

Für eine lange Zeit...

- wurde die sequentielle Ausführung schneller (Instruction Level Parallelism, Pipelining, Höhere Frequenzen)
- mehr und kleinere Transistoren = mehr Performance
- Programmierer warteten auf die nächste schnellere Generation

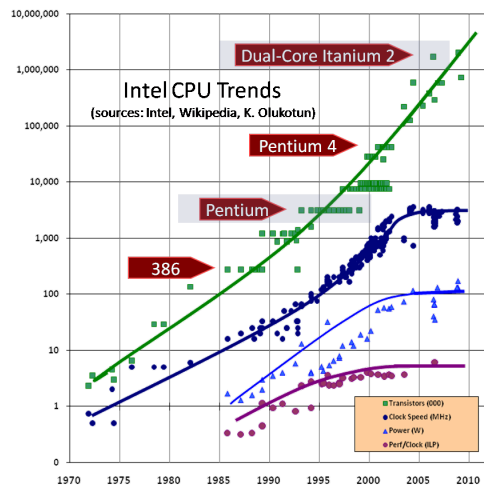
Heute

- steigt die Frequenz der Prozessoren kaum mehr an (Kühlproblematik)
- steigt die Instruction-Level Parallelität kaum mehr an
- ist die Ausführungsgeschwindigkeit in vielen Fällen dominiert von Speicherzugriffszeiten (Caches werden aber immer noch grösser und schneller)

775

776

Trends



<http://www.govt.ca/publications/concurrency-ddj.htm>

777

Multicore

- Verwende die Transistoren für mehr Rechenkerne
- Parallelität in der Software
- Implikation: Programmierer müssen parallele Programme schreiben, um die neue Hardware vollständig ausnutzen zu können

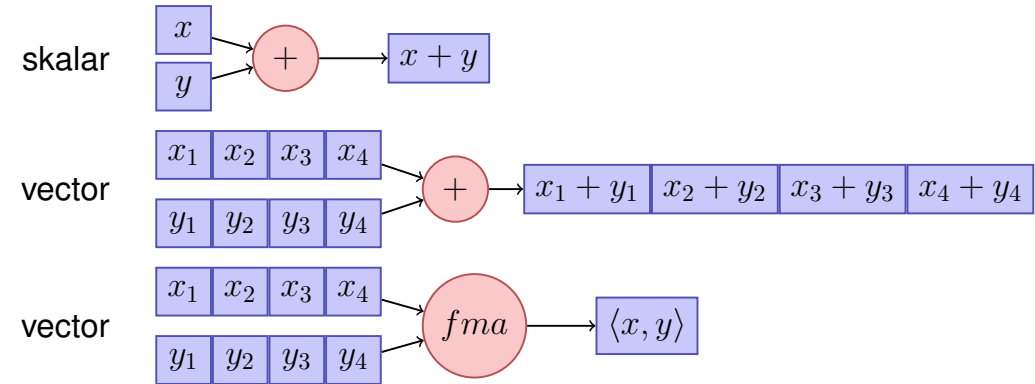
778

Formen der Parallelen Ausführung

- Vektorisierung
- Pipelining
- Instruction Level Parallelism
- Multicore / Multiprocessing
- Verteiltes Rechnen

Vektorisierung

Parallele Ausführung derselben Operation auf Elementen eines Vektor(Register)s



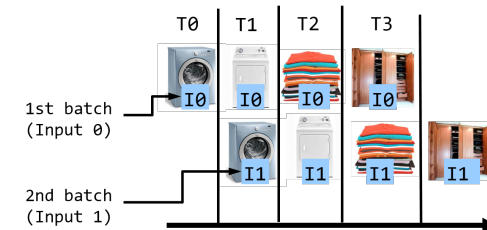
779

780

Hausarbeit



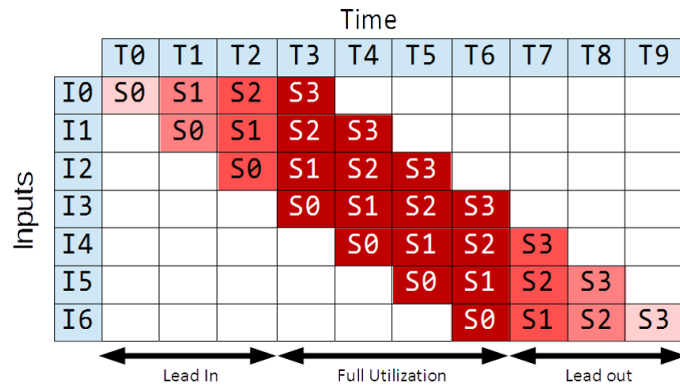
Effizienter



781

782

Pipeline



Throughput (Durchsatz)

