

BP 1	Transaktionen: Allgemeines
6 6.1	<p><b>Transaktionen</b></p> <p><b>Allgemeines</b></p> <ul style="list-style-type: none"> <li>Problembeispiel: Überweisung von \$5 von A an B           <ul style="list-style-type: none"> <li>read account A (obtaining \$10);</li> <li>read account B (obtaining \$15);</li> <li>write \$5 to account A;</li> <li>write \$20 to account B;</li> </ul> </li> <li>Vorstellung: Kunde A hat sein Konto bei der Bank A Kunde B hat sein Konto bei der Bank B Er tätigt die Anweisung bei der Bank C Jede Bank hat ihren eigenen Rechner usw.</li> <li>Problem: Systemzusammenbruch während oder zwischen den Schreibanweisungen Festplatte mit den Daten defekt Nachrichtenverlust oder Nachrichtenverdoppelung</li> <li>Ziel: Alle Schreibanweisungen oder keine werden durchgeführt und das Ergebnis ist bekannt</li> </ul>
26.11.01	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           6.1-1

BP 1	Transaktionen: Allgemeines
	<ul style="list-style-type: none"> <li>Transaktionen bei Betriebssystemen:           <ul style="list-style-type: none"> <li>auf die Dauerhaftigkeit kann oft verzichtet werden</li> <li>kein Logbuch notwendig</li> </ul> </li> <li>Betrachtung am Modell von Lampson           <ul style="list-style-type: none"> <li>die Probleme werden in dem Artikel detailliert dargestellt</li> <li>prinzipielle Lösungsmöglichkeiten werden vorgestellt</li> <li>effektive Implementierung</li> <li>keine effiziente Implementierung</li> <li>einige der Annahmen sind zu hinterfragen</li> </ul> </li> </ul>
26.11.01	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           6.1-2

BP 1	Transaktionen: System-Überblick
6.2	<p><b>System-Überblick</b></p> <pre> graph TD     subgraph "Das stabile System (stable system)"         T[Transaktionen (transactions)]         AI[Absichten (intentions)]         S[Sperren (locks)]         ZAT[Zusammengesetzte atomare Aktionen (compound atomic transactions)]         SM[Stabile Mengen (stable sets)]         F[Fernaufrufe (remote procedures)]         EID[Eindeutige Bezeichner (unique IDs)]         SP[Stabile Prozessoren (stable processors)]     end      subgraph "Das physikalische System (physical system)"         KA[Wiederholbare Aktionen (restartable actions)]         SS[Stabiler Speicher (stable storage)]         K[Kommunikationssystem (communications)]         P[Prozessoren (processors)]         PS[Plattenspeicher (disk storage)]     end      T --&gt; AI     T --&gt; S     AI --&gt; ZAT     S --&gt; SM     ZAT --&gt; F     SM --&gt; EID     F --&gt; KA     EID --&gt; SP     SP --&gt; SS     KA --&gt; K     SP --&gt; P     SS --&gt; PS     K --&gt; P     P --&gt; PS   </pre>
26.11.01	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           6.2-3

BP 1	Transaktionen: Konsistenz und Transaktionen
6.3	<p><b>Konsistenz und Transaktionen</b></p> <ul style="list-style-type: none"> <li><b>Konsistenz</b> <ul style="list-style-type: none"> <li>System erfüllt in jedem "beobachtbaren" Zustand ein bestimmtes Prädikat, die sogenannte Invariante.</li> </ul> </li> <li><b>Systemverhalten definiert durch Übergangsfunktion F</b> <ul style="list-style-type: none"> <li><math>state := F(state)</math></li> </ul> </li> <li><b>Ausgabemenge von F</b> <ul style="list-style-type: none"> <li>Modifizierbare Zustandsanteile.</li> </ul> </li> <li><b>Eingabemenge von F</b> <ul style="list-style-type: none"> <li>Zustandsanteile, die das Ergebnis beeinflussen.</li> </ul> </li> </ul>
26.11.01	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           6.3-4

