

**ARBEITSBERICHTE
DES INSTITUTS FÜR MATHEMATISCHE
MASCHINEN UND DATENVERARBEITUNG
(INFORMATIK)**

**FRIEDRICH ALEXANDER UNIVERSITÄT ERLANGEN NÜRNBERG
HERAUSGEBER: H. Billing, W. Händler, U. Herzog,
F. Hofmann, K. Leeb, P. Mertens, H. Müller, G. Nees, H. Niemann,
D. Seitzer, B. Schmidt, H.J. Schneider, H. Wedekind**

**WORKSHOP
FEHLERTOLERANTE MEHRPROZESSOR-
UND MEHRRECHNERSYSTEME
ERLANGEN, 14. OKTOBER 1983**

herausgegeben von
Erik Maehle
und
Ernst Schmitter

VERANSTALTET VON DER FACHGRUPPE 3.1.1
'FEHLERTOLERIERENDE RECHENSYSTEME' DER
GESELLSCHAFT FÜR INFORMATIK E.V. (GI)

BAND 16 · NUMMER 11 · ERLANGEN · DEZEMBER 1983

KONZEPTE DES FEHLERTOLERANTEN
ARBEITSPLATZRECHNERS ATTEMPTO

R.Brause, E.Ammann, M.Dal Cin,
E.Dilger, J.Lutz, T.Risse

Zusammenfassung: In diesem Artikel werden Konzepte beschrieben, welche die Entwicklung des Betriebs-Systems für den fehlertoleranten Arbeitsplatz-Rechner ATTEMPTO bestimmen. Besondere Bedeutung kommt dabei einer konsistenten Kooperation von Prozessoren zu.

1. Einführung

ATTEMPTO *) ist ein experimenteller, fehlertoleranter Arbeitsplatz-Rechner, der zur Zeit in Tübingen entwickelt wird.

Folgende Ziele werden verfolgt:

- 1) Entwurf eines Rechners, der dem Benutzer die Möglichkeit bietet, ein seinen jeweiligen Erfordernissen angemessenes Verhältnis zwischen Fehlertoleranz und Rechenleistung zu wählen.
- 2) Bei der Implementierung soll gezeigt werden, daß es möglich ist,
 - a) mit kommerzieller Hardware (Single-Board-Computer SBC) ohne nennenswerte Eigenentwicklung,

*) A Testable Experimental MultiProcessor system with fault-Tolerance: ATTEMPTO, Motto des Gründers der Universität Tübingen

- b) in einer höheren Programmiersprache und
- c) mit einem üblichen Betriebssystemkern einen fehlertoleranten Rechner aufzubauen.

3) Der Rechner soll auch als Testbett für die Erprobung von Diagnosealgorithmen und Rekonfigurationsstrategien dienen. Deshalb müssen die Fehlertoleranzmechanismen modular und hardwareunabhängig gewählt werden können.

2. Die Benutzersicht von ATTEMPTO

Der Benutzer sieht den Arbeitsplatzrechner als Single-User, Multi-Tasking-System; die Realisierung als Multiprozessor-system ist ihm verborgen. Das System bietet die Möglichkeit, die Fehlertoleranzeigenschaften im Gegensatz zu (3), (7), (8) individuell für jede Anwendung zu wählen.

Dies geschieht durch Festlegen eines Fehlertoleranzindex. Die Angabe des Fehlertoleranzindex t bedeutet, daß bei der Ausführung eines Programms transiente oder permanente Fehler in maximal t SBCs toleriert werden, in dem Sinne, daß keine fehlerhaften Ergebnisse ausgegeben werden. Syntaktisch ist der Index Teil des Programmnamens. Beispielsweise bedeutet ein Eintippen von 'MEINPROGRAMM (2)', daß 'MEINPROGRAMM' mit Toleranzgrad $t=2$ ausgeführt werden soll.

Der Benutzer kann das System so viele Programme gleichzeitig ausführen lassen, wie es die Systemkapazität erlaubt: mehrere Programme mit geringem oder wenige mit hohem Fehlertoleranzindex.

3. Die Hardwarekonfiguration des Systems

Die Hardware-Struktur der Implementierung besteht aus handelsüblichen Single-Board-Computern (Intel iSBC 86/12A) mit Dual-Port-Memory, die durch einen Multi-Master-Bus zur Inter-Prozessor-Kommunikation miteinander verbunden sind. Das Benutzer-Terminal ist über eine RS232-Schnittstelle mit allen SBCs verbunden.