- Throughput = Rate der ein- oder ausgehenden Daten
- Anzahl Operationen pro Zeiteinheit
- Je grösser, desto besser
- Approximation

$$\text{throughput} = \frac{1}{\max(\text{Berechnungszeit}(\text{Stufen}))}$$

ignoriert lead-in und lead-out Zeiten

783

784

Latenz

- Zeit zum Ausführen einer Berechnung
- Pipeline-Latenz ist nur konstant, wenn die Pipeline balanciert ist: Summe aller Operationen über die Stufen
- Unbalancierte Pipeline
 - Erster Durchlauf wie bei der balancierten Pipeline
 - Balancierte Version, Latenz = #stufen · max(Berechnungszeit(Stufen))

Beispiel Hausarbeit

Waschen $T_0 = 1h$, Trocknen $T_1 = 2h$, Bügeln $T_2 = 1h$, Versorgen $T_3 = 0.5h$

- Latenz Erster Durchlauf: $L = T_0 + T_1 + T_2 + T_3 = 4.5h$
- Latenz Zweiter Durchlauf: $L = T_1 + T_1 + T_2 + T_3 = 5.5h$
- Langfristiger Durchsatz: 1 Ladung alle $2h$ ($0.5/h$).

785

786

Throughput vs. Latency

- Erhöhen des Throughputs kann Latenz erhöhen
- Stufen der Pipeline müssen kommunizieren und synchronisieren: Overhead

787

Pipelines in CPUs

Fetch

Decode

Execute

Data Fetch

Writeback

Mehrere Stufen

- Jede Instruktion dauert 5 Zeiteinheiten (Zyklen)
- Im besten Fall: 1 Instruktion pro Zyklus, nicht immer möglich ("stalls")

Parallelität (mehrere funktionale Einheiten) führt zu *schnellerer Ausführung*.

788

ILP – Instruction Level Parallelism

Moderne CPUs führen unabhängige Instruktionen intern auf mehreren Einheiten parallel aus

- Pipelining
- Superscalar CPUs (multiple instructions per cycle)
- Out-Of-Order Execution (Programmer observes the sequential execution)
- Speculative Execution

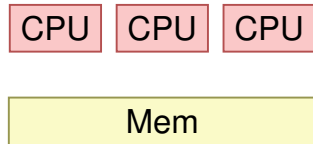
789

27.2 Hardware Architekturen

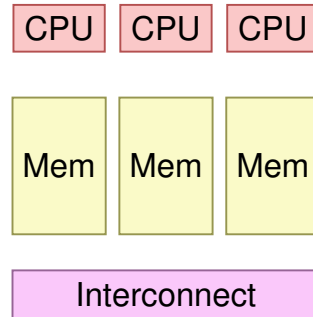
790

Gemeinsamer vs. verteilter Speicher

Gemeinsamer Speicher



Verteilter Speicher



791

Shared vs. Distributed Memory Programming

■ Kategorien des Programmierinterfaces

- Kommunikation via Message Passing
- Kommunikation via geteiltem Speicher

■ Es ist möglich:

- Systeme mit gemeinsamen Speicher als verteilte Systeme zu programmieren (z.B. mit Message Passing Interface MPI)
- Systeme mit verteiltem Speicher als System mit gemeinsamen Speicher zu programmieren (z.B. Partitioned Global Address Space PGAS)