<b>BP 1</b>	<b>Transaktionen: Konsistenz und Transaktionen</b>
	<ul style="list-style-type: none"> <li>◆ <b>Transaktionen (atomare Aktionen)</b> Eine Folge von Lese-/Schreibaufrufen mit der Eigenschaft, daß entweder alle oder keiner der Schreibaufrufe sichtbar sind. Sie werden durch Klammerung mit den (System-)Funktionsaufrufen Begin und End gekennzeichnet.</li> <li>◆ <b>Datenbefehle:</b> Read, Write</li> <li>◆ <b>Kontrollbefehle:</b> Begin, End, Abort</li> </ul>
26.11.01	<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>

6.3-5

<b>BP 1</b>	<b>Transaktionen: Das physikalische System</b>
<b>6.4</b>  <b>6.4.1</b>	<b>Das physikalische System</b>  <b>Der Plattenspeicher</b>  <ul style="list-style-type: none"> <li>◆ <b>Operationen:</b>             procedure put(at: address; data: datablock)             procedure get(at: address)                returns (status: (good, looksBad); data: datablock)         </li> </ul>
26.11.01	<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>

6.4-6

<b>BP 1</b>	<b>Transaktionen: Das physikalische System</b>
	<ul style="list-style-type: none"> <li>◆ <b>get</b> <ul style="list-style-type: none"> <li>• erwünscht: Seite ist (good, d) und get liefert (good, d) Seite ist (bad, ...) und get liefert (looksBad, ...)</li> <li>• Fehler: Seite ist (good, d) und get liefert (looksBad, ...) (Wiederholung, wie oft?)</li> <li>• Ruin: Seite ist (good, d) und get liefert <math>n_R</math>-mal (looksBad, ...) (Ist hier eingeordnet wegen der noch zu erläuternden Auswirkungen auf die Überlegungen zur Verlustverwandtschaft) Seite ist (bad, ...) und get liefert (good, ...) Seite ist (good, d) und get liefert (good, d') mit <math>d \neq d'</math></li> </ul> </li> </ul>
26.11.01	<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>

6.4-7

<b>BP 1</b>	<b>Transaktionen: Das physikalische System</b>
	<ul style="list-style-type: none"> <li>◆ <b>put</b> <ul style="list-style-type: none"> <li>• erwünscht: auf Plattenspeicher (good, d)</li> <li>• Fehler: auf Plattenspeicher unverändert auf Plattenspeicher (bad, ...) spontaner Datenverlust des Plattenspeichers Verlustmenge: Menge aller Seiten, die durch eine Verlustursache betroffen werden können, z. B. Spur durch Ausfall eines Schreib-/Lesekopfes Zylinder durch Justierungsfehler Oberfläche durch Kopfabsturz Einheit durch Fehler in der Steuerelektronik Verlustintervall <math>T_D</math>: Innerhalb des Verlustintervalls tritt höchstens eine Verlustursache auf</li> </ul> </li> </ul>
26.11.01	<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>

6.4-8

BP 1	Transaktionen: Das physikalische System
	<p>♦ Zwei Seiten heißen genau dann verlustverwandt, wenn es eine Verlustmenge gibt, der beide angehören.</p> <ul style="list-style-type: none"> <li>• Fehler: (good, d) geht selten nach (bad, d'), d. h. höchstens ein Verlust in <math>T_D</math> (bad, d) geht über in (good, d)</li> <li>• Ruin: (good, d) geht häufig nach (bad, d), d. h. mehrere Verluste in <math>T_D</math> (good, d) geht nach (good, d') mit <math>d \neq d'</math></li> </ul>
26.11.01	<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> 6.4-9

