

Effiziente Nutzung von Hochleistungsrechnern in der numerischen Strömungsmechanik

Dr. Georg Hager

georg.hager@rrze.uni-erlangen.de

HPC Services
Regionales RechenZentrum Erlangen

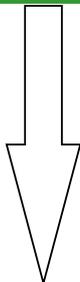


- Moderne Prozessoren für numerische Anwendungen
 - Allgemeine Optimierungsrichtlinien für CFD
- Parallelrechner
 - Möglichkeiten und Grenzen der Parallelität
 - Designprinzipien
 - effiziente Programmierung für CFD-Probleme
- Beispiel: SIPSolver nach Stone (Finite-Volumen)
- Zusammenfassung

Moderne Prozessoren für numerische Anwendungen

Vektorprozessor

Haupt-
speicher

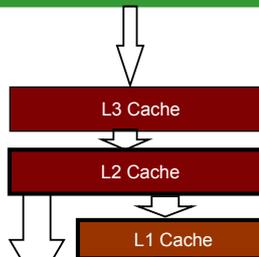


Vektorregister

Rechen-
werke

Klass. RISC Prozessor

Hauptspeicher

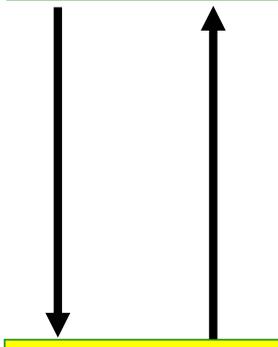


Floating Point
Register

Rechenwerke

CFD-Programm

Daten



Programm

Kurzlehrgang NUMET	Kennzahlen von Prozessoren	5 / 11
<ul style="list-style-type: none"> • Rechenleistung in Fließkommaoperationen pro Sekunde: Flop/s <ul style="list-style-type: none"> – i.A. Multiplikationen und Additionen mit doppelter Genauigkeit – Divisionen, SQRT etc. sind sehr langsam – max. FLOP/s-Zahl ist üblicherweise das Doppelte oder Vierfache der Prozessor-Taktfrequenz • Bandbreite (BW), d.h. Geschwindigkeit des Datentransfers: GB/s <ul style="list-style-type: none"> – in CFD-Anwendungen meist der limitierende Faktor: Bandbreite zwischen CPU und Hauptspeicher • Latenz: Zeitdauer vom Anstoßen eines Datentransfers bis zum Eintreffen des ersten Bytes <ul style="list-style-type: none"> – CPU↔CPU, CPU↔Speicher, Rechenwerke↔Cache – dominant, wenn pro Transfer nur wenig Daten übertragen werden 		

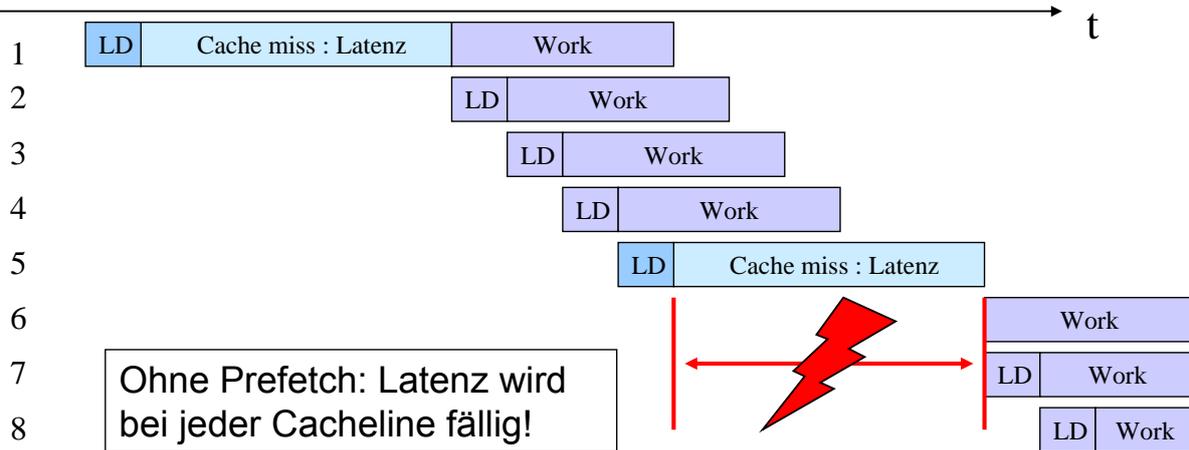
Kurzlehrgang NUMET	Leistungszahlen moderner Prozessoren				6 / 11
	Intel Xeon	AMD Opteron	Intel Itanium2	NEC SX8	
Spitzenleistung Taktfreq.	6.4 GFlop/s 3.2 GHz	4.8 GFlop/s 2.4 GHz	4.0 GFlop/s 1.0 GHz	16 GFlop/s 2.0 GHz	
# FP Registers	16/32	16/32	128	8 x 256 (Vector)	
L1	Size	16 kB	64 KB	16 KB (kein FP)	---
	BW	102 GB/s	51.2 GB/s	32 GB/s	---
	Latenz	12 Takte	3 Takte	1 Takt	---
L2	Size	1.0 MB	1.0 MB	256 KB	---
	BW	102 GB/s	51.2 GB/s	32 GB/s	---
	Latenz	18 Takte	13 Takte	5-6 Takte	---
L3	Size	---	---	3 MB	---
	BW	---	---	32 GB/s	---
	Latenz	---	---	12-13 Takte	---
Mem.	BW	6.4 GB/s r/w	6.4 GB/s r/w	6.4 GB/s r/w	64 GB/s r/w
	Latenz	~150 ns	~80 ns	~200 ns	Vektorisierung