792

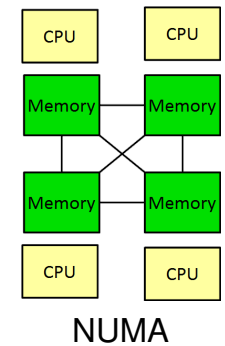
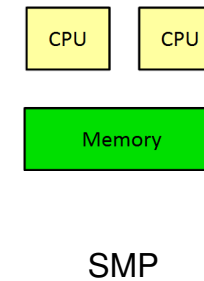
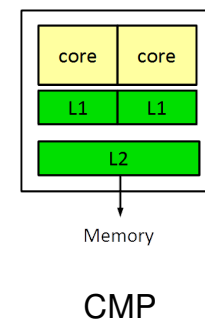
Architekturen mit gemeinsamen Speicher

- Multicore (Chip Multiprocessor - CMP)
- Symmetric Multiprocessor Systems (SMP)
- Simultaneous Multithreading (SMT = Hyperthreading)
 - nur ein physischer Kern, Mehrere Instruktionsströme/Threads: mehrere virtuelle Kerne
 - Zwischen ILP (mehrere Units für einen Strom) und Multicore (mehrere Units für mehrere Ströme). Limitierte parallele Performance
- Non-Uniform Memory Access (NUMA)

Gleiches Programmierinterface!

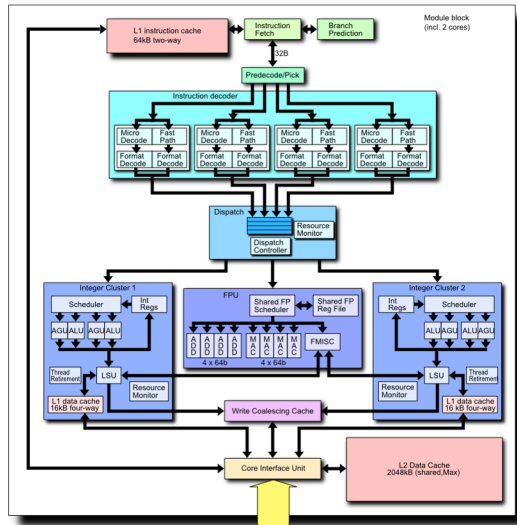
793

Übersicht



794

Ein Beispiel

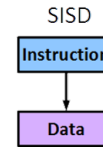


AMD Bulldozer: Zwischen CMP und SMT

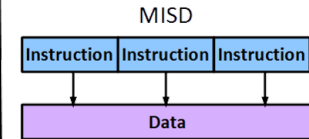
- 2x integer core
- 1x floating point core

Klassifikation nach Flynn

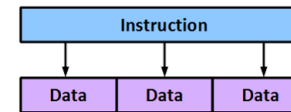
Single-Core



Fault-Tolerance

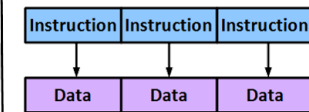


SIMD



Vector Computing / GPU

MIMD



Multi-Core

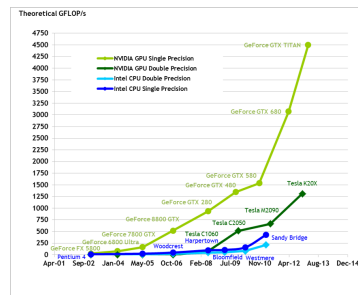
Wikipedia
795

796

Massiv Parallele Hardware

[General Purpose] Graphical Processing Units ([GP]GPUs)

- Revolution im High Performance Computing
 - Calculation 4.5 TFlops vs. 500 GFlops
 - Memory Bandwidth 170 GB/s vs. 40 GB/s
- SIMD
 - Hohe Datenparallelität
 - Benötigt eigenes Programmiermodell. Z.B. CUDA / OpenCL



27.3 Multi-Threading, Parallelität und Nebenläufigkeit

797

798

Prozesse und Threads

- Prozess: Instanz eines Programmes
 - jeder Prozess hat seinen eigenen Kontext, sogar eigenen Adressraum
 - OS verwaltet Prozesse (Ressourcenkontrolle, Scheduling, Synchronisierung)
- Threads: Ausführungsfäden eines Programmes
 - Threads teilen sich einen Adressraum
 - Schneller Kontextwechsel zwischen Threads

