

# Rechnerstrukturen

Vorlesung im Sommersemester 2008

Prof. Dr. Wolfgang Karl

Universität Karlsruhe (TH)

Fakultät für Informatik

Institut für Technische Informatik



- **Kapitel 3: Multiprozessoren – Parallelismus auf Prozess/Thread-Ebene**

## 3.3: Quantitative Maßzahlen

- **Ausführungszeit**

- Die (Gesamt-)Ausführungszeit  $T$  (Execution Time) eines parallelen Programms ist die Zeit zwischen dem Starten der Programmausführung auf einem der Prozessoren bis zu dem Zeitpunkt, an dem der letzte Prozessor die Arbeit an dem Programm beendet hat.
- Zu beachten:
  - Prozessorzustand:
    - Während der Programmausführung sind alle Prozessoren in einem der drei Zustände:
      - » rechnend
      - » kommunizierend
      - » untätig

- **Ausführungszeit**

- Die Ausführungszeit eines parallelen Programms auf einem dediziert zugeordneten Parallelrechner setzt sich zusammen aus:

- **Berechnungszeit  $T_{\text{comp}}$  (Computation Time)**

- Die Zeit, die für die Rechenoperationen verwendet wird

- **Kommunikationszeit  $T_{\text{comm}}$  (Communication Time)**

- Die Zeit, die für Sende- und Empfangsoperationen verwendet wird

- **Untätigkeitszeit  $T_{\text{idle}}$  (Idle Time)**

- Die Zeit, die mit Warten (auf zu empfangene Nachrichten oder und zu versendende) verbraucht wird

- **Es gilt:**

- $T = T_{\text{comp}} + T_{\text{comm}} + T_{\text{idle}}$

- **Ausführungszeit**

- Übertragungszeit einer Nachricht  $T_{msg}$

- die Zeit, die für das Verschicken einer Nachricht von einer bestimmten Länge zwischen zwei Prozessoren benötigt wird

- Die Übertragungszeit setzt sich zusammen aus:

- der Startzeit  $t_s$  (Message Startup Time):

- » Die Zeit, die benötigt wird, um die Kommunikation zu initiieren

- Transferzeit  $t_w$  pro übertragenem Datenwort:

- » hängt von der physikalischen Bandbreite des Kommunikationsmediums ab.

- Voraussetzung:

- » Verbindungsnetz ist konfliktfrei, da sonst die Übertragungszeit nicht fest berechnet werden kann

- Ausführungszeit

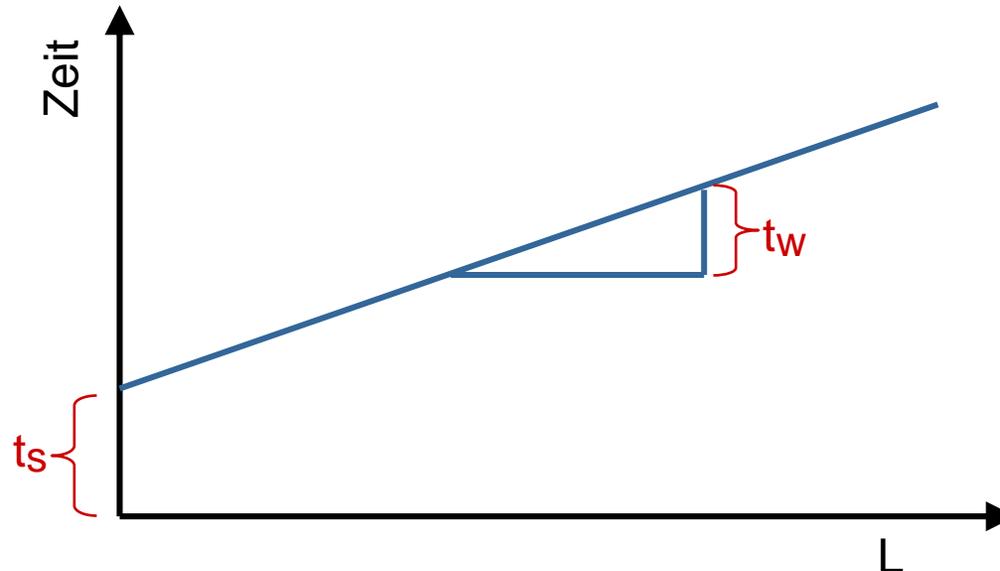
- Übertragungszeit einer Nachricht  $T_{msg}$

- Formel:

$$- T_{msg} = t_s + t_w * L,$$

mit L: Anzahl der Datenwörter

- Graphische Darstellung:



- **Parallelitätsprofil**

- misst die entstehende Parallelität in einem parallelen Programm bzw. bei der Ausführung auf einem Parallelrechner.

- Gibt eine Vorstellung von der inhärenten Parallelität eines Algorithmus/Programms und deren Nutzung auf einem realen oder ideellen Parallelrechner

- **Grafische Darstellung:**

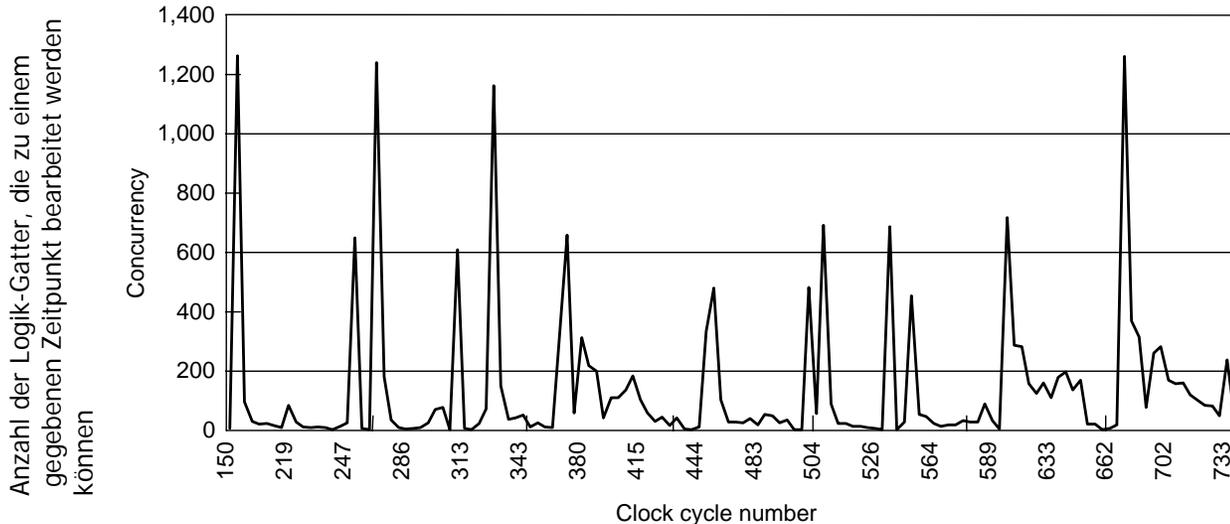
- Auf der x-Achse wird die Zeit und auf der y-Achse die Anzahl paralleler Aktivitäten angetragen.

- ➔ Perioden von Berechnungs- Kommunikations- und Untätigkeitszeiten sind erkennbar.

- **Parallelitätsprofil**

- Zeigt an, wie viele Tasks einer Anwendung zu einem Zeitpunkt parallel ausgeführt werden können
- **Parallelitätsgrad  $PG(t)$ :**
  - Anzahl der Tasks, die zu einem Zeitpunkt parallel bearbeitet werden können

**Beispiel: Parallele ereignisgesteuerte Simulation der Logik-Synthese**



Quelle: D. Culler: Parallel Computer Architecture. Morgan Kaufmann Publishers, 1999, p.87

## • Parallelitätsprofil

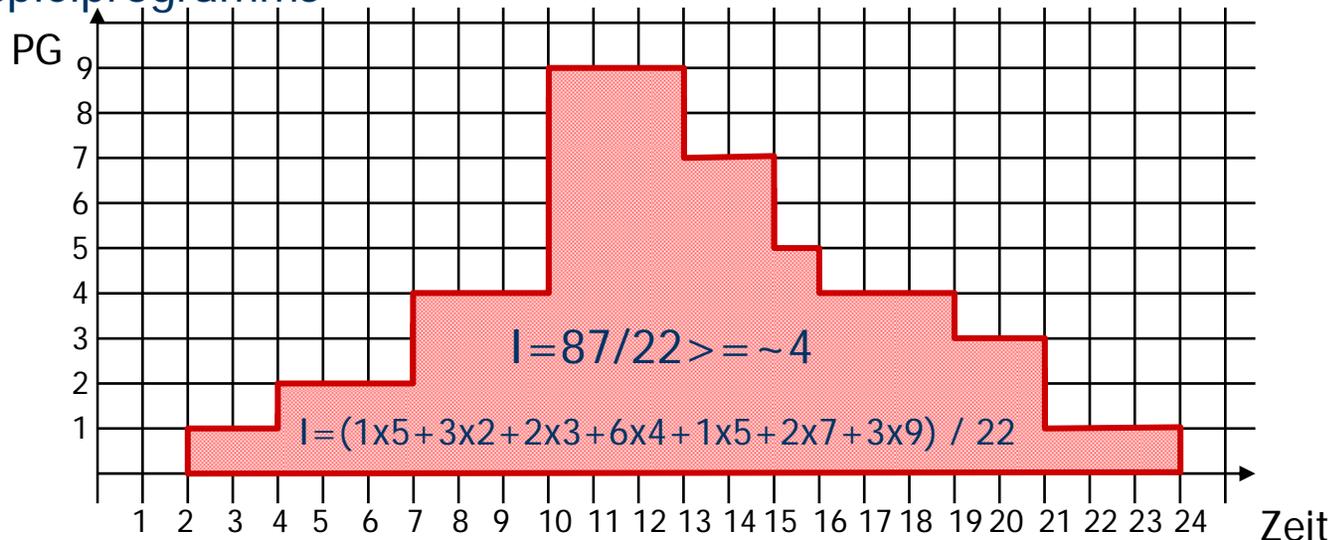
- Parallelindex I (Mittlerer Grad des Parallelismus):

$$I = \frac{1}{t_2 - t_1} \int_{t_1}^{t_2} PG(t) dt$$

$$I = \frac{\left( \sum_{i=1}^m i * t_i \right)}{\left( \sum_{i=1}^m t_i \right)}$$

Parallelitätsprofil eines  
Beispielprogramms

PG Bereich      Ausführungszeit



- **Quantitative Maßzahlen**
  - Leistungsangaben zu Multiprozessorsystemen werden mit Leistungsangaben zu Einprozessorsystemen in Beziehung gesetzt
  - **Notwendig:**
    - Programm das auf beiden zu vergleichenden Systemen ablaufen kann

- Vergleich von Multiprozessorsystemen zu Einprozessorsystemen

- Definitionen:

- $P(1)$ : Anzahl der auszuführenden (Einheits-) Operationen (Tasks) des Programms auf einem Einprozessorsystem.
- $P(n)$ : Anzahl der auszuführenden (Einheits-) Operationen (Tasks) des Programms auf einem Multiprozessorsystem mit  $n$  Prozessoren.
- $T(1)$ : Ausführungszeit auf einem Einprozessorsystem in Schritten (oder Takten).
- $T(n)$ : Ausführungszeit auf einem Multiprozessorsystem mit  $n$  Prozessoren in Schritten (oder Takten).

- Vereinfachende Voraussetzungen:

- $T(1) = P(1)$ ,
  - da in einem Einprozessorsystem (Annahme: einfacher Prozessor) jede (Einheits-) Operation in genau einem Schritt ausgeführt werden kann.
- $T(n) \leq P(n)$ ,
  - da in einem Multiprozessorsystem mit  $n$  Prozessoren ( $n \geq 2$ ) in einem Schritt mehr als eine (Einheits-)Operation ausgeführt werden kann.