4. Systemsoftware

Das Betriebssystem von ATTEMPTO besteht aus identischen, au-

tonomen, lokalen Betriebssystemen ATOS (ATTEMPTOs local Operating System), eines auf jedem SBC. Den Kern von ATOS bildet ein UNIX-ähnliches Einprozessor-Mehrprozeß-Betriebssystem (OS). Darauf baut eine Zwischenschicht von Systemfunktionen auf, welche die Fehlertoleranzeigenschaften bereitstellt. An das Betriebssystem wird für die Inter-Prozessor-Koordination nur die Anforderung gestellt, daß Nachrichten in der gleichen zeitlichen Reihenfolge an die Fehlertoleranzschicht weitergegeben werden, in der sie an das OS von den Device-Handlern übergeben werden.

Zum jetzigen Zeitpunkt kommunizieren die lokalen Betriebssysteme nur zum Zwecke der Fehlertoleranz und der Ressourcenverwaltung.

Um Fehler tolerieren zu können, werden Kopien des Benutzerprogramms asynchron auf mehreren SBCs parallel abgearbeitet. Die Ergebnisse werden anschließend verglichen.

Um aufgetretene Fehler nicht wirksam werden zu lassen, werden verschiedene Verfahren benutzt, z.B. Selbsttests, Rollback oder Rekonfiguration. Im Unterschied zu diesen Methoden basiert die Fehlertoleranz von ATTEMPTO auf einer Fehlermaskierung, und, unsichtbar für den Benutzer, einer Fehlerdiagnose mittels Vergleichstests. Dieser Ansatz weist z.B. gegenüber Selbsttests oder Rollback entscheidende Vorteile auf:

- Jedes Subsystem wird durch eine externe Instanz getestet - während Selbsttests ausfallsichere Hardware innerhalb des Subsystems voraussetzen.
- Die Methode der Vergleichstests ist konzeptionell einfach und unabhängig von der Hardware-Struktur der Subsysteme.
- Sie ist flexibel und übertragbar.
- Sie erfordert keine spezialisierte Hardware, wenig Verwaltungsaufwand und keinen Kontextwechsel für die Testphase. Deshalb ist ein funktionsbegleitendes Testen möglich.
- Sie ist außerdem weitgehend von der Art der vorliegenden Fehler unabhängig, allerdings wird eine Fehlerlokalisierung bis auf Bauteilebene nicht erreicht.

Die Zwischenschicht besteht aus der Kommunikationsinstanz (CI) und der Fehlertoleranzinstanz (FTI).

Diese Zwischenschicht ist transparent.

Als Systemprogrammiersprache wird Modula-2 verwendet (9).

4.1 Die Kommunikationsinstanz (CI)

Wenn ein Benutzerprogramm eine Dienstleistung des System anfordert (z.B. Eingabe/Ausgabe), so geschieht dies durch einen sogenannten CI-Call. Die Kommunikationsinstanz überprüft die Syntax und die Zulässigkeit der Anforderungen und schützt dadurch das System. Sie stellt auch die Möglichkeiten bereit zur Kommunikation zwischen Betriebssystemkern (OS) und FTI sowie zwischen zwei FTI auf verschiedenen SBCs.

4.2 Die Fehlertoleranzinstanz (FTI)

Die lokale FTI ist verantwortlich für folgende Betriebssystemaufgaben:

- lokale Interpretation der Benutzerkommandos an ATTEMPTO
- Management der systemweiten Ressourcen (z.B. Terminal)
- Kontrolle der Daten von und zu den Ein- und Ausgabegegeräten
- Systemweite Konsistenz der Systemtafeln

Außerdem werden von der FTI folgende Funktionen für Fehlertoleranzzwecke bereitgestellt:

- dezentrales Dispatching der Benutzerprogramme
- Vergleich der Ausgabedaten der lokalen Kopien eines Benutzerprogramms und anschließende Diagnose bei Nichtübereinstimmung
- Überwachung bei der Ausgabe der Daten
- Fehlerinterpretation und -behandlung

Aus Gründen der Fehlertoleranz besitzt jede FTI ihre eigene Systemtafel, den sogenannten Systemkontrollblock (SCB).