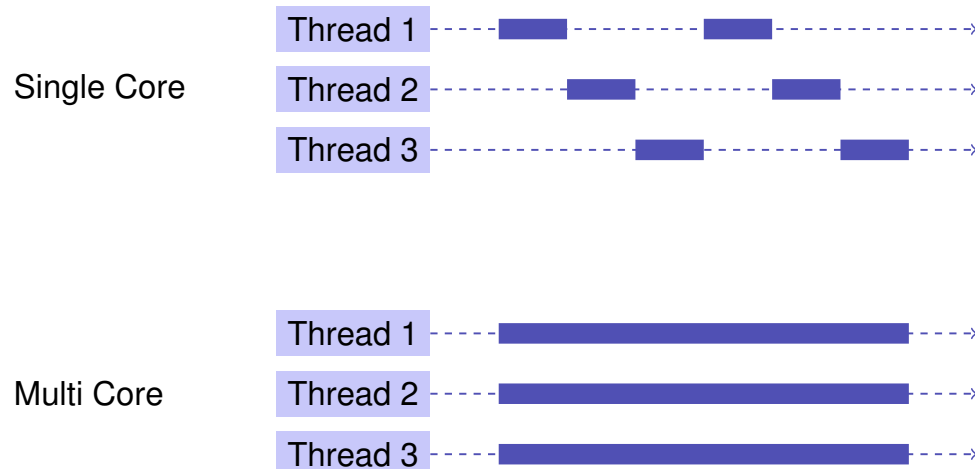
799

Warum Multithreading?

- Verhinderung vom "Polling" auf Ressourcen (Files, Netzwerkzugriff, Tastatur)
- Interaktivität (z.B. Responsivität von GUI Programmen)
- Mehrere Applikationen / Clients gleichzeitig instanzierbar
- Parallelität (Performanz!)

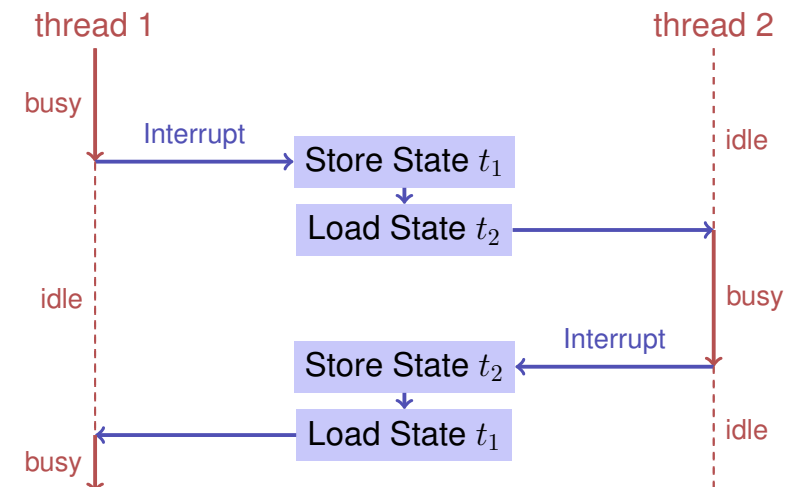
800

Multithreading konzeptuell



801

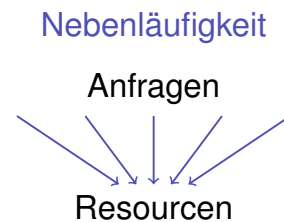
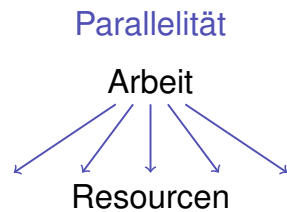
Threadwechsel auf einem Core (Preemption)



802

Parallelität vs. Nebenläufigkeit (Concurrency)

- **Parallelität:** Verwende zusätzliche Ressourcen (z.B. CPUs), um ein Problem schneller zu lösen
- **Nebenläufigkeit:** Verwalte gemeinsam genutzte Ressourcen (z.B. Speicher) korrekt und effizient
- Begriffe überlappen offensichtlich. Bei parallelen Berechnungen besteht fast immer Synchronisierungsbedarf.