Kurzlehrgang NUMET	Organisation von Caches	7 / 11
<ul style="list-style-type: none"> • Caches bestehen aus Cachelines, die nur als Ganzes gelesen/geschrieben werden <ul style="list-style-type: none"> – typische Längen: 8/16/32 Worte (je 8 Byte) • Cache Miss: Verwenden eines Datums, das noch nicht im Cache liegt <ul style="list-style-type: none"> – Folge: Cacheline wird aus Memory geladen (kostet Latenz + Zeit für Datentransfer) – alle weiteren Zugriffe auf Cacheline sind schnell • Cachelines werden nach einem bestimmten Algorithmus wieder aus dem Cache entfernt (zurückgeschrieben und/oder invalidiert) • Worst Case: Laden einer Cacheline, von der nur ein einziges Element benutzt wird <ul style="list-style-type: none"> – tritt auf bei irregulärem Speicherzugriff (z.B. indirekt) 		

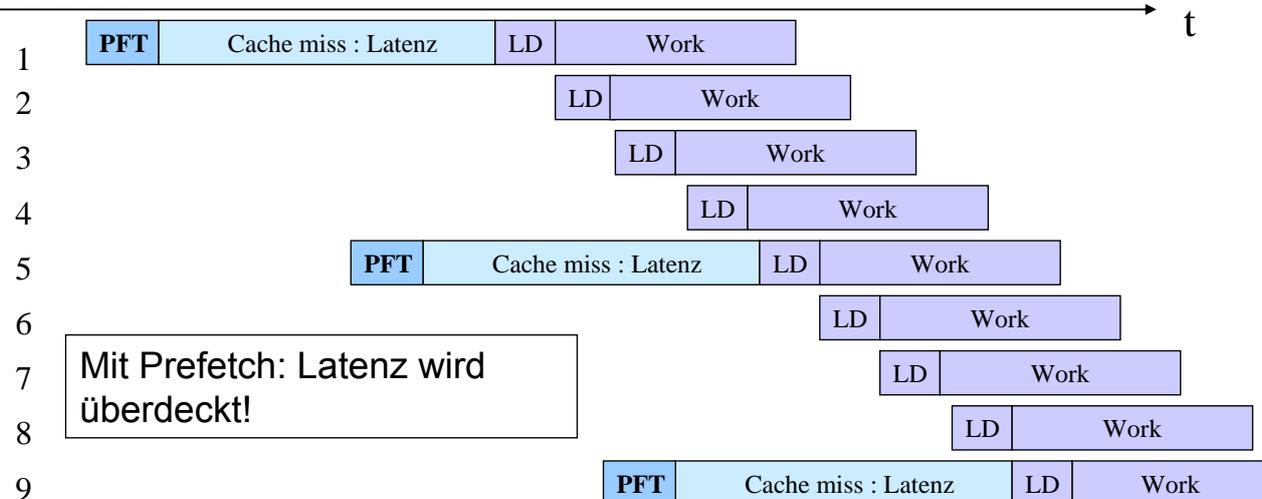
Kurzlehrgang NUMET	Optimierung für Cache (1)	8 / 11
<ul style="list-style-type: none"> • Caches sind schnell und klein, Speicher ist groß und langsam! • Folge: Langsame Datenpfade sind zu vermeiden! 		
<pre> DO I=1,N A(I)=B(I)+C(I) ENDDO DO I=1,N A(I)=A(I)*D(I) ENDDO </pre>		
<p style="text-align: center;">  Optimierung </p>		
<pre> DO I=1,N A(I)=(B(I)+C(I))*D(I) ENDDO </pre>		
		
Effiziente Nutzung schneller Datenpfade ist absolut vorrangig!		

- Wenn sich "Streaming" nicht vermeiden lässt (üblich in CFD), muss es wenigstens effizient geschehen: **Prefetching!**
 - **Prefetch** = Überdecken der Speicherlatenz durch rechtzeitiges Absetzen von Anfragen für Cachelines

Iteration



Iteration



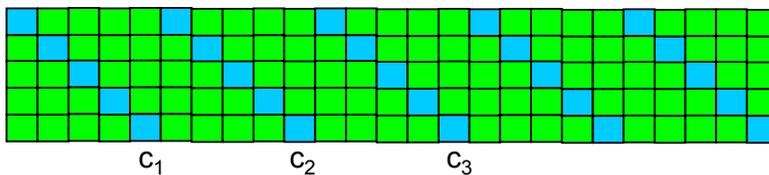
- Prefetching wird i.A. vom Compiler generiert (**überprüfen!**)
 - eingreifen u.U. im Quellcode möglich durch Compilerdirektiven
- Xeon, Opteron, Power4/5: zusätzlich **hardwarebasierter Prefetch**

- Analog RISC: Pipelining ist das A und O
- Aufteilung einer komplexen Operation (hier Subtraktion) in:

- compare exponent
- shift mantissa
- add mantissa
- normalize exponent

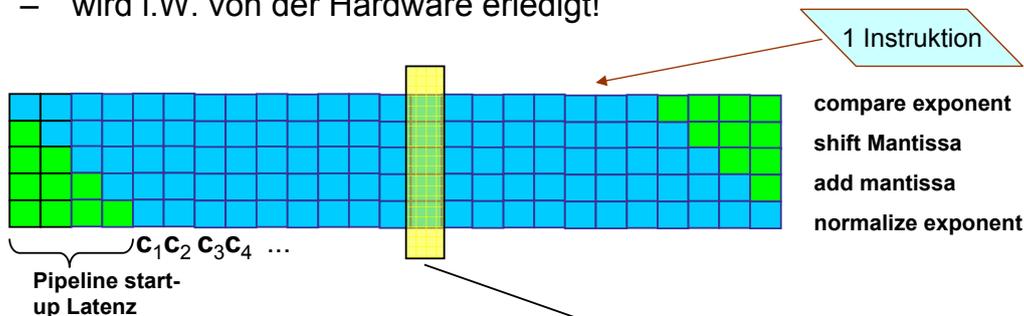
$$\begin{array}{r}
 1.77e4 - 9.3e3 \\
 1.77e4 - 0.93e4 \\
 0.84e4 \\
 8.40e3
 \end{array}$$

- Ablaufschema für Array-Operation ($C = A + B$)
- hier: ohne Pipelining



compare exponent
shift Mantissa
add mantissa
normalize exponent

- Wie kann man das verbessern?
- Pipelining: von 1 Ergebnis / 5 Takte zu einem Ergebnis / 1 Takt
- wird i.W. von der Hardware erledigt!