- Vergleich von Multiprozessorsystemen zu Einprozessorsystemen

- Beschleunigung  $S(n)$  (Speedup):

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

- Gibt die Verbesserung in der Verarbeitungsgeschwindigkeit an
- Wert bezieht sich auf das jeweils bearbeitete Programm oder kann als Mittelwert eine Menge von Programmen angesehen werden
- Üblicherweise gilt:  $1 \leq S(n) \leq n$

- Vergleich von Multiprozessorsystemen zu Einprozessorsystemen

- Effizienz  $E(n)$

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

- Gibt die relative Verbesserung in der Verarbeitungsgeschwindigkeit an
- Leistungssteigerung wird mit der Anzahl der Prozessoren  $n$  normiert
- Üblicherweise gilt:  $\frac{1}{n} \leq E(n) \leq 1$

- Vergleich von Multiprozessorsystemen zu Einprozessorsystemen
  - Beschleunigung (Speed-Up), Effizienz:
    - Algorithmenunabhängige Definition
      - Man setzt den besten bekannten sequentiellen Algorithmus für das Einprozessorsystem in Beziehung zum vergleichbaren parallelen Algorithmus für das Multiprozessorsystem.
    - ➔ Absolute Beschleunigung
    - ➔ Absolute Effizienz

- Vergleich von Multiprozessorsystemen zu Einprozessorsystemen
  - Beschleunigung (Speed-Up), Effizienz:
    - Algorithmenabhängige Definition
      - Man benutzt den parallelen Algorithmus so, als sei er sequentiell, und misst dessen Laufzeit auf einem Einprozessorsystem.
      - Der für die Parallelisierung erforderliche Zusatzaufwand an Kommunikation und Synchronisation kommt „ungerechterweise“ auch für den sequentiellen Algorithmus zum Tragen.
      - ➔ Absolute Beschleunigung
      - ➔ Absolute Effizienz

- Vergleich von Multiprozessorsystemen zu Einprozessorsystemen
  - Mehraufwand  $R(n)$  für die Parallelisierung:

$$R(n) = \frac{P(n)}{P(1)}$$

- Beschreibt den bei einem Multiprozessorsystem erforderlichen Mehraufwand für die Organisation, Synchronisation und Kommunikation der Prozessoren
- Es gilt:  $1 \leq R(n)$ 
  - Anzahl der auszuführenden Operationen eines parallelen Programms größer ist als diejenige des vergleichbaren sequentiellen Programms

- Vergleich von Multiprozessorsystemen zu Einprozessorsystemen

– Auslastung  $U(n)$ :

$$U(n) = \frac{I(n)}{n} = R(n) \cdot E(n) = \frac{P(n)}{n \cdot T(n)}$$

- Entspricht dem normierten Parallelindex
- Gibt an, wie viele Operationen (Tasks) jeder Prozessor im Durchschnitt pro Zeiteinheit ausgeführt hat

- Vergleich von Multiprozessorsystemen zu Einprozessorsystemen

- Folgerungen

- Alle definierten Ausdrücke haben für  $n = 1$  den Wert 1.
- Der Parallelindex gibt eine obere Schranke für die Leistungssteigerung:

$$1 \leq S(n) \leq I(n) \leq n$$

- Die Auslastung ist eine obere Schranke für die Effizienz:

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

## • Parallelitätsprofil

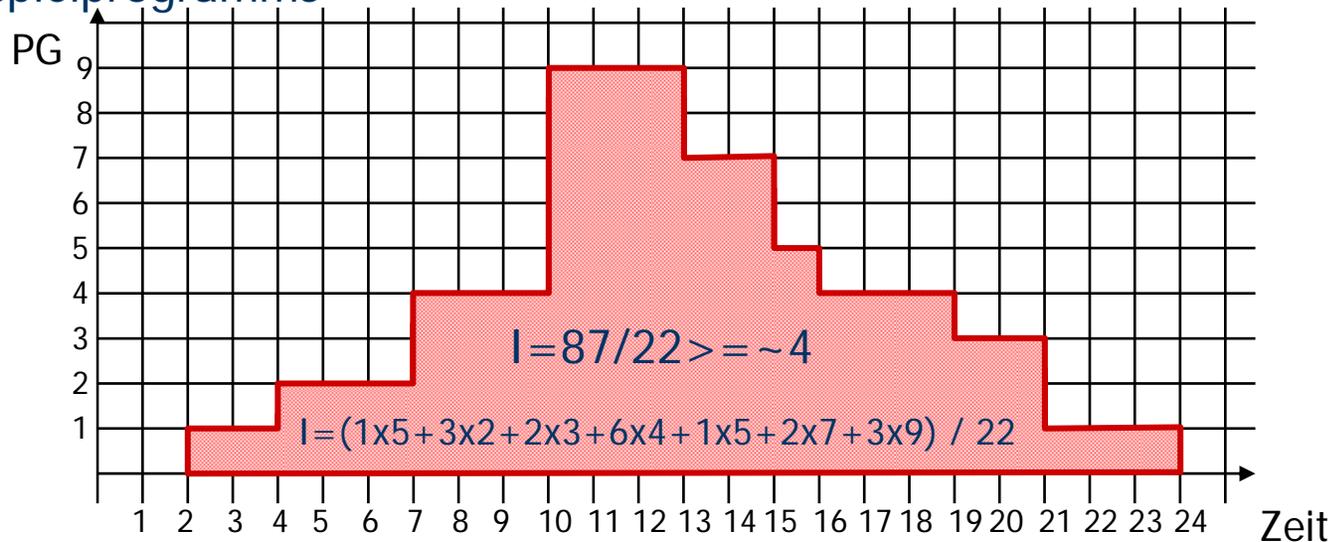
- Parallelindex I (Mittlerer Grad des Parallelismus):

$$I = \frac{1}{t_2 - t_1} \int_{t_1}^{t_2} PG(t) dt$$

$$I = \frac{\left( \sum_{i=1}^m i * t_i \right)}{\left( \sum_{i=1}^m t_i \right)}$$

Parallelitätsprofil eines  
Beispielprogramms

PG Bereich      Ausführungszeit



- Vergleich von Multiprozessorsystemen zu Einprozessorsystemen

- Parallelindex  $I(n)$ :

- Unter den vereinfachenden Voraussetzungen (siehe Folie 12-19) gilt:

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

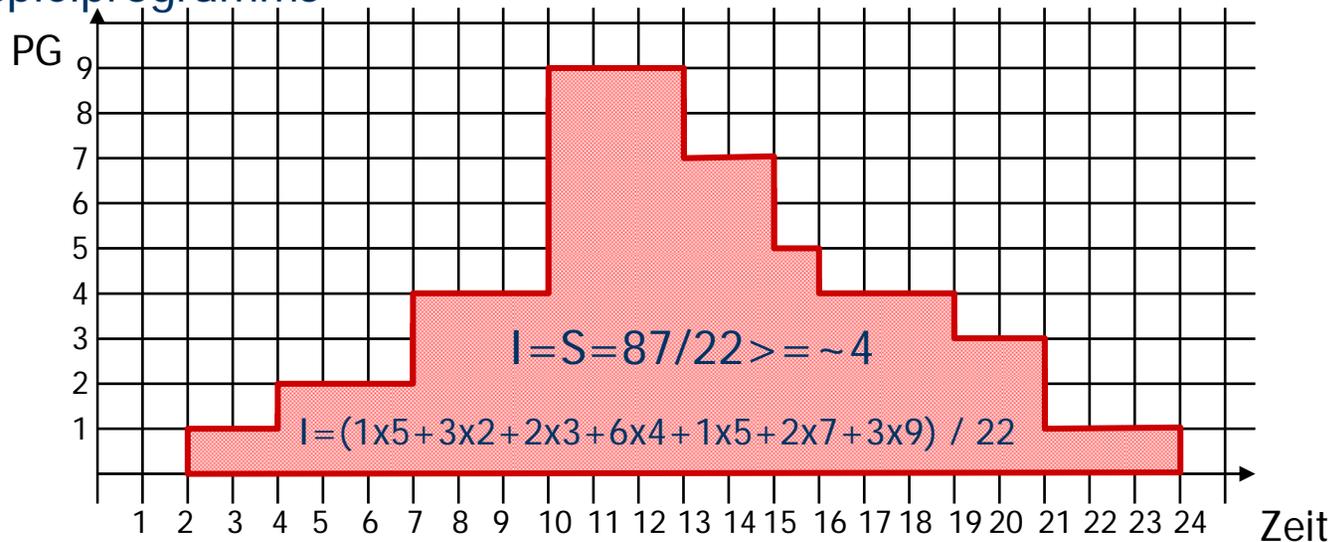
- und damit:  $I(n) = S(n)$

- Vergleich von Multiprozessorsystemen zu Einprozessorsystemen

- Beziehung Parallelindex und Speedup

- Unter den vereinfachenden Voraussetzungen (Folie 12-19)

Parallelitätsprofil eines Beispielprogramms



- Vergleich von Multiprozessorsystemen zu Einprozessorsystemen

- Zahlenbeispiel:

- Ein Einprozessorsystem benötige für die Ausführung von 1000 Operationen 1000 Schritte.
- Ein Multiprozessorsystem mit 4 Prozessoren benötige dafür 1200 Operationen, die aber in 400 Schritten ausgeführt werden können.

- Damit gilt:

$$P(1) = T(1) = 1000, P(4) = 1200, T(4) = 400$$

- Daraus ergibt sich:

$$S(4) = 2,5 \text{ und } E(4) = 0,625$$

- Die Leistungssteigerung verteilt sich als im Mittel zu 62,5% auf alle Prozessoren

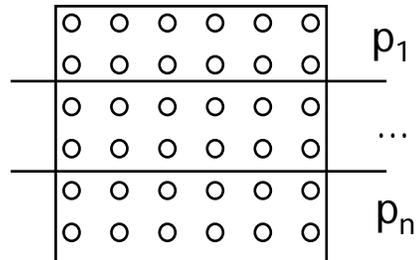
- Vergleich von Multiprozessorsystemen zu Einprozessorsystemen
  - Zahlenbeispiel:
    - Parallelindex und Auslastung:  
 $I(4) = 3$  und  $U(4) = 0,75$ 
      - Es sind im Mittel drei Prozessoren gleichzeitig tätig, d.h., jeder Prozessor ist nur zu 75% der Zeit aktiv.
    - Mehraufwand:  
 $R(4) = 1,2$ 
      - Bei Ausführung auf dem Multiprozessorsystem sind 20% mehr Operationen als bei Ausführung auf einem Einprozessorsystem notwendig.

- **Skalierbarkeit eines Parallelrechners**
  - das Hinzufügen von weiteren Verarbeitungselementen führt zu einer kürzeren Gesamtausführungszeit, ohne dass das Programm geändert werden muss.
  - Insbesondere meint man damit eine lineare Steigerung der Beschleunigung mit einer Effizienz nahe bei Eins.
  - Wichtig für die Skalierbarkeit ist eine angemessene Problemgröße.
  - Bei fester Problemgröße und steigender Prozessorzahl wird ab einer bestimmten Prozessorzahl eine Sättigung eintreten. Die Skalierbarkeit ist in jedem Fall beschränkt.
  - Skaliert man mit der Anzahl der Prozessoren auch die Problemgröße (scaled problem analysis), so tritt dieser Effekt bei gut skalierenden Hardware- oder Software-Systemen nicht auf.

- Gesetz von Amdahl

- Beispiel:

- Gegeben: ein zweidimensionales  $k \times k$ -Gitter
  - 1. Berechnungsphase: Ausführung einer Operation auf allen Gitterpunkten
    - » Annahme: keine Abhängigkeiten zwischen den Gitterpunkten
    - » Parallele Berechnung auf  $n$  Prozessoren



- 2. Berechnungsphase: Berechnung der Summe der  $k^2$  berechneten Werte der Gittersumme
  - » Jeder Prozessor addiert seine  $k^2/n$  berechneten Werte zur globalen Summe

- Gesetz von Amdahl

- Beispiel:

- Problem:

- Akkumulation der globalen Summe muss serialisiert werden!
- 2. Phase benötigt  $k^2$  Zeiteinheiten unabhängig von  $n$
- Ausführungszeit des parallelen Programms:  $k^2/n + k^2$
- Ausführungszeit des sequentiellen Programms:  $2k^2$
- Möglicher Speedup  $S$ :

$$\frac{2k^2}{\frac{k^2}{n} + k^2} = \frac{2n}{n+1}$$

- Selbst bei einer hohen Anzahl Prozessoren nicht mehr als 2!