Der SCB enthält eine doppelt verzeigerte Liste (JCBqueue) von Referenzen auf Job-Kontrollblöcke (JCB). Jeder Job-

Kontrollblock charakterisiert ein Benutzerprogramm mit seinem Namen und der Angabe der SBCs, die es ebenfalls ausführen, den sogenannten 'Kollegen'.

Die Prozesse der Zwischenschicht werden durch ein eigenes Prozeßmanagement verwaltet, das betriebsystemunabhängig konzipiert wurde.

4.2.1 Dispatching der Benutzerprogramme

Es findet kein Pre-Scheduling wie bei (8) und (3) statt. Stattdessen wird Dispatching nach dem Prinzip der Anziehung (attraction-principle) während der Laufzeit dezentral durchgeführt. Jeder freie Prozessor bewirbt sich um das in seiner Systemtafel als 'noch zu bearbeiten' gekennzeichnete in der Eingabereihenfolge nächste Benutzerprogramm. Wenn ein Prozessor ein Benutzerprogramm beginnt, ist er in allen Systemtafeln mit Hilfe des Synchronisations-Mechanismus von Abschnitt 5 vermerkt.

4.2.2 Vermeidung von fehlerhafter Ausgabe

Die FTI garantiert folgendermaßen, daß nur korrekte Ergebnisse ausgegeben werden:

Vor jeder Ausgabeoperation vergleichen alle Prozessoren, die dasselbe Programm bearbeiten, die zur Ausgabe anstehenden Daten. Die Daten werden dazu auf eine genormte Länge komprimiert (Signatur-Bildung). Jedem Fehlertoleranzindex ist ein sogenannter Testgraph zugeordnet, durch den festgelegt ist, welche Signaturen zu vergleichen sind. Die Diagnose basiert auf den Ergebnissen dieser Vergleiche. Damit gewinnt jeder SBC Information über den Zustand der Kollegen.

Die eigentliche Ausgabeoperation wird daraufhin von dem Prozessor durchgeführt, der als erster genügend Bestätigung erhalten hat. Die anderen intakten Kollegen überwachen diese Ausgabe.

Die Mechanismen zur Diagnose und Festlegung des ausgebenden Prozessors sind detaillierter in (1), (2) beschrieben.

Da wir primär transiente Fehler annehmen, ist es nicht ratsam, bei jedem auftretenden Fehler den SBC auszutauschen. Stattdessen führt jeder SBC eine eigene Fehlerfrequenzliste, in der auftretende Nichtübereinstimmungen der Resultate notiert werden. Für Service-Zwecke kann der Benutzer ein Rekonfigurationsprogramm starten, das auf allen SBCs des Systems abgearbeitet wird und zu einer Aktualisierung der dezentralen, systemweiten Systemtafeln mittels der lokalen, eventuell unterschiedlichen Fehlerfrequenzlisten führt.

5. Inter-Prozessor Koordination

In manchen Multiprozessorsystemen werden die Einzelaktivitäten der Prozessoren durch eine einzige, meist in einem 'common memory' gelagerte Systemtafel synchronisiert. Diese Lösung hat den gravierenden Nachteil, daß bei einer fehlerhaften Systemtafel (defekter Speicher) oder fehlerhaften Zugriffswegen (defekter Bus) das System zusammenbricht.

5.1 Konsistenz der Systemtafeln

Um solch einen Systemzusammenbruch zu verhindern, hat in ATTEMPTO jeder SBC seine eigene Systemtafel, die er durch Nachrichtenaustausch mit denen der anderen SBCs konsistent hält. Dabei tritt aber das Problem auf, daß die Verarbeitung einer Nachricht selbst wieder vom Zustand der Systemtafel abhängt; so kann beispielsweise ein Prozessor nur dann für einen Job in die Systemtafel eingetragen werden, wenn die Kollegenliste noch nicht vollständig ist. Eine Rückfrage führt in diesem Fall zu erhöhter Kommunikation und belastet das System zusätzlich. Deshalb wählten wir eine andere Lösung, nämlich die zeitliche Reihenfolge der Nachrichten zur Veränderung der Systemtafeln bei allen SBCs gleich zu halten. Damit sind bei gleichem initialem Ausgangszustand die Systemtafeln auch ohne Rückantwort jederzeit konsistent.

