem einzigen Schritt in die andere überführt werden kann.
- Es existieren Nachbarn  $C_0, C_1 \in Q$  derart, daß  $D_i = e(C_i)$   $i$ -valent ist. Die Indizierung der  $F_i$  sei so gewählt, daß die 0-valenten kleiner Indizes haben als die 1-valenten.



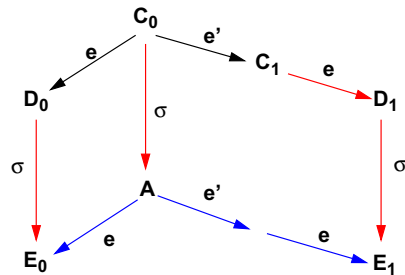
## BP 1 Einigungsprotokolle: Unlösbarkeitsaussagen

- Sei  $C_1 = e'(C_0)$  mit  $e' = (p', m')$ .
- Ist  $p' \neq p$ , dann gilt  $D_1 = e'(D_0)$ , was unmöglich ist, da jeder Nachfolger einer 0-valenten Konfiguration wieder 0-valent ist.



## BP 1 Einigungsprotokolle: Unlösbarkeitsaussagen

- Wenn  $p = p'$  ist, betrachte man irgendeinen Entscheidungslauf von  $C_0$  aus, an dem sich der Prozessor  $P_p$  nicht beteiligt.
  - Sei  $\sigma$  ein entsprechender Ablaufplan und  $A = \sigma(C_0)$ .
  - $\sigma$  ist auch auf  $D_i$  anwendbar und führt zu einer  $i$ -valenten Konfiguration  $E_i = \sigma(D_i)$ .
  - Wegen  $e(A) = E_0$  und  $e(e'(A)) = E_1$  muß  $A$  bivalent sein im Widerspruch zur Annahme, daß  $A$  ein bis zur Entscheidung geführter Lauf ist.



Also führen beide Fälle zum Widerspruch und somit ist die Annahme falsch.

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

12.2-17

## BP 1 Einigungsprotokolle: Unlösbarkeitsaussagen

### Beweis des Satzes (Fortsetzung)

Durch Konstruktion eines nicht-entscheidenden, beliebig langen Laufs

- Prozesse werden in einer Warteschlange geführt.
- Nachrichten werden in einem Nachrichtenpuffer gesammelt.
- Ein Abschnitt besteht aus einem oder mehreren Schritten. Er endet, wenn der erste Prozeß der Warteschlange einen Schritt ausgeführt hat, in dem die älteste, an ihn gerichtete und noch nicht bearbeitete Nachricht angenommen wurde (vorausgesetzt wird, daß es eine solche Nachricht gibt.)
- Der Prozeß wird dann wieder an das Ende der Warteschlange verwiesen.
- In einer endlosen Folge von Abschnitten führt jeder Prozeß unendlich viele Schritte aus und empfängt jede an ihn gerichtete Nachricht.

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

12.2-18

## BP 1 Einigungsprotokolle: Unlösbarkeitsaussagen

- Sei  $C_0$  eine bivalente Anfangskonfiguration. Dann existieren Abläufe derart, daß jeder Abschnitt wieder mit einer bivalenten Konfiguration endet.
  - $C$  sei eine bivalente Konfiguration und  $P_p$  der erste Prozeß in der Warteschlange.  $m$  oder  $\phi$  sei die älteste Nachricht an  $P_p$  in  $C$ .
  - Sei  $e = (p, m)$ . Nach Lemma 3 gibt es eine bivalente Konfiguration  $C'$ , die von  $C$  aus mit einem Ablaufplan erreichbar ist, in dem die Anwendung von  $e$  der letzte Schritt ist. Diese Folge definiert einen Abschnitt.
  - Da alle Abschnitte mit einer bivalenten Konfiguration enden, wird in dem so konstruierten Ablaufplan nie eine Entscheidung erreicht.

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

12.2-19

## BP 1 Einigungsprotokolle: Unlösbarkeitsaussagen

### Unlösbarkeit bei zu vielen fehlerhaften Prozessoren

#### Modell eines verteilten Systems

- Kommunikationsgraph:
  - Gerichteter Graph  $G$  mit Knotenmenge  $k(G)$  und Verbindungsmenge  $v(G)$
  - Verbindungen treten nur paarweise auf, d. h.  $v(k_1, k_2) \in v(G) \Leftrightarrow v(k_2, k_1) \in v(G)$
  - $(k_1, k_2)$  ist Ausgang von  $k_1$  und Eingang von  $k_2$ .

Ein Untergraph  $G_U$  von  $G$  ist die Menge aller Knoten aus  $U$  zusammen mit allen Verbindungen von  $v(G)$ , die Eingang oder Ausgang wenigstens eines Knotens aus  $U$  sind.

Die Eingangsschnittstelle eines Untergraphen  $U$  ist die Menge aller Eingänge von Knoten aus  $U$ , die Ausgänge von Knoten aus  $G - U$  sind.

17.01.02

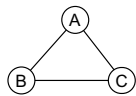
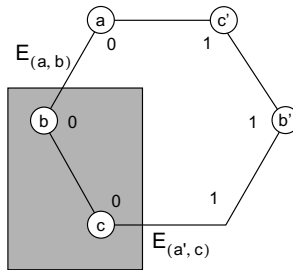
Universität Erlangen-Nürnberg, Lehrstuhl für Informatik 4 (Verteilte Systeme und Betriebssysteme), F. Hofmann  
Reproduktion jeder Art oder Verwendung dieser Unterlage zu Lehrzwecken außerhalb der Universität Erlangen-Nürnberg  
ist ohne Genehmigung des Autors unzulässig

12.2-20

BP 1	Einigungsprotokolle: Unlösbarkeitsaussagen
	<ul style="list-style-type: none"> <li>• <b>System:</b>  Ein System <math>S</math> ist ein Kommunikationsgraph, bei dem jedem Knoten eine <i>Einrichtung</i>, die ein Verhalten des Knotens erzeugt, und eine <i>Eingabe</i> zugeordnet ist.  Ein System <math>S</math> hat ein <i>Systemverhalten</i> <math>E</math>, das sich als Tupel des <i>Verhaltens</i> aller Knoten und Verbindungen ergibt.  Die Restriktion des Systemverhaltens <math>E</math> auf Knoten und Verbindungen eines Untergraphen <math>G_U</math> bildet das Szenario <math>E_U</math> von <math>G_U</math>.</li> </ul>
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 <div>12.2-21</div>

BP 1	Einigungsprotokolle: Unlösbarkeitsaussagen
	<ul style="list-style-type: none"> <li>• <b>Lokalitätsaxiom:</b>  Seien <math>S</math> bzw. <math>S'</math> Systeme mit Verhalten <math>E</math> bzw. <math>E'</math> und isomorphen Untersystemen <math>U</math> bzw. <math>U'</math>. Wenn das Verhalten korrespondierender Eingänge von <math>U</math> und <math>U'</math> in den Szenarien <math>E</math> und <math>E'</math> gleich ist, dann sind auch die Szenarien <math>E_U</math> und <math>E'_{U'}</math> gleich.  (D. h. insbesondere, daß das Verhalten eines Teilsystems durch das Verhalten seiner Knoten und inneren Verbindungen sowie dem seiner Eingänge bestimmt ist.)</li> <li>• <b>Fehleraxiom:</b>  Sei <math>A</math> eine Einrichtung. Weiter seien <math>(E_1, \dots, E_d)</math> Verbindungsverhalten derart, daß <math>E_i</math> in einem Systemverhalten <math>E_i</math> das Verhalten des <math>i</math>-ten Ausgangs eines Knotens mit zugeordneter Einrichtung <math>A</math> ist. Außerdem sei <math>u</math> ein Knoten mit <math>d</math> Ausgängen <math>(u, v_1), \dots, (u, v_d)</math>.  Dann existiert eine Einrichtung <math>F</math> derart, daß in jedem System, in dem <math>u</math> die Einrichtung <math>F</math> zugeordnet ist, das Verhalten der Ausgangskanten <math>(u, v_i)</math> gleich <math>E_i</math> ist.  Um diesen Sachverhalt deutlich zu machen, wird für <math>F</math> in diesem Fall <math>F_A(E_1, \dots, E_d)</math> geschrieben.</li> </ul>
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 <div>12.2-22</div>

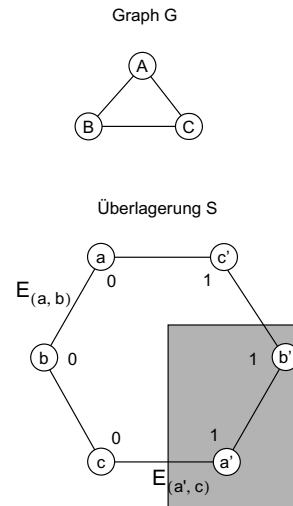
BP 1	Einigungsprotokolle: Unlösbarkeitsaussagen
	<ul style="list-style-type: none"> <li>• <b>Graphentheoretische Begriffe:</b>  Ein Graph <math>S</math> überlagert einen Graphen <math>G</math>, wenn eine Abbildung der Knoten von <math>S</math> auf die von <math>G</math> existiert mit folgender Eigenschaft:  Wenn Knoten <math>u</math> aus <math>S</math> mit <math>d</math> Nachbarknoten <math>v_1, \dots, v_d</math> verbunden ist und auf <math>w</math> aus <math>G</math> abgebildet wird, dann hat <math>w</math> ebenfalls <math>d</math> Nachbarn <math>x_1, \dots, x_d</math>, wobei für <math>1 \leq i \leq d</math> Knoten <math>v_i</math> auf Knoten <math>x_i</math> abgebildet wird.  (D. h. <math>S</math> sieht lokal wie <math>G</math> aus.)</li> </ul> <p><b>S12.2 Satz</b></p> <p>In einem System aus <math>n</math> Prozessoren, von denen maximal <math>m</math> fehlerhaft sind, können keine Protokolle zur Lösung des Problems der Byzantinischen Einigung (Byzantine consensus) existieren, wenn <math>n \leq 3m</math> ist.</p> <p><b>Beweis:</b></p> <p>1. Sei <math>n = 3</math> und <math>m = 1</math></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 <div>12.2-23</div>