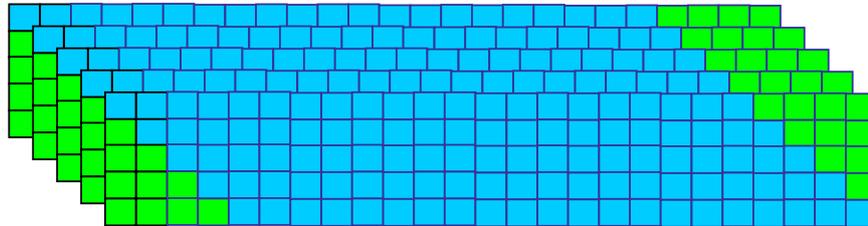


gleichzeitige Arbeit an
Elementen
14, 13, 12, 11, 10

- Vektor-CPU's haben die notwendige Speicherbandbreite, um Pipelines ohne Verzögerung zu "füttern"
- Prefetch erfolgt automatisch

Kurzlehrgang NUMET	Prinzipien der Vektorarchitektur (3): Multi-Track Pipelines	13 / 11
-----------------------	--	---------

- Noch besser:
 - mehrere parallele Pipelines
 - N_p "Tracks" $\rightarrow N_p$ Ergebnisse pro Takt



- Voraussetzungen für gute Performance:
 - Pipes beschäftigt halten!
 - Compiler muss Vektorstrukturen im Code erkennen können
 - Unabhängigkeit aufeinander folgender Operationen

Kurzlehrgang NUMET	Abschätzung von Performancezahlen	14 / 11
-----------------------	-----------------------------------	---------

- Woher weiß man, ob ein Code die Ressourcen des Systems effizient nutzt?
- In vielen Fällen ist eine Abschätzung der zu erwartenden Performance aus den Eckdaten der Architektur und der Codestruktur möglich

- Randbedingungen der Architektur (maschinenabhängig):

Speicherbandbreite GWorte/s
 Fließkommaleistung GFlop/s

daraus abgeleitet:

Maschinenbalance

$$B_m = \frac{\text{Speicherbandbreite [Worte / s]}}{\text{Fließkommaleistung [Flop / s]}}$$