5.2 Nachrichten-Austausch in ATTEMPTO

Eine Möglichkeit, die zeitliche Reihenfolge der Nachrichten auf allen SBCs identisch zu halten, bietet ein Broadcast-Bus. Bei dieser Lösung wird jede Nachricht von allen SBCs im System empfangen. Da dieser Bus von verschiedenen Prozes-

soren nur nacheinander benutzt werden kann, ist damit eine eindeutige Reihenfolge der Nachrichten gegeben.

Vom Standpunkt der Fehlertoleranz ist ein solcher Broadcast-Bus aber nicht unbedenklich. Es ist dabei nicht möglich, zwischen zwei Prozessoren Nachrichten auszutauschen, ohne daß ein im System befindlicher, defekter Prozessor diese lesen oder verfälschen kann. Es wäre deshalb besser, ein Kommunikationssystem zu verwenden, bei dem der Sender nur dem eine Nachricht übermittelt, der sie auch erhalten soll.

Da außerdem im Bussystem unserer Implementation (MULTIBUS) für eine Broadcast-Eigenschaft eine zusätzliche, nicht-triviale Hardware-Änderung auf jeder SBC-Platine nötig wäre, entschieden wir uns im Einklang mit der Vorgabe, nur Standard-Hardware zu verwenden, für ein anderes physikalisches Kommunikationsprotokoll. Dazu benutzen wir die Eindeutigkeit der Adresskennung auf dem MULTIBUS (s. Abb. 1).

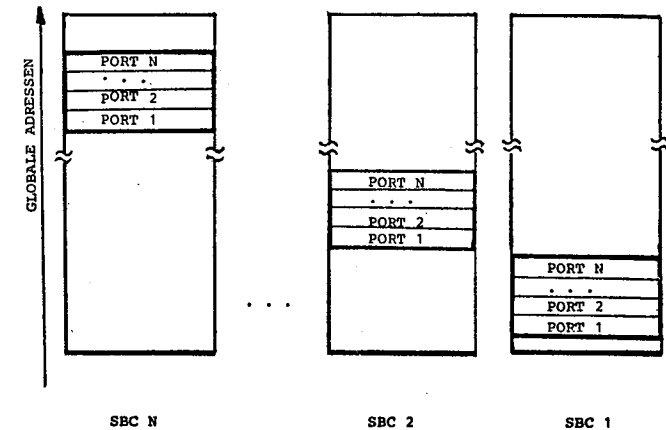


Abb. 1 Aufteilung des globalen Adressbereichs
Die Interprozessor-Synchronisation und Kommunikation basiert in ATTEMPTO auf einem Nachrichtenaustausch, der

über 'ports' (spezielle Speicherbereiche des Dual-Port RAM) abgewickelt wird. Auf jedem SBC existiert ein port für den Empfang von Nachrichten von jedem SBC des Systems. Auf die Bedeutung des ports für Nachrichten an sich selbst ('pseudo-port') wurde bereits in 4.2.1 hingewiesen.

Der sicherere Austausch von Daten wird somit durch eine höhere Busbelastung erkauft - nämlich dann, wenn eine Nachricht an mehrere zu versenden ist. Da jeder Prozessor sofort feststellen kann, ob er eine neue Nachricht erhalten hat, entfällt dafür im Unterschied zu einem Broadcast-Bus die Notwendigkeit, jede einzelne Nachricht zu lesen, um den Empfänger festzustellen.

5.3 Koordination des Nachrichtenaustauschs

Allerdings stellt sich ein weiteres Problem: wenn eine Nachricht von einem Prozessor an alle anderen geht, darf das Versenden einer Nachricht an mehrere Empfänger nicht durch einen anderen Prozessors unterbrochen werden, da sonst die Empfangsreihenfolge der beiden Nachrichten auf allen SBCs nicht mehr identisch ist.

Eine übliche, aber nicht fehlertolerante Möglichkeit, dieses Problem zu lösen, ist die Sperrung des Busses während der gesamten Sendedauer einer Nachricht an alle Empfänger (ein defekter Prozessor könnte damit den Bus dauerhaft blockieren).

Stattdessen entschieden wir uns für folgende Lösung:

Jedem Prozessor am Kommunikationsbus ist eine Interrupt-Signalleitung zugeordnet, die er aktiviert, nachdem er eine Nachricht zu allen Empfängern gesendet hat.

Dabei wird das folgende logische Protokoll einer asynchronen Nachrichtenübertragung benutzt:

Der Interrupt aktiviert die Port-Handler aller SBCs.