803

Thread-Sicherheit

Thread-Sicherheit bedeutet, dass in der nebenläufigen Anwendung eines Programmes dieses sich immer wie gefordert verhält.

Viele Optimierungen (Hardware, Compiler) sind darauf ausgerichtet, dass sich ein *sequentielles* Programm korrekt verhält.

Nebenläufige Programme benötigen für ihre Synchronisierungen auch eine Annotation, welche gewisse Optimierungen selektiv abschaltet

804

Beispiel: Caches

- Speicherzugriff auf Register schneller als auf den gemeinsamen Speicher
- Prinzip der Lokalität
- Verwendung von Caches (transparent für den Programmierer)

Ob und wie weit die Cache-Kohärenz sichergestellt wird ist vom eingesetzten System abhängig.



805

27.4 Skalierbarkeit: Amdahl und Gustafson

806

Skalierbarkeit

In der parallelen Programmierung:

- Geschwindigkeitssteigerung bei wachsender Anzahl p Prozessoren
- Was passiert, wenn $p \rightarrow \infty$?
- Linear skalierendes Programm: Linearer Speedup

Parallele Performanz

Gegeben fixierte Rechenarbeit W (Anzahl Rechenschritte)

Sequentielle Ausführungszeit sei T_1

Parallele Ausführungszeit T_p auf p CPUs

- Perfektion: $T_p = T_1/p$
- Performanzverlust: $T_p > T_1/p$ (üblicher Fall)
- Hexerei: $T_p < T_1/p$

807

808

Paralleler Speedup

Paralleler Speedup S_p auf p CPUs:

$$S_p = \frac{W/T_p}{W/T_1} = \frac{T_1}{T_p}$$

- Perfektion: Linearer Speedup $S_p = p$
- Verlust: sublinearer Speedup $T_p > T_1/p$ (der übliche Fall)
- Hexerei: superlinearer Speedup $T_p < T_1/p$

Effizienz: $E_p = S_p/p$

Erreichbarer Speedup?

Paralleles Programm

Paralleler Teil	Seq. Teil
80%	20%

$$T_1 = 10$$

$$T_8 = ?$$

$$T_8 = \frac{10 \cdot 0.8}{8} + 10 \cdot 0.2 = 1 + 2 = 3$$

$$S_8 = \frac{T_1}{T_8} = \frac{10}{3} = 3.33$$

809

810

Amdahl's Law: Zutaten

Zu Leistende Rechenarbeit W fällt in zwei Kategorien

- Parallelisierbarer Teil W_p
- Nicht parallelisierbarer, sequentieller Teil W_s

Annahme: W kann mit einem Prozessor in W Zeiteinheiten sequentiell erledigt werden ($T_1 = W$):

$$T_1 = W_s + W_p$$
$$T_p \geq W_s + W_p/p$$

Amdahl's Law

$$S_p = \frac{T_1}{T_p} \leq \frac{W_s + W_p}{W_s + \frac{W_p}{p}}$$

811

812

Amdahl's Law

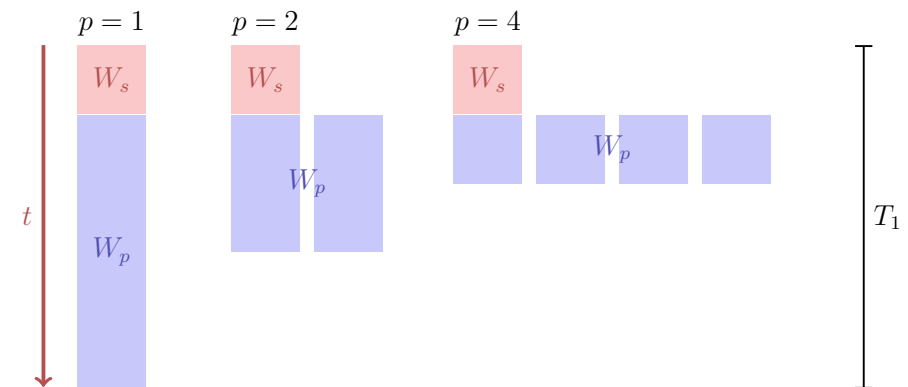
Mit seriellem, nicht parallelisierbarem Anteil λ : $W_s = \lambda W$,
 $W_p = (1 - \lambda)W$:

$$S_p \leq \frac{1}{\lambda + \frac{1-\lambda}{p}}$$

Somit

$$S_\infty \leq \frac{1}{\lambda}$$

Illustration Amdahl's Law



813

814

Amdahl's Law ist keine gute Nachricht

Alle nicht parallelisierbaren Teile können Problem bereiten und stehen der Skalierbarkeit entgegen.