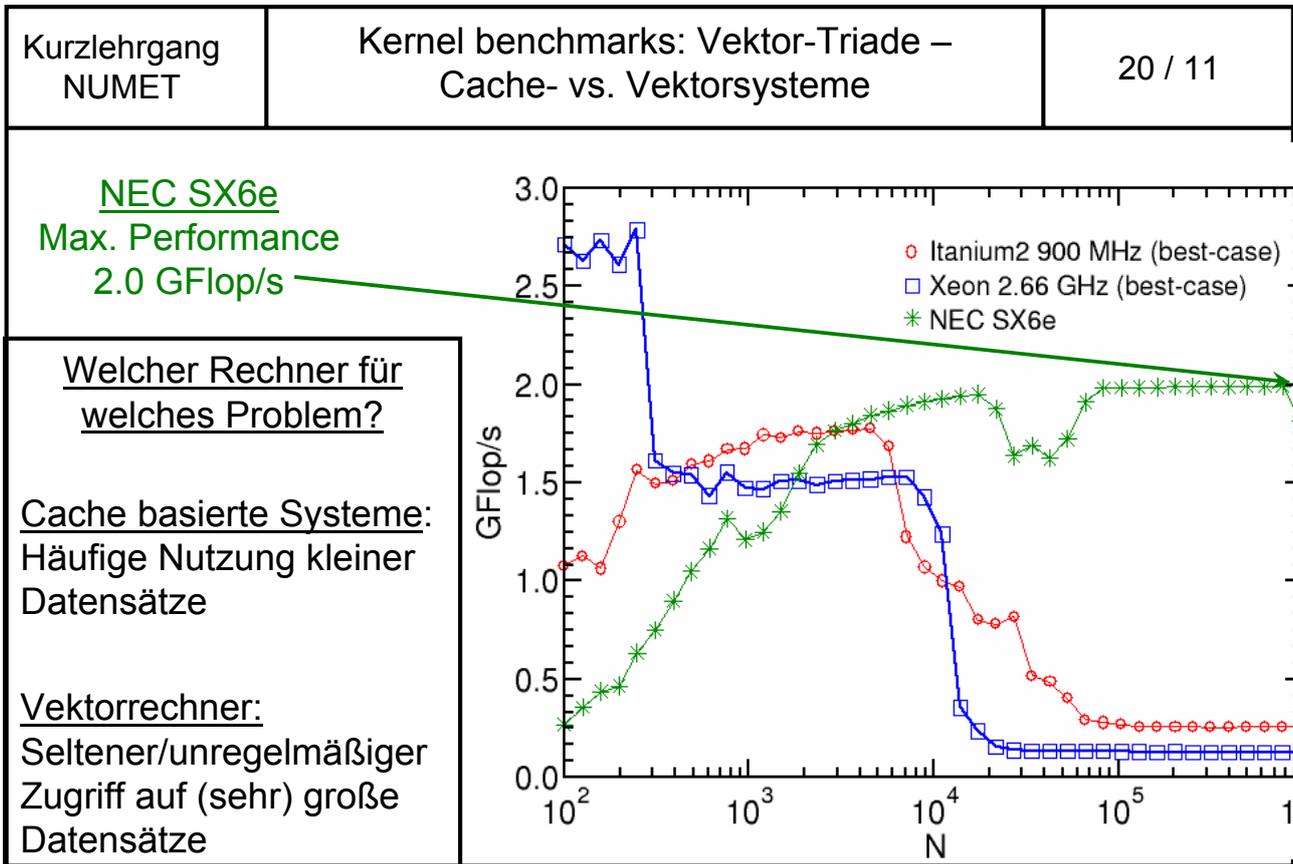
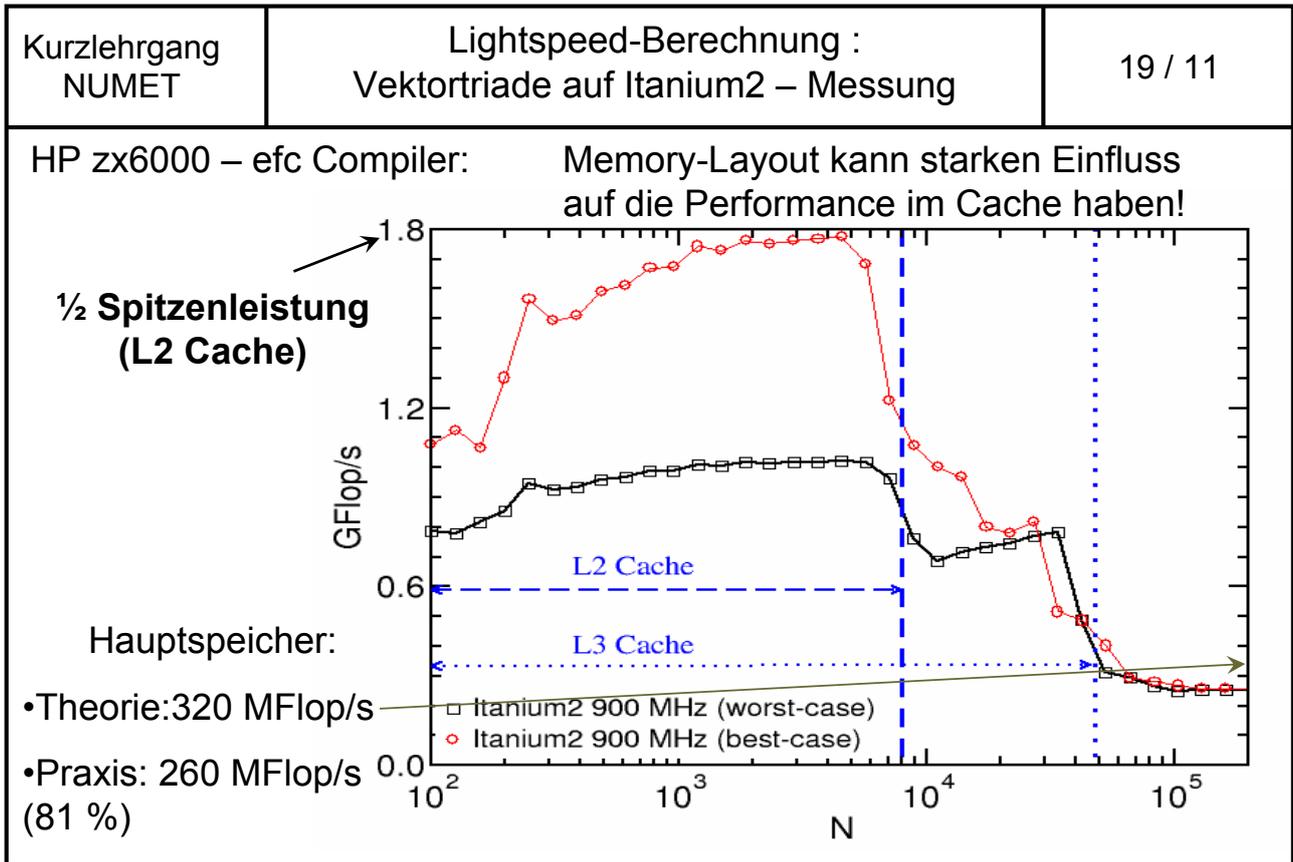
- Typ. Werte (Hauptspeicher): 0.2 W/Flop (Itanium2 1.0 GHz)
 0.125 W/Flop (Xeon 3.2 GHz),
 0.5 W/Flop (NEC SX8)

Kurzlehrgang NUMET	Abschätzung von Performancezahlen	15 / 11
<ul style="list-style-type: none"> Zu erwartende Performance auf Schleifenebene? Codebalance: $B_c = \frac{\text{Datentransfer (LD/ST) [Worte]}}{\text{arithmetische Operationen [Flops]}}$ erwarteter Bruchteil der Peak Performance ("Lightspeed"): $l = \frac{B_m}{B_c}$ <ul style="list-style-type: none"> Beispiel? 		