BP 1	Einigungsprotokolle: Unlösbarkeitsaussagen
	<p>Angenommen die Einrichtungen <math>A</math>, <math>B</math> und <math>C</math> könnten bei der gezeichneten Zuordnung in <math>G</math> immer Einigung erzielen.</p> <p>In der Überlagerung <math>S</math> ordne man den Knoten <math>a</math> und <math>a'</math> die Einrichtung <math>A</math>, den Knoten <math>b</math> und <math>b'</math> die Einrichtung <math>B</math> und den Knoten <math>c</math> und <math>c'</math> die Einrichtung <math>C</math> zu.</p> <p>Nach Fehleraxiom existiert eine Einrichtung <math>F_A(E_{(a,b)}, E_{(a',c)})</math> die für <math>E_{(a,b)}</math> und <math>E_{(a',c)}</math> das gleiche Eingangsverhalten erzeugt, wie das System <math>S</math>.</p> <p>Ersetzt man in <math>G</math> die Einrichtung <math>A</math> durch <math>F_A(E_{(a,b)}, E_{(a',c)})</math>, so ist nach Lokalitätsaxiom <math>E_{b,c} = E_{B,C}</math>.</p> <p>Auf Grund der Gültigkeitsforderung müssen sich <math>B</math> und <math>C</math>, und damit auch <math>b</math> und <math>c</math> auf den Wert 0 einigen.</p> <p>Eine analoge Betrachtung für <math>E_{c,a'}</math> führt zusammen mit der Tatsache, daß <math>c</math> zum Wert 0 kommen muß dazu, daß auch <math>a'</math> zum Ergebnis 0 kommen muß.</p> <div> <div>Graph G</div>  <div>Überlagerung S</div>  </div>
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 <div>12.2-24</div>

## BP 1 Einigungsprotokolle: Unlösbarkeitsaussagen

Wiederholung dieser Betrachtung für  $a'$  und  $b'$  zeigt, daß auch  $b'$  zum Ergebnis 0 kommen muß. Andererseits schließt man analog zur anfänglichen Überlegung, daß sich die Knoten  $a'$  und  $b'$  auf den Wert 1 einigen müßten.

Dieser Widerspruch zeigt, daß die Annahme falsch ist.



## BP 1 Einigungsprotokolle: Unlösbarkeitsaussagen

### 2. Allgemeiner Fall $n \leq 3m$

Es werde wiederum angenommen, das Einigungsproblem sei lösbar.

Man partitioniere die Menge der Knoten in Teilmengen  $a$ ,  $b$  und  $c$  derart, daß jede Teilmenge wenigstens ein und höchstens  $m$  Elemente enthält.

Die Vereinigung von je zwei Teilmengen enthält mindestens  $n-m$  Elemente.

$A(B, C)$  sei die Menge der Einrichtungen von Knoten in  $a(b, c)$ , die das Einigungsproblem lösen.

## BP 1 Einigungsprotokolle: Unlösbarkeitsaussagen

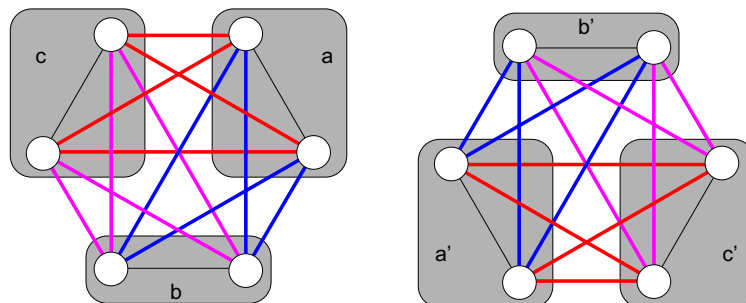
Man konstruiere eine Überlagerung durch Verdopplung des Graphen (Bezeichnungen des Doppels entstehen durch Anfügen von Apostroph) und Kantenersetzungen nach folgender Vorschrift

$$x \in a \wedge y \in c \wedge (x, y) \in v(G) \Rightarrow (x, y') \in v(S)$$

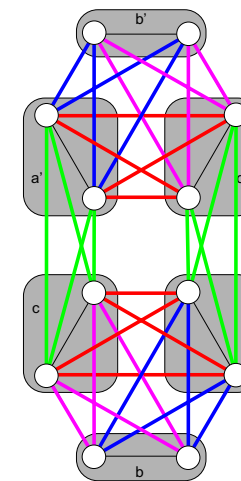
$$x \in a \wedge y \in c \wedge (y, x) \in v(G) \Rightarrow (y', x) \in v(S)$$

$$x' \in a' \wedge y' \in c' \wedge (x', y') \in v(G) \Rightarrow (x', y) \in v(S)$$

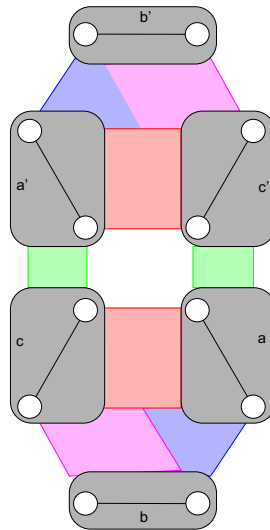
$$x' \in a' \wedge y' \in c' \wedge (y', x') \in v(G) \Rightarrow (y, x') \in v(S)$$



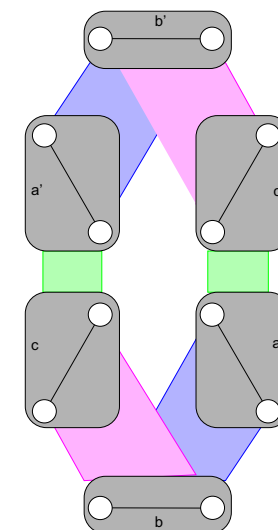
## BP 1 Einigungsprotokolle: Unlösbarkeitsaussagen



## BP 1 Einigungsprotokolle: Unlösbarkeitsaussagen



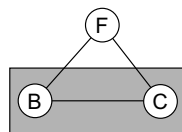
## BP 1 Einigungsprotokolle: Unlösbarkeitsaussagen



## BP 1 Einigungsprotokolle: Unlösbarkeitsaussagen

Ersetzt man in der Argumentation von Fall 1, die dortigen Knoten durch die entsprechenden Knotenmengen des jetzt betrachteten Falls, so kann wieder A durch eine Einrichtung F mit Fehlverhalten ersetzt werden.

Da B und C zusammen wenigstens  $n-m$  fehlerfreie Knoten enthalten, muß G korrektes Verhalten aufweisen, d. h. alle korrekten Knoten in b und c müssen sich für den Wert 0 entscheiden.



Damit läßt sich die Argumentation analog dem Fall 1 fortsetzen und führt in gleicher Weise zum Widerspruch.

## BP 1 Einigungsprotokolle: Byzantinische Verständigung

### 12.3 Byzantinische Verständigung

#### 12.3.1 Byzantinische Verständigung (Lamport-Shostak-Pease, 1982) mit unbeglaubigten Nachrichten

#### S12.3 Satz: Lamport-Shostak-Pease, 1982

In Systemen mit synchronisierbaren Runden löst nachfolgender rekursive Algorithmus das Verständigungsproblem auf der Basis unbeglaubigter Nachrichten, wenn  $n > 3m$  ist.



## BP 1 Einigungsprotokolle: OM-Algorithmus von Lamport-Shostak-Pease

Algorithmus:

$\text{Mehrheit}(v_1, v_2, \dots, v_{n-1}) = v$ , wenn wenigstens die Hälfte der  $v_i$  den Wert  $v$  hat.

Algorithmus OM(0)

1. Der Anführer sendet seinen Wert an alle anderen.
2. Jeder, außer dem Anführer, nimmt als Ergebnis den Wert, den er erhalten hat.

Algorithmus OM(m),  $m > 0$

1. Der Anführer sendet seinen Wert an alle anderen.
2. Für jedes  $i$  sei  $v_i$  der Wert den Prozessor  $i$  vom Anführer erhalten hat, bzw. 1 wenn er keine Nachricht bekommt. Prozessor  $i$  handelt als Anführer im Algorithmus OM(m-1) mit den übrigen  $n-2$  Prozessoren als Teilnehmern.
3. Für jedes  $i$  und jedes  $j \neq i$  sei  $v_j$  der Wert den Prozessor  $i$  im Schritt 2 von Prozessor  $j$  (unter Verwendung von OM(m-1)) erhalten hat oder 1, wenn er einen solchen nicht erlangt hat. Prozessor  $i$  benutzt den Wert  $\text{Mehrheit}(v_1, v_2, \dots, v_{n-1})$ .

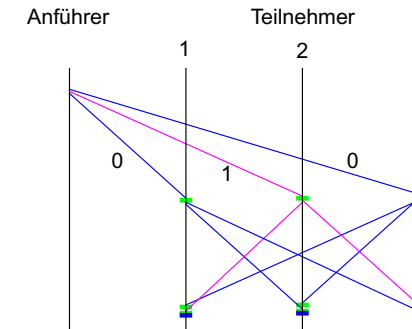
17.01.02

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

12.3-33

## BP 1 Einigungsprotokolle: Einigungsprotokolle: Lamport OM

Beispiel1: Der Algorithmus OM(1) bei fehlerhaftem Anführer



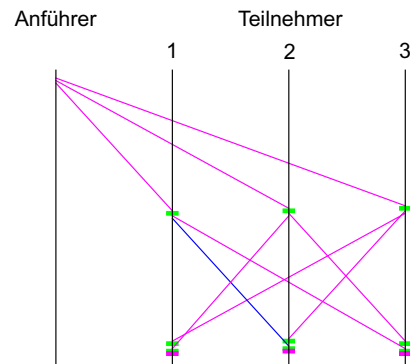
17.01.02

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

12.3-34

## BP 1 Einigungsprotokolle: Einigungsprotokolle: Lamport OM

Beispiel 2: Der Algorithmus OM(1) bei fehlerhaftem Teilnehmer 1



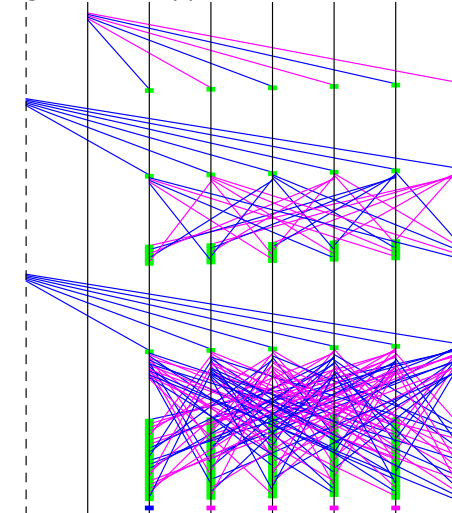
17.01.02

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

12.3-35

## BP 1 Einigungsprotokolle: Einigungsprotokolle: Lamport OM

Beispiel 3: Der Algorithmus OM(2), Teilnehmer 0 und 1 fehlerhaft



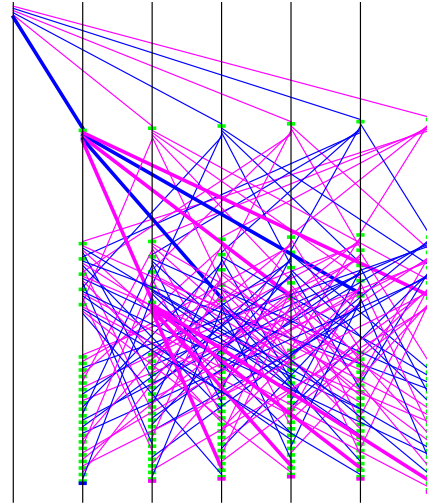
17.01.02

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

12.3-36

## BP 1 Einigungsprotokolle: Einigungsprotokolle: Lamport OM

Beispiel 4: Der Algorithmus OM(2) von Lamport, Teilnehmer 0 und 1 fehlerhaft, Synchronisation durch Timeouts



## BP 1 Einigungsprotokolle: Einigungsprotokolle: Lamport OM

Angenommen: Richtig für  $m - 1$  mit  $m > 0$

Im 1. Schritt  $v$  an alle anderen

Im 2. Schritt benutzen alle OM( $m-1$ ) für  $n-1$  Prozessoren.