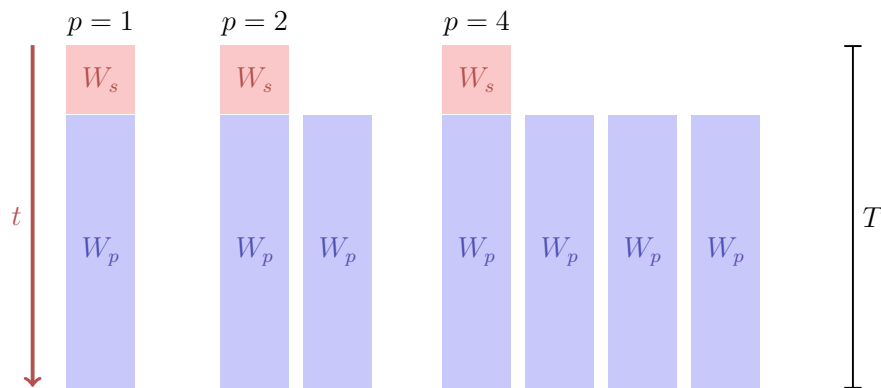
Gustafson's Law

- Halte die Ausführungszeit fest.
- Variiere die Problemgröße.
- Annahme: Der sequentielle Teil bleibt konstant, der parallele Teil wird grösser.

815

816

Illustration Gustafson's Law



817

Gustafson's Law

Arbeit, die mit einem Prozessor in der Zeit T erledigt werden kann:

$$W_s + W_p = T$$

Arbeit, die mit p Prozessoren in der Zeit T erledigt werden kann:

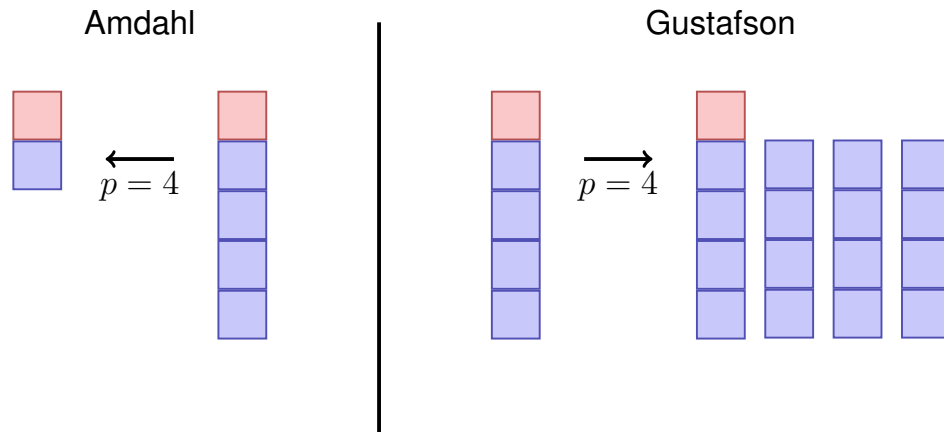
$$W_s + p \cdot W_p = \lambda \cdot T + p \cdot (1 - \lambda) \cdot T$$

Speedup:

$$\begin{aligned} S_p &= \frac{W_s + p \cdot W_p}{W_s + W_p} = p \cdot (1 - \lambda) + \lambda \\ &= p - \lambda(p - 1) \end{aligned}$$

818

Amdahl vs. Gustafson



27.5 Task- und Datenparallelität

819

820

Paradigmen der Parallelen Programmierung

- **Task Parallel:** Programmierer legt parallele Tasks explizit fest.
- **Daten-Parallel:** Operationen gleichzeitig auf einer Menge von individuellen Datenobjekten.