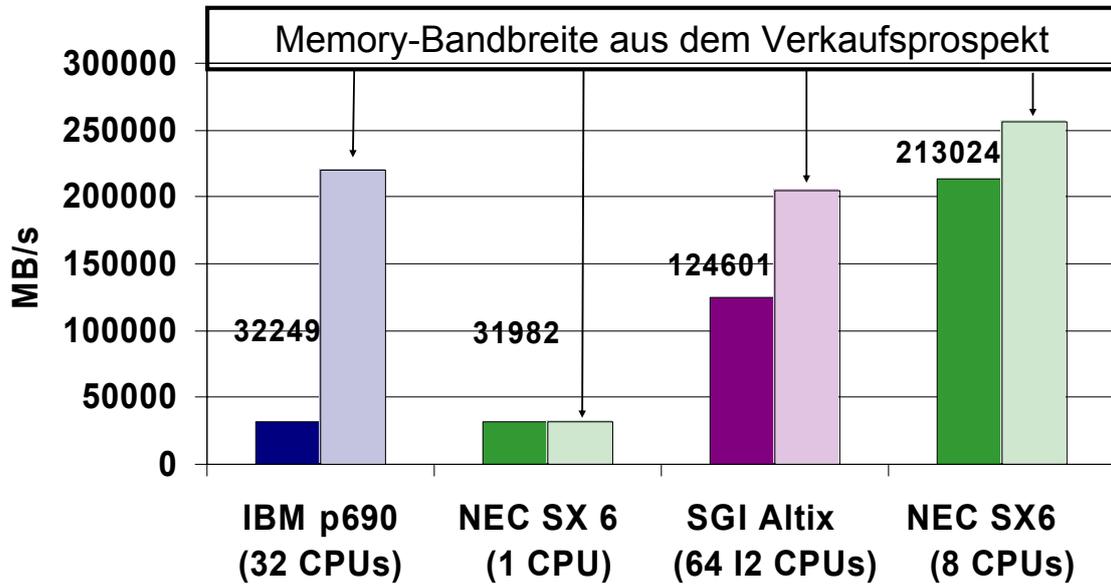
Kurzlehrgang NUMET	Beispiel für Lightspeed-Berechnung: Vektortriade	16 / 11
<ul style="list-style-type: none"> Kernel benchmarks: <ul style="list-style-type: none"> – Charakterisierung des Prozessors – Resultate können i.A. leicht verstanden und verglichen werden – Erlauben oft eine obere Schranke für die Leistungsfähigkeit eines Systems für eine spezifische Anwendergruppe abzuschätzen – Beispiele: <i>streams</i>, <i>cachebench</i>,... Vektortriade: Charakterisierung der erzielbaren Speicherbandbreiten 		
<pre> REAL*8 (SIZE): A,B,C,D DO ITER=1,NITER DO i=1,N A(i) = B(i) + C(i) * D(i) ENDDO <OBSCURE> ENDDO </pre>	<ul style="list-style-type: none"> Codebalance: <ul style="list-style-type: none"> • Recheneinheiten: 2 Flops / i-Iteration • Bandbreite: (3 LD & 1 ST) / i-Iteration • $B_c = 4/2 = 2 \text{ W/Flop}$ 	

Kurzlehrgang NUMET	Lightspeed-Berechnung: Vektortriade auf Itanium2 (1)	17 / 11
<ul style="list-style-type: none"> Funktionseinheiten: Performance allein durch die arithm. Einheiten beschränkt: 2 MULTIPLY-ADD (2 MADD) Pipelines <u>Vektortriade:</u> 2 Iterationen \longleftrightarrow 1 Takt 4 Flops / Takt \longleftrightarrow Peak Performance ➡ Die Daten liegen jedoch i.A. nicht in den Funktionseinheiten/Registern, sondern im L2 Cache: Bandbreite von 2 LD und (2 LD oder 2 ST) pro Takt (=32 GB/s bei 1 GHz) Maschinenbalance: 4/4 = 1 W/Flop <u>Vektortriade:</u> $B_m/B_c = 0.5$ ➡ 1/2 Peak Performance (2 GFlop/s) 		

Kurzlehrgang NUMET	Lightspeed-Berechnung: Vektortriade auf Itanium2 (2)	18 / 11
<ul style="list-style-type: none"> L3 Cache: komplexe Situation, weil immer aus L2 geladen wird Abschätzung: halbe Performance wie aus L2 Speicher: Bandbreite von <ul style="list-style-type: none"> – 2 LD oder 2 ST mit 400 MHz (Frequenz des Speicherbuses!) Maschinenbalance: 0.8/4 = 0.2 W/Flop – Achtung: "Read for Ownership" benötigt einen zusätzlichen LD vor jedem ST, d.h. effektive Codebalance: $B_c = 2.5$ W/Flop <u>Vektortriade:</u> $B_m/B_c = 0.08$ ➡ 8% der Peak Performance (320 MFlop/s) 		



Streams Triade ($B_c=1.5$ W/Flop)



<http://www.cs.virginia.edu/stream/standard/Bandwidth.html>

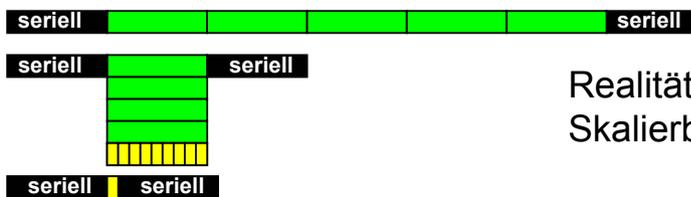
Möglichkeiten und Grenzen der Parallelität

"Parallelrechnen" = mehr als eine CPU zur Lösung eines numerischen Problems einsetzen

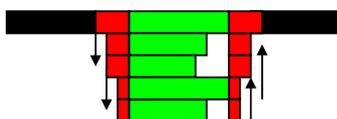
- Begriffe
 - **Skalierbarkeit** $S(N)$: Wieviel mal schneller wird mein Problem gelöst, wenn es auf N CPUs läuft statt auf einer?
 - **Performance** $P(N)$: Wieviele Rechenoperationen pro Sekunde macht das Programm mit N CPUs im Vergleich mit einer CPU?
 - **parallele Effizienz** $\epsilon(N)=S(N)/N$: Welchen Bruchteil jeder CPU nutzt das parallele Programm?
- Viele numerische Aufgaben sind im Prinzip parallelisierbar
 - Lösen linearer Gleichungssysteme: **$Ax = b$**
 - Eigenwertberechnungen: **$Ax = \lambda x$**
 - Zeitpropagation: **$p(x,t+1)=F(p(x,t))$**