$n > 2k + m \rightarrow n - 1 > 2k + (m - 1)$ , d. h. Ind.-Vor. ist erfüllt.

Also  $v_j = v$  für jeden fehlerfreien Prozessor  $j$ .

Höchstens  $k$  fehlerhafte und  $n - 1 > 2k + (m - 1) \geq 2k$

$\rightarrow$  Mehrheit von  $n-1$  Prozessoren fehlerfrei

$\rightarrow$  Mehrheit ermittelt unter Verwendung von OM( $m-1$ )  
den Wert  $v$

$\rightarrow$  Mehrheit( $v_1, v_2, \dots, v_{n-1}$ ) =  $v$

## BP 1 Einigungsprotokolle: Einigungsprotokolle: Lamport OM

Beweis

Zu erfüllende Bedingungen

IC1: Alle fehlerfreien Prozessoren ermitteln den gleichen Wert.

IC2: Wenn der Anführer fehlerfrei ist, ermitteln alle fehlerfreien Prozessoren seinen Wert.

### L12.4 Lemma

Für jedes  $m$  und  $k$ , erfüllt OM( $m$ ) die Bedingung IC2, wenn  $n > 2k + m$  ist und höchstens  $k$  Prozessoren fehlerhaft sind.

Beweis:

Induktion über  $m$

OM(0): trivial

## BP 1 Einigungsprotokolle: Einigungsprotokolle: Lamport OM

Beweis des Hauptsatzes durch Induktion über  $m$

OM(0) trivial

Annahme OM( $m-1$ ) erfüllt Satz.

1. Anführer fehlerfrei

$\rightarrow$  mit  $k = m$  folgt die Behauptung aus Lemma 1.

2. Anführer fehlerhaft

$\rightarrow m$  fehlerhafte, darunter Anführer

$n > 3m \rightarrow n-1 > 3m-1 > 3(m-1)$

$\rightarrow$  Voraussetzung für OM( $m-1$ ) im 2. Schritt erfüllt

$\rightarrow$  für jeden Prozessor  $j$  ermitteln alle fehlerfreien den gleichen Wert und damit den gleichen Wertvektor

$\rightarrow$  alle fehlerfreien ermitteln im 3. Schritt den gleichen Wert  $v$

## BP 1 Einigungsprotokolle: SM-Algorithmus von Lamport, Shostak und Pease

### S12.4 Satz: Lamport-Shostak-Pease, 1982

Es sei  $V$  eine Menge von Werten. Weiter sei  $\text{choice}(V)$  eine Funktion mit folgenden Eigenschaften:

- Im weiteren bezeichne  $x_i$  den Wert  $x$  beglaubigt von Prozessor  $i$ .  
Es bezeichnet also  $v_{j:i}$  den Wert  $v$  beglaubigt von Prozessor  $j$  und dann  $v_{j:i}$  beglaubigt von Prozessor  $i$ .

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

12.3-41

## BP 1 Einigungsprotokolle: SM-Algorithmus von Lamport, Shostak und Pease

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

12.3-42

The diagram shows a quantum circuit with 7 qubits, labeled 0 through 6 at the top. The circuit consists of 6 CNOT gates, represented by green squares with a cross. The gates are arranged in a grid-like pattern. The first gate is between qubits 0 and 1. The second gate is between qubits 0 and 2. The third gate is between qubits 0 and 3. The fourth gate is between qubits 0 and 4. The fifth gate is between qubits 0 and 5. The sixth gate is between qubits 0 and 6. The qubits are represented by vertical lines, and the gates are represented by green squares with a cross. The qubits are labeled 0, 1, 2, 3, 4, 5, 6 at the top. The gates are arranged in a grid-like pattern. The first gate is between qubits 0 and 1. The second gate is between qubits 0 and 2. The third gate is between qubits 0 and 3. The fourth gate is between qubits 0 and 4. The fifth gate is between qubits 0 and 5. The sixth gate is between qubits 0 and 6. The qubits are represented by vertical lines, and the gates are represented by green squares with a cross. The qubits are labeled 0, 1, 2, 3, 4, 5, 6 at the top.

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

12.3-43

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

12.3-44

BP 1	Einigungsprotokolle: SM-Algorithmus von Lamport, Shostak und Pease
	<p>Es genügt zu zeigen, daß <math>V_i = V_j</math> ist, wenn die Prozessoren <math>i</math> und <math>j</math> fehlerfrei sind.</p> <p>Prozessor <math>i</math> nehme in Schritt 2 <math>v</math> in <math>V</math> auf.</p> <p>Macht er dies in Schritt 2A, dann sendet er ihn auch an <math>j</math>.</p> <p>Macht er dies in Schritt 2B, dann muß er eine Nachricht <math>v:0:j_1: \dots :j_k</math> erhalten haben.</p> <p>Es sind zwei Unterfälle zu betrachten:</p> <ol style="list-style-type: none"> <li>1. <math>k &lt; m</math>: In diesem Fall sendet er die Nachricht <math>v:0:j_1: \dots :j_k:i</math> an <math>j</math>, so daß auch Prozessor <math>j</math> diesen Wert erhält.</li> <li>2. <math>k = m</math>: Anführer fehlerhaft <math>\rightarrow</math> höchstens <math>m - 1</math> der Prozessoren <math>j_1, \dots, j_m</math> fehlerhaft. Also muß wenigstens einer fehlerfrei sein und den Wert auch an Prozessor <math>j</math> gesendet haben.</li> </ol>
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           12.3-45

BP 1	Einigungsprotokolle: Algorithmus von Dolev
	<ul style="list-style-type: none"> <li>• <b>Nachrichtenarten</b> <ul style="list-style-type: none"> <li>*-Nachricht, die als Wert des Anführers 1 zusichert.</li> <li>k-Nachricht, die mitteilt, daß <math>k</math> überzeugt ist, eine zu recht bestehende *-Nachricht erhalten zu haben.</li> </ul> </li> <li>• <b>Lokale Zustände</b> <math>q^j \in Q</math> mit  <math>Q = \wp((\{*\} \cup \{1, \dots, n\}) \times \{1, \dots, n\}) \times \mathbb{N}</math>            Paare (Nachrichtenkennung, Sender), Runde</li> <li>• <math>LOW = m + 1</math> und <math>HIGH = 2m + 1</math>            Wenn eine bestimmte Aussage von wenigstens <math>LOW</math> Prozessen erhalten wurde, so muß sie von wenigstens einem korrekten gemacht worden sein.            Wenn eine bestimmte Aussage von wenigstens <math>HIGH</math> Prozessen erhalten wurde, so muß sie von wenigstens <math>LOW</math> korrekten Prozessen gemacht worden sein.</li> </ul>
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           12.3-47

BP 1	Einigungsprotokolle: Algorithmus von Dolev
12.3.3	<b>Byzantinische Verständigung (Dolev)</b> Reduktion der Nachrichtenzahl durch Erhöhung der maximalen Rundenzahl  Es sei $n = 3m + 1$ : Definitionen: <ul style="list-style-type: none"> <li>• <b>Runde</b>                Algorithmus arbeitet in Runden  <math>q_r = (q_r^1, \dots, q_r^n)</math> Zustand am Anfang der <math>r</math>-ten Runde  <math>q_r^i</math> enthält u. a. alle bis zu Beginn der <math>r</math>-ten Runde empfangenen Nachrichten</li> <li>• <b>Abarbeitung einer Runde in 2 Schritten:</b> <ol style="list-style-type: none"> <li>1. Aussenden von Nachrichten in Abhängigkeit vom Zustand am Anfang der Runde</li> <li>2. Empfang aller in der Runde gesandten Nachrichten und Ermittlung des Anfangszustandes für die nächste Runde</li> </ol> </li> </ul>
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           12.3-46

BP 1	Einigungsprotokolle: Algorithmus von Dolev
	<ul style="list-style-type: none"> <li>• <math>W_x^j(q) = \{i   i \in \{1, \dots, n\} \wedge (x, i) \in \pi_1(q^j)\}</math>            Menge aller Prozessoren, die an <math>j</math> eine <math>x</math>-Nachricht gesandt haben (<math>x \in \{*, 1, \dots, n\}</math>).</li> <li>• <b><math>j</math> unterstützt <math>i</math> direkt</b>, wenn <math>j</math> von <math>i</math> eine *-Nachricht erhielt (d. h. <math>i \in W_x^j(q)</math>).</li> <li>• <b><math>j</math> unterstützt <math>i</math> indirekt</b>, wenn <math> W_i^j(q)  \geq LOW</math> ist, d. h. <math>j</math> weiß infolge der Mitteilung wenigstens eines korrekten Prozessors, daß <math>i</math> eine *-Nachricht erhalten hat.</li> <li>• <b><math>j</math> ist überzeugt</b>, daß <math>i</math> eine *-Nachricht ausgesandt hat, wenn <math> W_i^j  \geq HIGH</math> ist.</li> <li>• <math>C^j(q) = \{k   k \in \{1, \dots, n\} \wedge k \neq s \wedge  W_k^j(q)  \geq HIGH\}</math> ist die Menge der <b>von <math>j</math> bestätigten Prozessoren</b> (ohne den Anführer <math>s</math>),            d. h. <math>j</math> ist zu Recht überzeugt, daß eine Mehrheit korrekter Prozessoren eine *-Nachricht versandt hat.</li> </ul>
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           12.3-48

