

Überblick

- Einleitung
 - Lit., Motivation, Geschichte, v. Neumann-Modell, VHDL
- Befehlsschnittstelle
- Mikroarchitektur
- Speicherarchitektur
- Ein-/Ausgabe
- Multiprozessorsysteme, ...

JR - RA - SS2002

Kap. 6

5/1

Kap. 6 Multiprozessorsysteme

- 6.1 Einleitung/Klassifikation**
- 6.2 Verbindungsstrukturen**
- 6.3 Beispiele**
- 6.4 Leistungsbewertung**
- 6.5 Cache-Kohärenz**

Kap. 6.1 Einleitung/Klassifikation

Einsatz der zusätzlich verfügbaren Chipfläche

- Parallelität auf Bitebene: bis etwa 1985
 - **Kombinatorische Addierer** und **Multiplizierer**, etc.
 - → *wachsende Wortbreite auf 64 Bit*
- Parallelität auf Instruktionsebene: 1985 bis heute
 - **Pipelining** der Instruktionverarbeitung
 - **Mehrere Funktionseinheiten** (→ *superskalare Prozessoren*)
- Integration von Caches und Hauptspeicher auf die Chipfläche: 1990 bis heute
 - Zur Verringerung der mittleren Zugriffszeiten
 - → *DEC Alpha 21164: 77% der Fläche für Caches*
- Parallelität auf Prozessorebene

JR - RA - SS2002

Kap. 6

5/4

Performanz von Rechnern läßt sich durch Parallelverarbeitung steigern.

- **Verteilten Systemen**
Diese bestehen aus mehreren Prozessoren ohne gemeinsamen Speicher. Jeder Prozessor hat einen eigenen Speicher (local memory). Die Prozessoren tauschen Daten durch Nachrichten aus.
Nachteil: Evtl. hoher Aufwand für Nachrichtenaustausch.
- **Parallelrechnern**
Diese bestehen ebenfalls aus mehreren Prozessoren und haben einen gemeinsamen Speicher (shared memory). Die Kommunikation erfolgt über den gemeinsamen Speicher.
Nachteil: Evtl. Performanz-Probleme, wenn viele Prozessoren gleichzeitig auf den Speicher zugreifen wollen. Konsistenzprobleme beim Einsatz von lokalen Caches

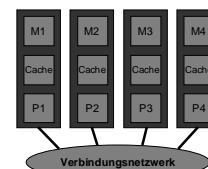
JR - RA - SS2002

Kap. 6

5/5

Architektur von verteilten/parallelen Systemen

- Systeme mit Verbindungsnetzwerk und lokalem Speicher



- Zugriff auf fremden lokalen Speicher nicht möglich
- Kommunikation durch Austausch von Nachrichten
- gute physikalische Skalierbarkeit

JR - RA - SS2002

Kap. 6

5/6

Architektur von verteilten/parallelen Systemen ff

- Symmetrische Multiprozessoren (SMP = shared memory processing)
 - gemeinsamer nichtverteilter Speicher

- meist busbasiert mit physikalisch gemeinsamem Speicher
- Reduktion der Speicherzugriffslatenz durch lokale Caches
- UMA (Uniform-Memory-Access-Modell)
- Problem:** Cache-Kohärenz
- Beispiele:** SGI Challenge, Sun Enterprise

JR - RA - SS2002 Kap. 6 5/7

Architektur von verteilten/parallelen Systemen ff

- Verteilter gemeinsamer Speicher (DSM=distributed shared memory)
 - lokale Speichermodule
 - gemeinsamer Adreßraum
 - NUMA (Non-Uniform-Memory-Access-Modell)

Beispiel: Cray T3D

JR - RA - SS2002 Kap. 6 5/8

Architektur von verteilten/parallelen Systemen ff

- Verteilter gemeinsamer Speicher mit lokalem Cache
 - lokale Speichermodule
 - gemeinsamer Adreßraum
 - Zugriff der Prozessoren über prozessor-eigene Cache
 - cc-NUMA (cache coherent NUMA)

Beispiele: HP Convex SPP, Sequent NUMA-Q

JR - RA - SS2002 Kap. 6 5/9

Architektur von verteilten/parallelen Systemen ff

- Cache-Only-Memory-Access-Modell (COMA) Maschinen
 - es gibt keinen Hauptspeicher, sondern nur Cachespeicher

Beispiele: DDM, KSR

Problem: Cache-Kohärenz, Auffinden von Daten in fremden Caches

JR - RA - SS2002 Kap. 6 5/10

Unterteilung von paralleln/verteilten Systemen

nach dem verwendeten Speichermodell:

- verteilter Speicher, getrennte Adressräume
- verteilter Speicher, gemeinsamer Adressraum
- Nichtverteilter Speicher, gemeinsamer Adressraum

nach der Homogenität der Prozessoren

- homogene Parallelrechner:** alle Prozessoren sind gleich
- heterogene Parallelrechner:** Prozessoren dürfen sich hardwaremäßig unterscheiden

Nach der Hierarchie zwischen den Prozessoren

- symmetrische Parallelrechner:** die Prozessoren sind bzgl. ihrer Rolle im System untereinander austauschbar
- nichtsymmetrische Parallelrechner:** es gibt Masters und Slaves

Nach der Eigenständigkeit der Prozessoren

- lose gekoppelte Parallelrechner:** Netzwerk von eigenständigen Rechnern
- eng gekoppelte Parallelrechner:** physikalisch ein Rechner

JR - RA - SS2002 Kap. 6 5/11

Bsp.: Flynn's Schema '66

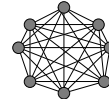
- Einteilung von Rechnern in Klassen:**
 - SISD (Single Instruction stream, Single Data stream)**
 - Personal Computer
 - SIMD (Single Instruction stream, Multiple Data streams)**
 - Vektorrechner z.B. Cray-1 (1976), Cray-2 (1985)
 - ein Befehl wird auf ein Feld (Vektor) von Daten angewendet (Vektorpipeline, komplexe Simulationen → viele FP-Berechnungen)
 - Feldrechner z.B. Illiac IV (1972, 64 PEs), CM-2 (1987, 65.536 PEs)
 - Die gleiche Instruktion wird parallel auf ein Feld von Verarbeitungseinheiten angewendet.
 - MISD (Multiple Instruction streams, Single Data stream)**
 - Diese Klasse ist leer!
 - MIMD (Multiple Instruction streams, Multiple Data streams)**
 - INMOS Transputer, CONVEX SPP, CRAY T3D/T3E, IBM SP2

JR - RA - SS2002 Kap. 6 5/12

Kap. 6.2 Verbindungsstrukturen

Verbindungsstrukturen

- Der **Kommunikationsaufwand** zwischen den Prozessoren ist einer der Hauptpunkte für die Leistung des parallelen/verteilten Systems.
⇒ nicht jedes Problem ist für Parallelisierung geeignet.
- Verschiedene Kommunikationsstrukturen unterscheiden sich hinsichtlich ihrer Kosten und ihrer Leistung.



Beispiel:
Vollständiger
Verbindungsgraph

JR - RA - SS2002

Kap. 6

5/14

Modellierung von Verbindungen

- Die Topologie eines Parallelrechners/Netzwerkes wird durch einen abstrakten Graphen $G=(V,E)$ dargestellt, mit
 - $V = \{ 1, \dots, n \}$ Menge der Knoten, d.h. der Prozessoren bzw. Schaltelemente
 - $E \subseteq \{ \{a,b\}; a,b \in V \}$, die Menge der Kanten, d.h. der Verbindungen

JR - RA - SS2002

Kap. 6

5/15

Charakteristika einer Verbindungsstruktur

- Komplexität oder Kosten:**
Hardware- Aufwand für das Verbindungsnetz gemessen in der Anzahl und der Art der Schaltelemente und Verbindungsleitungen.
- Verbindungsgrad**
eines Knotens ist definiert als die Anzahl der Verbindungen, die von dem Knoten zu anderen Knoten bestehen.
- Diameter oder Durchmesser:**
maximale Distanz für die Kommunikation zweier Prozessoren, also die Anzahl der Verbindungen, die durchlaufen werden müssen. Man spricht auch von der maximalen Pfadlänge zwischen zwei Knoten.