Idealfall: Arbeit ist komplett parallelisierbar



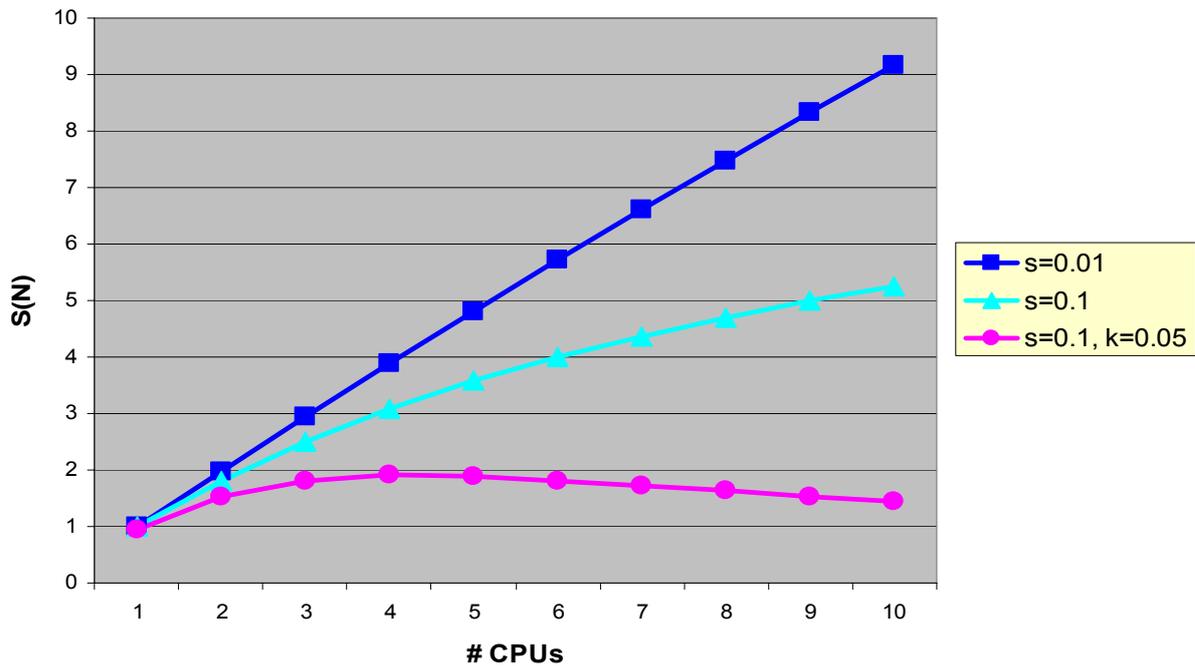
Realität: rein serielle Anteile begrenzen Skalierbarkeit



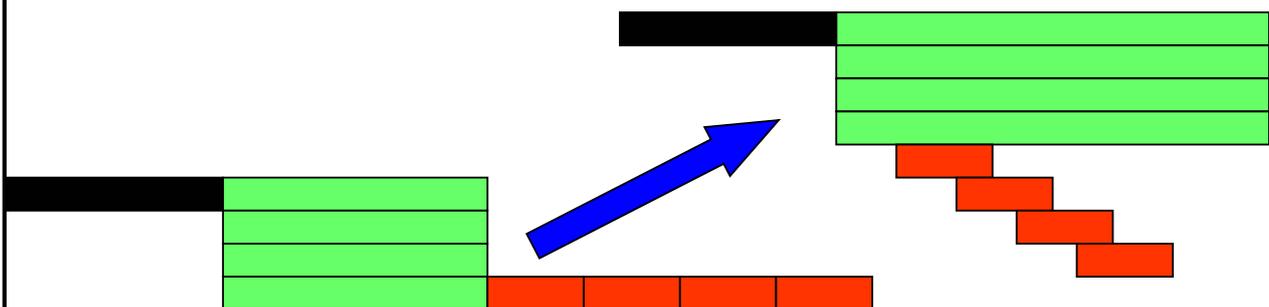
Kommunikation verschlimmert die Situation noch weiter

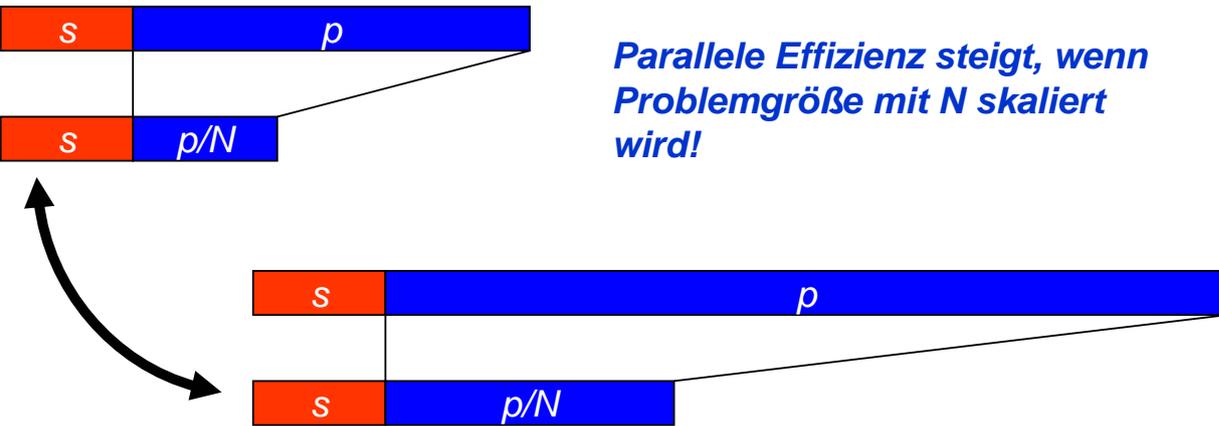
Kurzlehrgang NUMET	Grenzen der Parallelität: ein Modell	25 / 11
<p style="text-align: center;">$T(1) = s+p$ = serielle Rechenzeit</p> <p style="text-align: center;">parallelisierbarer Anteil: $p = 1-s$ rein serieller Anteil s</p> <p>parallel: $T(N) = s+p/N+Nk$</p> <p style="text-align: right;">Bruchteil k für Kommunikation zwischen je 2 Prozessoren</p> <p>Allgemeine Formel für Skalierbarkeit (worst case): Fall $k=0$: "Amdahl's Law" (strong scaling)</p> <div style="border: 1px solid blue; padding: 5px; display: inline-block;"> $S_p^k(N) = \frac{T(1)}{T(N)} = \frac{1}{s + \frac{1-s}{N} + Nk}$ </div>		