BP 1

Einigungsprotokolle: Algorithmus von Dolev

• Ein Prozessor  $j$  ist im Zustand  $q_r$  **aktiviert**, d. h. er sendet in der  $r$ -ten Runde an alle einschließlich sich selbst eine \*-Nachricht,

wenn er schon in einer früheren Runde aktiviert war,  
d. h.  $j \in W^j_r(q_r)$  ist, (I1)

oder wenn  $|C^j_r(q_r)| \geq \text{LOW} + \max(0, \lceil r/2 \rceil - 2)$  ist (I2)

oder wenn  $r = 2$  ist und er in der ersten Runde vom Anführer eine \*-Nachricht erhielt, d. h.  $s \in W^j_1(q_2)$  ist. (I3)

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

12.3-49

BP 1

Einigungsprotokolle: Algorithmus von Dolev

**Beispiel 6:** Fehlerhafter Anführer ist  $P_0$ , Teilnehmer  $P_2$ ,  $P_4$  und  $P_6$  erhalten anfänglich keine Nachricht

— \*-Nachricht  
— k-Nachricht

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

12.3-51

BP 1

Einigungsprotokolle: Algorithmus von Dolev

**Algorithmus**

- In der ersten Runde sendet ein korrekter Anführer eine \*-Nachricht an alle einschließlich sich selbst, falls sein Wert 1 ist.
- In der  $k$ -ten Runde ( $k > 1$ ) senden die Prozessoren an alle anderen (einschließlich sich selbst)
  - eine \*-Nachricht, wenn sie aktiviert sind, (M1)
  - eine  $j$ -Nachricht, wenn sie  $j$  direkt unterstützen (d. h.  $j \in W^j_k(q)$ ), (weitergeleitete Aussage:  $j$  hat \* mitgeteilt)
  - eine  $j$ -Nachricht, wenn sie  $j$  indirekt unterstützen (d. h.  $|W^j_k(q)| \geq \text{LOW}$ ). (M3) (weitergeleitete Aussage: wenigstens ein Korrekter hat \* von  $j$  erhalten)
- Prozessoren, die von wenigstens HIGH Prozessoren überzeugt sind, d. h.  $|\{i \mid |W^i_k| \geq \text{HIGH}\}| \geq \text{HIGH}$ , entscheiden sich für den Wert 1.
- Wenn ein Prozessor sich nach der  $(2m + 3)$ -ten Runde für 1 entschieden hat, so nimmt er den Wert 1 als Ergebnis, sonst den Wert 0.

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

12.3-50

BP 1

Einigungsprotokolle: Algorithmus von Dolev

**Veranschaulichung**

$i \rightarrow j \rightarrow k$ :  $P_i$  hat an  $P_j$  eine \*-Nachricht gesandt und  $P_j$  hat dies an  $P_k$  weitergemeldet

$k \in W^j_i$ :  $j \rightarrow k \rightarrow i$

$k \in W^i_j$ :  $k \rightarrow i$

$k \in C^i$ : Es gibt mindestens  $m+1$  korrekte vom Anführer verschiedene Prozessoren  $j_v$  mit  $k \rightarrow j_v \rightarrow i$ . (Eine Mehrheit Korrekter hat eine \*-Nachricht bekommen)  
Prozessor  $k$  wird von  $i$  **bestätigt**, d. h. er verhält sich so, als ob er von  $k$  eine \*-Nachricht bekommen hätte.

$i$  zeugt für  $j$ :  $j \rightarrow i$

$i$  zeugt indirekt für  $j$ : Für mindestens einen korrekten Prozessor  $k$  gilt  $j \rightarrow k \rightarrow i$

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

12.3-52

## BP 1 Einigungsprotokolle: Algorithmus von Dolev

### Aktivierung von i

(I1)  $i \rightarrow i$  in vorangehender Runde

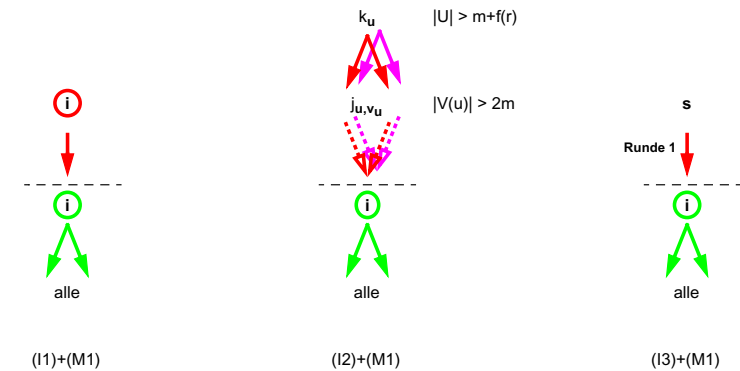
(I2)  $k_u \rightarrow j_u, v_u \rightarrow i$  mit  $u \in U \wedge |U| \geq m+1 + f(r)$  und  $v_u \in V_u \wedge |V_u| \geq 2m+1$   
(also  $k_u \in C^i$ )

(I3)  $s \rightarrow i$  in Runde 1

## BP 1 Einigungsprotokolle: Algorithmus von Dolev

### Senden:

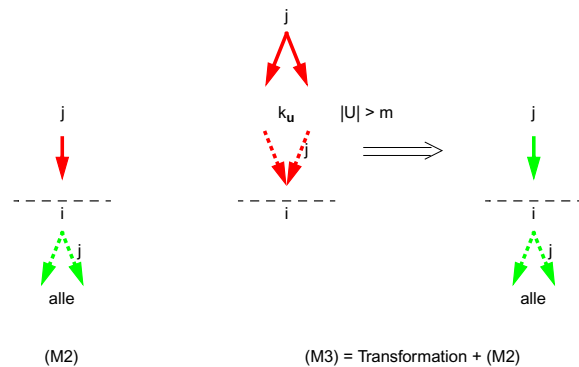
(M1)  $(I1) \vee (I2) \vee (I3) \rightarrow i \rightarrow \text{alle}$



## BP 1 Einigungsprotokolle: Algorithmus von Dolev

(M2)  $j \rightarrow i \rightarrow j \rightarrow i \rightarrow \text{alle}$

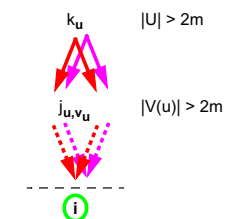
(M3) Für alle  $u \in U$  mit  $|U| \geq m+1$ :  $j \rightarrow k_u \rightarrow i \rightarrow j \rightarrow i \rightarrow \text{alle}$



## BP 1 Einigungsprotokolle: Algorithmus von Dolev

### Entscheidung von i für Wert 1:

$k_u \rightarrow j_u, v_u \rightarrow i$  mit  $u \in U \wedge |U| \geq 2m+1$  und  $v_u \in V_u \wedge |V_u| \geq 2m+1$



## BP 1 Einigungsprotokolle: Algorithmus von Dolev

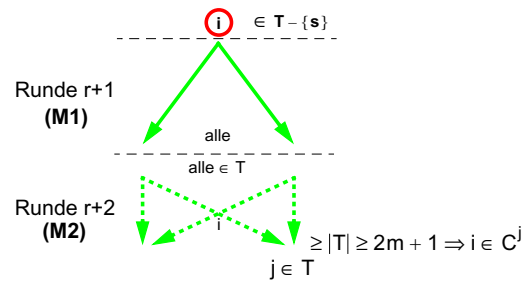
**Beweis für das geforderte Verhalten**

**L12.5** Für korrekte Prozessoren sind die Mengen  $W$  und  $C$  mit der Rundenzahl monoton wachsend.

**Beweis: Definitionen**

**L12.6**  $i, j \in T, i \neq s$ :  $i$  aktiviert in  $q_r$   $\rightarrow$   $i$  von  $j$  bestätigt in  $q_{r+2}$ , d. h.  $i \in C^j(q_{r+2})$

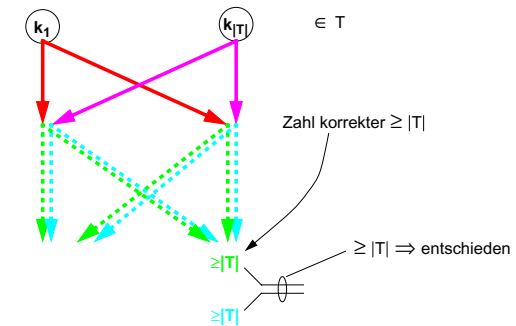
**Beweis:**



## BP 1 Einigungsprotokolle: Algorithmus von Dolev

**L12.7** Alle  $i \in T$  aktiviert in  $q_r \rightarrow$  alle  $i \in T$  entschieden in  $q_{r+2}$

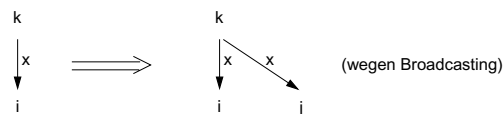
**Beweis**



## BP 1 Einigungsprotokolle: Algorithmus von Dolev

**L12.8**  $i, j \in T$ :  $W_x^i(q_r) \cap T = W_x^j(q_r) \cap T$

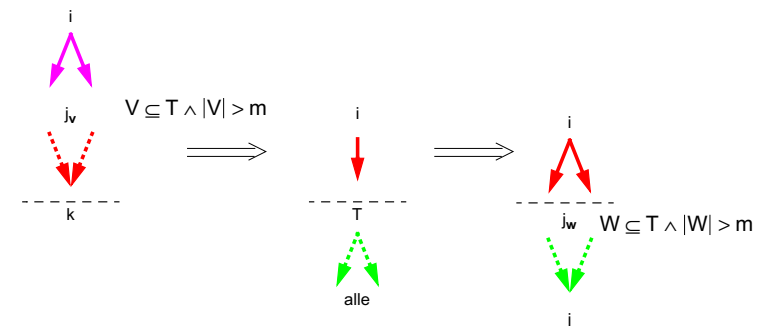
**Beweis**



## BP 1 Einigungsprotokolle: Algorithmus von Dolev

**L12.9**  $i, j \in T$ : ex.  $k \in T$  mit  $i \in C^k(q_r) \rightarrow$  für alle  $j \in T$  ist  $i \in C^j(q_{r+1})$