- Gesetz von Amdahl
  - Gesamtausführungszeit  $T(n)$

$$T(n) = T(1) \cdot \frac{1-a}{n} + T(1) \cdot a$$

Ausführungszeit  
des sequentiell  
ausführbaren  
Programmteils a

Ausführungszeit  
des parallel  
ausführbaren  
Programmteils 1-a

a: Anteil des Programmteils,  
der nur sequentiell  
ausgeführt werden kann

- Beschleunigung

$$S(n) = \frac{T(1)}{T(n)} = \frac{T(1)}{T(1) \cdot \frac{1-a}{n} + T(1) \cdot a} = \frac{1}{\frac{1-a}{n} + a}$$

- Für  $n \rightarrow \infty$  :  $S(n) = \frac{1}{a}$

- Gesetz von Amdahl

- Beispiel:

- Erhöhung der Parallelität

- Aufteilung der 2. Berechnungsphase in zwei weiteren Teilphasen:

- » 1. Teilphase: Jeder Prozessor berechnet die Summe seiner berechneten Werte

- » Kann vollständig parallel abgearbeitet werden

- » 2. Teilphase: Akkumulation der Teilsummen

- » Weiterhin seriell!

- Ausführungszeit  $T(n) = k^2/n + k^2/n + n$

- Beschleunigung  $S(n) = n * 2k^2 / (2k^2 + n^2)$

- Wenn  $n$  groß genug, dann nahezu linear!

- Gesetz von Amdahl

- Diskussion

- Amdahls Gesetz zufolge kann eine kleine Anzahl von sequentiellen Operationen die mit einem Parallelrechner erreichbare Beschleunigung signifikant begrenzen.
- Beispiel:  $a = 1/10$  des parallelen Programms kann nur sequenziell ausgeführt werden, → das gesamte Programm kann maximal zehnmal schneller als ein vergleichbares, rein sequenzielles Programm sein.
- Jedoch: viele parallele Programme haben einen sehr geringen sequenziellen Anteil ( $a \ll 1$ )

- Synergetischer Effekt und superlinearer Speedup
  - Theorie : einen „superlinearen Speedup“ kann es nicht geben:
    - Jeder parallele Algorithmus lässt sich auf einem Einprozessorsystem simulieren, indem in einer Schleife jeweils der nächste Schritt jedes Prozessors der parallelen Maschine emuliert wird.

- Synergetischer Effekt und superlinearer Speedup
  - Ein „superlinearer Speed-up“ kann real beobachtet werden bei
    - parallelem Backtracking (*depth-first search*)
    - Beim Programmmlauf auf einem Rechner passen die Daten nicht in den Hauptspeicher des Rechners (häufiger Seitenwechsel), aber: bei Verteilung auf die Knoten des Multiprozessors können die parallelen Programme vollständig in den Cache- und Hauptspeichern der einzelnen Knoten ablaufen.

- Weitere grundsätzliche Probleme bei Multiprozessoren
  - Verwaltungsaufwand (Overhead)
    - Steigt mit der Zahl der zu verwaltenden Prozessoren
  - Möglichkeit von Systemverklemmungen (*deadlocks*)
  - Möglichkeit von Sättigungserscheinungen
    - können durch Systemengpässe (*bottlenecks*) verursacht werden.

- **Kapitel 3: Multiprozessoren – Parallelismus auf Prozess/Thread-Ebene**

## 3.4: Verbindungsstrukturen

- **Verbindungsnetze in Multiprozessoren**
  - Ermöglichen die Kommunikation und Kooperation zwischen den Verarbeitungselementen (Knoten)
    - Zuverlässiger Austausch von Informationen
  - Einsatz eines Verbindungsnetzwerks
    - Multiprozessor mit verteiltem Speicher (nachrichtenorientierter Multiprozessor)
      - Verbinden physikalisch jeden Knoten für das Versenden von Nachrichten
      - Direkte Send/Receive-Kommunikation zwischen den Knoten
    - Multiprozessor mit gemeinsamem Speicher
      - Ermöglicht den Zugriff aller Knoten auf den gemeinsamen Speicher
      - Kommunikation durch Lesen und Schreiben auf gemeinsame Daten

- Charakterisierung von Verbindungsnetzwerken

- Latenz

- Übertragungszeit einer Nachricht  $T_{\text{msg}}$ 
  - die Zeit, die für das Verschicken einer Nachricht von einer bestimmten Länge zwischen zwei Prozessoren benötigt wird
- Die Übertragungszeit setzt sich zusammen aus:
  - der Startzeit  $t_s$  (Message Startup Time):
    - » Die Zeit, die benötigt wird, um die Kommunikation zu initiieren
  - Transferzeit  $t_w$  pro übertragenem Datenwort:
    - » hängt von der physikalischen Bandbreite des Kommunikationsmediums ab.
  - Voraussetzung:
    - » Verbindungsnetz ist konfliktfrei, da sonst die Übertragungszeit nicht fest berechnet werden kann

- Charakterisierung von Verbindungsnetzwerken

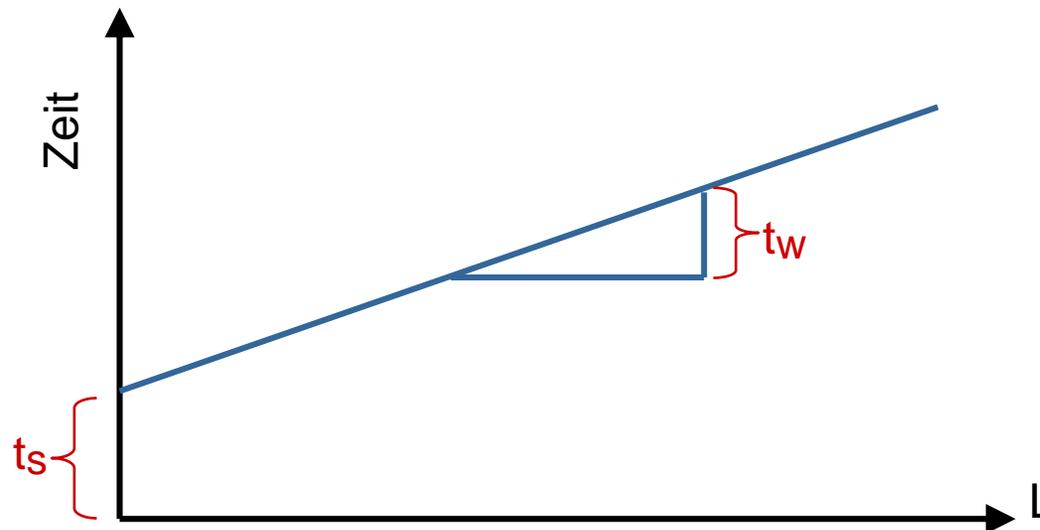
- Latenz

- Formel:

$$- T_{msg} = t_s + t_w * L,$$

mit L: Anzahl der Datenwörter

- Graphische Darstellung:



- Charakterisierung von Verbindungsnetzwerken
  - Latenz (Übertragungszeit einer Nachricht, latency)
    - Software-Overhead
      - in Verbindung mit dem Senden und Empfangen von Nachrichten

- Charakterisierung von Verbindungsnetzwerken
  - Beispiel für Software-Overhead

Erzeugerprozess:

`send(proci, processi, @sbuffer, num_bytes)`

**Sender**

Systemaufruf  
 Prüfe Schutzbed.  
 DMA Init.  
 DMA nach NI

Verbraucherprozess:

`receive(@rbuffer, max_bytes)`

**Empfänger**

DMA vom Netzwerk in den Puffer  
 BS Interrupt und Dekodierung der Nachr.  
 BS kopiert Systempuffer in Userpuffer  
 Reschedule Benutzerprozess  
 Lesen der Nachricht

Zeit

NI: Netzwerkschnittstelle  
 DMA: Direktspeicherzugriff

- Charakterisierung von Verbindungsnetzwerken
  - Latenz (Übertragungszeit einer Nachricht, latency)
    - Kanalverzögerung (channel delay)
      - Dauer Belegung eines Kommunikationskanals durch eine Nachricht
      - Kanal:
        - » Physikalische Verbindung zwischen Schalterelementen oder Knoten mit einem Puffer zum Halten der Daten während ihrer Übertragung
      - Verbindung (link)
        - » Menge von Leitungen

- Charakterisierung von Verbindungsnetzwerken
  - Latenz (Übertragungszeit einer Nachricht, latency)
    - Schaltverzögerung, Routing-Verzögerung (switching delay, routing delay)
      - Zeit, einen Weg zwischen zwei Knoten aufzubauen
      - Pfadberechnung oder Wegefindung (Routing)
        - » die Art, wie der Weg einer Nachricht vom Sender zum Zielknoten berechnet wird
        - » Zu einer Verbindungsstruktur kann es mehrere Wegefindungsalgorithmen geben
        - » einfache Implementierung in Verbindungselementen mit Hilfe eines schnellen Hardware-Algorithmus

- Charakterisierung von Verbindungsnetzwerken
  - Latenz (Übertragungszeit einer Nachricht, latency)
    - Blockierungszeit (contention time)
      - Wird verursacht, wenn zu einem Zeitpunkt mehr als eine Nachricht auf eine Netzwerkressource zugreifen
    - Blockierung (contention)
      - Ein Verbindungsnetzwerk heisst blockierungsfrei, falls jede gewünschte Verbindung zwischen Prozessoren oder zwischen Prozessoren und Speichern unabhängig von schon bestehenden Verbindungen hergestellt werden kann

- Charakterisierung von Verbindungsnetzwerken
  - Durchsatz oder Übertragungsbandbreite (bandwidth)
    - Maximale Übertragungsleistung des Verbindungsnetzwerkes oder einzelner Verbindungen, meist in Megabits pro Sekunde (MBit/s) oder Megabytes pro Sekunde (MB/s)
  - Bisektionsbandbreite (bisection bandwidth)
    - Maximale Anzahl von Megabytes pro Sekunde, die das Netzwerk über die Bisektionslinie, die das Netzwerk in zwei gleiche Hälften teilt, transportieren kann

- **Charakterisierung von Verbindungsnetzwerken**
  - **Diameter oder Durchmesser  $r$  (*diameter*):**
    - 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.
  - **Verbindungsgrad eines Knotens  $P$  (node degree, connectivity)**
    - ist definiert als die Anzahl der direkten Verbindungen, die von einem Knoten zu anderen Knoten bestehen.
  - **Mittlere Distanz  $d_a$  (average distance) zwischen zwei Knoten**
    - Anzahl der Links auf dem kürzesten Pfad zwischen zwei Knoten
    - $P/d_a$  ist die maximale Anzahl neuer Nachrichten, die von jedem Knoten in einem Zyklus in das Netzwerk eingebracht werden können

- **Charakterisierung von Verbindungsnetzwerken**
  - **Komplexität oder Kosten:**
    - Kosten für die Implementierung einer Hardware
    - Aufwand für das Verbindungsnetz gemessen in der Anzahl und der Art der Schaltelemente und Verbindungsleitungen.
  - **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.

- Charakterisierung von Verbindungsnetzwerken
  - Ausfallstoleranz 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.

- Charakterisierung von Verbindungsnetzwerken
  - Art der Adressierung:
    - 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.

- Charakterisierung von Verbindungsnetzwerken

- Art des Datentransfers:

- Durchschalte- oder Leitungsvermittlung (circuit switching):

- Eigenschaft eines Netzes eine direkte Verbindung zwischen zwei oder mehreren Knoten eines Netzes zu schalten. Die physikalische Verbindung bleibt für die gesamte Dauer der Informationsübertragung bestehen.