JR - RA - SS2002

Kap. 6

5/16

Charakteristika einer Verbindungsstruktur

- Regelmäßigkeit** des Verbindungsmusters:
Ein regelmäßiges Verbindungsmuster lässt sich meist besser implementieren.
- Notwendige **Leitungslängen**:
Kurze Leitungslängen für alle Verbindungen eines Verbindungsnetzes sind vorteilhaft.
- Blockierung**:
Ein Verbindungsnetz heißt **blockierungsfrei**, falls jede gewünschte Verbindung zwischen den Prozessoren oder zwischen den Prozessoren und Speichern unabhängig von schon bestehenden Verbindungen hergestellt werden kann.

JR - RA - SS2002

Kap. 6

5/17

Charakteristika einer Verbindungsstruktur

- Erweiterbarkeit**:
Multiprozessoren können begrenzt, stufenlos oder nur durch Verdopplung der Anzahl der Prozessoren erweiterbar sein.
- Skalierbarkeit**:
Fähigkeit, die wesentlichen Eigenschaften des Verbindungsnetzes auch bei beliebiger Erhöhung der Knotenzahl beizubehalten.

JR - RA - SS2002

Kap. 6

5/18

Charakteristika einer Verbindungsstruktur

- **Ausfalltoleranz** oder **Redundanz** :
Verbindungen zwischen Knoten sind selbst dann noch zu schalten, wenn einzelne Elemente des Netzes (Schaltelemente, Leitungen) ausfallen. Ein fehlertolerantes Netz muss also zwischen jedem Paar von Knoten mindestens einen zweiten, redundanten Weg bereitstellen. Die Eigenschaft eines Systems, bei Ausfall einzelner Komponenten unter deren Umgehung funktionstüchtig zu bleiben, wenn auch mit verminderter Leistung, wird als *Graceful degradation* bezeichnet.

JR - RA - SS2002 Kap. 6 5/19

Charakteristika einer Verbindungsstruktur

- **Durchsatz** oder **Übertragungsbandsbreite**:
Die maximale Übertragungsleistung des Verbindungsnetzes oder einzelner Verbindungen, meist in Megabits pro Sekunde (MBit/ s).
- **Komplexität der Pfadberechnung** oder **Wegefindung**:
die Art, wie der Weg einer Nachricht vom Sender- zum Zielknoten berechnet wird. Die Wegefindung sollte einfach sein, um mittels eines schnellen Hardware- Algorithmus in jedem Verbindungselement implementierbar zu sein. Zu einer Verbindungsstruktur kann es mehrere Wegefindungsalgorithmen geben. Unterscheidungsmerkmale

JR - RA - SS2002 Kap. 6 5/20

Charakteristika einer Verbindungsstruktur

- Bei **statischen Netzen** existieren fest installierte Verbindungen zwischen Paaren von Netzknoten.
- **Dynamische Netze** enthalten eine Komponente „Schaltnetz“, an die alle Knoten über Ein- und Ausgänge angeschlossen sind. Direkte fest installierte Verbindungen zwischen den Knoten existieren nicht.

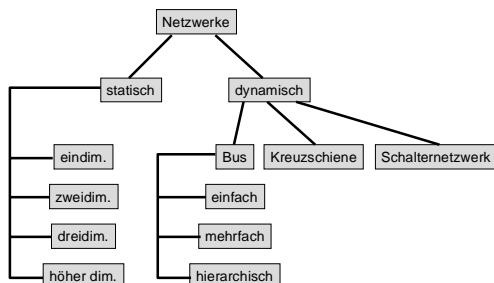
JR - RA - SS2002 Kap. 6 5/21

Charakteristika einer Verbindungsstruktur

- Durchmesser(G)
 $= \max_{v,w \in V} \{ \text{Länge des kürzesten Pfades von } v \text{ nach } w \}$
- $\text{Grad}(G) = \max_{v \in V} | \{ w \in V; \{v,w\} \in E \} |$
- Anzahl der physikalischen Verbindungen(G) = | E |
- minimale Bisektionsbreite:
Teilt man das Netz in zwei (annähernd) gleich große Hälften A und B, so daß die Anzahl der Kanten e zwischen A und B minimal ist, so wird e als minimale Bisektionsbreite bezeichnet
- Diskonnektivität $d = n/e$
n = Anzahl der Prozessorknoten, e = minimale Bisektionsbreite
Schlimster Fall: alle Knoten in A senden eine Nachricht an Knoten in B und umgekehrt
- Kosteneffektivität $K(G) = \frac{\text{Grad}(G)}{\max(\text{Durchmesser}(G), \text{Diskonnektivität}(G))}$

JR - RA - SS2002 Kap. 6 5/22

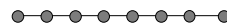
Klassifikation von Verbindungsstrukturen



JR - RA - SS2002 Kap. 6 5/23

Statische Verbindungsstrukturen, 1-dim

■ **Kette**



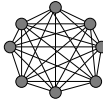
Großer Durchmesser
sehr fehleranfällig

- Durchmesser = n-1
- Grad = 2
- |Verbindungen| = n-1
- e = 1
- d = n
- K = 2n

JR - RA - SS2002 Kap. 6 5/24

Statische Verbindungsstrukturen, 2-dim

■ Vollständiger Verbindungsgraph

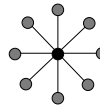


Zu teuer wegen dem großen Fanout

- Durchmesser = 1
- Grad = n-1
- |Verbindungen| = n(n-1)/2, keine „Kollisionen“
- e = (n/2)²
- d = 4/n
- K = (n-1)*max(1,4/n)=n-1 n>4

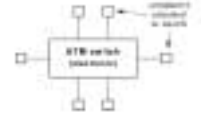
Statische Verbindungsstrukturen, 2-dim

■ Stern



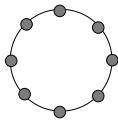
Zentraler Knoten ist Flaschenhals

- Durchmesser = 2
- Grad = n-1
- |Verbindungen| = n-1
- e = n/2
- d = 2
- K = (n-1)*max(2,2)=2*(n-1)



Statische Verbindungsstrukturen, 2-dim

■ Ring



Beispiel für Ring: Token-Ring
 Es kreist ein sogenanntes **Token** (spezielles Paket).
 Ein Rechner darf nur dann senden, wenn er das Token besitzt.

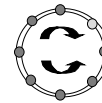
- Durchmesser = n/2
- Grad = 2
- |Verbindungen| = n
- e = 2
- d = n/2
- K = 2*max(n/2,n/2)=n

Beispiel für Ringtopologie : CDDI/FDDI-Ring

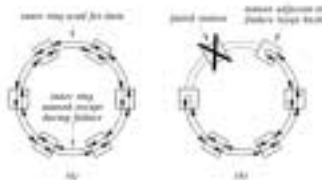
[Copper / Fiber Distributed Data Interconnect]

Charakteristika

- Ring-Topologie
- Besteht aus zwei gegenläufigen Ringen
 → **Fehlertolerantes Netz**
- ... ansonsten wie beim Token-Ring

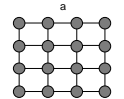


CDDI/FDDI-Ring ist fehlertolerant!