BP 1	Transaktionen: Das physikalische System
6.4.2	<b>Prozessoren</b> <ul style="list-style-type: none"> <li>• Jeder Prozessor besitzt feste Zahl von Prozessen</li> <li>• Ein Prozeß liefert Zeit für Messung von Verlustintervallen <math>T_D</math> (Timeout)</li> <li>• Fehler: Zusammenbruch und Wiederaufsetzen mit irgendeinem Standardzustand, d. h. (virtueller) Speicher ist flüchtig</li> <li>• Annahme: Alle anderen Fehlfunktionen werden entdeckt und führen zu einem Zusammenbruch, bevor Plattenspeicher oder das Kommunikationssystem betroffen werden. Zusammenbrüche werden entdeckt.</li> </ul>
26.11.01	<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> 6.4-10

BP 1	Transaktionen: Das physikalische System
6.4.3	<b>Kommunikationssystem</b> <ul style="list-style-type: none"> <li>• Operationen: <ul style="list-style-type: none"> <li>send(to: processor; data: messageblock)</li> <li>receive() returns ( status: (good, bad), data: messageblock)</li> </ul> </li> <li>• receive <ul style="list-style-type: none"> <li>• erwünscht: <ul style="list-style-type: none"> <li>es existiert (good, d, p) und es wird (good, d) abgeliefert.</li> <li>es existiert (bad, d, p) und es wird (bad, ...) abgeliefert.</li> </ul> </li> </ul> </li> <li>• send <ul style="list-style-type: none"> <li>• erwünscht: <ul style="list-style-type: none"> <li>nach Ausführung existiert (good, d, p)</li> </ul> </li> </ul> </li> </ul>
26.11.01	<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> 6.4-11

BP 1	Transaktionen: Das physikalische System
	<ul style="list-style-type: none"> <li>• spontane Ereignisse</li> <li>• Fehler: <ul style="list-style-type: none"> <li>Nachricht geht verloren</li> <li>Nachricht wird vervielfacht</li> <li>(good, d, p) geht nach (bad, ..., ...)</li> </ul> </li> <li>• Ruin: <ul style="list-style-type: none"> <li>(bad, d, p) geht nach (good, ..., ...)</li> <li>(good, d, p) geht nach (good, d', p') mit <math>d \neq d'</math> oder <math>p \neq p'</math>.</li> </ul> </li> </ul>
26.11.01	<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> 6.4-12

BP 1	Transaktionen: Das physikalische System
6.4.4	<p><b>Einfache, zusammengesetzte und wiederholbare Aktionen</b></p> <ul style="list-style-type: none"> <li>• <b>Einfache Aktionen</b> <ul style="list-style-type: none"> <li>• <b>atomar:</b> <ul style="list-style-type: none"> <li>- <b>Ganzheit des Effekts</b> (alles oder nichts beobachtbar) <ul style="list-style-type: none"> <li>kein Zusammenbruch: <math>\{P\} S \{Q\}</math></li> <li>Zusammenbruch: <math>\{P\} S \{P \vee Q\}</math></li> </ul> </li> <li>- <b>Serialisierbarkeit</b></li> <li>-</li> </ul> </li> </ul> </li> <li>• <b>Zusammengesetzte Aktionen</b> <ul style="list-style-type: none"> <li>• <b>Nicht alle Aktionen aller Abstraktionsebenen können atomar gemacht werden.</b></li> </ul> </li> <li>• <b>Allgemeine Eigenschaften:</b> <ul style="list-style-type: none"> <li>- <b>kein Zusammenbruch:</b> <math>\{P\} S \{Q_{ok}\}</math></li> <li>- <b>Zusammenbruch:</b> <math>\{P\} S \{Q_{crash}\}</math></li> </ul> <p>Die Gültigkeit von <math>Q_{crash}</math> wird beim Wiederanlauf (recovery) hergestellt.</p> </li> </ul>
26.11.01	<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>

6.4-13

BP 1	Transaktionen: Das physikalische System
26.11.01	<p>• <b>Im Vergleich zur Atomarität abgeschwächte Forderungen:</b></p> <p>a) <math>Q_{ok} \Rightarrow Q_{crash}</math> Zusammenbruch unmittelbar vor return möglich</p> <p>b) <math>Q_{crash} \Leftrightarrow P \vee Q_{ok}</math> Gleichwertig zu Ganzheit</p> <p>c) <math>Q_{crash} \Rightarrow P</math> Wiederholbarkeit (Idempotenz): wird notwendig, wenn Atomarität nicht gewährleistet werden kann</p>
26.11.01	<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>