Falls ein SBC zu den Empfängern einer Nachricht gehört, leitet dessen Handler sie an das OS weiter, sendet als ACK-Signal eine Nachricht mit der Signatur der empfangenen Nachricht an deren Absender und löscht den Header der Nachricht zum Zeichen, sie gelesen zu haben. Bekommt der Absender keine Empfangsbestätigung (timeout) oder eine solche mit falscher Signatur, so versucht er es mehrmals. Bei Mißerfolg erkennt er auf Bus/Prozessorfehler und kommuniziert nicht mehr mit der betreffenden Einheit (über diesen Bus).

Dies wird in Abb. 2 verdeutlicht. In der Skizze wird die Situation gezeigt, in der die FTI auf SBC 1 erkannt hat, daß ein Job zu bearbeiten und sie selbst frei ist.

Sie sendet diese Anfrage, wie in 4.2.1 beschrieben, mit dem erläuterten Kommunikationsprotokoll an alle SBCs des Systems, einschließlich an sich selbst. Da die Reihenfolge der Nachrichten bei den Port-handlern auf allen SBC und in allen Systempuffern die gleiche ist, ist auch die konsistente Veränderung der Systemtafeln durch die FTI sichergestellt.

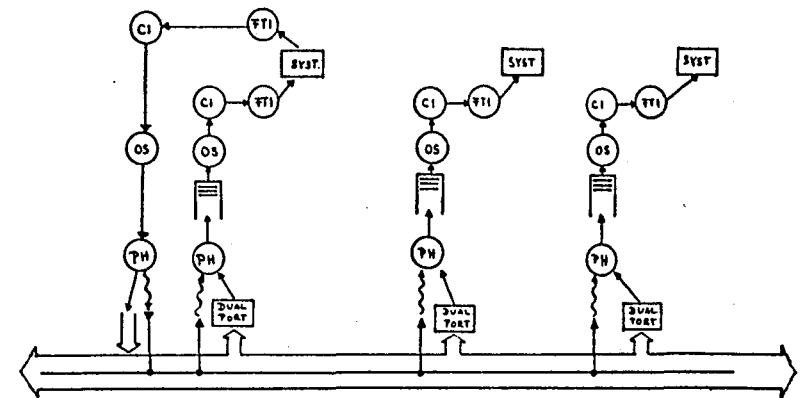


Abb. 2 Koordination der Prozessoren

Die Reihenfolge, in der die Nachrichten vom lokalen ATOS registriert werden, ist somit nicht vom Eintreffen der Nach-

richten auf den verschiedenen SBCs bestimmt, sondern von der Reihenfolge der dazugehörigen Interrupts. Im Vergleich zu den üblichen Broadcast-Systemen ist für den Nachrichtenaustausch eine Interruptleitung für jeden SBC nötig.

Die Forderung, auf allen SBCs die gleiche zeitliche Reihenfolge der Nachrichten zu garantieren, wird damit zu der Forderung, auf allen SBCs die Bearbeitung der Interrupts in der gleichen Reihenfolge sicherzustellen. Im Folgenden werden vier Probleme geschildert, die bei der Realisierung dieser Forderung auftreten, und Lösungsmöglichkeiten angegeben.

a) Fehlende Interruptsequenz-Hardware

Jede Interrupt-Service-Routine (ISR) benötigt eine gewisse Mindestdauer, z.B. die Zeit zum Umladen der Register und Initialisieren der Routine. Erfolgen innerhalb dieser Zeitspanne erneut Interrupts, so können diese nicht sofort bearbeitet werden, sondern werden in einem Interruptregister als Ereignis gespeichert.

Diese Register erlauben aber keine Aussage mehr über das zeitliche Eintreffen der Ereignisse.

Würde man statt der konventionellen nun spezielle Hardware entwickeln, die auch die zeitliche Reihenfolge der Interrupt-Ereignisse beachtet (Aufbau einer FIFO), so wäre es trotzdem nicht möglich, die Unterschiede der fertigungsmäßig und thermisch bedingten Offset-Spannungen der Interrupteingänge zu beseitigen. Dies bedeutet für zwei gleichzeitig generierte Interrupts, daß sie je nach Offset früher oder später registriert und damit unterschiedlich eingeordnet werden.