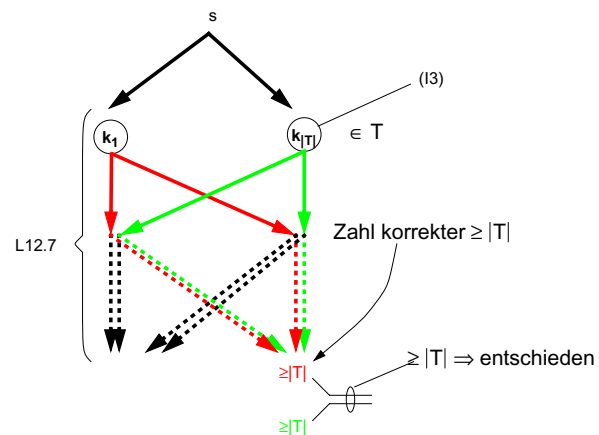
**Beweis**



## BP 1 Einigungsprotokolle: Algorithmus von Dolev

**L12.10**  $s \in T \wedge s_0 = 1 \rightarrow$  alle  $j \in T$  entschieden in  $q_4$

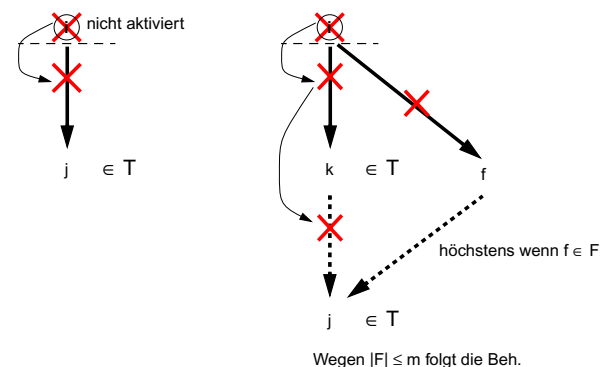
**Beweis**



## BP 1 Einigungsprotokolle: Algorithmus von Dolev

**L12.11**  $i, j \in T : i$  nicht aktiviert in  $q_r \rightarrow i \notin W^j(q_{r+1})$  und  $|W^j(q_{r+2})| < \text{LOW}$

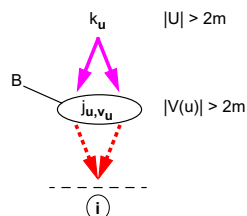
**Beweis**



## BP 1 Einigungsprotokolle: Algorithmus von Dolev

**L12.12**  $i \in T : i$  entschieden in  $q_r \rightarrow$  ex.  $B \subseteq T$  mit  $|B| = \text{LOW}$  und alle  $j \in B$  aktiviert in  $q_{r-1}$

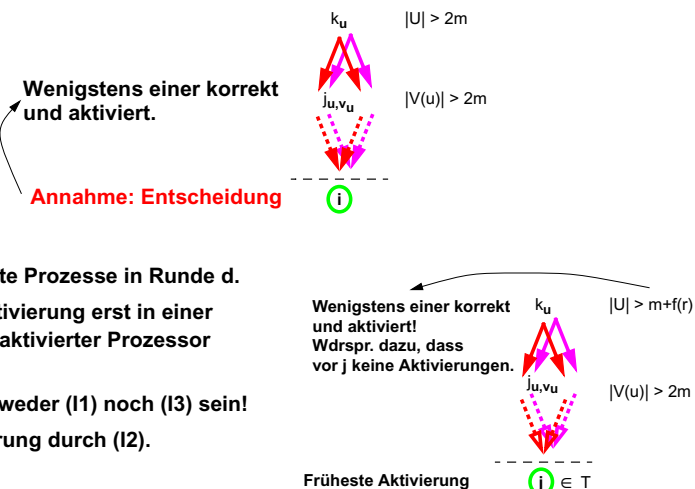
**Beweis**



## BP 1 Einigungsprotokolle: Algorithmus von Dolev

**L12.13**  $s \in T \wedge s_0 = 0 \rightarrow$  kein  $i \in T$  (aktiviert oder entschieden in  $q_{d+1}$ )

**Beweis**



Also aktivierte Prozesse in Runde d.

Früheste Aktivierung erst in einer Runde  $\geq 2$ , aktivierter Prozessor sei j.

Grund kann weder (I1) noch (I3) sein!

Also Aktivierung durch (I2).



BP 1 Einigungsprotokolle: Algorithmus von Dolev	
<b>L12.14</b>	<p><math>A \subseteq T - \{s\} \wedge  A  = \text{LOW}</math> und alle <math>i \in A</math> aktiviert in <math>q_r</math>  <math>\implies</math> alle <math>j \in T</math> entschieden in <math>q_{r+4}</math></p> <p><b>Beweis</b></p> <p>a) <math>s \in T</math>: Wegen L12.13 muß <math>s_0 = 1</math> sein. Dann folgt Behauptung aus (I3) und L12.7.</p> <p>b) <math>s \notin T</math>: Es seien alle <math>i \in A</math> aktiviert in <math>q_r</math> und <math>r'</math> minimal.</p> <p>Nach L12.6 alle <math>i \in A</math> von allen <math>j \in T</math> bestätigt in <math>q_{r+2}</math></p> <p>b1) <math>r' = 1</math>: Kann nach Definition nicht sein.</p> <p>b2) <math>r' = 2</math>: Wegen L12.6 alle <math>i \in A</math> von allen <math>j \in T</math> bestätigt in <math>q_4</math>. Da <math>f(4) = 0</math> also auch entschieden.</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> 12.3-65

BP 1 Einigungsprotokolle: Algorithmus von Dolev	
<b>L12.15</b>	<p><math>i \in T</math>: <math>i</math> entschieden in <math>q_{d+1}</math> <math>\implies</math> alle <math>j \in T</math> entschieden in <math>q_{d+1}</math></p> <p><b>Beweis</b></p> <p>a) <math>(m = 0) \vee s \in T</math>: L12.10 und L12.13</p> <p>b) <math>m &gt; 0 \wedge s \notin T</math>: <math>\implies d &gt; 4</math></p> <p><math>i \in T</math> habe sich in <math>q_{d+1}</math> entschieden  <math>\implies \exists A (A \subseteq T \wedge ( A  = \text{LOW}) \wedge \forall j (j \in A \implies j \text{ aktiviert\_in } q_d))</math> wegen L12.12.</p> <p><math>r</math> sei erste Runde in der ein derartiges <math>A</math> existiert.</p> <p>b1) <math>r \leq d - 3 \implies</math> Beh. wegen L12.14.</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> 12.3-67

BP 1 Einigungsprotokolle: Algorithmus von Dolev	
<b>b3) <math>r' &gt; 2</math>:</b>	<p>Es muß <math>k</math> existieren, das in <math>q_r</math> erstmalig aktiviert wird. Grund kann nur I2 sein.</p> <p><math>\implies  C^k(q_r)  \geq \text{LOW} + \max(0, \lceil r'/2 \rceil - 2)</math></p> <p><math>r'</math> minimal <math>\implies k \notin C^k(q_r)</math>, denn <math>k \in C^k(q_r) \implies  W_k^k(q_r)  \geq \text{HIGH}</math> nach Def.  <math>\implies k</math> aktiviert in <math>q_{r-2}</math> wegen L12.11, Wdrspr.</p> <p><math>\forall j (j \in T \implies C^j(q_{r+2}) \supseteq C^k(q_r))</math> wegen L12.5 und L12.9.</p> <p><math>\forall j (j \in T \implies k \in C^j(q_{r+2}))</math> wegen L12.6, da <math>k</math> in <math>q_r</math> aktiviert ist.</p> <p>Wegen der letzten beiden Aussagen zusammen mit <math>k \notin C^k(q_r)</math></p> <p><math>\implies  C^j(q_{r+2})  \geq \text{LOW} + \max(0, \lceil \frac{r'}{2} \rceil - 2) + 1 = \text{LOW} + \max(0, \lceil \frac{r'+2}{2} \rceil - 2)</math></p> <p>also <math>j</math> aktiviert in <math>q_{r+2}</math> wegen (I2) und wegen L12.6 entschieden in <math>q_{r+4}</math>.</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> 12.3-66

BP 1 Einigungsprotokolle: Algorithmus von Dolev	
<b>b2) <math>r \geq d - 2 = 2m + 1 &gt; 2</math></b>	<p><math>\implies \exists k (k \in A \wedge k \text{ aktiviert\_in } q_r \wedge \neg(k \text{ aktiviert\_in } q_{r-1}))</math> nach Def. von <math>r</math> und LOW</p> <p>Wegen (I2) <math>\implies  C^k(q_r)  \geq \text{LOW} + \max(0, \lceil \frac{r}{2} \rceil - 2) \geq \text{LOW} + m - 1 = \text{HIGH} - 1</math>.</p> <p>Es ist <math>s \notin T \wedge s \notin C^k(q_r)</math> nach Voraussetzung und Def. von <math>C^k(q_r)</math>,</p> <p>also enthält <math>C^k(q_r)</math> höchstens <math>m-1</math> fehlerhafte Prozessoren.</p> <p><math display="block">C^k(q_r) = \overline{\{j   j \in N - \{s\} \wedge  W_j^k(q_r)  \geq \text{HIGH}\}}</math></p> <p><math>\implies  \{j   j \in T \wedge  W_j^k(q_r)  \geq \text{HIGH}\}  \geq \text{LOW}</math></p> <p><math>\implies \exists A' (A' \subseteq T \cap C^k(q_r) \wedge ( A'  = \text{LOW})</math> wegen L12.11 und L12.5.  <math>\wedge \forall j (j \in A' \implies j \text{ aktiviert\_in } q_{r-1}))</math></p> <p>Widerspruch zur Wahl von <math>r</math>, also ist dieser Fall nicht möglich.</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> 12.3-68

**BP 1** Einigungsprotokolle: Algorithmus von Dolev

**L12.16**  $s \in T$ :  $s_0 = 1 \rightarrow$  alle  $j \in T$  entschieden in  $q_{d+1}$   
 $s_0 = 0 \rightarrow$  kein  $j \in T$  entschieden in  $q_{d+1}$

**Beweis**

1. Zeile folgt unmittelbar aus L12.7.
2. Zeile bewiesen als L12.13.

**BP 1** Einigungsprotokolle: Algorithmus von Burns und Neiger**12.3.4** Algorithmus von Burns und Neiger (Einigungs-Algorithmus, 1994)**Bezeichnungen**

- $n$  Zahl der Prozessoren
- $m$  Maximalzahl fehlerhafter Prozessoren
- $f$  Zahl tatsächlich fehlerhafter Prozessoren

**Bedingungen**

$$f \leq m \text{ und } n \geq m^2 + 4m + 1$$

**BP 1** Einigungsprotokolle: Algorithmus von Burns und Neiger

**Partitionierung von  $m^2 + 4m + 1$  Prozessoren in  $m + 1$  disjunkte Gruppen  $G_i$**

$G_1$  enthält  $2m + 1$  Prozessoren.

Für  $2 \leq i \leq m + 1$  besteht  $G_i$  aus  $2(m + 2 - i) + 1$  Prozessoren.

d. h. es enthält

$G_1$   $2m + 1$  Prozessoren,

$G_2$   $2m + 1$  Prozessoren,

$G_3$   $2m - 1$  Prozessoren,

.

.

.

$G_{m+1}$  3 Prozessoren.