6.4-14

BP 1	Transaktionen: Das stabile System
6.5	<p><b>Das stabile System</b></p> <p>6.5.1 <b>Stabiler Speicher</b></p> <ul style="list-style-type: none"> <li>• <b>Höhere Abstraktionen des Plattenspeichers:</b> <ul style="list-style-type: none"> <li>• <b>realer Plattenspeicher</b> <math>\rightarrow</math> <b>sorgfältiger Plattenspeicher</b> <math>\rightarrow</math> <b>stabiler Plattenspeicher</b></li> </ul> </li> <li>• <b>Sorgfältiger Plattenspeicher:</b> <ul style="list-style-type: none"> <li>• <b>sorgfältiges Lesen:</b> Wiederholung bis Zustand good oder <math>n_R</math>-mal</li> </ul> </li> <li>• <b>sorgfältiges Schreiben:</b> Evtl. wiederholtes Schreiben und jeweils anschließendes einfaches Lesen, bis Leseergebnis good</li> </ul>
26.11.01	<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>

6.5-15

BP 1	Transaktionen: Das stabile System
26.11.01	<p>• <b>Stabiler Plattenspeicher:</b></p> <ul style="list-style-type: none"> <li>• <b>stabile Seite:</b> Seite und nicht verlustverwandte Schattenseite</li> <li>• <b>stabiles Lesen:</b> Seite sorgfältig lesen; falls Zustand bad Schattenseite sorgfältig lesen</li> <li>• <b>stabiles Schreiben:</b> Seite sorgfältig schreiben, anschließend Schattenseite sorgfältig schreiben</li> </ul>
26.11.01	<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>

6.5-16

BP 1	Transaktionen: Das stabile System
	<ul style="list-style-type: none"> <li>• <b>Bereinigung:</b> <ul style="list-style-type: none"> <li>- Seite und Schattenseite sorgfältig lesen; Beide im Zustand good, dann fertig; falls eine davon im Zustand bad, dann andere sorgfältig lesen und die "bad"-Seite damit sorgfältig beschreiben; falls beide im Zustand good, aber mit verschiedenen Daten, dann eine der beiden auswählen und in die andere sorgfältig schreiben.</li> </ul> </li> <li>• Bereinigung bei Systemstart, nach Zusammenbruch und mindestens alle <math>T_D</math> Zeiteinheiten.</li> <li>• Es muß eine Zuordnungsfunktion vorgesehen werden, <ul style="list-style-type: none"> <li>- mit deren Hilfe alle stabilen Seiten aufgezählt werden können und</li> <li>- die zu jeder stabilen Seite die sie darstellende Seite und Schattenseite zu ermitteln gestattet.</li> </ul> </li> </ul>
26.11.01	<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> 6.5-17

BP 1	Transaktionen: Das stabile System
6.5.2	<b>Stabile Prozessoren</b> <ul style="list-style-type: none"> <li>• <b>Plattenspeicher als verlängerter Arbeitsspeicher</b> <math>\Rightarrow</math> flüchtig</li> <li>• <b>Stabiler Speicher zur Sicherung von Prozeßzuständen;</b> <ul style="list-style-type: none"> <li>• <b>save:</b>      Speichert PCB, Aufrufparameter, ...</li> <li>• <b>reset</b></li> <li>• <b>restart:</b>    Durch diese Operation kann nach Zusammenbruch ein Prozeß in seinem zuletzt gesicherten Zustand wieder aufgesetzt werden.</li> </ul> </li> <li>• <b>Stabile Monitore</b> <ul style="list-style-type: none"> <li>• Lokale Daten, ausgenommen die Sperre zum gegenseitigen Ausschluß, werden im stabilen Speicher geführt.</li> <li>• Modifikation der lokalen Variablen durch eine Monitorprozedur erfolgt mit einem einzigen stabilen Schreibaufruf.</li> <li>• Prozeßzustände werden in Monitoren nicht gesichert.</li> <li>• Nur Monitore können Daten enthalten, die nicht prozeßlokal sind.</li> <li>• Höchstens ein stabiles Put (also atomar!)</li> </ul> </li> </ul>
26.11.01	<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> 6.5-18