Statische Verbindungsstrukturen, 2-dim

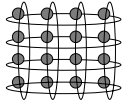
■ 2D Gitter



- Durchmesser = a+b-2
- Grad = 4
- |Verbindungen| = a(b-1)+b(a-1) = 2ab - a - b
- e = min(a,b)
- d = ab/min(a,b)
- K = 4*(a+b-2) a>1, b>1

Statische Verbindungsstrukturen, 2-dim

■ **MESH** (torusähnliches Gitter)



Typischer Vertreter war das Transputer-Netz von INMOS

- Durchmesser = $n^{1/2}$
- Grad = 4
- |Verbindungen| = $2n$
- $e = 2n^{1/2}$
- $d = n^{1/2}/2$
- $K = 4 \cdot n^{1/2}$

Statische Verbindungsstrukturen, 3-dim

■ **Hypercube** (d-dimensionaler Würfel)



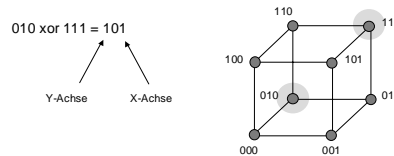
- Durchmesser = $\log n$
- Grad = $\log n$
- |Verbindungen| = $(n \log n) / 2$
- $e = n/2$
- $d = 2$
- $K = (\log n)^2$

e-Cube-Routing

- Die Knotennummern werden als Binärzahlen geschrieben, dadurch unterscheiden sich benachbarte Knoten in genau einer Stelle, die zudem die Richtung der Verbindung angeben kann.
- Eine einfache Wegewahl: die Bits in Start- und Zieladresse werden mittels einer XOR-Verbindung verknüpft und das Resultat bestimmt die möglichen Wege.

e-Cube-Routing

- „higher dimensions of channels until the destination is reached“
Dimension eines Kanals = Bitposition von (Knoten# XOR Knoten#)
- Beispiel: A = (010) und B = (111)



Statische Verbindungsstrukturen, 3-dim

■ **Cube Connected Cycle (CCC)**



- d = Dimension
- r = Anzahl Knoten in Ringen
- häufig: $r=d$
- |Knoten| = $2^{dr} = n$

- Durchmesser = $(r/2) \cdot d$
- Grad = 3
- |Verbindungen| = $(d \cdot d^r) / 2 + 2^{dr}$
- $e = 2^{dr}/2$
- $d = 2r$
- $K = 3 \cdot (r/2) \cdot d$ für $d > 4$, $6r$ für $d < 4$

Dynamische Verbindungsstrukturen

■ **Bus**



Topologie versagt bei den heutigen Technologien bei grossem Datentransfer zwischen den Prozessoren

- Ein Bus lässt sich als Stern modellieren, wobei der zentrale Knoten aber kein Prozessor ist, sondern der zentrale Bus.
- Gleiches „Flaschenhalsproblem“ wie bei Stern.

Beispiel für Bustopologie: Ethernet

Charakteristika

- Bus-Topologie
- 10 - 100 Mbit / Sekunde
- Paket-Versand mit Paketgrößen von 64-1518 Bytes
- Nicht abhörsicher: *alle hören mit!*
- (Lokales) Rechnernetz über eine **Bridge** mit der Aussenwelt verbunden, die die Nachrichten filtert und verstärkt

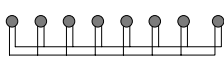
Übertragungsvorgang

- Nachrichten werden in Pakete fester Länge zerteilt. Jedes Paket enthält **Headerinformation** mit *Zielfresse* und *Sequenznummer*
- Jeder Rechner horcht am Bus und empfängt die Pakete, die seine Adresse tragen
- Kollisionen von mehreren Sendern werden erkannt.
- Falls Kollision, dann später erneuter Sendeversuch

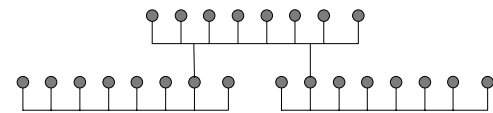
JR - RA - SS2002 Kap. 6 5/37

Alternative Busstrukturen

■ **Mehrfachbus**



■ **Hierarchischer Bus**

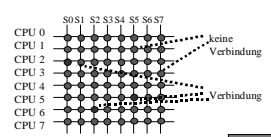


JR - RA - SS2002 Kap. 6 5/38

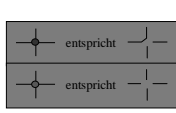
Dynamische Verbindungsstrukturen ff

■ **Crossbar Switch**

meist verwendete Struktur bei Parallelrechnern mit gemeinsamen "nichtverteilten" Speicher




Nachteil: n * m Crosspoints



JR - RA - SS2002 Kap. 6 5/39

Permutationsnetzwerke

- Vermeidung des hohen Aufwandes von Kreuzschienen
- Aufgebaut aus Zweierschaltern (2 Zustände = 1 Bit):

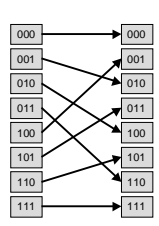


- Permutation: p Eingänge werden parallel auf p Ausgänge geschaltet (Eingänge werden permutiert)
- Einstufige, mehrstufige oder rückgekoppelte Netze
- Reguläre/irreguläre Permutationsnetzwerke

JR - RA - SS2002 Kap. 6 5/40

Permutationsnetzwerke

■ **Mischpermutation (perfect shuffle)**



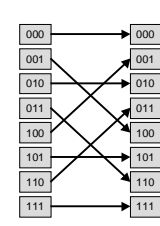
Perfect shuffle:
 $PS(n) = 2 * n$ for $n < N/2$
 $PS(n) = 2 * n - N + 1$ for $n \geq N/2$

$$M(a_n, \dots, a_1) = a_{n-1}, \dots, a_1, a_n$$

JR - RA - SS2002 Kap. 6 5/41

Permutationsnetzwerke

■ **Kreuzpermutation**

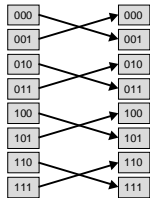


$$K(a_n, \dots, a_1) = a_1, a_{n-1}, \dots, a_2, a_n$$

JR - RA - SS2002 Kap. 6 5/42

Permutationsnetzwerke

■ Tauschpermutation



$$T(a_n, \dots, a_1) = a_n, a_{n-1}, \dots, a_2, \bar{a}_1$$

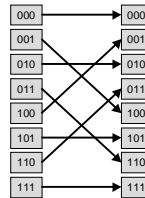
JR - RA - SS2002

Kap. 6

5/43

Permutationsnetzwerke

■ Umkehrpermutation (Butterfly)



$$U(a_n, \dots, a_1) = a_1, a_2, \dots, a_{n-1}, a_n$$

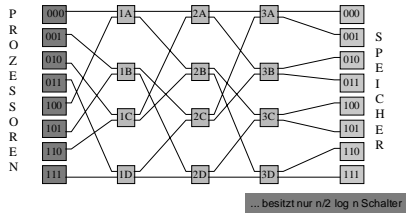
JR - RA - SS2002

Kap. 6

5/44

Dynamische Verbindungsstrukturen ff

■ Omega Netzwerk (log n Stufen von Mischpermutationen)



...besitzt nur n/2 log n Schalter

JR - RA - SS2002

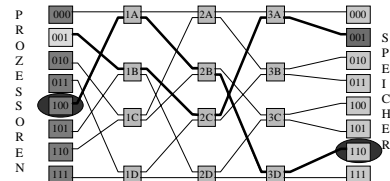
Kap. 6

5/45

Dynamische Verbindungsstrukturen ff

■ Omega Netzwerk

- der obere Ausgang eines Schalters ist der 0-Ausgang
- der untere Ausgang eines Schalters ist der 1-Ausgang
- ein Schalter der Stufe i schaltet gemäß dem i-ten Bit der Zieladresse



JR - RA - SS2002

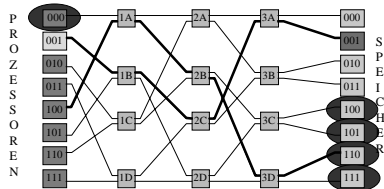
Kap. 6

5/46

Dynamische Verbindungsstrukturen ff

■ Omega Netzwerk

- nicht jede Kommunikation ist gleichzeitig möglich, auch wenn alle Zieladressen paarweise verschieden sind (**blocking network**)

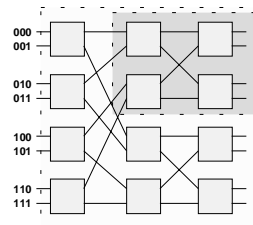


JR - RA - SS2002

Kap. 6

5/47

Delta-Netzwerk



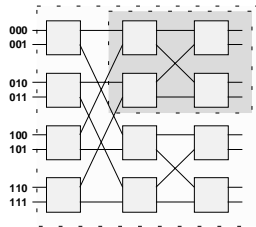
- Log n Stufen
- Mischpermutation zwischen den Stufen
- Zwischen den Stufen i und i+1 wird die Permutation nur auf $2^{(\log n - i + 1)}$ Ein-/Ausgänge angewandt
- Eingänge werden ohne Permutation mit der ersten Stufe verbunden