**BP 1** Einigungsprotokolle: Algorithmus von Burns und Neiger**Algorithmus**

Der Algorithmus läuft in Runden.

- Jeder Prozessor besitzt einen Wert, der anfänglich den Anfangswert darstellt.
- In jeder Runde  $r$  senden die Prozessoren der Gruppe  $G_r$  ihren Wert an alle anderen.
- Am Rundenende ermittelt jeder Prozessor der Runde seinen Wert durch Mehrheitsbildung. Für ausgebliebene Nachrichten wird der eigene Wert benutzt.
- Wenn ein Prozessor von allen Prozessoren der Gruppe  $r$  den gleichen Wert erhält, entscheidet er sich für diesen Wert.
- Prozessoren, die sich bis zum Ende der  $(m+1)$ -ten Runde noch nicht entschieden haben, entscheiden sich für ihren Wert.

**BP 1** Einigungsprotokolle: Algorithmus von Burns und Neiger

Jeder Prozessor  $p$  führt folgenden Algorithmus (in synchronisierten Runden) aus:

```
for (i = 1; i <= m + 1; i++) {  
  // Ein Schleifendurchlauf stellt jeweils eine Runde dar.  
  if (G[i].contains(p))  
    for (j = 1; j <= n; j++) send(j, p, v);  
  tmp = 0;  
  for (pix = G[i].first(); pix != 0; G[i].next(pix)) {  
    receive(pix, &m); tmp += m;  
  }  
  if (!decided) {  
    if (tmp == 0 || tmp == G[i].length()) {  
      v = (tmp ? 1 : 0); decided = TRUE; break;  
    } else {  
      v = 2 * tmp > G[i].length() ? 1 : 0;  
    }  
  }  
}
```

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

12.3-73

**BP 1** Einigungsprotokolle: Algorithmus von Burns und Neiger**Grundgedanke**

- Da es  $m + 1$  Gruppen gibt, muß wenigstens eine fehlerfreie existieren (d. h. alle ihre Mitglieder sind fehlerfrei).  
Wenn eine fehlerfreie Gruppe sendet, erhalten alle die gleiche Multimenge von Werten.  
Wenn einmal alle fehlerfreien Prozessoren den gleichen Wert besitzen, so behalten sie ihn bei.
- Falls  $f < m$  müssen unter  $G_1, G_2, \dots, G_{f+2}$  wenigstens zwei fehlerfrei sein.
- Die Mitglieder der zweiten fehlerfreien Gruppe senden alle den gleichen Wert, also haben sich spätestens nach der Runde  $f + 2$  alle fehlerfreien entschieden.

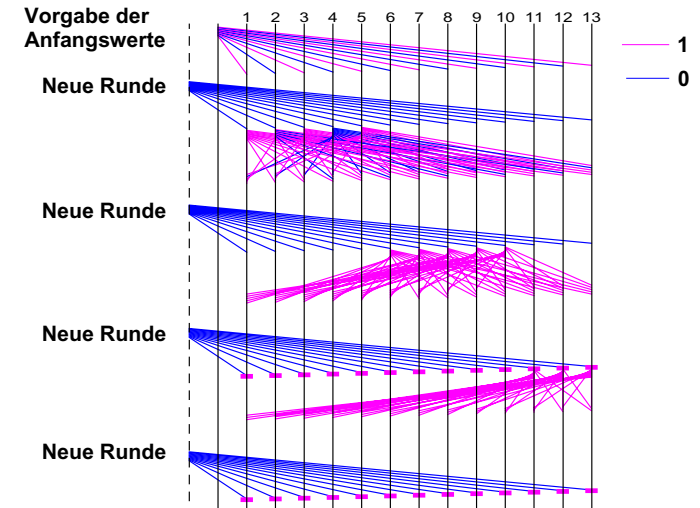
17.01.02

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

12.3-75

**BP 1** Einigungsprotokolle: Algorithmus von Burns und Neiger

Beispiel: Einigungsprotokoll von Burns und Neiger ( $f = m = 2$ )



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

12.3-74

**BP 1** Einigungsprotokolle: Algorithmus von Burns und Neiger

- Wenn unter  $G_1, \dots, G_i$  wenigstens eine fehlerfreie Gruppe ist, beenden alle fehlerfreien Prozessoren die  $i$ -te Runde mit dem gleichen Wert.
- Wenn unter  $G_1, \dots, G_i$  wenigstens zwei fehlerfreie Gruppen sind, haben sich nach der  $i$ -ten Runde alle fehlerfreien Prozessoren entschieden.
- Wenn alle fehlerfreien Prozessoren, die  $i$ -te Runde mit dem gleichen Wert beginnen, beenden sie die Runde mit eben diesem Wert.
- Wenn sich ein fehlerfreier Prozessor am Ende der  $i$ -ten Runde entschieden hat, dann beenden alle fehlerfreien Prozessoren die Runde mit dem gleichen Wert.

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

12.3-76

**BP 1** Einigungsprotokolle: Algorithmus von Burns und Neiger**Beweis durch vollständige Induktion****Induktionsanfang:  $i = 1$** 

- (a)  $G_1$  fehlerfrei ---> alle erhalten gleiche Wertemenge  
---> alle fehlerfreien entscheiden sich für gleichen Wert
- (b) entfällt
- (c) Folgt aus a) und Gewinnung des Wertes durch Mehrheitsbestimmung
- (d) Entschieden ---> alle erhaltenen Werte gleich  
---> von den fehlerfreien gleichen Wert erhalten  
---> bilden die Mehrheit, da  $G_1$  mindestens  $2m+1$  Prozessoren enthält  
---> bei allen liefert die Mehrheit den gleichen Wert

**BP 1** Einigungsprotokolle: Algorithmus von Burns und Neiger**Induktionsannahme: Aussage gültig bis einschließlich  $i$** **Induktionsschluß: Fall  $i+1$** 

- (a) 1.  $G_1, \dots, G_i$  fehlerhaft  
--->  $G_{i+1}$  fehlerfrei  
---> alle erhalten gleiche Wertemenge  
---> alle fehlerfreien ermitteln gleichen Wert
- 2. Wenigstens eine Gruppe  $G_j$  fehlerfrei,  $j$  sei maximal  
---> alle fehlerfreien haben am Ende der  $j$ -ten Runden den gleichen Wert  
---> alle fehlerfreien haben am Ende der folgenden Runden den gleichen Wert (Eigenschaft c)

**BP 1** Einigungsprotokolle: Algorithmus von Burns und Neiger

- (b) 1.  $G_1, \dots, G_i$  enthalten nur eine fehlerfreie Gruppe  
--->  $G_{i+1}$  fehlerfrei  
---> alle erhalten in der  $(i+1)$ -ten Runde lauter gleiche Werte  
---> alle fehlerfreien entscheiden sich für diesen Wert
- 2. Wenigstens zwei Gruppen aus  $G_1, \dots, G_i$  sind fehlerfrei  
---> alle fehlerfreien haben sich schon zu Beginn der  $(i+1)$ -ten Runde entschieden
- (c) 1.  $G_1, \dots, G_i$  enthalten wenigstens eine fehlerfreie  
---> Behauptung wegen a) und Algorithmus
- 2.  $G_1, \dots, G_i$  fehlerhaft  
--->  $G_{i+1}$  enthält höchstens  $m - i < 2(m + 2 - (i + 1)) + 1$  fehlerhafte  
---> die fehlerfreien sind in der Mehrzahl  
---> alle fehlerfreien entscheiden sich für den gemeinsamen Anfangswert

**BP 1** Einigungsprotokolle: Algorithmus von Burns und Neiger

- (d)  $G_p$  entscheide sich am Ende der  $(i+1)$ -ten Runde  
--->  $G_p$  hat von allen aus  $G_{i+1}$  den gleichen Wert bekommen und unter  $G_1, \dots, G_i$  ist höchstens eine fehlerfreie Gruppe  
--->  $G_{i+1}$  enthält höchstens  $m - i + 1$  fehlerhafte Prozessoren
- Da  $G_{i+1}$  insgesamt  $2(m + 2 - (i + 1)) + 1 > 2(m + 1 - i)$  Prozessoren enthält, ist die Mehrheit fehlerfrei  
---> alle fehlerfreien haben am Ende der Runde den gleichen Wert

## BP 1 Einigungsprotokolle: Algorithmus von Burns und Neiger

**S12.5** Satz: Der Algorithmus von Burns und Neiger löst das Einigungsproblem. Spätestens nach  $\min(f + 2, m + 1)$  Runden haben sich alle Prozessoren entschieden.

Beweis:

Terminierung: trivial

Richtigkeit: Unmittelbare Folgerung aus c)

**Einigung:** Da es  $m + 1$  Gruppen gibt, aber höchstens  $m$  fehlerhafte Prozessoren, muß wenigstens eine Gruppe fehlerfrei sein. Damit folgt aus a) die Einigung spätestens in der  $(m+1)$ -ten Runde.

Ist  $f < m$ , so enthält  $G_1, \dots, G_{f+2}$  wenigstens zwei fehlerfreie Gruppen, woraus nach b) die Behauptung folgt.

## BP 1 Einigungsprotokolle: PAXOS

Synchronisation

- Die Ausführung von Zustandsübergängen erfolgt atomar
- L maximale Ausführungsdauer der (atomaren) Durchführung eines Zustandsübergangs
- D maximale Übertragungsdauer von Nachrichten

Architektur

Multi-Paxos (Folge von Einigungsvorgängen)

Basic Paxos (ein einzelner Einigungsvorgang)

Wahl eines Anführers

Ausfallerkennung

Nachrichtensystem

Hardware

## BP 1 Einigungsprotokolle: PAXOS

### 12.4 PAXOS

#### 12.4.1 Aufgabenstellung

Modellannahmen

• Prozessoren

- Mögliches Fehlverhalten: Fail Stop  
d. h. Prozessoren sind entweder *lebendig* oder *ausgefallen*
- Wiederanlauf erlaubt

• Kommunikationssystem

- vollständig vermascht
- jeder Prozessor kann auch an sich selbst senden
- zuverlässige FIFO-Verbindungen

Aufgabenstellung

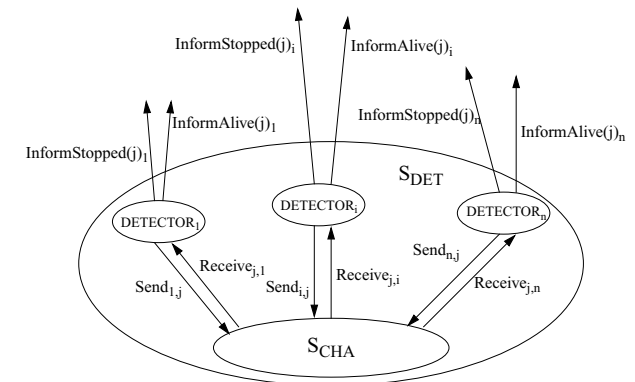
• Einigung

- Wertevorrat für Anfangswerte beliebig
- Einigung auf den Anfangswert eines der Prozessoren