Kurzlehrgang NUMET	Grenzen der Parallelität	26 / 11
<ul style="list-style-type: none"> Limes großer Prozessorzahlen: <ul style="list-style-type: none"> bei $k=0$: Amdahl's Law! <div style="border: 1px solid black; padding: 5px; display: inline-block;"> $\lim_{N \rightarrow \infty} S_p^0(N) = \frac{1}{s}$ </div> <p style="text-align: center;">unabhängig von N</p> bei $k \neq 0$: einfaches Kommunikationsmodell liefert: <div style="border: 1px solid black; padding: 5px; display: inline-block;"> $S_p^k(N) \xrightarrow{N \gg 1} \frac{1}{Nk}$ </div> Die Realität ist noch viel schlimmer: <ul style="list-style-type: none"> Load Imbalance Overhead beim Starten paralleler Programmteile 		



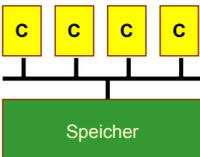
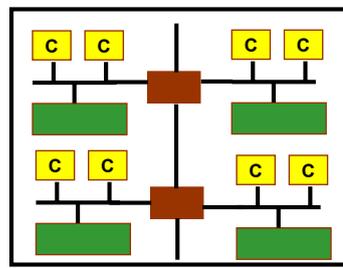
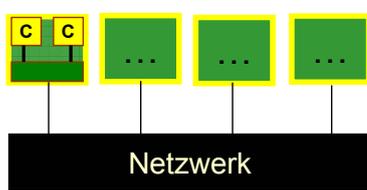
- Kommunikation ist nicht immer rein seriell
 - nichtblockierende Netzwerke können gleichzeitig verschiedene Kommunikationsverbindungen schalten – Faktor Nk im Nenner wird zu $\log k$ oder k (technische Maßnahme)
 - u.U. kann Kommunikation mit nützlicher Arbeit überlappt werden (Implementierung, Algorithmus):



Kurzlehrgang NUMET	Grenzen der Parallelität: Erhöhen der parallelen Effizienz	29 / 11
<ul style="list-style-type: none"> • Vergrößern des Problems hat oft eine Vergrößerung von p zur Folge <ul style="list-style-type: none"> – p skaliert, während s konstant bleibt – Anteil von s an der Gesamtzeit sinkt – Nebeneffekt: Kommunikationsaufwand sinkt eventuell  <p data-bbox="831 546 1369 674"><i>Parallele Effizienz steigt, wenn Problemgröße mit N skaliert wird!</i></p>		

Kurzlehrgang NUMET	Grenzen der Parallelität: Erhöhen der parallelen Effizienz	30 / 11
<ul style="list-style-type: none"> • Quantifizierung? Betrachte zunächst $k=0!$ <ul style="list-style-type: none"> – Sei wie vorher $s+p=1$ – Perfekte Skalierung: Laufzeit bleibt konstant bei Vergrößerung von N – Skalierbarkeit ist keine relevante Größe mehr! – „Parallele Performance“ = $\frac{\text{Arbeit/Zeit für Problemgröße } N \text{ mit } N \text{ CPUs}}{\text{Arbeit/Zeit für Problemgröße } 1 \text{ mit } 1 \text{ CPU}}$ <div style="border: 1px solid black; padding: 10px; width: fit-content; margin: 10px auto;"> $P_s(N) = \frac{s + pN}{s + p} = s + pN = s + (1 - s)N$ </div> <p style="text-align: center;">Gustafsson's Law ("weak scaling")</p> <ul style="list-style-type: none"> – Linear in N – aber die Definition der "Arbeit" spielt eine große Rolle! 		

Kurzlehrgang NUMET		31 / 11
<h2>Designprinzipien moderner Parallelrechner</h2>		

Kurzlehrgang NUMET	Moderne Parallelrechner Schematischer Aufbau	32 / 11
(1) Bus-basierte SMP-Systeme		<ul style="list-style-type: none"> • 2-4 Prozessoren mit einem Speicherkanal • Itanium / Xeon / Power4 / NEC SX
(2) Skalierbare „Shared-memory“ Systeme (ccNUMA)		<ul style="list-style-type: none"> • Speicherbandbreite skaliert • Sehr schnelle interne Verbindungen (0.8 GB/s – 6.4 GB/s) • SGI Origin / Altix, IBM p690 (ähnlich), Opteron-Knoten • 4-256 Prozessoren
(3) Cluster von Rechenknoten (Architektur (1) oder (2))		<ul style="list-style-type: none"> • Speicherbandbreite skaliert, aber verteilter Speicher • Netzwerk: 50 – 1000 MByte/s • 4 – 100000 Prozessoren

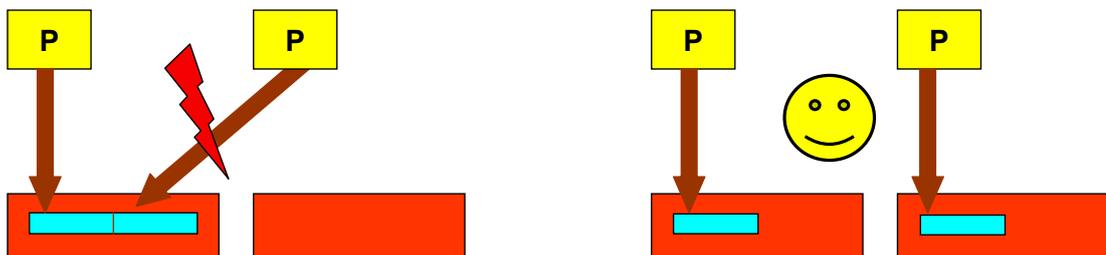
Kurzlehrgang NUMET	Moderne Parallelrechner	33 / 11
<ul style="list-style-type: none">• Bei allen 3 Varianten (heute): Dual-Core möglich<ul style="list-style-type: none">– d.h.: in einem "Sockel" steckt ein Chip mit 2 oder mehr (etwas langsameren) CPUs auf einem Stück Silizium– Vorteil: potenziell höhere Rechenleistung pro Sockel bei nahezu gleicher Leistungsaufnahme ("Macho-FLOPs")– Nachteil: Speicheranbindung der einzelnen CPU wird schlechter• Dominant bei Parallelrechnern: Cluster mit SMP-Knoten<ul style="list-style-type: none">– Kompromiss bzgl. Speicherbandbreite pro CPU– skalierend bis mehrere 10000 CPUs (akt. Rekord: 131072)• Spezialfall: "Constellation"-Systeme<ul style="list-style-type: none">– CPUs/Knoten > # Knoten– SGI Altix (HLRB II): ccNUMA-Architektur		