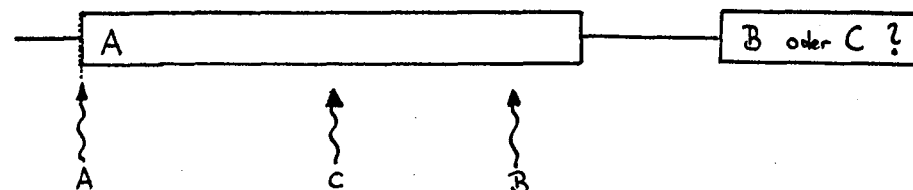


Abb. 3 Keine Hardware-FIFO der Interrupts

Das Problem der identischen zeitlichen Reihenfolge läßt sich stattdessen mit Hilfe der Standard-Hardware befriedigend lösen, indem man den Interruptleitungen der Kommunikation Prioritäten zuordnet. Dies induziert eine feste, andere Ordnung der Abarbeitung der Interrupts als die der zeitlichen Reihenfolge. Da diese Ordnung auf allen SBCs identisch ist, ist die Konsistenz der Systemtafeln gewahrt. Durch den Nachrichten-Rückantwortmechanismus sind zwei Interrupts in kurzem Abstand für den selben Empfänger auf der selben Leitung (vom selben Sender) ausgeschlossen, so daß auch keine Interrupts von intakten Prozessoren verloren gehen können.

b) Unterschiedliche ISR-Abarbeitungszeiten

Wenn beim Bearbeiten der Interrupts Routinen benutzt werden, die verschiedene zeitliche Längen haben (z.B. wenn ein Prozessor eine Nachricht erhält, der andere aber nicht), so

ist eine korrekte Reihenfolge ebenfalls nicht gewährleistet.

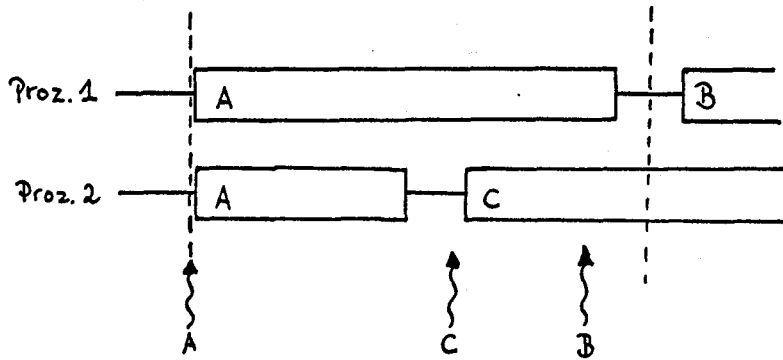


Abb. 4 Unterschiedliche zeitl. Länge der ISR

Nach der Registrierung eines Interrupt zur Kommunikation darf während einer Zeitspanne, die zum Entnehmen einer Nachricht aus einem port ausreicht, kein weiterer Interrupt generiert werden. Vor dem Auslösen eines Interrupts ermittelt deshalb jeder Prozessor, ob diese Zeitspanne seit dem letzten Interrupt schon vergangen ist.

c) Interrupts bei der Zeitabfrage

Ein Interrupt kann gerade in der Situation erfolgen, in der ein Prozessor vor dem Senden ermittelnt hat, daß die oben geforderte Zeitspanne verstrichen ist. Würde er nach dem Abarbeiten der ISR an dieser Stelle das Programm weiter bearbeiten, so würde- er ohne weitere Zeitabfrage den Inter-

rupt für's Senden sofort (im Gegensatz zu b) auslösen.

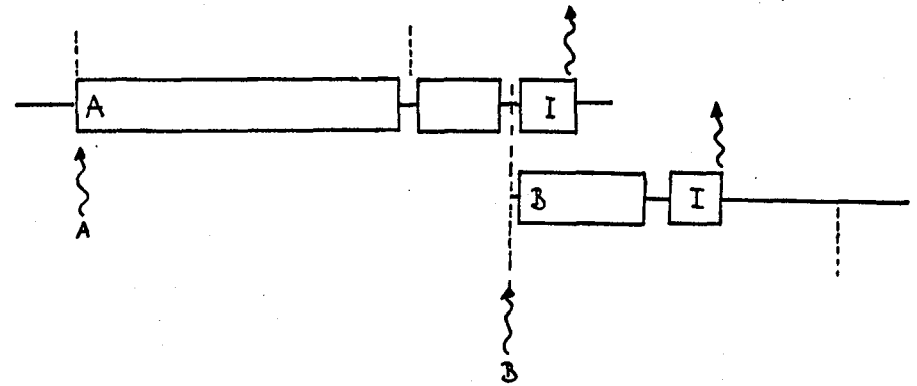


Abb. 5 Interrupts bei der Zeitabfrage

Dieses Problem kann dadurch vermieden werden, daß die Instruktionen zwischen der Zeitabfrage und dem Auslösen des Interrupts ununterbrechbar (clear interrupt) abgearbeitet werden.