## BP 1 Einigungsprotokolle: PAXOS

### 12.4.2 Ausfallerkennung

Systemstruktur beschrieben mit synchron gekoppelten E/A-Automaten



**BP 1** Einigungsprotokolle: PAXOS

## Formale Beschreibung der Komponenten

**DETECTOR**(z, c)<sub>i</sub>

Signature:

Input: Receive(m)<sub>j,i</sub>, Stop<sub>i</sub>, Recover<sub>i</sub>Internal: Check(j)<sub>i</sub>Output: InformStopped(j)<sub>i</sub>, InformAlive(j)<sub>i</sub>, Send(m)<sub>i,j</sub>

State:

Clock  $\in \mathbb{R}$  initially arbitraryStatus  $\in \{\text{alive}, \text{stopped}\}$  initially aliveAlive  $\in 2^I$  initially Ifor all  $j \in I$ Prevrec(j)  $\in \mathbb{R}^{\geq 0}$  initially arbitraryLastinform(j)  $\in \mathbb{R}^{\geq 0}$  initially ClockLastsend(j)  $\in \mathbb{R}^{\geq 0}$  initially ClockLastcheck(j)  $\in \mathbb{R}^{\geq 0}$  initially Clock**BP 1** Einigungsprotokolle: PAXOS

## Actions:

**input Stop<sub>i</sub>**

Eff: Status = stopped;

**input Recover<sub>i</sub>**

Eff: Status = alive;

**input Receive("Alive")<sub>j,i</sub>**

Eff: if (Status == alive) {

Prevrec(j) = Clock;

if (j  $\notin$  Alive) {Alive = Alive  $\cup \{j\}$ ;

Lastcheck(j) = Clock + c;

}

}

**BP 1** Einigungsprotokolle: PAXOS**internal Check(j)<sub>i</sub>**Pre: Status == alive && j  $\in$  Alive

Eff: Lastcheck(j) = Clock + c;

if (Clock &gt; Prevrec(j) + z + D) { Alive = Alive - {j}; }

**output Send("Alive")<sub>i,j</sub>**

Pre: Status == alive

Eff: Lastsend(j) = Clock + z;

**output InformStopped(j)<sub>i</sub>**Pre: Status == alive && j  $\notin$  Alive

Eff: Lastinform(j) = Clock + L;

**output InformAlive(j)<sub>i</sub>**Pre: Status == alive && j  $\in$  Alive

Eff: Lastinform(j) = Clock + L;

**BP 1** Einigungsprotokolle: PAXOS**time\_passage v(t)**

Pre: Status == alive

Eff: Let t' be such that {

 $\forall j, \text{Clock} + t' \leq \text{Lastinform}(j);$  $\forall j, \text{Clock} + t' \leq \text{Lastsend}(j);$  $\forall j, \text{Clock} + t' \leq \text{Lastcheck}(j);$ 

}

Clock = Clock + t'

## BP 1 Einigungsprotokolle: PAXOS

### 12.4.3 Wahl eines Anführers

Anführer fungiert als Leiter von Einigungsrunden

- Zur Vereinfachung wird ein Verfahren ausgewählt, das eventuell zu mehreren Anführern führt
- Durch Zuteilung von logischen Zeitstempeln kann gewährleistet werden, daß sich schließlich einer durchsetzt (ähnlich Wahl bei allgemeinen Kommunikationsgraphen)

## BP 1 Einigungsprotokolle: PAXOS

Formale Beschreibung

LEADERELECTOR<sub>i</sub>

Signature:

Input: InformStopped(j)<sub>i</sub>, InformAlive<sub>i</sub>, Stop<sub>i</sub>, Recover<sub>i</sub>

Output: Leader<sub>i</sub>, NotLeader<sub>i</sub>

State:

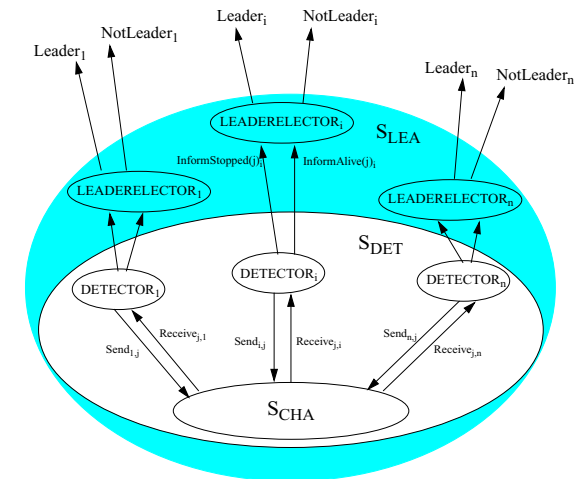
Status ∈ {alive, stopped} initially alive

Pool ∈ 2<sup>I</sup> initially {i}

Leader ∈ I initially i

## BP 1 Einigungsprotokolle: PAXOS

Systemstruktur



## BP 1 Einigungsprotokolle: PAXOS

Actions:

**input Stop<sub>i</sub>**

Eff: Status = stopped;

**input Recover<sub>i</sub>**

Eff: Status = alive;

**input InformStopped(j)<sub>i</sub>**

Eff: if (Status == alive) {

Pool = Pool - {j};

Leader = max(Pool);

}

**input InformAlive(j)<sub>i</sub>**

Eff: if (Status == alive) {

Pool = Pool ∪ {j};

Leader = max(Pool);

}

BP 1	Einigungsprotokolle: PAXOS
	<p><b>output Leader<sub>i</sub></b>  Pre: Status == alive &amp;&amp; i == Leader  Eff: skip;</p> <p><b>output NotLeader<sub>i</sub></b>  Pre: Status == alive &amp;&amp; i != Leader  Eff: none;</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
	12.4-93

BP 1	Einigungsprotokolle: PAXOS
12.4.4	<p><b>Basic Paxos</b></p> <p><b>Modularisierung</b></p> <p><b>BASICPAXOS</b> versucht eine Entscheidungsrunde durchzuführen, an deren Ende sich alle Teilnehmer der Runde entschieden haben.</p> <p>Runden können möglicherweise wegen Teilausfällen nicht beendet werden. Entscheidungsrunde wird deshalb von <b>STARTERALG</b> überwacht. Diese Komponente startet nach einer gewissen Zeit eine neue Runde mit dem gleichen Anführer.</p> <p><b>BPSUCCESS</b> verbreitet das Ergebnis einer erfolgreich beendeten Entscheidungsrunde.</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
	12.4-94

BP 1	Einigungsprotokolle: PAXOS
	<p><b>Systemstruktur</b></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
	12.4-95

BP 1	Einigungsprotokolle: PAXOS
	<p><b>BASICPAXOS</b></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
	12.4-96



BP 1	Einigungsprotokolle: PAXOS
12.4.5	<p><b>BASICPAXOS</b></p> <p><b>Informelle Beschreibung</b></p> <p><b>"Collect"-Nachrichten (r, "Collect"):</b> Werden vom Rundenanführer als Zeichen des Beginns der Runde r an alle versandt.</p> <p><b>"Last"-Nachrichten (r, "Last", r', v'):</b> Werden bei Empfang einer "Collect"-Nachricht der Runde r an den Rundenanführer gesandt. Dabei ist r' die letzte Runde, in der der Sender einen Wert, nämlich v', als zustimmungsfähig akzeptiert hat. Existiert keine solche Runde wird r' = (0,i) gesetzt und als Wert der Anfangswert zurückgegeben bzw. nil.</p> <p><b>"Begin"-Nachricht (r, "Begin", v):</b> Wird vom Rundenanführer als Vorschlag für den Einigungswert v versandt, der anhand der "Last"-Nachrichten ermittelt wird.</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> <p>12.4-97</p>

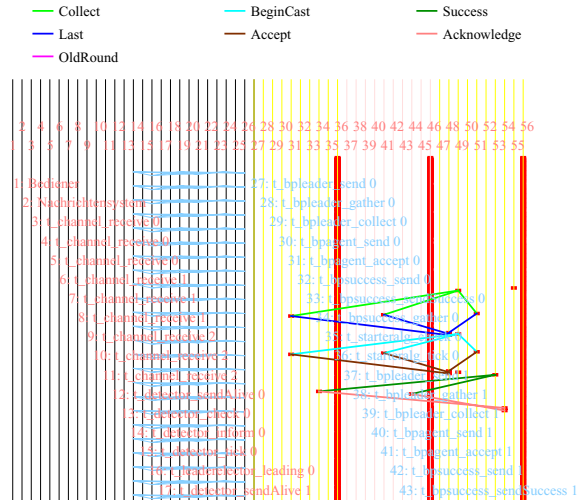
BP 1		Einigungsprotokolle: PAXOS	
rndV	value of this round	<b>Anführer<sub>i</sub></b>	<b>Teilnehmer<sub>j</sub></b>
rndVFrom	round from which rndV is taken	rndV = initV curRnd = new(max) rndInfQuo = { } rndAccQuo = { } acked = { }	
curRnd	number of current round	(r = curRnd, Collect)	
rndInfQuo	set of processes for which "Last" has been received		r > commit commit = r r' = lastR v' = lastV
rndAccQuo	set of processes for which "Accept" has been received	rndInfQuo += {j} rndVFrom < r' rndV = v' rndVFrom = r'	(r, Last, r', v')
acked	set of processes for which "Acknowledge" has been received	rndInfQuo  > N/2 (r = curRnd, BeginCast, v = rndV)	
commit	number of committed round		r ≥ commit: lastR = r lastV = v
lastR	number of latest round	(r, Accept)	
lastV	value for round lastR	r == curRnd rndAccQuo += {j}	
decision	value for which processor decided	rndAccQuo  > N/2 (Success, v = rndV) decision = rndV	
		(Acknowledge)	decision = v
		acked += {j}	
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.		12.4-99

BP 1	Einigungsprotokolle: PAXOS
	<p><b>"Accept"-Nachricht (r, "Accept"):</b> Wird als Antwort auf "Begin"-Nachricht zurückgesandt zum Zeichen, daß der vorgeschlagene Wert akzeptiert wird.</p> <p><b>"OldRound"-Nachricht (r, "OldRound", r'):</b> Wird als Antwort auf eine "Collect"- oder "Begin"-Nachricht gesandt zum Zeichen, daß der Sender in der späteren Runde r' einen Einigungswert akzeptiert hat.</p> <p><b>"Success"-Nachricht ("Success", v):</b> Wird vom Rundenanführer am Ende einer erfolgreichen Runde als Entscheidungswert gesandt.</p> <p><b>"Ack"-Nachricht ("Ack"):</b> Wird an den Rundenanführer gesandt, als Zeichen des Empfangs einer "Success"-Nachricht, d. h. der Kenntnis des Entscheidungswertes.</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> <p>12.4-98</p>