JR - RA - SS2002

Kap. 6

5/48

Butterfly-Netzwerk



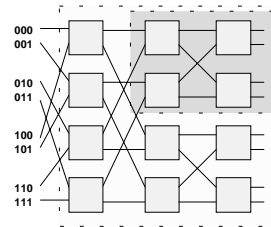
- Log n Stufen
- Umkehrpermutation zwischen den Stufen
- Zwischen den Stufen i und $i+1$ wird die Permutation nur auf $2^{(\log n - i + 1)}$ Ein-/Ausgänge angewandt
- Eingänge werden ohne Permutation mit der ersten Stufe verbunden

JR - RA - SS2002

Kap. 6

5/49

Banyan-Netzwerk



- log n Stufen
- Umkehrpermutation zwischen den Stufen
- Zwischen den Stufen i und $i+1$ wird die Permutation nur auf $2^{(\log n - i + 1)}$ Ein-/Ausgänge angewandt
- Eingänge werden per Mischpermutation mit der ersten Stufe verbunden

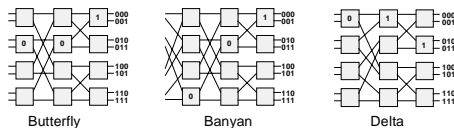
JR - RA - SS2002

Kap. 6

5/50

Self-Routing Property

- Zieladresse liefert Weeginformation
- jede Stufe betrachtet das korrespondierende Bit der Zieladresse
 - '0' bedeutet, Nachricht an den oberen Ausgang
 - '1' bedeutet, Nachricht an den unteren Ausgang



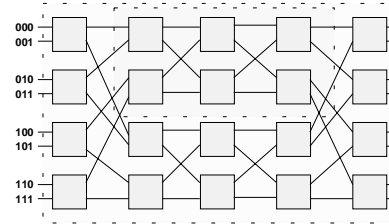
JR - RA - SS2002

Kap. 6

5/51

Beneš Netzwerk

- $2 \log n - 1$ Stufen
- Zusammengesetzt aus zwei gespiegelte Butterfly-Netzwerken



JR - RA - SS2002

Kap. 6

5/52

Eigenschaften des Beneš Netzwerkes

- Vorteil
 - Jede Permutation kann konfliktfrei geschaltet werden
 - Achtung: zwei Butterfly-Netzwerke in Serie haben nicht diese Eigenschaft!
- Nachteil:
 - die Pfadbestimmung ist sehr komplex und muß deshalb off-line durchgeführt werden

JR - RA - SS2002

Kap. 6

5/53

Vermittlungsarten

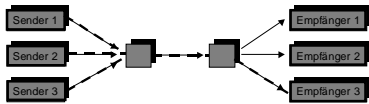
- **Paket-vermittelnde Vermittlung**
 - Nachricht wird in ein oder mehrere Pakete verpackt
 - Jedes der Pakete enthält die Adresse des Empfängers
 - Es wird kein Pfad vom Sender zum Empfänger freigeschaltet
 - Das Paket wird in Abhängigkeit der Empfängeradresse immer nur zu einem direkten Nachbarn geschickt.
- **Leitungsvermittelnde Vermittlung**
 - Es wird im Netzwerk ein Pfad vom Sender zum Empfänger geschaltet, über den alle Nachrichten geschickt werden (Bsp: Telefonverbindung)

JR - RA - SS2002

Kap. 6

5/54

Leitungsvermittelnde Vermittlung



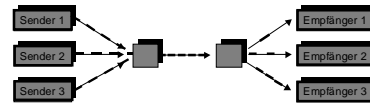
- Schnelle Übertragung grosser Datenmengen
- Geschalteter Pfad blockiert andere Verbindungen

JR - RA - SS2002

Kap. 6

5/55

Paket-vermittelnde Vermittlung



- Keine Verbindung muss lange warten
- Unterbrechungen während einer Übertragung möglich

JR - RA - SS2002

Kap. 6

5/56

Adressierungsarten

- **zielbasiert** (*destination-based routing*):
Kopfteil eines Pakets (oder einer Nachricht) wird mit einer systemweit eindeutigen Empfängeradresse versehen, die bei der Wegefindung von jedem Zielknoten zur Auswahl eines Übertragungskanals genutzt wird
- **quellenbasiert** (*source-based routing*):
Paket wird mit allen Informationen versehen, um über die Zwischenknoten zum Empfänger zu gelangen. Für jeden Zwischenknoten wird im voraus die Abzweigung bestimmt, die das Paket nehmen muss.

JR - RA - SS2002

Kap. 6

5/57

Wegewahl

- **Deterministische Wegewahl:**
Nachrichten zwischen zwei Knoten müssen immer den gleichen Weg gehen
 - Vorteil: einfachen Pfadberechnung
 - Nachteile: erhöhten Möglichkeit zu Blockierungen, Mangel an Fehlertoleranz.
- **Adaptive Wegewahl:**
Möglichkeit auch andere Wege zu nehmen,
 - Nachteil: Etwas höheren Hardware-Aufwand.
 - Vorteil: Blockierte Strecken und ausgefallene Knoten oder Schaltelemente können umgangen werden.

JR - RA - SS2002

Kap. 6

5/58

Kommunikationskontrolle

- deterministische Kontrolle
der Weg eines jeden Pakets ist reproduzierbar
- randomisierte Kontrolle
an bestimmten Stellen des Algorithmus werden zufällige Entscheidungen getroffen
 - Beispiel: Valiant-Paradigma
 - Route zu einer zufälligen Zwischenadresse
 - Route dann erst zum Ziel

JR - RA - SS2002

Kap. 6

5/59

Phit und Flusskontrolle

- Eine Nachricht selbst wird in eine Anzahl von **Übertragungseinheiten** (*phits – physical transfer units* – oder auch *flits – flow control digits* – genannt) zerlegt.
- Ein *Phit* ist dabei die Datenportion, die zu einem Zeitpunkt zwischen zwei Knoten übertragen werden kann.
- Bei der Nachrichtenübertragung zwischen nicht benachbarten Sender- und Empfängerknoten sind Puffer nötig.

JR - RA - SS2002

Kap. 6

5/60

Phit und Flusskontrolle

- Bei der Zwischenspeicherung muss der Knoten dafür sorgen, dass sein Puffer nicht überläuft.
- Dies geschieht durch eine zusätzliche **Fluss-Steuerung (flow control)** z.B. über
 - spezielle Leitungen (z. B. *ack* - oder *strobe*-Leitung) vom empfangenden (Zwischen-)Knoten zum sendenden (Zwischen-)Knoten

JR - RA - SS2002

Kap. 6

5/61

Übertragungsmodi

- **Store-and-forward-Modus:**
Nachricht wird von jedem Zwischenknoten in Empfang genommen, vollständig zwischengespeichert und dann weiter übertragen

JR - RA - SS2002

Kap. 6

5/62

Übertragungsmodi

- **Virtual-cut-through-Modus:**
 - Nachricht wird als Kette von Phits transportiert.
 - Der Kopfteil der Nachricht enthält die Empfängeradresse und bestimmt den einzuschlagenden Weg.
 - Bei Ankunft des ersten Phits an einem Zwischenknoten wird diese sofort an den nächsten Knoten weitergeleitet, falls der Pfad dorthin frei ist.
 - Alle anderen folgen ähnlich wie bei einer Pipeline- Verarbeitung.
 - ankommende Daten werden *nur im Konfliktfall* im Knoten vollständig zwischengespeichert
 - In jedem Knoten werden Puffer bereit gehalten, die auch ein maximal großes Nachrichtenpaket zwischenspeichern können

JR - RA - SS2002

Kap. 6

5/63

Übertragungsmodi

- **Wormhole-routing-Modus:**
 - solange keine Übertragungskanäle blockiert sind, mit den Virtual-cut-through-Modus identisch.
 - Falls der Kopfteil der Nachricht auf einen Kanal trifft, der gerade belegt ist, wird er abgeblockt.
 - Alle nachfolgenden Übertragungseinheiten der Nachricht verharren dann ebenfalls an ihrer augenblicklichen Position, bis die Blockierung aufgehoben ist.
 - Durch das Verharren werden die Puffer nachfolgender Kanäle auch für weitere Nachrichten blockiert.