BP 1	Transaktionen: Zusammengesetzte Aktionen
6.5.3	<b>Fernaufrufe</b> <ul style="list-style-type: none"> <li>• Kommunikation zwischen Maschinen erfolgt nur durch Fernaufrufe.</li> <li>• At-least-once-Semantik reicht wegen der Wiederholbarkeit</li> </ul>
6.6	<b>Zusammengesetzte Aktionen</b> <ul style="list-style-type: none"> <li>• Bisher bereitgestellte Basis: <ul style="list-style-type: none"> <li>• Stabiler Speicher</li> <li>• Stabile Prozesse</li> <li>• Stabile Monitore</li> </ul> </li> </ul> Keine flüchtigen Daten
26.11.01	<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> 6.6-19

BP 1	Transaktionen: Zusammengesetzte Aktionen
6.6.1	<b>Stabile Mengen</b> <ul style="list-style-type: none"> <li>• Eine stabile Menge besteht aus <ul style="list-style-type: none"> <li>• Verbunden, die kleiner sind als eine stabile Seite,</li> <li>• einem eindeutigen Bezeichner.</li> </ul> </li> <li>• Alle Operationen zur Manipulation stabiler Mengen sind wiederholbar.</li> <li>• Alle Operationen liefern eine Fehlermeldung, wenn sie nicht zwischen Create und dem nächsten Erase ausgeführt werden (Create und Erase ausgenommen).</li> </ul>
26.11.01	<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> 6.6-20

BP 1	Transaktionen: Zusammengesetzte Aktionen
	<ul style="list-style-type: none"> <li>• <b>Atomare Operationen:</b> <ul style="list-style-type: none"> <li>• <b>Create(ID i)</b> Erzeugt eine neue stabile Menge mit Bezeichner i. Falls sie schon existiert, ist die Operation wirkungslos.</li> <li>• <b>Insert(StableSet s, StableSet t, Record new)</b> Setzt voraus, daß new noch nicht in s oder t enthalten ist, und nimmt es in s und t auf. Einer der beiden Parameter s oder t kann auch nil sein.</li> <li>• <b>Replace(StableSet s, StableSet t, Record old, Record new)</b> Setzt voraus, daß old in s und t enthalten ist. Old wird entfernt und new aufgenommen.</li> <li>• <b>IsEmpty(StableSet s)</b> Liefert true bzw. false, je nachdem ob s leer ist oder nicht.</li> <li>• <b>IsMember(StableSet s, Record r)</b> Liefert true, wenn r in s enthalten ist, sonst false.</li> </ul> </li> </ul>
26.11.01	<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> 6.6-21

BP 1	Transaktionen: Zusammengesetzte Aktionen
	<ul style="list-style-type: none"> <li>• <b>Nicht atomare Operationen:</b> <ul style="list-style-type: none"> <li>• <b>Enumerate(StableSet s, procedure p)</b> Ruft p für jedes Element in s auf. In den Programmen geschrieben als: for r in s do p;</li> <li>• <b>Erase(StableSet s)</b> Entfernt alle Elemente aus s. Wurde ein Element bei seiner Einfügung mittels Insert oder Replace gleichzeitig in eine weitere stabile Menge eingefügt, so wird es auch dort entfernt.</li> </ul> </li> <li>• <b>Implementierungsmöglichkeiten</b> <ul style="list-style-type: none"> <li>• <b>Pool</b> lokal auf einem Prozessor</li> <li>• <b>Wide Stable Set</b> Root auf einem Processor kennt die Verteilung der Blätter auf andere Prozessoren</li> </ul> </li> </ul>
26.11.01	<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> 6.6-22