- » Blockierungsfreie Kommunikation

- » Kurze Latenz

- » Gut geeignet für lange Nachrichten, da die Zeit zum Aufsetzen einer Nachricht im Verhältnis zur Übertragungszeit kurz ist

- » Übertragungszeit einer Nachricht der Länge  $L$  über eine Distanz  $d$  beträgt:  $L/b + d\delta$ , mit individueller Schaltverzögerung  $\delta$  und der Bandbreite  $b$  eines Kanals

- Charakterisierung von Verbindungsnetzwerken
  - Art des Datentransfers:
    - Paketvermittlung (packet switching):
      - Datenpakete fester Länge oder Nachrichten variabler Länge werden entsprechend einem Wegefindungsalgorithmus (routing) vom Absender zum Empfänger geschickt
      - Nachrichten mit Adresse und Daten werden durch das Netzwerk verschickt
        - » Adresse wird in jedem Knoten gelesen und die Nachricht wird zum nächsten Knoten weitergeleitet, bis die Nachricht das Ziel erreicht
        - » Günstig für kurze Nachrichten

- Charakterisierung von Verbindungsnetzwerken
  - Art des Datentransfers:
    - Paketvermittlung (packet switching):
      - Übertragungsmodi: Store-and-forward-Modus
        - » Jeder Knoten enthält einen Puffer zum Aufnehmen der vollständigen Nachricht
        - » Nachricht wird von jedem Zwischenknoten in Empfang genommen, vollständig zwischengespeichert und dann weiter übertragen
        - » Nachfolgende Pakete werden nacheinander verschickt
        - » Gegenüber Circuit Switching: höhere Bandbreite, aber auch höhere Latenz
        - » Übertragungszeit einer Nachricht der Länge  $L$  über eine Distanz  $d$  von einer Quelle zum Ziel beträgt  $d(L/b + \delta)$

- Charakterisierung von Verbindungsnetzwerken
  - Art des Datentransfers:
    - Paketvermittlung (packet switching):
      - Übertragungsmodi: Cut-through oder wormhole
      - 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.

- Charakterisierung von Verbindungsnetzwerken
  - Art des Datentransfers:
    - Paketvermittlung (packet switching):
      - Übertragungsmodi: Virtual-cut-through-Modus:
        - » Nachricht wird aufgeteilt in Zellen (Flits) fester Größe
        - » Der Kopfteil der Nachricht enthält die Empfängeradresse und bestimmt den einzuschlagenden Weg. Flits mit Daten folgen dem Kopf auf dem Pfad von der Quelle zum Ziel
        - » Bei Ankunft des Kopfs einer Nachricht wird dieser dekodiert . Nachfolgende Flits werden automatisch an den nächsten Knoten auf dem ausgewählten Pfad weitergeleitet , ein Flit pro Zeiteinheit gemäß einer Pipeline-Verarbeitung
        - » Kopf-Information wird festgehalten bis letztes Flit angekommen ist.
        - » 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

- Charakterisierung von Verbindungsnetzwerken
  - Art des Datentransfers:
    - Paketvermittlung (packet switching):
      - 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.

- Charakterisierung von Verbindungsnetzwerken
  - Art des Datentransfers:
    - Paketvermittlung (packet switching):
      - Übertragungsmodi: Buffered wormhole routing:
        - » Kompromisslösung zwischen Virtual-cut-through- und Wormhole-routing-Modus eingesetzt,
        - » begrenzter Puffer zur Aufnahme kleinerer Pakete vorhanden
        - » größere Pakete werden im Blockierungsfall – ähnlich dem Wormhole-routing-Modus – in den Puffern mehrerer Knoten zwischengespeichert.

- **Topologie**

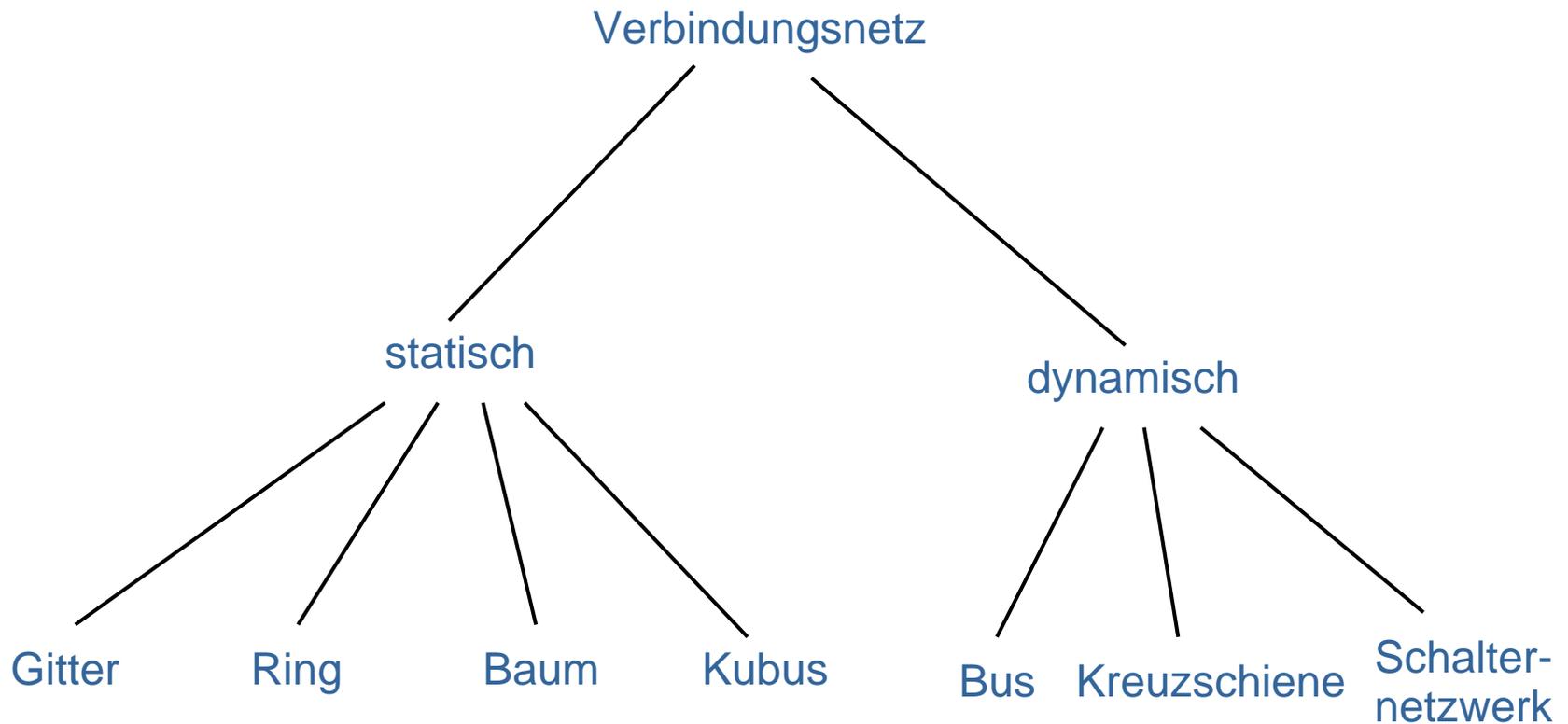
- **Statische Verbindungsnetzwerke**

- Die Prozessoren sind mit jedem anderen Prozessor direkt verbunden
- Fest installierte Verbindungen zwischen Paaren von Netzknoten
  - Die Knoten sind mit Punkt-zu-Punkt-Verbindungen verbunden, die sich während der Programmausführung nicht ändern

- **Dynamische Verbindungsnetze**

- Die Knoten sind über Schaltelemente miteinander verbunden, die konfiguriert werden können, um den Kommunikationsanforderungen des ausführenden Programms zu genügen
  - Direkte fest installierte Verbindungen zwischen den Knoten existieren nicht

- Topologie: Klassifizierung

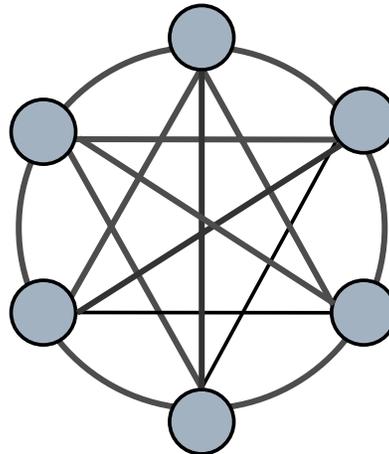


- **Statische Verbindungsnetze**
  - Nach Aufbau des Verbindungsnetzes bleiben die Verbindungen fest
    - Gute Leistung für Probleme mit vorhersagbaren Kommunikationsmustern zwischen benachbarten Knoten

- **Statische Verbindungsnetze**

- **Vollständige Verbindung**

- Jeder Knoten ist mit jedem anderen Knoten verbunden
- Höchste Leistungsfähigkeit
  - Arbeitet für alle Anwendungen mit allen Arten von Kommunikationsmustern effizient
- Aber: nicht praktikabel in Parallelrechnern
  - Netzwerkkosten steigen quadratisch mit der Anzahl der Prozessoren



- Statische Verbindungsnetze

- Gitterstrukturen

- 1-dimensionales Gitter (lineares Feld, Kette)
  - Verbindet  $N$  Knoten mit  $(N-1)$  Verbindungen
  - Endknoten haben den Grad 1, Zwischenknoten den Grad 2 und sind mit benachbarten Knoten verbunden
  - Diameter  $r$  ist  $N-1$
  - Disjunkte Bereiche des linearen Netzwerkes können gleichzeitig genutzt werden, aber es sind mehrere Schritte notwendig, um eine Nachricht zwischen zwei nicht benachbarte Knoten zu verschicken



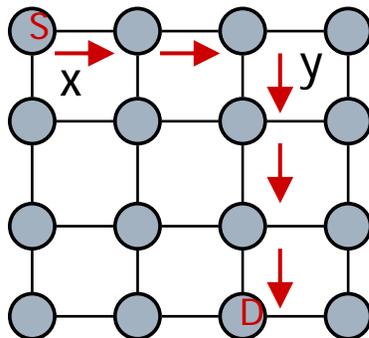
- Statische Verbindungsnetze

- Gitterstrukturen

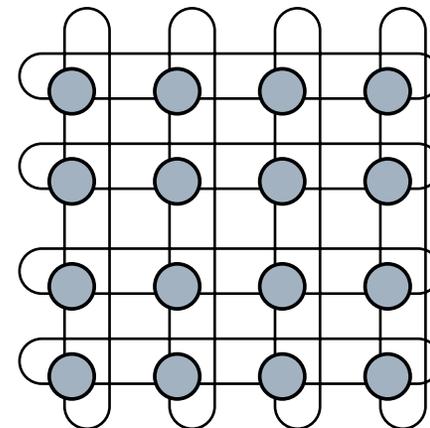
- k-dimensionales Gitter mit N Knoten

- Innere Knoten haben den Grad  $2k$ , wobei die  $2k$  benachbarten Knoten miteinander verbunden sind
- In einem k-dimensionalen Netzwerk mit  $\sqrt[k]{N}$  Knoten in jeder Dimension beträgt der Diameter  $k(\sqrt[k]{N} - 1)$

2-dim. Gitter



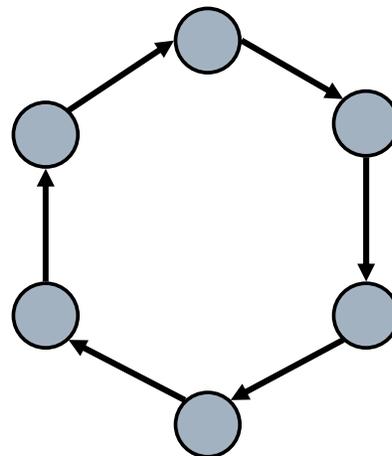
2-dim. Torus



- **Statische Verbindungsnetze**

- **Ring**

- Erhält man, wenn man die Endknoten eines linearen Feldes miteinander verbindet
- **Unidirektionaler Ring mit N Knoten**
  - Nachrichten werden in einer Richtung vom Quellknoten zum Zielknoten verschickt
  - Diameter  $r$  ist  $N-1$
  - Bei Ausfall einer Verbindung bricht die Kommunikation zusammen

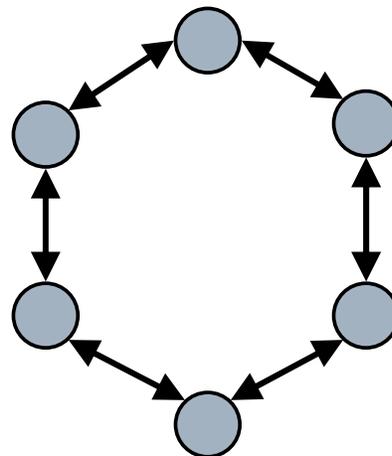


- Statische Verbindungsnetze

- Ring

- Bidirektionaler Ring mit N Knoten

- symmetrisches Netzwerk
- Der längste Pfad, den eine Nachricht nehmen muss, ist nicht länger als  $N/2$
- Bei Ausfall einer Verbindung bricht die Kommunikation noch nicht zusammen, während zwei Ausfälle von Verbindungen das Netzwerk in zwei disjunkte Teilnetzwerke aufteilen