Beispiel Data Parallel (OMP)

```
double sum = 0, A[MAX];  
#pragma omp parallel for reduction (+:ave)  
for (int i = 0; i < MAX; ++i)  
    sum += A[i];  
return sum;
```

821

822

Beispiel Task Parallel (C++11 Threads/Futures)

```
double sum(Iterator from, Iterator to)
{
    auto len = from - to;
    if (len > threshold){
        auto future = std::async(sum, from, from + len / 2);
        return sumS(from + len / 2, to) + future.get();
    }
    else
        return sumS(from, to);
}
```

823

Partitionierung und Scheduling

- Aufteilung der Arbeit in parallele Tasks (Programmierer oder System)
 - Ein Task ist eine Arbeitseinheit
 - Frage: welche Granularität?
- Scheduling (Laufzeitsystem)
 - Zuweisung der Tasks zu Prozessoren
 - Ziel: volle Ressourcennutzung bei wenig Overhead

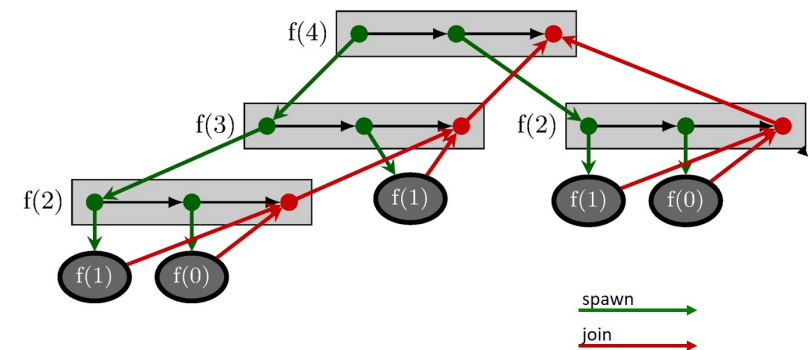
824

Beispiel: Fibonacci P-Fib

```
if  $n \leq 1$  then
    return  $n$ 
else
     $x \leftarrow$  spawn P-Fib( $n - 1$ )
     $y \leftarrow$  spawn P-Fib( $n - 2$ )
    sync
    return  $x + y$ ;
```

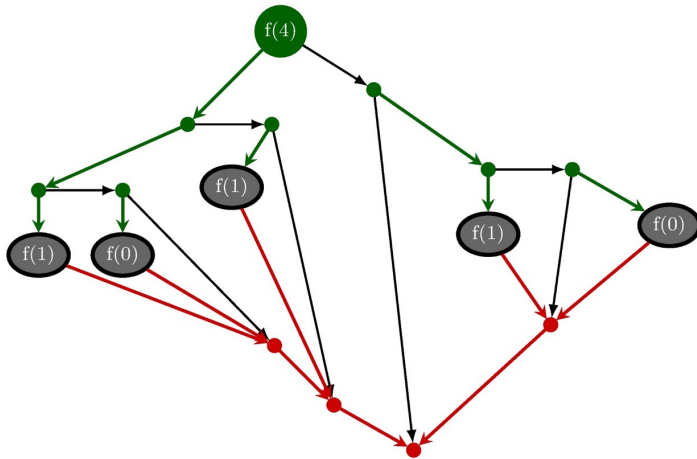
825

P-Fib Task Graph



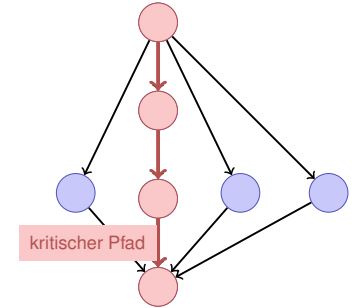
826

P-Fib Task Graph



Frage

- Jeder Knoten (Task) benötigt 1 Zeiteinheit.
- Pfeile bezeichnen Abhängigkeiten.
- Minimale Ausführungseinheit wenn Anzahl Prozessoren = ∞ ?