BP 1	Transaktionen: Zusammengesetzte Aktionen
6.6.2	<b>Zusammengesetzte atomare Aktionen</b>  <b>Gestalt einer Methode:</b> A(...) { R; }  <b>Implementierung (nur für wiederholbare Prozeduren):</b>  <pre> A(...) {   Save /* Process State */;   R;   Reset /* Process State */; } </pre>
26.11.01	<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> 6.6-23

BP 1	Transaktionen: Transaktionen
6.7	<b>Transaktionen</b> <ul style="list-style-type: none"> <li>• <b>Eindeutiger Transaktionsbezeichner:</b> typedef UniqueId TI</li> <li>• <b>Seitenadresse:</b> class SA {    FileName Dateiname;                   int Seitennummer;                   };</li> <li>• <b>Aktion:</b> enum RW {read, write};</li> <li>• <b>Absicht:</b> class Intention {    TI    t;                           SA    p;                           RW    a;                           char    d[...]; // zu modif. Daten;                   };</li> </ul>
26.11.01	<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> 6.7-24

BP 1

Transaktionen: Transaktionen

- Phase:
 

```
enum Phase
{nonexistent, running, committed, aborted};
```
- Transaktionsmarke:
 

```
class TM {  TI      t;
           Phase  ph;
};
```

26.11.01

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

6.7-25

BP 1

Transaktionen: Transaktionen

```
entry Phase GetPhase(TI t)
{
  for r in s.flags do {
    if (r.t == t)
      return r.ph;
  }
  return nonexistent;
}
```

26.11.01

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

6.7-26

BP 1

Transaktionen: Transaktionen

```
entry SetPhase(TI t, Phase desiredPhase)
/* desiredPhase may not be nonexistent */
{ switch (GetPhase(t)) {
  case committed:
  case aborted: break; /* do nothing */
  case running: overwrite(s.flags, TM(t, desiredPhase));
                 break;
  case nonexistent: if (desiredPhase == running)
                     replace(s.flags, nil,
                             TM(t, GetPhase(t)),
                             TM(t, running));
                     break;
  }
};
```

```

graph LR
    nonexistent --> running
    running --> committed
    running --> aborted
  
```

26.11.01

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

6.7-27

BP 1

Transaktionen: Transaktionen

- Für jede Dateiseite auf dem Server-Rechner:
  - Stabile Menge p.locks (mit Absichten als Elementen)
  - Stabiler Monitor zur Bearbeitung dieser Menge
  - Bedingungsvariable p.lockFreed, an der auf Sperren gewartet wird.
- Für jede Transaktion beim Koordinator:
  - Wurzel für verteilte stabile Menge
  - t.intentions (mit Absichten als Elementen; sie sind zu schreiben, wenn die Transaktion mit End beendet wird.)
  - Blatt bei jedem betroffenen Server
- Bei jedem Server:
  - Stabile Menge s.flags (mit Transaktionsmarken als Elementen)
  - Stabiler Monitor zur Bearbeitung dieser Menge.

26.11.01

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

6.7-28

BP 1	Transaktionen: Transaktionen
	<pre> entry function Read(t: TI; p: SA): Data;   var noConflict: Boolean;   var i, i0: intention;   begin     repeat noConflict := true; i0.t = 0;       for i in p.locks do         if i.a=write           then if i.t≠t             then beginWait(p.lockFreed);                   noConflict := false end;             else i0 = i;           until noConflict;         if i0.t = t           then return(i0.d);           else begin Insert(p.locks, t.intentions, &lt;read, nil&gt;);                   return StableGet(p); end         end;       end;     </pre>
26.11.01	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           6.7-29