JR - RA - SS2002

Kap. 6

5/64

Übertragungsmodi

- **Buffered wormhole routing:**
 - Kompromisslösung zwischen Virtual-cut-through- und Wormhole-routing-Modus
 - begrenzter Puffer zur Aufnahme kleinerer Pakete vorhanden
 - größere Pakete werden im Blockierungsfall ähnlich dem Wormhole-routing-Modus in den Puffern mehrerer Knoten zwischengespeichert.

JR - RA - SS2002

Kap. 6

5/65

Kap. 6.3 Multiprozessorsysteme: Beispiele

Intel Pentium Pro Quad

- Bis zu 4 Prozessoren
- 66 MHz Busverbindung mit bis zu 528 MB/sek
- Kohärenz- und Multiprozessorlogik im Prozessor integriert

JR - RA - SS2002 Kap. 6 5/67

SUN Enterprise

- Bis zu 16 Einsteckkarten: (bis zu 4) Prozessor + Speicher oder I/O-Karte
- 100 MHz Bus mit bis zu 2GB/sek

JR - RA - SS2002 Kap. 6 5/68

Cray T3E

- Bis zu 2048 Prozessoren, 480MB/s Verbindungen
- DSM: Memory controller generiert Kommunikationsanforderungen für nichtlokale Speicherzugriffe
- Spitzen-Bisektionsbandbreite: 166GB/sek (512 Prozessoren)

JR - RA - SS2002 Kap. 6 5/69

IBM SP-2

- 2 bis 512 Knoten
- Jeder Knoten ist eine vollständige RS6000 Workstation
- Verbindungsnetz:
 - 4x4 Crossbars
 - 8 Crossbars = 16x16 Verbindungsnetz

JR - RA - SS2002 Kap. 6 5/70

Intel Paragon

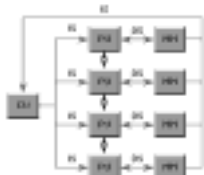
- 2D Gitternetzwerk
- Wormhole routing
- größtes System mit 1984 Knoten

JR - RA - SS2002 Kap. 6 5/71

Connection Machine: CM5

- 32 - 16384 Prozessoren
- message-passing, distributed memory
- Control-, Data- und Diagnostic-Network
- MIMD, SIMD möglich durch Control-Network

JR - RA - SS2002 Kap. 6 5/72

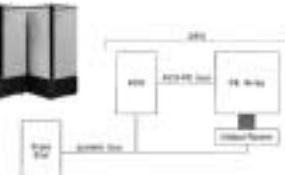
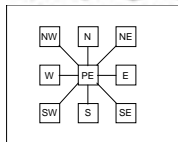
ILLIAC 4 (SIMD)

- 1 "control unit" (CU) mit Zugriff auf den Speicher steuert 64 Prozessoreinheiten (PUs)
- jede PE hat lokalen Speicher
- jede PE teilt Speicher mit benachbarten PUs
- 2D Gitterstruktur

JR - RA - SS2002

Kap. 6

5/73

MasPar MP2 (SIMD)

- Bis zu 16.384 Prozessorelemente
- Jeder Knoten ist mit seinen 8 Nachbarn verbunden

JR - RA - SS2002

Kap. 6

5/74

Kap. 6.4 Leistungsbewertung

Leistungsbewertung

- Leistung eines Rechners läßt sich unterschiedlich definieren.
 - ┆ Antwortzeit
 - ┆ Zeit für ein Programm.
 - ┆ Durchsatz
 - ┆ Zahl der gleichzeitig verarbeiteten Programme.
- Sehr einfaches Leistungsmaß sind MIPS.
- Andere Form der Leistungsmessung ist der Leistungsvergleich
 - ┆ Betrachtung eines konkreten Programms
 - ┆ Benchmarks

JR - RA - SS2002

Kap. 6

5/76

Leistungsbewertung durch MIPS

- Vorteil
 - ┆ sehr einfaches und leicht verständliches Maß.
- Nachteil
 - ┆ Schnellere Rechner können niedrigere MIPS-Zahl haben.
 - ┆ Hängt vom Befehlssatz des Rechners ab.
 - ┆ Speicher wird nicht berücksichtigt.
 - ┆ Betriebssystem wird nicht berücksichtigt.

JR - RA - SS2002

Kap. 6

5/77

Leistungsbewertung durch Benchmarks

- Es werden reale Programme auf den Rechnern ausgeführt.
- Dadurch werden Speicher, I/O, Betriebssystem, Compiler, etc. berücksichtigt.
- Um aussagekräftige Vergleiche erhalten zu können, müssen typische Anwendungen verglichen werden.
- Solche typische Anwendungen werden in Benchmarks zusammengefaßt:
 - ┆ SPEC-Benchmarks (www.specbench.org)
 - ┆ TPC-Benchmarks
 - ┆ DIN-Benchmarks
 - ┆ usw.

JR - RA - SS2002

Kap. 6

5/78

Vergleich einiger Architekturen

	DEC Alpha21262	MIPS R10000	IBM PPC150	SUN UltraSparc	PII Klamath	PII Yocco	AMD K6	AMD K6-3D
Taktfrequenz	667	250	400	333	300	400	300	350
SPECint95	44	14.7	17.6	14.2	11.9	16.5	?	?
SPECfp95	66	24.5	12.2	16.9	8.6	13.7	?	?
Transistoren (Mio)	15.2	6.8	6.35	5.4	7.5	7.5	8.8	9.3
Leistungsverbrauch (W)	72	>30	5.7	<30	?	23.3	?	?

JR - RA - SS2002

Kap. 6

Quelle: Rosenstiel, 1999

5/79

Vergleich einiger Architekturen

Hersteller Prozessor	Intel P4	AMD Athlon	Sun UltraSparc III
Taktfrequenz in GHz	2,53	1,733	0,9
SPECint2000	922	749	533
SPECfp2000	901	660	713

Quelle: www.specbench.org

JR - RA - SS2002

Kap. 6

5/80

Leistungsbewertung von Parallelen Programmen

- MIPS: gibt die Leistung des Rechners an
- Benchmarks: Leistung einer Programmsuit auf einem Rechner/Speicher/Betriebssystem
- Wie mißt man den Gewinn durch den Einsatz von Parallelrechnern?

JR - RA - SS2002

Kap. 6

5/81

Speed-Up

- Definitionen: Gegeben ein Programm
 - $P(1)$ = Zahl der auszuführenden Operationen auf Monoprozessorssystem
 - $P(n)$ = Zahl der auszuführenden Operationen auf Multiprozessorssystem
 - $T(1)$ = Ausführungszeit auf Monoprozessorssystem in Schritten
 - $T(n)$ = Ausführungszeit auf Multiprozessorssystem in Schritten
- Speedup nach Lee (Leistungssteigerung):

$$S(n) = \frac{T(1)}{T(n)}$$

nach Lee gilt: $P(1) = T(1)$

JR - RA - SS2002

Kap. 6

5/82

Was ist eine gute Beschleunigung?

- Hoffentlich: $S(n) > 1$
- Linear speedup:
 - $S(n) = n$
 - Programm skaliert perfekt mit der Zahl der Prozessoren
- Superlinear speedup:
 - $S(n) > n$
 - Ist das möglich?
- Realität (Lee'sches Prinzip): $1 \leq S(n) \leq n$

JR - RA - SS2002

Kap. 6

5/83

Alternative Definitionen

- Schnellster bekannter Algorithmus der auf einer Einprozessormaschine ausgeführt wird

$$S(n) = \frac{T(1)}{T(n)}$$

- Schnellster bekannter Algorithmus der auf einem Prozessor der Multiprozessormaschine ausgeführt wird (Gustafson)

$$S(n) = \frac{T(n)}{T(1)}$$

JR - RA - SS2002

Kap. 6

5/84

Kann Speedup superlinear sein?