- Statische Verbindungsnetze

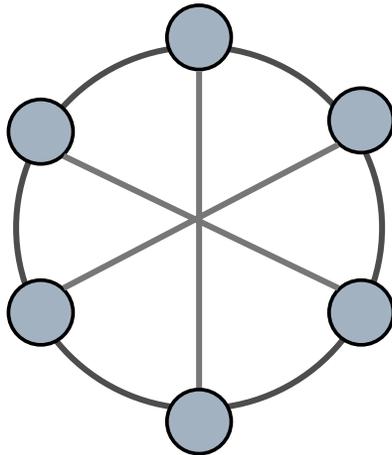
- Ring

- Chordialer Ring

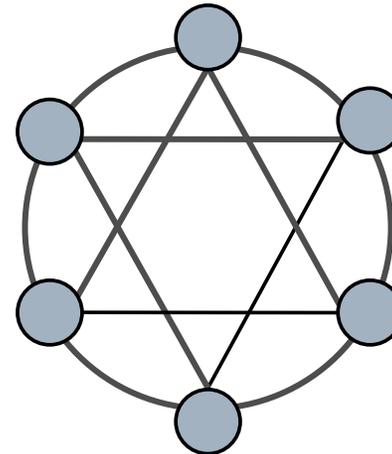
- Hinzufügen redundanter Verbindungen

- » erhöht Fehlertoleranzeigenschaft des Verbindungsnetzwerkes

- » Höherer Knotengrad und kleinerer Diameter gegenüber Ring



Chordaler Ring mit Knotengrad 3



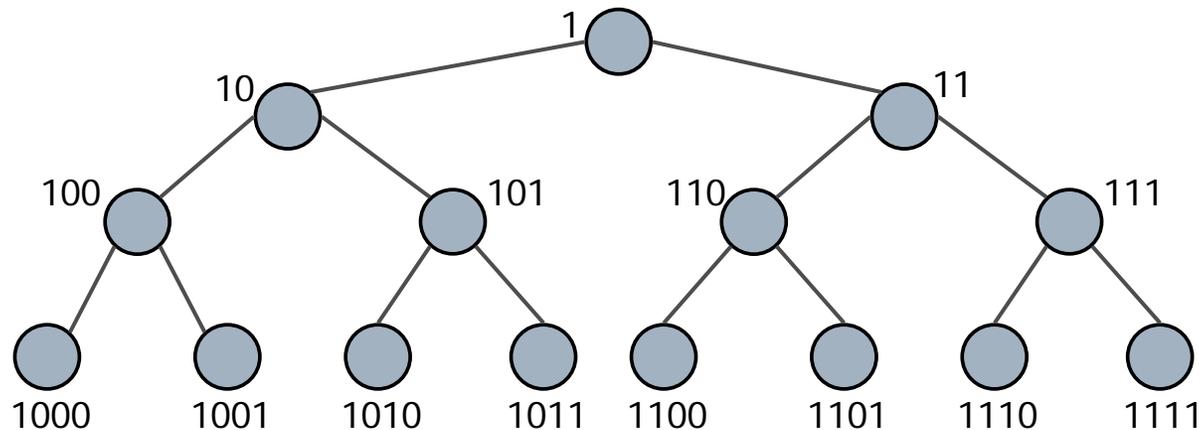
Chordaler Ring mit Knotengrad 4

- **Statische Verbindungsnetze**

- **Baumstrukturen**

- **Binärer Baum mit m-Ebenen:**

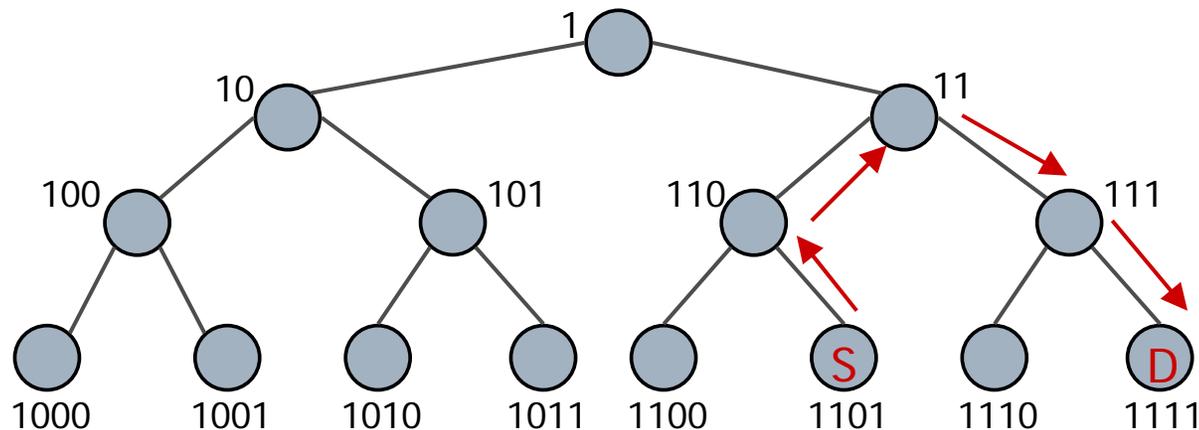
- Auf Ebene m:  $N=2^m-1$  Knoten
- Diameter:  $2(m-1)$
- Adressierung der Knoten:
  - » Die Knotennummer auf Ebene m besteht aus m Bits
  - » Der Wurzelknoten hat die Nummer 1
  - » Die Nummer des linken Kindknoten erhält man durch Hinzufügen einer 0 an die niederwertige Stelle der Adresse des Elternknoten
  - » Die Nummer des rechten Kindknoten erhält man durch Hinzufügen einer 1 an die niederwertige Stelle der Adresse des Elternknoten



- Statische Verbindungsnetze

- Baumstrukturen

- Routing:
- Finde gemeinsamen Elternknoten P von S und D
- Gehe von S nach P und von P nach D



- Statische Verbindungsnetze

- Baumstrukturen

- Routing:

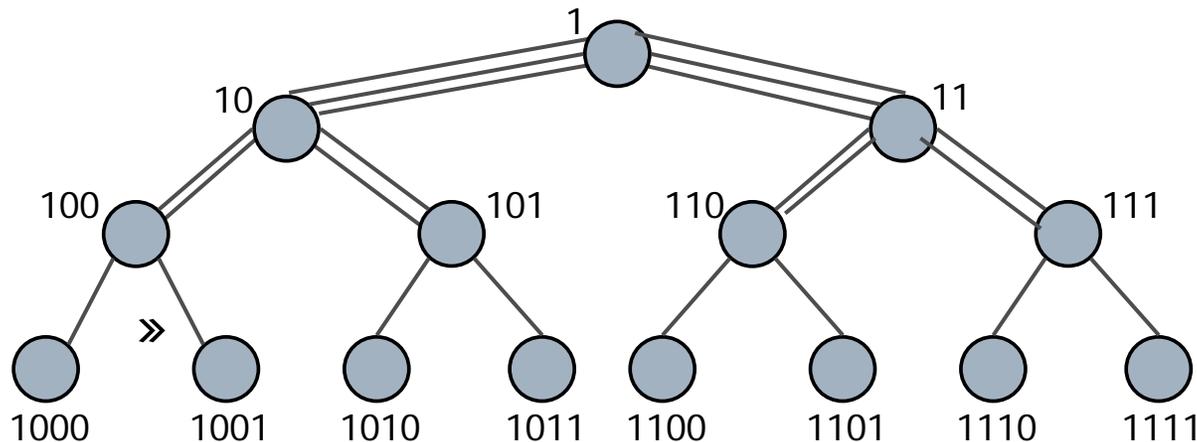
- » Die Binärdarstellung der Adresse eines Quellknotens  $S$  auf Ebene  $i$  sei  $S_i S_{i-1} \dots S_1$  und die der Adresse des Zielknotens  $D$  auf Ebene  $j$  sei  $D_j D_{j-1} \dots D_1$
- » Finde die gemeinsamen höchstwertigen Bits von  $S$  und  $D$ , so dass die Adresse des Elternknotens  $P$  gleich  $D_j D_{j-1} \dots D_x = S_i S_{i-1} \dots S_{(i-j+x)}$  ist
- » Steige von  $S$   $(i-j+x)$  Ebenen auf nach  $P$
- »           for  $k=x-1$  step 1 until 0  
              {steige nach links ab, falls  $D_x=0$   
              steige nach rechts ab, falls  $D_x=1$ }

- **Statische Verbindungsnetze**

- Baumstrukturen

- Fat Tree:

- Lösung des Blockierungsproblems in Richtung Wurzel
- Kommunikationskanäle werden größer, je näher man sich der Wurzel nähert



- **Statische Verbindungsnetze**
  - Baumstrukturen
    - **Dynamic Fat Tree:**
      - Interne Knoten sind Schalter
      - Beispiel: Quadrics QSnet
        - » Los Alamos Lab (LANL): ASCI Q



Images Courtesy of LANL, LLNL, PNNL, PSC, CEA  
Quelle: <http://www.c3.lanl.gov/~fabrizio/quadrics.html>

- **Statische Verbindungsnetze**
  - Los Alamos Lab (LANL): ASCI Q  
AlphaServer SC45, 1.25 GHz
    - 8125 Prozessoren Alpha 1250 MHz
    - Verbindungsnetzwerk: Quadrics
    - Linpack: 13880 GFLOPS (TOP500: 12, 06/2005)

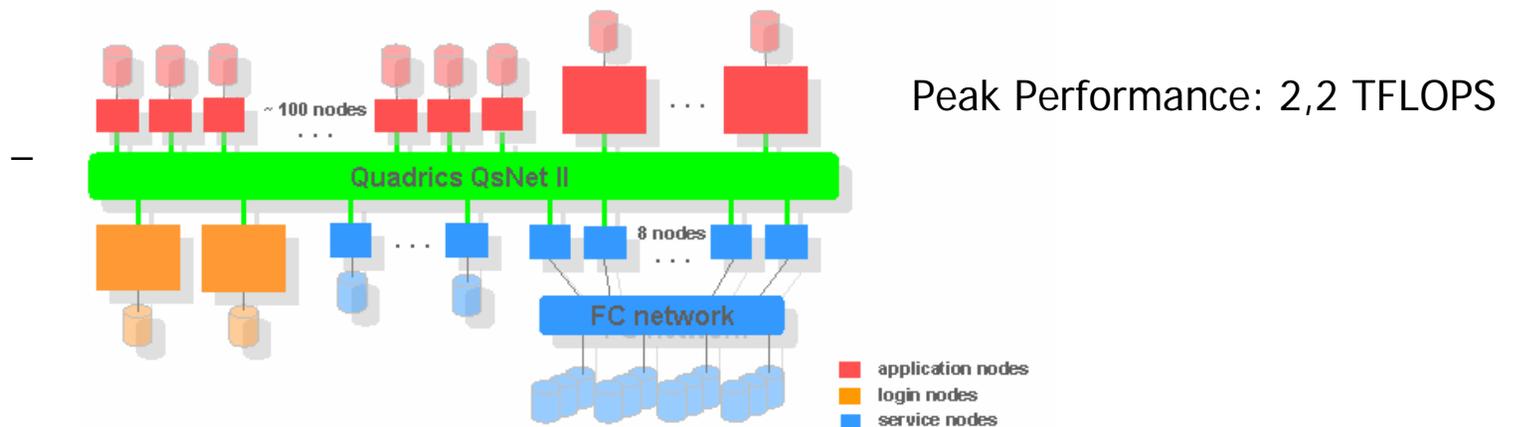
- Statische Verbindungsnetze
  - Baumstrukturen
    - Fat Tree:
      - Beispiel: Quadrics QSnet
        - » Universität Karlsruhe (TH),  
Landeshöchstleistungsrechner, HP XC 6000

»