BP 1	Transaktionen: Transaktionen
	<pre> entry procedure Write(t: TI; p: SA; d: Data);   var noConflict: Boolean;   var d': Data;   begin     repeat noConflict := true; d' := d;       for i in p.locks do /* enumerate(p.locks, q) */         /* The next 5 lines correspond to procedure q */         if i.t &lt;&gt; t           then beginWait(p.lockFreed);                 noConflict := false end             else if i.a = write then d' := i.d;           until noConflict;         if d'≠d           then Replace(p.locks, t.intentions, &lt;write,d'&gt;, &lt;write,d&gt;)           else             Insert(p.locks, t.intentions, &lt;write,d&gt;           end;         end;       end;     </pre>
26.11.01	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           6.7-30

BP 1	Transaktionen: Transaktionen
	<pre> procedure Begin(): TI;   const t = UniqueID();   begin     SetPhase(t, running);     CreateWideStableSet(t);     return t   end;  function End(t: TI): Phase;   return Complete(t, committed);  function Abort(t: TI): Phase;   return Complete(t, aborted);   </pre>
26.11.01	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           6.7-31

BP 1	Transaktionen: Transaktionen
	<pre> function Complete(t: TI; desiredResult: Phase   /* committed or aborted only */): Phase;   begin     if GetPhase(t) == nonexistent then return nonexistent;     Save /* process state */;     SetPhase(t,desiredResult);     /* now the transaction is committed or aborted */     if GetPhase(t) == committed       then for i in t.intentions do         if i.a == write then StablePut(i.p, i.d);       Erase(t.intentions)       /* also erases all corresponding entries in all p.locks and         signals all p.lockFreed conditions */;     Reset /* process state */;     return GetPhase(t)   end;   </pre>
26.11.01	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           6.7-32



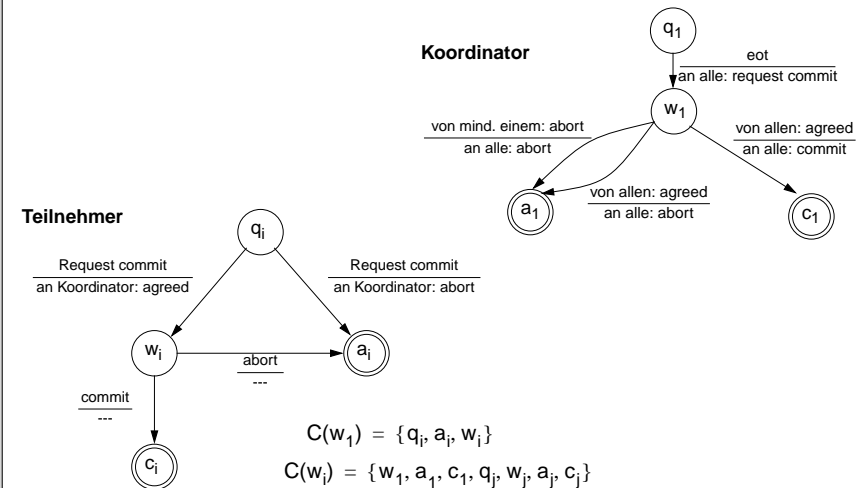
## BP 1 Transaktionen: 2-Phasen-Commit

### 6.8 Commit-Protokolle

- **Systemzustand:** Tupel der lokalen Zustände
- **Atomare Aktionen:**
  - Übergabe einer Nachricht an das Kommunikationssystem
  - Entnahme einer Nachricht aus dem Kommunikationssystem
  - Änderung des lokalen Zustandes
- Nur atomare Aktionen bewirken eine Änderung des Systemzustandes.
- Wenn die beteiligten Prozesse fehlerfrei sind, treten nur durch das Protokoll bedingte Übergänge zwischen Systemzuständen auf.
- Die Nebenläufigkeitsmenge  $C(s)$  eines lokalen Zustandes  $s$  besteht aus allen lokalen Zuständen, die zusammen mit  $s$  in einem erreichbaren Systemzustand vorkommen können.
- Die Sendermenge  $S(s)$  eines lokalen Zustandes  $s$  ist die Menge aller lokalen Zustände, die der Nebenläufigkeitsmenge von  $s$  angehören und von denen Nachrichten an  $s$  gerichtet sein können.