- Antwort: Ja und Nein, es hängt von der Definition ab
- 1. Antwort: NEIN, Speedup kann nicht superlinear sein!
 - Sei M eine parallele Maschine mit n Prozessoren
 - Sei $T(x)$ die Zeit um das Problem mit x Prozessoren auf M zu lösen
 - Speedup definition: $S(n) = \frac{T(1)}{T(n)}$
 - Nimm an, daß ein paralleler Algorithmus A eine Instanz des Problems in t Zeitschritten löst
 - Dann kann A das gleiche Problem in nt Zeitschritten durch Zeitmultiplexing lösen
 - Die beste serielle Zeit kann nicht größer sein als nt
 - Also kann der Speedup nicht größer sein als n

$$S(n) = \frac{T(1)}{T(n)} \leq \frac{nt}{t} = n$$

JR - RA - SS2002

Kap. 6

5/85

Kann Speedup superlinear sein?

- Antwort 2: Speedup kann superlinear sein!
 - Sei M eine parallele Maschine mit n Prozessoren
 - Sei $T(x)$ die Zeit um das Problem mit x Prozessoren auf M zu lösen
 - Speedup definition: $S(n) = \frac{T_s}{T(n)}$
 - Die serielle Version kann mehr Overhead erzeugen als die parallele Version
 - Z.B. $A=B+C$ auf SIMD Maschine mit A,B,C Matrizen, versus Schleifenoverhead einer seriellen Maschine
 - unterschiedliche Algorithmen, z.B. Tiefensuche bei serieller Maschine, Breitensuche bei Parallelermaschine
 - Es werden unterschiedliche Algorithmen verglichen

JR - RA - SS2002

Kap. 6

5/86

Grenzen des Speedups

- Software Overhead
z.B. zusätzliche Indexberechnungen um den Code auf mehrere Prozessoren zu verteilen
- Load Balancing
Kann die Ausgabe gleichmäßig auf mehrere Prozessoren verteilt werden?
- Communication Overhead
Wird Kommunikation vom Prozessor ausgeführt?
Wie groß sind die Latenzzeiten?

JR - RA - SS2002

Kap. 6

5/87

Grenzen des Speedup

- f = Anteil des Programms, der nicht parallelisiert werden kann
 - z.B. Filezugriffe auf eine Platte
- Amdahls Gesetz:

$$T(n) = fT(1) + (1-f) \frac{T(1)}{n}$$

$$S(n) = \frac{T(1)}{fT(1) + \frac{(1-f)T(1)}{n}} = \frac{n}{nf + 1 - f} = \frac{1}{\frac{(n-1)f + 1}{n}}$$

$$\lim_{n \rightarrow \infty} = \frac{1}{f}$$

JR - RA - SS2002

Kap. 6

5/88

Beispielanwendung von Amdahls Gesetz

- Algorithmus hat 4% seriellen Code, wie groß ist der maximal erreichbare Speedup mit 16 Prozessoren?
 - $16 / (1 + (16 - 1) * 0.04) = 10$
- Wie groß ist der maximale Speedup?
 - $1 / 0.04 = 25$

JR - RA - SS2002

Kap. 6

5/89

Andere Leistungskennzahlen

- Effizienz = relative Verbesserung in der Verarbeitungsgeschwindigkeit

$$E(n) = \frac{S(n)}{n} \quad \frac{1}{n} \leq E(n) \leq 1$$
- Redundanz = Mehraufwand (im Code) durch Parallelverarbeitung

$$R(n) = \frac{P(n)}{P(1)} \quad 1 \leq R(n)$$
- Parallelexponent = Anzahl der parallelen Operationen pro Zeiteinheit

$$I(n) = \frac{P(n)}{T(n)}$$

JR - RA - SS2002

Kap. 6

5/90

Andere Leistungskennzahlen

- Auslastung = Operationen pro Zeiteinheit in jedem Prozessor
- Qualität = qualitativ erzielte Leistungssteigerung
- Es gilt

$$U(n) = \frac{l(n)}{n}$$

$$Q(n) = S(n) \frac{E(n)}{R(n)}$$

$$1 \leq Q(n) \leq S(n) \leq l(n) \leq n$$

$$\frac{1}{n} \leq E(n) \leq U(n) \leq 1$$

Kap. 6.5 Cache-Kohärenzprotokolle

Prinzipien Standard-Protokolle MESI-Protokoll

Multiprozessor-Cache-Kohärenz

- Zur Sicherstellung der Kohärenz von Cache und Hauptspeicher gibt es sogenannte cache coherency protocols
- zwei Prinzipien:
 - directory based: Information über einen Block befindet sich nur an einer Stelle
 - Nachteil: Kohärenzinformation im Prinzip proportional zur Anzahl Blöcke im Hauptspeicher
 - snooping: Information über Blöcke befindet sich beim jeweiligen Block im jeweiligen Cache. Jeder Cache-Controller betreibt "Snooping" am Bus, um zu überwachen, was mit Block passiert
 - erfordert gemeinsamen Bus
 - üblich und wird im folgenden diskutiert

Snooping-Protokolle

- zwei Arten, die sich bei Write unterscheiden:
 - write invalidate: wenn ein Prozessor einen Block verändern will, macht er vorher alle anderen Kopien dieses Blocks ungültig
 - write broadcast: nach dem Ändern wird der geänderte Block an alle anderen Prozessoren geschickt, die ihre Kopie dann ggf. aktualisieren können
- beide Protokolle haben ihre Berechtigung und werden auch verwendet

„Write-Once“ Cache Protokoll

Prinzip:
Cache-Block wird beim Ersetzen in den Speicher zurückgeschrieben, falls er, nachdem er in den Cache geschrieben wurde, verändert wurde (=Zustand "Dirty", s.u.)

- 4 Zustände
- Invalid:** inkonsistente Cache-Kopie
- Valid:** Cache-Kopie mit Speicher konsistent
- Reserved:** Cache-Kopie mit Speicher konsistent und einzige Kopie und Daten zum ersten Mal geschrieben
- Dirty:** einzige Kopie und Daten mehr als einmal verändert



...Erläuterung zum "write once" (Standard-) Protokoll

- Übergänge durch
 - Lesen und Schreiben durch jeweiligen Prozessor (P-Read, P-Write)
 - (normale) Speicher-Lese/Schreib-Operationen (Read-Blk, Write-Blk)
 - Konsistenz-Operationen, die nicht durch den Prozessor, sondern durch den gemeinsamen Bus ausgelöst werden (gestrichelte Übergänge):
 - Write-Invalidate (Write-Inv): macht alle anderen Cache-Kopien eines Blocks ungültig (invalid)
 - Read-Invalidate (Read-Inv): liest einen Block und macht alle anderen Cache-Kopien ungültig (invalid)

...Erläuterung zum "write once" (Standard-) Protokoll

- **Treffer beim Lesen:**
 - l kann lokal im Cache abgewickelt werden
 - l keine Zustandsübergänge erforderlich
- **Miss beim Lesen:**
 - l keine Kopie im Dirty-Zustand:
 - l Block wird aus Speicher gelesen
 - l Zustand = Valid
 - l Kopie im Dirty-Zustand:
 - l Cache mit Dirty-Kopie verhindert, daß Speicher Kopie schickt, und schickt selbst seine Kopie zum anfordernden Cache
 - l Speicher wird aktualisiert, Zustände beider Kopien = Valid

JR - RA - SS2002

Kap. 6

...Erläuterung zum "write once" (Standard-) Protokoll

- **Treffer beim Schreiben:**
 - l **Kopie im Dirty- oder Reserved-Zustand:**
 - l Schreiben kann lokal erfolgen
 - l neuer Zustand = Dirty
 - l **Kopie im Valid-Zustand:**
 - l Write-Invalidate wird ausgelöst
 - l Speicher-Kopie wird aktualisiert
 - l neuer Zustand = Reserved