- **Statische Verbindungsnetze**

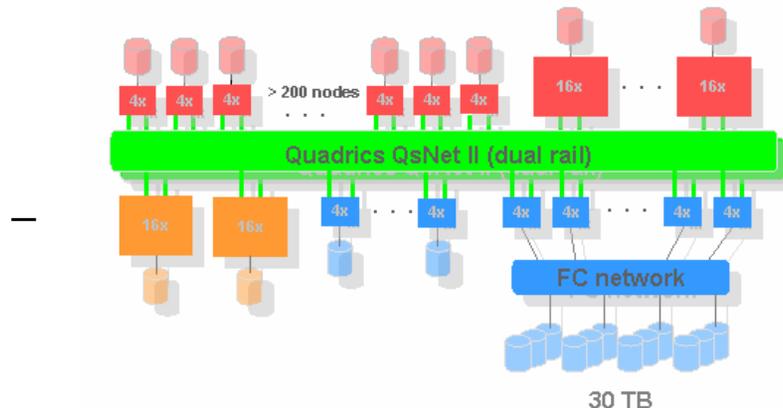
- Universität Karlsruhe (TH), Landeshöchstleistungsrechner,
  - HP XC 6000: Konfiguration
    - 108 Knoten mit jeweils 2 Intel Itanium2 Prozessoren (1,5 GHz) ,12 GB Hauptspeicher pro Knoten und 146 GB lokalem Plattenplatz,
    - 12 Knoten mit jeweils 8 Intel Itanium2 Prozessoren (1,6 GHz), 64 GB Hauptspeicher pro Knoten und ca. 500 GB lokalem Plattenplatz,
    - 8 2-Wege Fileserver-Knoten basierend auf Xeon Prozessoren mit angeschlossenen Platten in der Größe von 10 TB und
    - Quadrics QsNet II Interconnect
    - <http://www.rz.uni-karlsruhe.de/ssc/hpxc.php>



- **Statische Verbindungsnetze**

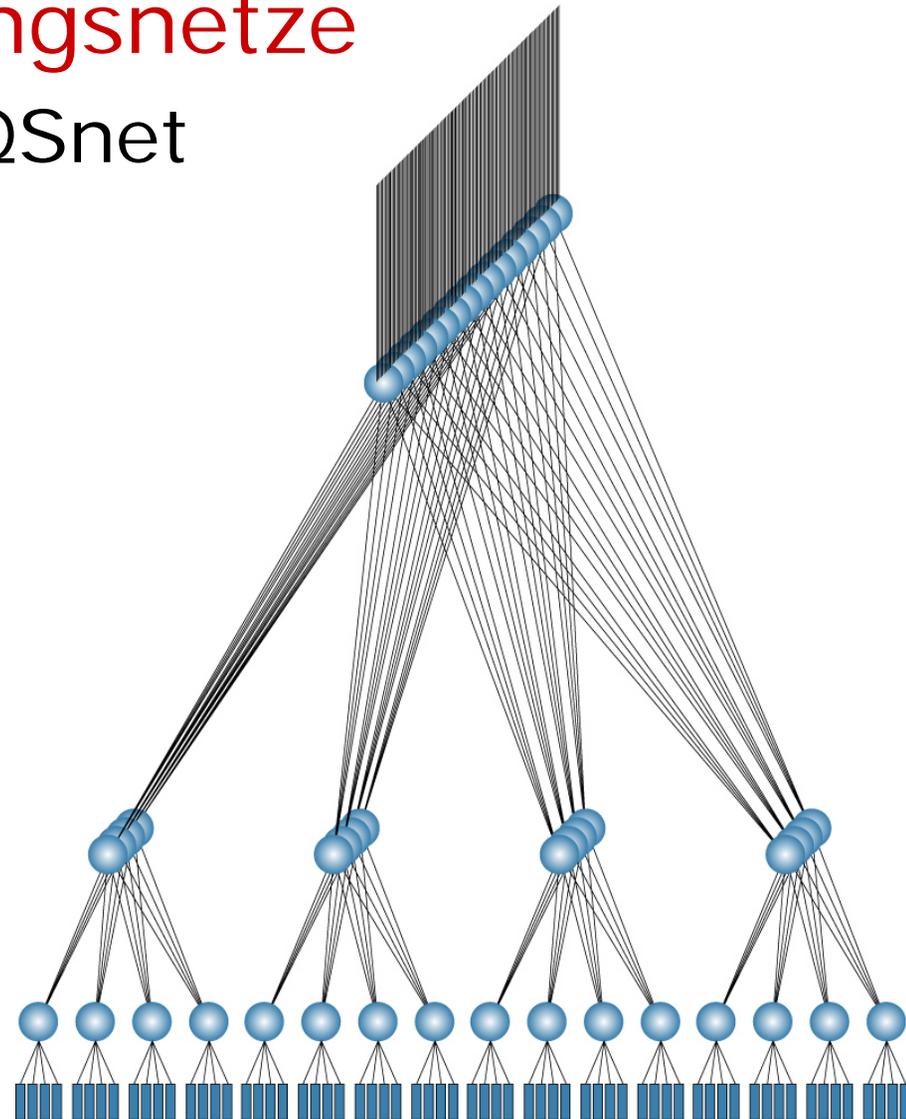
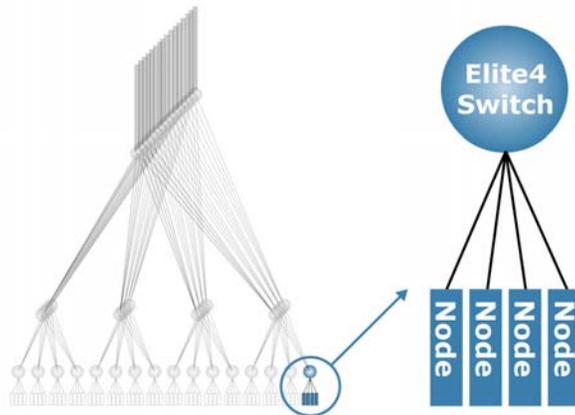
- Universität Karlsruhe (TH),  
Landeshöchstleistungsrechner, HP XC 6000

- Im 1. Quartal 2006 wird das Produktionssystem um 218 4-Wege Knoten erweitert mit
  - 2 "dual core" Intel Itanium2 Prozessoren,
  - 24 GB Hauptspeicher pro Knoten,
  - 146 GB lokalem Plattenplatz,
  - dual rail Quadrics QsNet II Interconnect und
  - 30 TB globaler Plattenplatz.



Peak Performance: 11 TFLOPS

- **Statische Verbindungsnetze**
  - Beispiel: Quadrics QSnet

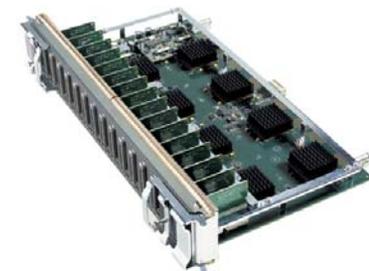


Quelle: D. Addison, J. Beecroft, D. Hewson, M. McLaren (Quadrics Ltd.),  
 Fabrizio Petrini (LANL): A network for Supercomputing Applications. Hot Chips,  
 August 2003, <http://www.c3.lanl.gov/~fabrizio/quadrics.html>

- **Statische Verbindungsnetze**
  - Beispiel: Quadrics QSnet
    - Elan 4 network interface card:



- Elite 4 Switch Component:

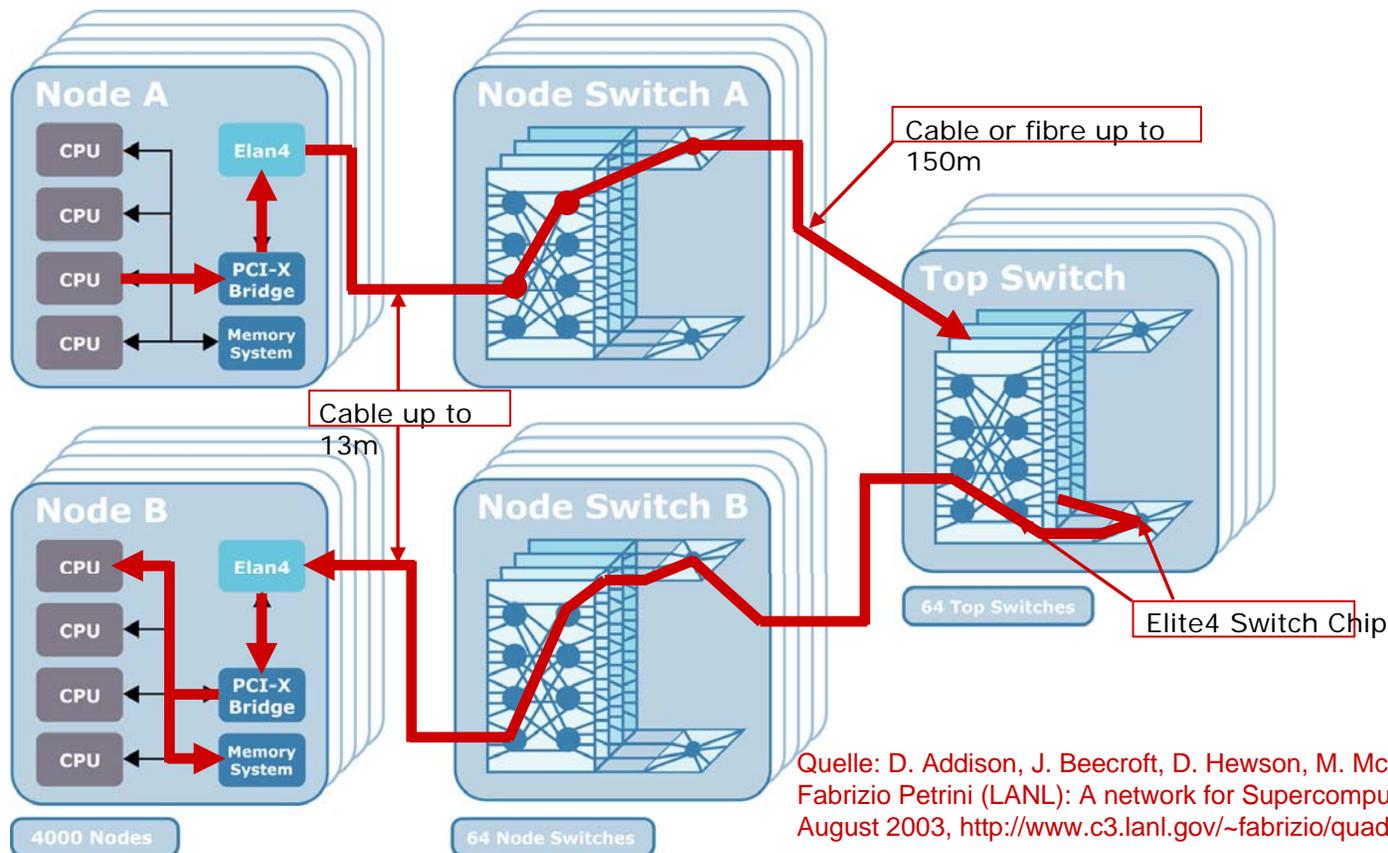


- QsNet II Switch:

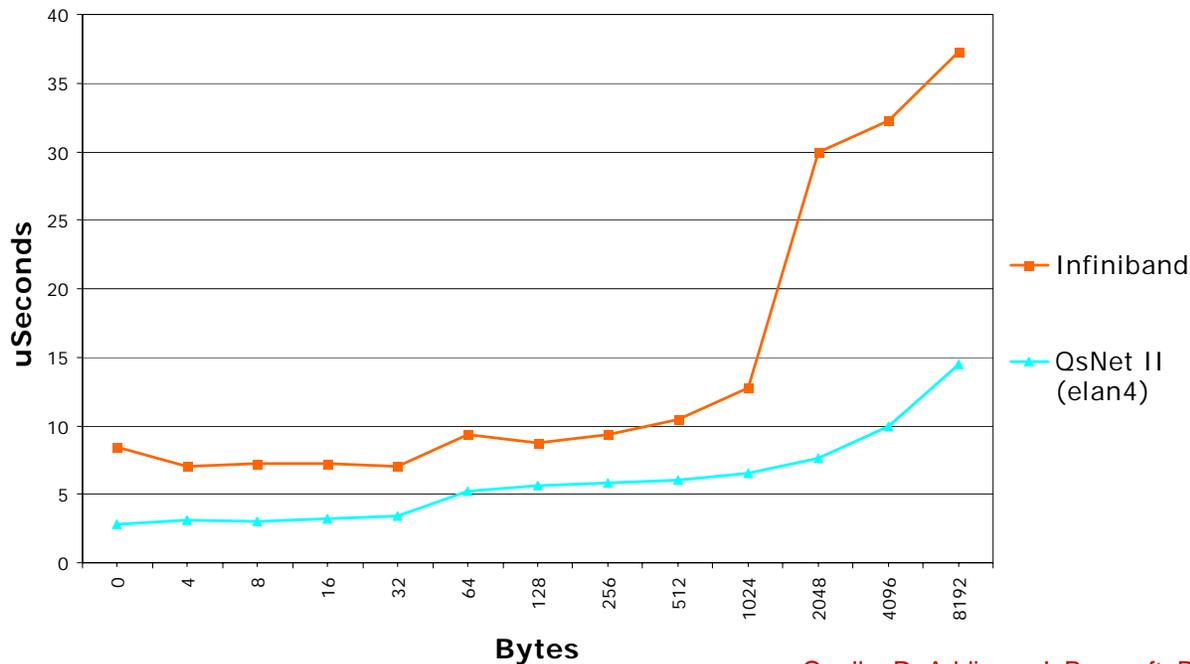


Quelle: D. Addison, J. Beecroft, D. Hewson, M. McLaren (Quadrics Ltd.),  
Fabrizio Petrini (LANL): A network for Supercomputing Applications. Hot Chips,  
August 2003, <http://www.c3.lanl.gov/~fabrizio/quadrics.html>

- Statische Verbindungsnetze
  - Beispiel: Quadrics QSnet

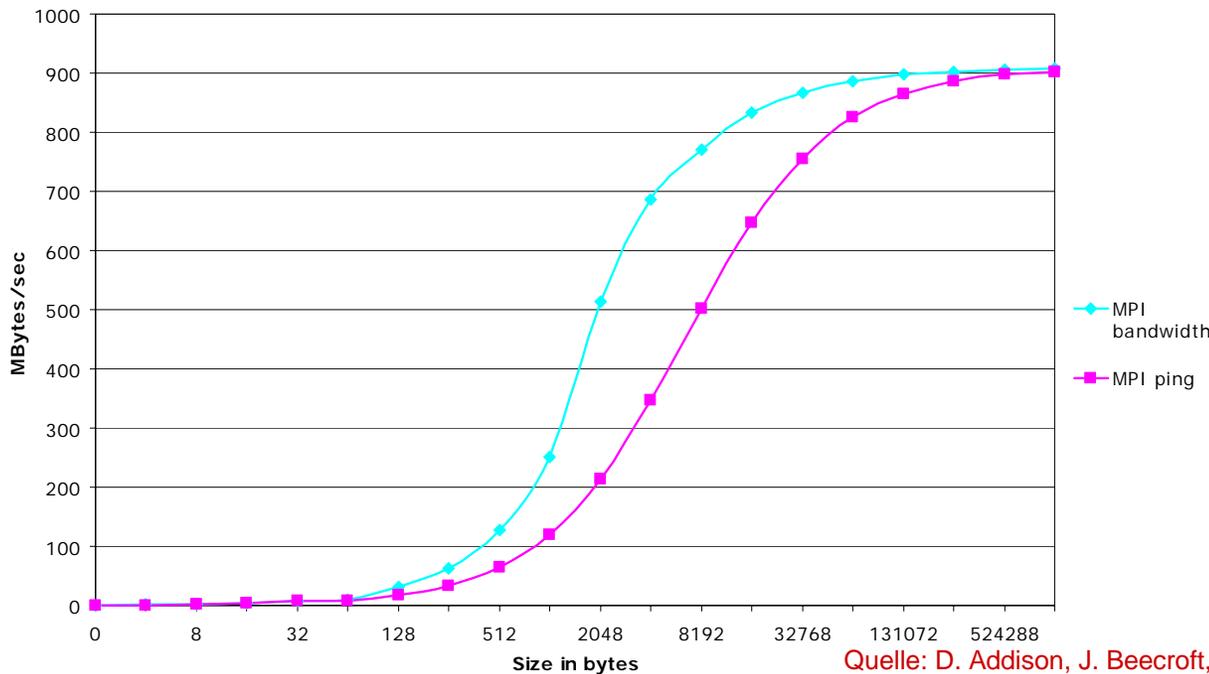


- **Statische Verbindungsnetze**
  - **Beispiel: Quadrics QSnet:**
    - MPI short message latency



Quelle: D. Addison, J. Beecroft, D. Hewson, M. McLaren (Quadrics Ltd.), Fabrizio Petrini (LANL): A network for Supercomputing Applications. Hot Chips, August 2003, <http://www.c3.lanl.gov/~fabrizio/quadrics.html>

- **Statische Verbindungsnetze**
  - **Beispiel: Quadrics QSnet:**
    - MPI Bandwidth



Quelle: D. Addison, J. Beecroft, D. Hewson, M. McLaren (Quadrics Ltd.),  
 Fabrizio Petrini (LANL): A network for Supercomputing Applications. Hot Chips,  
 August 2003, <http://www.c3.lanl.gov/~fabrizio/quadrics.html>

- **Statische Verbindungsnetze**

- **K-ärer n-Kubus (Cubes, Würfel)**

- Allgemeine Form eines Kubus-Verbindungsnetzwerkes
- Ringe, Gitter, oder Hyperkubi sind topologisch isomorph zu einer Familie von K-ären n-Kubus Netzwerken
  - n ist die Dimension
  - Radius K ist die Anzahl der Knoten, die einen Zyklus in einer Dimension bilden
- Enthält  $N=K^n$  Knoten
- Die Knoten werden über eine n-stellige binäre Radius K Zahl der Form  $a_0, a_1, \dots, a_{n-1}$  adressiert
  - Jede Stelle  $0 \leq a_i < K$  stellt die Position des Knotens in der entsprechenden i-ten Dimension dar, mit  $0 \leq i \leq n-1$
  - Ein Nachbarknoten in der i-ten Dimension zu einem Knoten mit Adresse  $a_0, a_1, \dots, a_{n-1}$  kann erreicht werden mit  $a_0, a_1, \dots, a_{(i \pm 1)} \bmod k \dots a_{n-1}$ .
- Knotengrad ist  $2n$  und der Diameter ist  $n \left\lceil \frac{k}{2} \right\rceil$

- **Statische Verbindungsnetzwerke:**
  - **Hyperkubus (Hypercubes)**
    - Verallgemeinerter Würfel:
      - die  $N = 2^n$  Prozessoren sind Ecken eines  $n$ -dimensionalen Würfels, wobei die Verbindungen dann die Kanten des Würfels darstellen.
    - Komplexität ist  $(N \cdot \log_2 N) / 2$ .
    - Diameter beträgt  $\log_2 N$ .
    - Lange Zeit häufigste Verbindungsstruktur bei den nachrichtengekoppelten Multiprozessoren, aber:
      - Skalierbarkeit:
        - » Jede Erweiterung benötigt mindestens die Verdopplung der Prozessorenanzahl.
        - » Aus räumlichen Anordnungsgründen begrenzt.

- **Statische Verbindungsnetzwerke:**
  - Hyperkubus
    - 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 (Hamming Distanz)
      - Eine einfache Wegewahl:  
die Bits in Start- und Zieladresse werden mittels einer XOR-Verbindung verknüpft und das Resultat bestimmt die möglichen Wege.

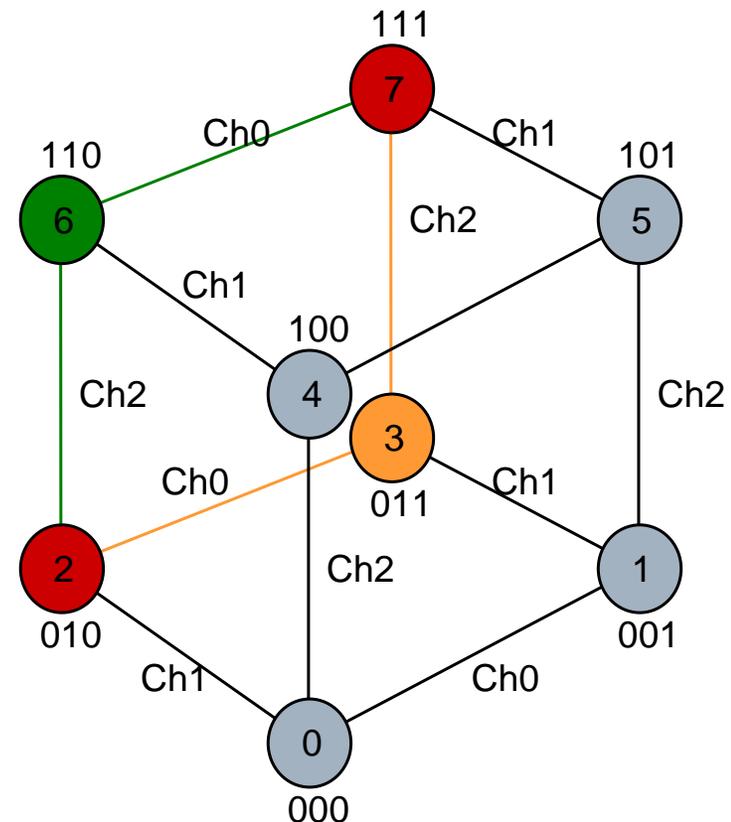
- Statische Verbindungsnetzwerke:
  - Hyperkubus
    - e-Cube Routing Algorithmus:
      - „messages are routed in increasingly higher dimensions of channels until the destination is reached“
      - Dimension eines Kanals = Bitposition von (Knoten# XOR Knoten#)

- **Statische Verbindungsnetzwerke:**

- Hyperkubus: e-cube Routing

- Beispiel:

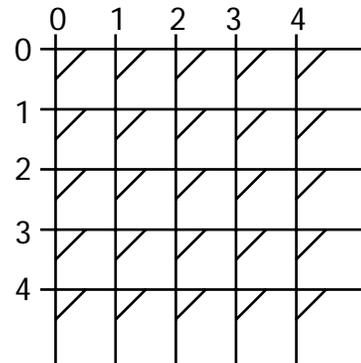
- $A = (010)$  und  $B = (111)$
- $W = A \text{ XOR } B = 101$
- $(010) \rightarrow (011) \rightarrow (111)$ ,
- $(010) \rightarrow (110) \rightarrow (111)$



- **Dynamische Verbindungsnetzwerke:**
  - Geeignet für Anwendungen mit variablen und nicht regulären Kommunikationsmustern
    - **Bus:**
      - Wird von den am Bus angeschlossenen Prozessoren gemeinsam benützt
      - Ein Datentransport zu einem Zeitpunkt
      - Nachricht von einer Quelle zu jedem Ziel in einem Schritt
      - Busbandbreite =  $w * f$ 
        - »  $w$ : Anzahl der Datenleitungen (Busbreite)
        - »  $F$ : Frequenz
        - » Bestimmt maximale Anzahl der Prozessoren, d. h. die Bandbreite muss mit dem Produkt der Anzahl der Prozessoren und ihrer Geschwindigkeit abgestimmt werden

- **Dynamische Verbindungsnetzwerke:**
  - **Bus:**
    - Reduzierung des Busverkehrs
      - » Verwendung von Cache-Speichern mit Cache-Kohärenz-Protokollen
    - Verwendung von sog. Split-Phase Busprotokollen
      - » Das Protokoll gibt den Bus nach der Übertragung einer Speichereferenzanforderung wieder frei
      - » Wenn der Speicher bereit ist, das Datum zu liefern, fordert dieser den Bus an und schickt die Daten als Antwort
      - » Ermöglicht, dass andere Prozessoren in der Zwischenzeit den Bus anfordern können, vorausgesetzt, dass ein verschränkter Speicher vorliegt oder Pipelining möglich ist

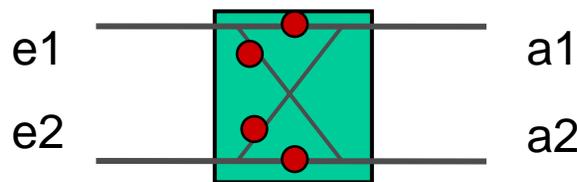
- **Dynamische Verbindungsnetzwerke:**
  - Kreuzschienenverteiler (Crossbar)
    - Vollständig vernetztes Verbindungswerk mit allen möglichen Permutationen der N Einheiten, die über das Netzwerk verbunden werden



- **Dynamische Verbindungsnetzwerke:**
  - **Kreuzschienenverteiler (Crossbar)**
    - Hardware-Einrichtung, die so geschaltet werden kann, dass in einer Menge von Prozessoren alle möglichen disjunkten Paare von Prozessoren gleichzeitig und blockierungsfrei miteinander kommunizieren können.
      - In Abhängigkeit vom Zustand der Schaltelemente im Kreuzschienenverteiler können dann je zwei beliebige Elemente aus den verschiedenen Mengen miteinander kommunizieren.
      - Alle  $N!$  Permutationen sind möglich
    - An den Kreuzungspunkten sitzen Schaltelemente: hoher Hardware-Aufwand
      - Kosten:  $N^2$  Schaltelemente (bei  $N$  Knoten pro Dimension)
      - Ein Schaltelement entspricht einem Paar von Quelle und Ziel, so dass die Darstellung einer Permutation als eine Liste solcher Paare direkt zu der korrekten Schaltung der Schalterelemente führt.