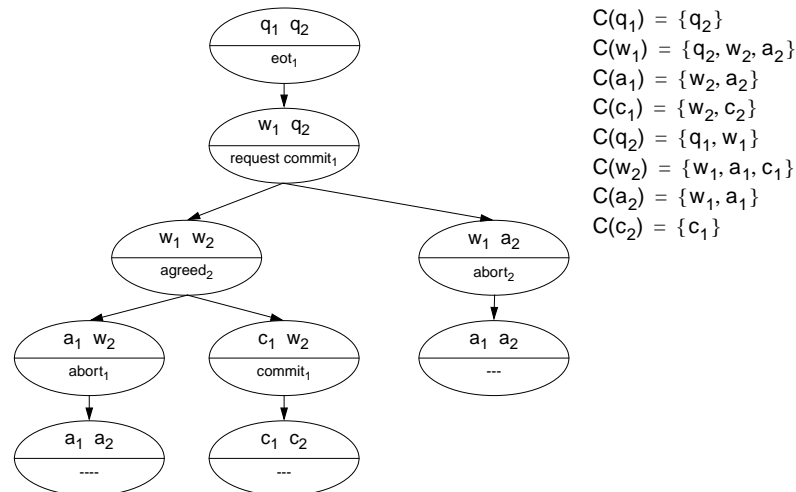
## BP 1 Transaktionen: 2-Phasen-Commit

### Das 2-Phasen-Commit Protokoll.



## BP 1 Transaktionen: 2-Phasen-Commit

### Erreichbare Zustände



## BP 1 Transaktionen: 2-Phasen-Commit

### S6.1 Satz

Besitzt ein Protokoll einen Zustand  $s$ , dessen Nebenläufigkeitsmenge sowohl 'Abort'- als auch 'Commit'-Zustände enthält, so kann es bei unabhängigem Wiederanlauf keinerlei Ausfälle - auch nur eines Prozessors - tolerieren.

### Fehler-Regel

Von jedem lokalen Zustand  $s$  aus, der kein Endzustand ist, wird eine F-Kante nach commit eingefügt, wenn  $C(s)$  einen Commit-Zustand enthält, sonst eine F-Kante nach abort.

### Timeout-Regel

Von jedem lokalen Zustand  $s$  aus, der kein Endzustand ist, wird eine T-Kante nach commit (abort) eingefügt, wenn ein  $t \in S(s)$  existiert, von dem eine F-Kante nach commit (abort) führt. Bei einer T-Kante nach abort wird an alle abort gesandt.

### 3-Phasen-Commit

[illegible]

## 6.2 Satz

### 6.3 Satz

26.11.01 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

26.11.01 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

6.8-40

Also muß ein kleinstes  $k$  mit folgender Eigenschaft existieren:

**Sequenz G**

**F, wenn  $i$  und  $j$  fehlerhaft**

$\langle \dots, s_{0,i}, \dots, s_{0,j}, \dots \rangle$

$\langle \dots, a_i, \dots, a_j, \dots \rangle$

$\langle \dots, s_{k-1,i}, \dots, s_{k-1,j}, \dots \rangle$

$\langle \dots, a_i, \dots, a_j, \dots \rangle$

$\langle \dots, s_{k,i}, \dots, s_{k,j}, \dots \rangle$

**entw.  $f(s_{k,i}) = c_i$  oder  $f(s_{k,j}) = c_j$**

$\langle \dots, s_{m,i}, \dots, s_{m,j}, \dots \rangle$

$\langle \dots, c_i, \dots, c_j, \dots \rangle$

**Zwei aufeinanderfolgende Zustände unterscheiden sich nämlich nach Definition in genau einer Komponente. Deshalb ist  $f(s_{k-1,i}) = f(s_{k,i})$  oder  $f(s_{k-1,j}) = f(s_{k,j})$ . Also muß einer ein Abort-Zustand sein.**

**Nach Annahme muß dann der andere ein commit-Zustand sein, also ist der erreichte Endzustand inkonsistent.**