JR - RA - SS2002

Kap. 6

...Erläuterung zum "write once" (Standard-) Protokoll

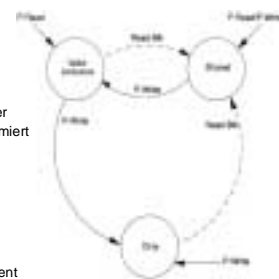
- **Miss beim Schreiben:**
 - l Kopie kommt aus Speicher oder von anderem Cache mit Kopie im Dirty-Zustand
 - l wird durch Read-Invalidate ausgelöst
 - l Speicher wird aktualisiert
 - l neuer Zustand = Dirty
- **Ersetzen eines Blocks:**
 - l Zurückschreiben in den Speicher falls im Dirty-Zustand
 - l sonst keine Aktion erforderlich

JR - RA - SS2002

Kap. 6

Firefly

Prinzip:
Firefly-Protokoll benutzt "copy-back" für private Blöcke und "write-through" für gemeinsame Blöcke, wobei über privat und gemeinsam erst zur Laufzeit entschieden wird zur Implementierung wird zusätzliche Busleitung "shared" benutzt, die während des Snooping-Mechanismus Schreiber über tatsächliche Existenz weiterer Kopien informiert



3 Zustände
Valid-exclusive: einzige Cache-Kopie mit Speicher konsistent
Shared: Cache-Kopie mit Speicher konsistent und es existieren weitere konsistente Kopien
Dirty: einzige Kopie mit Speicher inkonsistent

JR - RA - SS2002

Kap. 6

...Erläuterung zum "Firefly" Write-Update-Protokoll

- **Übergänge durch**
 - l Lesen und Schreiben durch jeweiligen Prozessor (P-Read, P-Write)
 - l (normale) Speicher-Lese/Schreib-Operationen (Read-Blk, Write-Blk)
 - l Konsistenz-Operation Write-Update:
 - l alle anderen Cache-Kopien einschließlich des Speichers werden aktualisiert
 - l falls Shared-Zustand bei anderen Kopien nicht mehr vorliegt, ist Shared-Busleitung nicht aktiv
 - wird z.B. ausgenutzt, um beim Schreiber Shared-Zustand durch Valid-exclusive-Zustand zu ersetzen

JR - RA - SS2002

Kap. 6

...Erläuterung zum "Firefly" Write-Update-Protokoll

- **Treffer beim Lesen:**
 - l kann lokal im Cache abgewickelt werden
 - l keine Zustandsübergänge erforderlich
- **Miss beim Lesen:**
 - l keine Cache-Kopie existiert:
 - l Block wird aus Speicher gelesen
 - l Zustand = Valid-exclusive
 - l eine Cache-Kopie im Dirty-Zustand:
 - l Cache mit Dirty-Kopie verhindert, daß Speicher Kopie schickt, und schickt selbst seine Kopie zum anfordernden Cache
 - l Speicher wird aktualisiert, Zustände beider Kopien = Shared

JR - RA - SS2002

Kap. 6

...Erläuterung zum "Firefly" Write-Update-Protokoll

- wenn mehrere Kopien im Shared-Zustand existieren:
 - ┆ diese Caches synchronisieren Übertragung und schicken die Kopie direkt zum anfordernden Cache
 - ┆ Zustand = Shared

JR - RA - SS2002

Kap. 6

...Erläuterung zum "Firefly" Write-Update-Protokoll

- Treffer beim Schreiben:
 - ┆ Cache-Kopie im Dirty- oder Valid-exclusive-Zustand:
 - ┆ Schreiben kann lokal erfolgen
 - ┆ neuer Zustand = Dirty
 - ┆ Kopie im Shared-Zustand:
 - ┆ Write-Update aktualisiert alle Kopien einschließlich Speicher
 - ┆ Zustand = Shared
 - ┆ Ausnahme: während Write-Update ist Shared-Busleitung nicht aktiv, der Shared-Zustand ist also bei den anderen Kopien zwischenzeitlich beendet worden
 - dann gilt: Zustand = Valid-exclusive

JR - RA - SS2002

Kap. 6

...Erläuterung zum "Firefly" Write-Update-Protokoll

- Miss beim Schreiben:
 - ┆ Kopie kommt aus Speicher oder von anderem Cache
 - ┆ falls Kopie vom Speicher kommt: neuer Zustand = Dirty
 - ┆ falls Kopie von anderem Cache kommt: Write-Update aktualisiert alle Kopien einschließlich Speicher: neuer Zustand = Shared
- Ersetzen eines Blocks:
 - ┆ Zurückschreiben in den Speicher falls im Dirty-Zustand
 - ┆ sonst keine Aktion erforderlich

JR - RA - SS2002

Kap. 6

"Dragon"-Protokoll

- Dragon-Protokoll für Dragon-Multiprozessor Workstation von Xerox vorgeschlagen
- verbessert Effizienz der Cache-Cache-Transfers dadurch, daß auf Speicher-Aktualisierung verzichtet wird
- Speicher wird erst beim Ersetzen eines Blocks aktualisiert

JR - RA - SS2002

Kap. 6

Pentium-MESI-Protokoll

- ist i.w. modifiziertes Write- Once-Protokoll
- Modified: einzige Cache-Kopie, nicht konsistent mit Hauptspeicher, lesen und schreiben lokal
- Exclusive: einzige Cache-Kopie, nicht modifiziert und daher konsistent mit Hauptspeicher, lesen und schreiben lokal, schreiben führt zu Übergang in Modified-Zustand
- Shared: Kopie in anderen Caches, lesen lokal, schreiben führt zu write-through und invalidate der anderen Kopien
- Invalid: Cache-Block nicht verfügbar, Lesen/Schreiben führt u.U. zum Nachladen und zwar in Abänderung des Write-Once Prot.:

JR - RA - SS2002

Kap. 6

...Pentium-MESI-Protokoll

- Read Miss: Mögliche Folgezustände:
 - ┆ S: falls Block Read-Only
 - ┆ S: falls CPU-Write-Through Attribut aktiv
 - ┆ S: falls Kopie in anderem Cache existiert
 - ┆ M: falls modifizierte Daten direkt von anderem Cache kommen, Speicher wird nicht aktualisiert
 - ┆ E: sonst
- Write Miss: 2 Möglichkeiten direktes Schreiben in Speicher mit evtl. Invalidate anderer Kopie
 - ┆ ohne Nachladen des Cache
 - ┆ mit (folgendem) Nachladen des Cache: WriteMiss with Allocation

JR - RA - SS2002

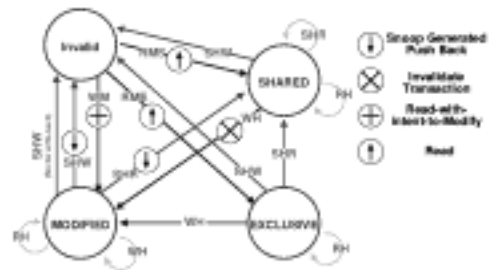
Kap. 6

...Pentium-MESI-Protokoll

- Ergänzung um explizite Befehle für Invalide-Zustand:
 - INVD: setzt interne und externe Caches in Invalide-Zustand
 - ┆ i.w. für Testzwecke oder ähnliches geeignet
 - ┆ modifizierte Cache-Blöcke werden nicht in Speicher zurückgeschrieben
 - WBINVD: schreibt modifizierte Cache-Blöcke zurück bzw. veranlaßt Zurückschreiben und setzt Blöcke in Invalide-Zustand

JR - RA - SS2002 Kap. 6

Zustandsübergänge beim MESI-Protokoll



JR - RA - SS2002 Kap. 6

Zustände des (Pentium-)MESI-Protokolls

Zustand	M modified	E exclusive	S shared	I Invalid
Cache-Block	gültig	gültig	gültig	ungültig
Speicher-Block	ungültig	gültig	gültig	-
weitere Cache-Kopien	nein	nein	möglich	möglich
Treffer beim Lesen	lokal	lokal	lokal	abh. von Cache Contr.
Treffer beim Schreiben	lokal	lokal	Update-Zyklus	geht zum Speicher