BP 1	Einigungsprotokolle: PAXOS
	<p><b>Abstrakte Darstellung</b></p> <p><b>Signature of BPLEADER</b></p> <p>Input: Receive(m)<sub>j,i</sub>, m ∈ {"Last", "Accept", "OldRound"}</p> <p>Internal: Collect<sub>i</sub>, GatherLast<sub>i</sub>, GatherAccept<sub>i</sub>, GatherOldRound<sub>i</sub></p> <p>Output: Send(m)<sub>i,j</sub>, m ∈ {"Collect", "Begin"}, BeginCast<sub>i</sub>, RndSuccess(v)<sub>i</sub></p> <p><b>Signature of BPAGENT</b></p> <p>Input: Receive(m)<sub>j,i</sub>, m ∈ {"Collect", "Begin"}</p> <p>Internal: LastAccept<sub>i</sub>, Accept<sub>i</sub></p> <p>Output: Send(m)<sub>i,j</sub>, m ∈ {"Last", "Accept", "OldRound"}</p> <p><b>Signature of BPSUCCESS</b></p> <p>Input: Receive(m)<sub>j,i</sub>, m ∈ {"Ack", "Success"}, RndSuccess(v)<sub>i</sub></p> <p>Internal: SendSuccess, GatherSuccess<sub>i</sub>, GatherAck<sub>i</sub></p> <p>Output: Send(m)<sub>i,j</sub>, m ∈ {"Success"}, Decide(v)<sub>i</sub></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> <p>12.4-100</p>

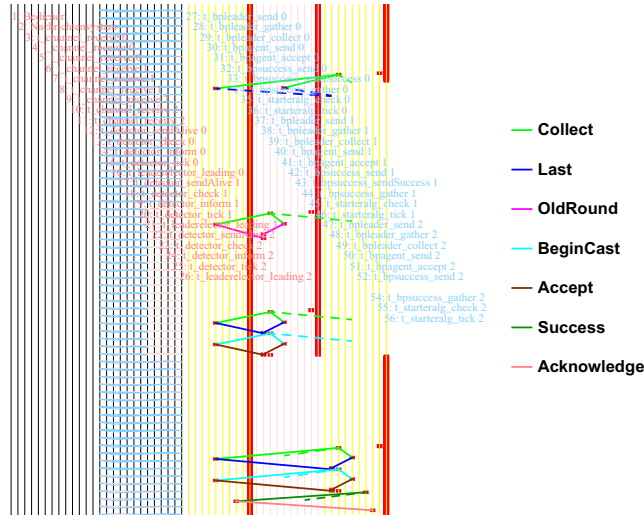
## BP 1 Einigungsprotokolle: PAXOS

### Ungestörter Ablauf



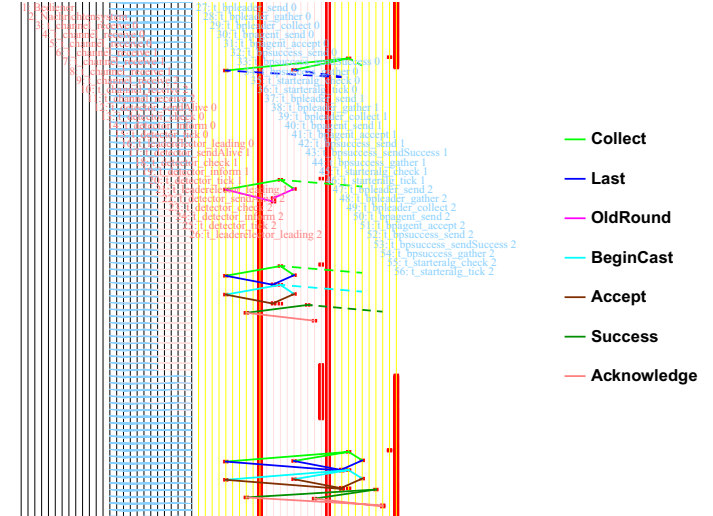
## BP 1 Einigungsprotokolle: PAXOS

### Zeitweiser Ausfall der Prozessoren 1 und 2



## BP 1 Einigungsprotokolle: PAXOS

### Zeitweiser Ausfall von Prozessor 2



## BP 1 Einigungsprotokolle: PAXOS

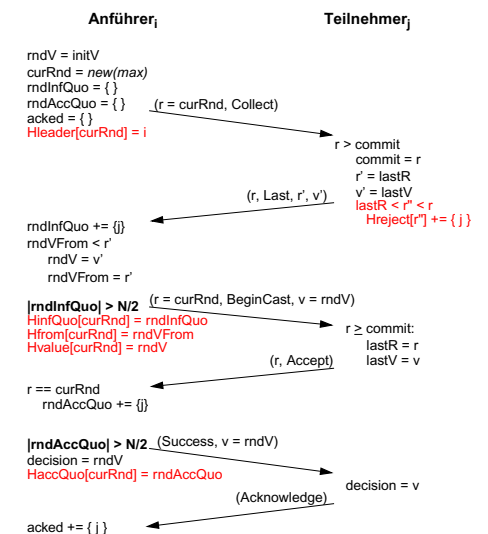
### Gültigkeit

#### Anreicherung mit "Historien"-Variablen

Invariante 1: In jedem Zustand der Ausführung gilt für jede Runde  $r$  mit  $Hleader[r] = nil$ :

$Hvalue[r] = nil$   
 $Hfrom[r] = nil$   
 $HinfQuo[r] = \{ \}$   
 $HaccQuo[r] = \{ \}$

**D11.1** Eine Runde  $r$  heißt tot, wenn  $|Hreject[r]| \geq n/2$  ist.



BP 1 Einigungsprotokolle: PAXOS	
D4.2	R bezeichne den Wertevorrat für Rundennummern. $R_S$ bezeichne die Menge $\{ r \in R \mid Hleader_i[r] \neq nil \}$ . $R_V$ bezeichne die Menge $\{ r \in R \mid Hvalue_i[r] \neq nil \}$ . Invariante 2: In jedem Systemzustand ist $R_V \subseteq R_S$ .
	D4.3 $r \in R_V$ heißt verankert, wenn für jede Runde $r' \in R_V$ mit $r' < r$ $r'$ tot oder $Hvalue[r'] = Hvalue[r]$ ist. Invariante 3: In jedem Zustand, in dem Prozessor j an Prozessor i die Nachricht (r, "Last", r", v) gesendet hat, besitzen alle $r'$ mit $r'' < r' < r$ die Eigenschaft $j \in Hreject(r')$ . Invariante 4: In jedem Zustand, in dem Prozessor $j \in rndInfQuo_i$ und $curRnd_i = r$ ist, gilt für alle Runden $r'$ mit $rndVFrom_i < r' < r$ die Beziehung $j \in Hreject(r')$ .
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 12.4-105

BP 1 Einigungsprotokolle: PAXOS	
	Terminierung
D4.4	Unter einem <i>Ausführungsfragment</i> wird ein Ausschnitt aus einer Ausführung des Algorithmus verstanden, der während eines (gewählten) Zeitintervalls stattfand. Ein Ausführungsfragment heißt <i>stabil</i> , wenn während des Zeitintervalls, in dem es ausgeführt wurde, Prozessoren weder ausgefallen noch wiederangelaufen sind. Ein Ausführungsfragment heißt <i>gut</i> , wenn es stabil ist und eine Mehrheit von Prozessoren während seiner Ausführung lebendig ist.
L12.17	In einer Runde, die in einem stabilen Ausführungsfragment enthalten ist, werden maximal 6n Nachrichten versandt.
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 12.4-107

BP 1 Einigungsprotokolle: PAXOS	
	Invariante 5: In jedem Zustand, in dem Prozessor $j \in HinfQuo(r)$ ist, gilt für alle Runden $r'$ mit $Hfrom(r) < r' < r$ die Beziehung $j \in Hreject(r')$ . Invariante 6: In jedem Zustand ist jede lebendige (nicht tote) Runde verankert. Invariante 7: In jedem Zustand besitzen alle Variablen $decision_i \neq nil$ den gleichen Wert. Dieser Wert ist Anfangswert mindestens eines Prozessors.
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 12.4-106

BP 1 Einigungsprotokolle: PAXOS	
L12.18	Ein Ausführungsfragment a, in dem noch kein Prozessor sich für einen Wert entschieden hat, gelte folgendes: <ol style="list-style-type: none"> <li>a ist stabil,</li> <li>In a existiert genau ein Anführer i,</li> <li>a dauert mindestens <math>4 * L + 2 * n * L + 2 * D</math> Zeiteinheiten,</li> <li>Prozeß i ist in a Anführer der (geeignet gewählten) Runde r und</li> <li>Runde r ist erfolgreich.</li> </ol> Dann wird $rndSuccess$ spätestens zum Zeitpunkt $7 * L + 4 * n * L + 4 * D$ ab Beginn des Ausführungsfragmentes ausgeführt.
L12.19	Wenn ein Ausführungsfragment stabil ist, wenigstens $3 * L + 2 * n * L + 2 * D$ Zeiteinheiten dauert und genau ein Anführer sich bereits vor Beginn des Fragmentes entschieden hat, dann weiß er, daß jeder lebendige Prozessor sich entschieden hat und es wurden höchstens 2n Nachrichten versandt.
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 12.4-108

**L12.20** Ein Ausführungsfragment  $a$  besitze folgende Eigenschaften:

1.  $a$  ist gut,
2. in  $a$  gibt es genau einen Anführer  $i$  und
3.  $a$  dauert mindestens  $16 * L + 8 * n * L + 9 * D$  Zeiteinheiten.

Dann hat der Anführer zum Zeitpunkt  $16 * L + 8 * n * L + 9 * D$  eine Entscheidung getroffen.

**L12.21** Es sei  $a$  ein gutes Ausführungsfragment, das mindestens  $24 * L + 10 * n * L + 13 * D$  Zeiteinheiten dauert. Dann gilt:

1. Der Anführer entscheidet sich nach höchstens  $21 * L + 8 * n * L + 11 * D$  Zeiteinheiten und es werden höchstens  $8n$  Nachrichten versandt.
2. Spätestens nach  $24 * L + 10 * n * L + 13 * D$  Zeiteinheiten ab Beginn des Ausführungsfragments entscheidet sich jeder lebendige Prozessor und es wurden maximal  $2n$  weitere Nachrichten versandt.