d) Verzögerte Interrupts

Führt ein Prozessor gerade selbst eine Instruktion in einem nichtunterbrechbaren Codestück aus und es tritt ein Interrupt auf, so kann der Prozessor den Interrupt erst verzögert bearbeiten. Erfolgt innerhalb dieser Zeit ein weiterer Interrupt, so bearbeitet der eine Prozessor die Interrupts gemäß ihren Prioritäten, die anderen aber möglicherweise nach der zeitlichen Reihenfolge. Dies führt zu Inkonsistenz

der Nachrichten.

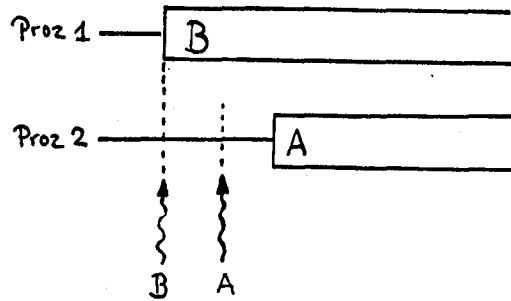


Abb. 6 Verzögerte Bearbeitung der Interrupts

Alle Prozessoren müssen vor Abarbeitung des aktuellen Interrupts eine gewisse 'Karenzzeit' abwarten. Damit ist sichergestellt, daß alle eventuellen Interrupts auch eingetroffen sind. Nach diesem Zeitabschnitt können keine Interrupts mehr von intakten Einheiten eintreffen, da sich alle Prozessoren im 'interrupted'-Zustand befinden. Treffen trotzdem zusätzliche Interrupts ein, so stammen sie von defekten (nichtsynchronisierten) Boards. Über ein Maskierungsbit im Interrupt-Controller kann der Interrupt eines als defekt erkannten SBC wirkungslos gemacht werden.

Da unsere Implementation nicht für zeitkritische Anwendungen gedacht ist, erlaubt der beschriebene, in Software ausgeführte Mechanismus zur Synchronisation im Unterschied zu SIFT (8) und FTMP (5) den Verzicht auf eine globale Systemuhr mit allen ihren Problemen (4).

Diese Arbeit wurde von der Deutschen Forschungsgemeinschaft gefördert.

Referenzen

- (1) E.Ammann, R.Brause, M.Dal Cin, E.Dilger, J.Lutz, T.Risse ATTEMPTO: A Fault-Tolerant Multiprocessor Working Station, Design and Concepts, Proc. FTCS-13, Milano (1983) 10-13.
- (2) M.Dal Cin, E.Dilger (Eds.), Self-Diagnosis and Fault-Tolerance, ATTEMPTO-Verlag, Tübingen 1981.
- (3) G.Färber, Task-Specific Implementation of Fault-Tolerance in Process Automation, in (2) 84-102.
- (4) S.G.Frison, J.H.Wensley, Interactive Consistency and its Impact on the Design of TMR Systems, Proc. FTCS-12, Santa Monica, (1982) 228-234.
- (5) A.L.Hopkins et al., FTMP: A Highly Reliable Fault-Tolerant Multiprocessor for Aircraft Control, Proc. IEEE, Vol 66/10 (1978), 1221-1239.
- (6) D.Katsuki et al., PLURIBUS- an Operational Fault-Tolerant Multiprocessor, Proc. IEEE, Vol 66/10 (1978) 1146-1159.
- (7) Y.Tohma, The SAFE-Project, Tokyo Institute of Technology, private Mitteilung.
- (8) J.H.Wensley et al., SIFT: Design and Analysis of a Fault-Tolerant Computer for Aircraft Control, Proc. IEEE, Vol.66, Oct.(1978), 1240-1255.
- (9) N.Wirth, Programming in Modula-2, Springer-Verlag (1982).

Anschrift der Autoren:
 Institut für Informationsverarbeitung
 Universität Tübingen
 Köstlinstr.6
 D-74 Tübingen