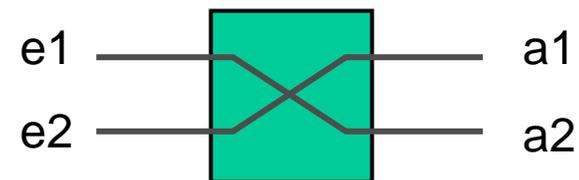
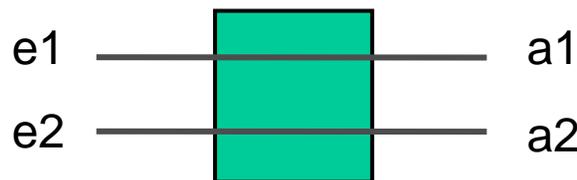
- **Dynamische Verbindungsnetzwerke:**
  - Schaltelemente (2x2 Kreuzschienenverteiler)
    - bestehen aus Zweierschaltern mit zwei Eingängen und zwei Ausgängen, die entweder durchschalten oder die Ein- und Ausgänge überkreuzen können

## → Schalernetzwerke



Durchschalten

● Kontakt, der geöffnet oder geschlossen werden kann



Vertauschen

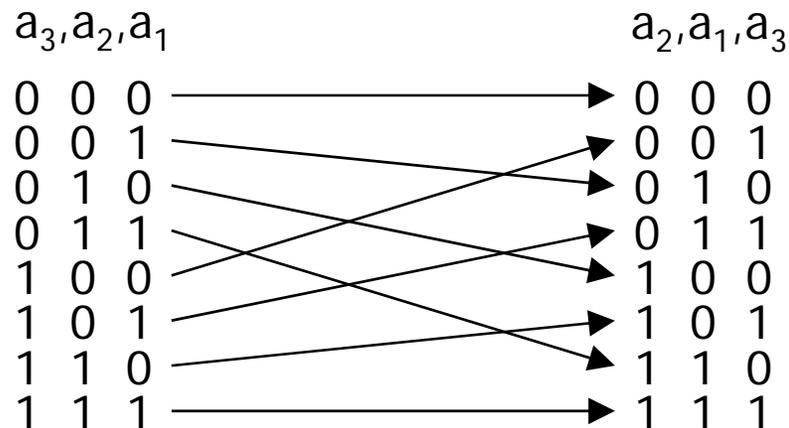
- **Dynamische Verbindungsnetzwerke:**
  - Mehrstufige Verbindungsnetzwerke (Schalternetzwerke, Permutationsnetzwerke)
    - Kompromiss zwischen der niedrigeren Leistungsfähigkeit von Bussen und hohem Hardware-Aufwand von Kreuzschienenverteilern
    - Oft 2 x 2 Kreuzschienenverteiler (Schalterelement) als Grundelement

- **Dynamische Verbindungsnetzwerke:**
  - **Permutationsnetze**
    - $p$  Eingänge des Netzes können gleichzeitig auf  $p$  Ausgänge geschaltet werden und somit wird eine Permutation der Eingänge erzeugt.
    - **Einstufige Permutationsnetze**
      - enthalten eine einzelne Spalte von Zweierschaltern,
    - **mehrstufige Permutationsnetze**
      - enthalten mehrere solcher Spalten
      - Spalten: Stufen des Permutationsnetzwerkes

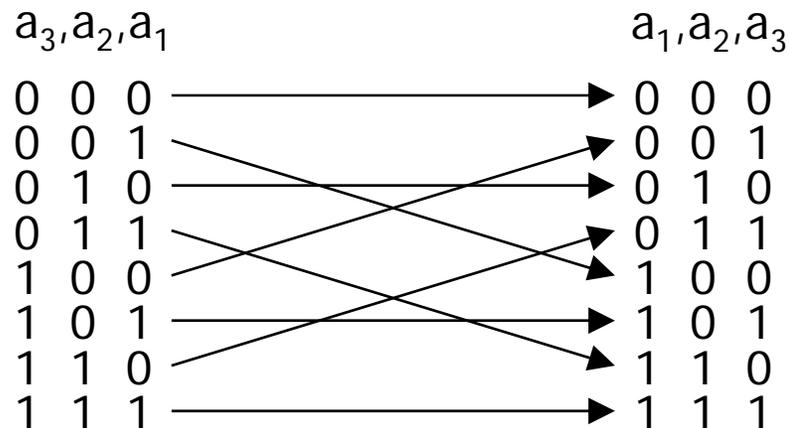
- **Dynamische Verbindungsnetzwerke:**
  - Permutationsnetze
    - reguläre Permutationsnetzwerke:
      - $p$  Eingänge,  $p$  Ausgänge und  $k$  Stufen mit jeweils  $p/2$  Zweierschaltern, wobei die Zahl  $p$  normalerweise eine Zweierpotenz ist.
    - Irreguläre Permutationsnetzwerke
      - weisen gegenüber der vollen regulären Struktur Lücken auf

- **Dynamische Verbindungsnetzwerke:**
  - **Permutationen**
    - eineindeutige (bijektive) Zuordnungen von Eingängen zu Ausgängen
    - Man stellt die Eingänge als binären Zahlenwert dar, d.h., man nummeriert die Eingänge beginnend mit 0 bis zum  $(2n-1)$ -ten Eingang durch.
    - Auf diese Weise ordnet man also jedem Eingang eine Art Adresse zu.
    - Die Permutation lässt sich nun durch eine Bitmanipulation dieser Adresse darstellen, so dass am Ausgang neue Bitmuster entstehen.

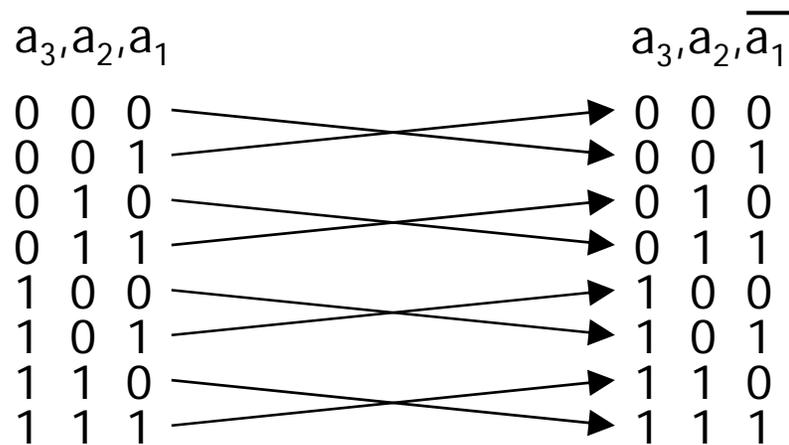
- **Dynamische Verbindungsnetzwerke:**
  - Permutationen
    - Mischpermutation M (Perfect Shuffle):
      - $M(a_n, a_{n-1}, \dots, a_2, a_1) = (a_{n-1}, \dots, a_2, a_1, a_n)$



- **Dynamische Verbindungsnetzwerke:**
  - Permutationen
    - Kreuzpermutation K (Butterfly):
      - $K(a_n, a_{n-1}, \dots, a_2, a_1) = (a_1, a_{n-1}, \dots, a_2, a_n)$

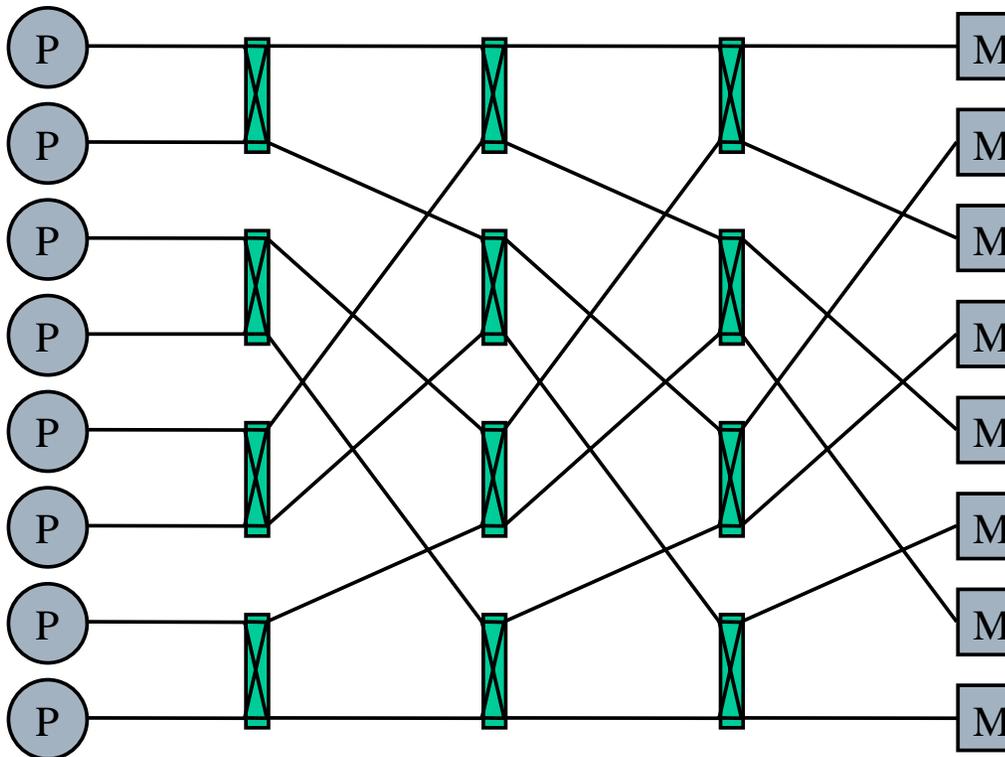


- **Dynamische Verbindungsnetzwerke:**
  - Permutationen
    - Tauschpermutation T (Butterfly):
      - Negation des niedrigwertigen Bits
      - $T(a_n, a_{n-1}, \dots, a_2, a_1) = (a_n, a_{n-1}, \dots, a_2, \bar{a}_1)$



- **Dynamische Verbindungsnetzwerke:**
  - **Omega-Netzwerk**
    - Das Netzwerk für  $p=2n$  Ein-/Ausgänge umfasst  $n = \lg p$  Stufen von Zweierschaltern, die untereinander jeweils nach dem Grundmuster der Mischpermutation verknüpft sind.
    - Gesamtzahl der Zweierschalter in einem Omega-Netzwerk mit  $p = 2n$  Ein-/Ausgängen beträgt  $(p/2) * \lg p$
    - Nicht blockierungsfrei

- **Dynamische Verbindungsnetzwerke:**
  - Speichergekoppeltes Omega-Netzwerk mit 8 Eingängen und 8 Ausgängen



# Speichergekoppeltes Switching-Banyan- Netzwerk mit 16 Ein- und 16 Ausgängen