JR - RA - SS2002 Kap. 6

Cache-Aktualisierung beim Pentium

- Pentium ist bzgl. Aktualisierungsstrategie programmierbar
- Cache mit und ohne Hauptspeicherkonsistenz abschaltbar
- im Betrieb mit Cache Write-Back oder WriteThrough wählbar:
 - ┆ Standard-Protokoll ist Write-Back, um gerade bei Multiprocessing Anzahl der Speicherzugriffe zu reduzieren
 - ┆ Write-Through kann manchmal erforderlich sein, z.B. für Frame Buffer, damit immer aktuelles Bild auf Schirm
 - ┆ seitenweise durch Software (PWT-Bit in Seitentabelle) oder Hardware (WB/WT#) kontrollierbar
- Besonderheit: I-Cache nicht beschreibbar

JR - RA - SS2002 Kap. 6

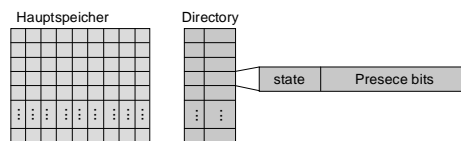
Directory-basierte Protokolle

- Snoopy-Protokolle
 - ┆ benötigen ein Broadcast Kommunikationsmedium (Bus)
 - ┆ beschränkte Zahl von Komponenten (bis zu 30)
- Directory-based Protokolle
 - ┆ werden in DSM-Systemen eingesetzt
 - ┆ anwendbar für große Systeme mit vielen Prozessoren

JR - RA - SS2002 Kap. 6 5/13

Bit-Vector-Protocol

- Für jede Cache-line existiert ein Eintrag in einer Tabelle (directory)
- Eintrag hat Information über
 - ┆ den Zustand der Cache-line (shared, clean, dirty)
 - ┆ die Knoten, die eine Kopie der Cache-line halten
 - ┆ als Bitvector, jedem Bit ist ein Knoten zugeordnet (presence bits)



JR - RA - SS2002 Kap. 6 5/14

Bit-Vector-Protocol

- Read-Miss
 - Anfordern der Daten beim „Home Node“ per GET
 - Cache-Line im Zustand *clean* oder *shared*
- Home-Node
 - schickt Daten übers Netz (PUT) und
 - setzt das Presence-Bit des Anfragenden Knotens
 - neuer Zustand der Cache-line ist *shared*

JR - RA - SS2002 Kap. 6 5/115

Bit-Vector-Protocol

- Read-Miss
 - Anfordern der Daten beim „Home Node“ per GET
 - Cache-Line im Zustand *dirty*
 - Home-Node
 - leitet Anfrage an Modifying-Node weiter
 - Modifying-Node
 - schickt die Daten an den Requesting-Node und ein SWB (shared writeback) packet mit der aktuellen Cache-line an den Home-Node
 - Home-Node
 - setzt den Zustand auf *shared*
 - aktualisiert das Presence-Bit
 - schreibt Cache-line in den Speicher

JR - RA - SS2002 Kap. 6 5/116

Bit-Vector-Protocol

- Write-Miss
 - GETX-Anfrage an den Home-Node
 - Cache-Line in keinem Cache:
 - Home-Node sendet ein PUTX mit der Cache-line
 - aktualisiert den Presence-Bit-Vector
 - setzt Cache-line in den Zustand *dirty*

JR - RA - SS2002 Kap. 6 5/117

Bit-Vector-Protocol

- Write-Miss
 - Cache-Line in mehreren Caches vorhanden:
 - Requesting Node sends GETX to Home Node.
 - Home Node sends INVAL to other Sharing Nodes.
 - Sharing Nodes respond with INVAL_ACK to Home Node.
 - Home Node sends PUT to Requesting Node.

JR - RA - SS2002 Kap. 6 5/118

Bit-Vector-Protocol

- Write-Miss
 - Cache-Line ist *dirty*:
 - Requesting Node sends GETX to Home Node.
 - Home Node forwards GETX to Modifying Node.
 - Modifying Node sends data to Requesting Node and ownership transfer to Home Node.
 - Home Node sends PUTX to Requesting Node.
 - Modifying Node sends OWN_ACK to Home Node.

JR - RA - SS2002 Kap. 6 5/119

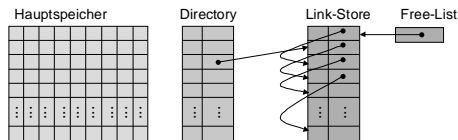
Bit-Vector-Protocol

- Pro Cache-Line im Hauptspeicher
 - 2 Bits für Zustand
 - 1 Bit pro Prozessorknoten im System
- großer Speicherbedarf
- Protokoll nur für kleine und mittelgroße Systeme
- Coarse-Vector
 - jedes Bit repräsentiert ein Cluster von Prozessoren
 - Cluster kann selbst z.B. ein SMP mit snoopy-Protokoll sein

JR - RA - SS2002 Kap. 6 5/120

Dynamic-Pointer-Allocation

- Prozessoren die eine Cache-Line teilen sind als verkettete Liste abgelegt



- Durch dynamische Allokierung auch für große Systeme mit vielen Prozessoren

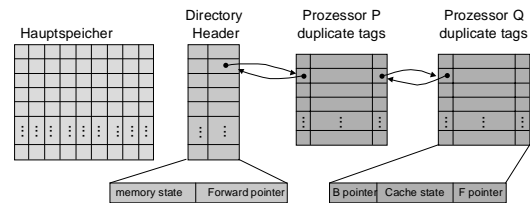
JR - RA - SS2002

Kap. 6

5/121

Scalable Coherent Interface

- SCI ist standardisiert IEEE 1596-1992
- Die Verkettete Liste wird über die Knoten verteilt (doppelverzeigert)



JR - RA - SS2002

Kap. 6

5/122

Software-basierte Cache-Coherency Verfahren

- Hardwareunterstützung gerade bei Netzwerk-Architekturen sehr teuer
- Softwarelösung durch Compiler und Betriebssystem interessante Alternative:
 - Cacheability Marking
 - Cache Coherency Enforcement
 - Indiscriminate Invalidation
 - Fast Selective Invalidation
 - Timestamp

JR - RA - SS2002

Kap. 6

Verschiedene Softwareverfahren

- Cacheability Marking:
 - Einteilung sogenannter "computational units" (z.B. Schleife etc.) in Art der Zugriffe
 - Abhängig von Zugriff kann Variable in Cache oder nicht
 - Beim Übergang zwischen computational units erfolgt Cache-Invalidation bzw. Speicher-Update
 - Compiler kann/muß:
 - computational unit identifizieren
 - Typ des Zugriffs aufgrund einer Datenflußanalyse ermitteln

JR - RA - SS2002

Kap. 6

Einteilung bzgl. Speicherung in Cache

- Arten von Zugriffen (Ann.: Prozesse auf versch. Prozessoren)
 - Read-Only für beliebig viele Prozesse (cacheable)
 - Read-Only für beliebig viele Prozesse und Read-Write für genau einen Prozeß (höchstens Read-Write Prozeß darf Variable in Cache halten, muß aber dann Konsistenz mit Hauptspeicher sicherstellen z.B. durch write through)
 - Read-Write für genau einen Prozeß (cacheable und copy back)
 - Read-Write für eine beliebige Zahl von Prozessen (non cacheable, gilt z.B. typischerweise für Semaphorevariable und andere Synchronisationspunkte)

JR - RA - SS2002

Kap. 6

Cache Coherency Enforcement

- indiscriminate invalidation: zwischen den computation units wird Cache ungültig und es erfolgt Hauptspeicher-Update, alle Variablen eines bestimmten Typs werden hierbei gleich behandelt (leicht zu implementieren aber zu pessimistisch)
- selective invalidation: nur diejenigen Cache-Variablen, die tatsächlich Inkonsistenz erzeugen können werden ungültig
- timestamp: verhindert, daß Cache-Inhalt ungültig wird, wenn gleicher Prozessor wiederholt (aber in verschiedenen computation units) auf dieselbe Variable zugreift

JR - RA - SS2002

Kap. 6