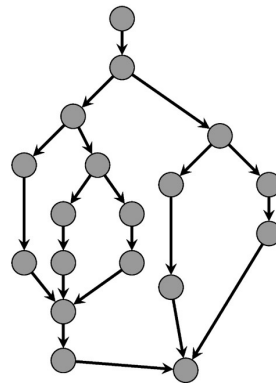


827

828

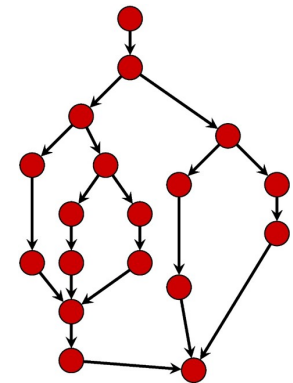
Performanzmodell

- p Prozessoren
- Dynamische Zuteilung
- T_p : Ausführungszeit auf p Prozessoren



Performanzmodell

- T_p : Ausführungszeit auf p Prozessoren
- T_1 : *Arbeit*: Zeit für die gesamte Berechnung auf einem Prozessor
- T_1/T_p : Speedup



829

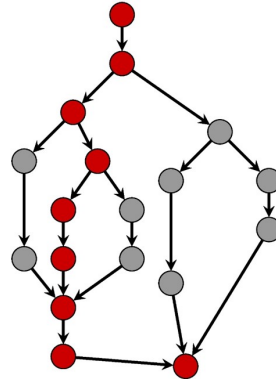
830

Performanzmodell

- T_∞ : **Zeitspanne**: Kritischer Pfad. Ausführungszeit auf ∞ Prozessoren. Längster Pfad von der Wurzel zur Senke.
- T_1/T_∞ : **Parallelität**: breiter ist besser
- Untere Grenzen

$$T_p \geq T_1/p \quad \text{Arbeitsgesetz}$$

$$T_p \geq T_\infty \quad \text{Zeitspannengesetz}$$



Greedy Scheduler

Greedy Scheduler: teilt zu jeder Zeit so viele Tasks zu Prozessoren zu wie möglich.

Theorem

Auf einem idealen Parallelrechner mit p Prozessoren führt ein Greedy-Scheduler eine mehrfädige Berechnung mit Arbeit T_1 und Zeitspanne T_∞ in Zeit

$$T_p \leq T_1/p + T_\infty$$

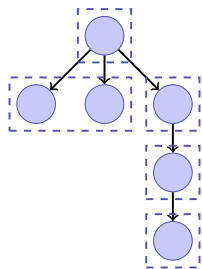
aus.

831

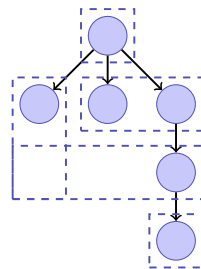
832

Beispiel

Annahme $p = 2$.



$$T_p = 5$$



$$T_p = 4$$

Beweis des Theorems

Annahme, dass alle Tasks gleich viel Arbeit aufweisen.

- Vollständiger Schritt: p Tasks stehen zur Berechnung bereit
- Unvollständiger Schritt: weniger als p Tasks bereit.

Annahme: Anzahl vollständige Schritte grösser als $\lfloor T_1/p \rfloor$.
 Ausgeführte Arbeit $\geq P \cdot (\lfloor T_1/p \rfloor \cdot p) = T_1 - T_1 \bmod p + p \geq T_1$.
 Widerspruch. Also maximal $\lfloor T_1/p \rfloor$ vollständige Schritte.

Jeder unvollständige Schritt führt zu jedem Zeitpunkt alle vorhandenen Tasks t mit $\deg^-(t) = 0$ aus und verringert die Länge der Zeitspanne. Andernfalls wäre die gewählte Zeitspanne nicht maximal. Anzahl unvollständige Schritte also maximal T_∞ .

833

834

Granularität: Wie viele Tasks?

Antwort: so viele Tasks wie möglich mit sequentiellem Cut-off, welcher den Overhead vernachlässigen lässt.

Beispiel: Parallelität von Mergesort

- Arbeit (sequentielle Laufzeit) von Mergesort $T_1(n) = \Theta(n \log n)$.
- Span $T_\infty(n) = \Theta(n)$
- Parallelität $T_1(n)/T_\infty(n) = \Theta(\log n)$ (Maximal erreichbarer Speedup mit $p = \infty$ Prozessoren)