Kurzlehrgang NUMET		34 / 11
<p style="text-align: center;">Parallele Programmierung in der Strömungsmechanik</p>		

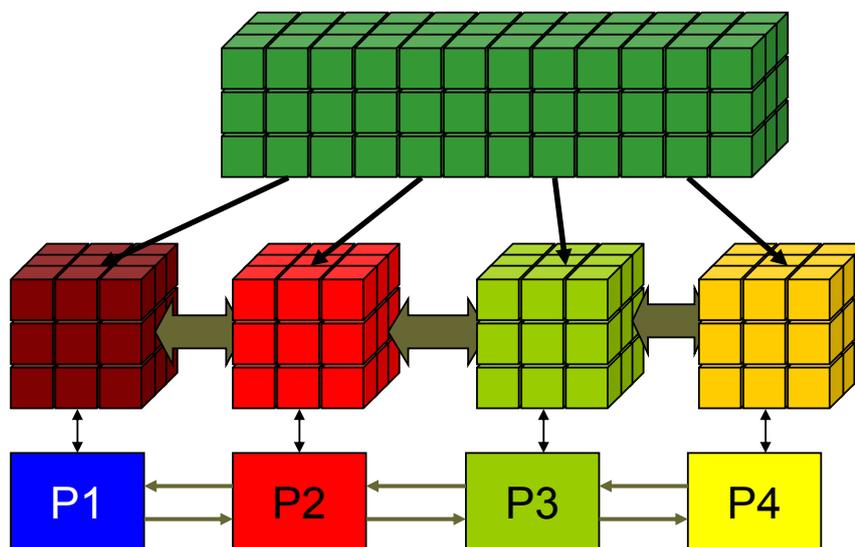
Kurzlehrgang NUMET	Programmierparadigmen	35 / 11
<ul style="list-style-type: none"> • Parallele Programmierung beginnt immer mit der Identifikation von Parallelität im Code <ul style="list-style-type: none"> – Datenparallelität: Welche Teile eines großen Datensatzes (Feld, Matrix, ...) können im Prinzip gleichzeitig bearbeitet werden? Beispiel: <pre style="margin-left: 20px;"> do i=1,N do j=1,N A(i,j)=A(i,j-1)*B(i,j) enddo enddo </pre> <div style="margin-left: 40px; border: 1px solid green; border-radius: 5px; padding: 2px; display: inline-block;">datenparallel auf i-Ebene</div> <div style="margin-left: 40px; border: 1px solid red; border-radius: 5px; padding: 2px; display: inline-block;">nicht datenparallel auf j-Ebene</div> – funktionelle Parallelität: Welche funktionalen Einheiten im Programm können im Prinzip gleichzeitig ablaufen? <ul style="list-style-type: none"> • in der Praxis selten benutzt 		

Kurzlehrgang NUMET	Programmierparadigmen	36 / 11
<ul style="list-style-type: none"> • Abhängig von der Rechnerarchitektur bieten sich 2 Zugänge an: <ul style="list-style-type: none"> – Distributed-Memory-Systeme <p>Ein paralleles Programm besteht aus mehreren Prozessen. Jeder Prozess hat seinen eigenen Adressbereich und keine Zugriff auf die Daten anderer Prozesse. Notwendige Kommunikation findet durch Übermittlung von Messages statt. Standard: MPI</p> – Shared-Memory-Systeme (auch ccNUMA) <p>Ein paralleles Programm besteht aus mehreren Threads. Alle Threads haben Zugriff auf den kompletten Datenbereich. Standard: OpenMP (Compiler-Direktiven)</p> 		

- Bei **ccNUMA-Systemen** (SGI Altix, Origin) sind Verbindungen zwischen Knoten i.A. langsamer als die lokale Speicherbandbreite
- Folge: Es muss darauf geachtet werden, dass jede CPU auf den benötigten Speicher **möglichst lokal** zugreifen kann
 - Problem dabei: **“First Touch”** Policy mappt Speicherseiten dort, wo sie das erste Mal angefasst werden (i.A. bei der Initialisierung)
 - Wenn die Initialisierung nicht in gleicher Weise parallelisiert ist wie die Rechenschleifen, folgt eine sehr ungünstige Zugriffsstruktur:



- Unabhängig von Programmiermodell: Parallelisierung von CFD-Codes läuft oft über **"Domain Decomposition"**
 - Zerlegung des Rechengebietes in N Teile – eines für jede CPU

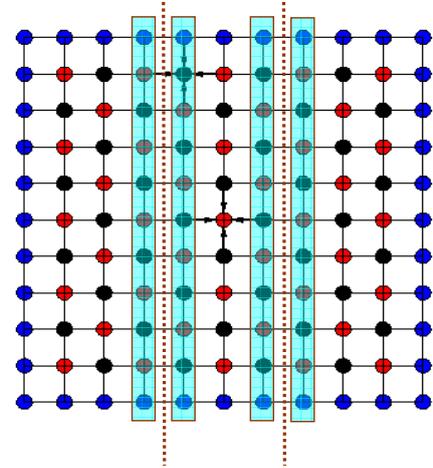


- Problem: Temperaturverteilung auf quadratischer Platte mit Randbedingungen

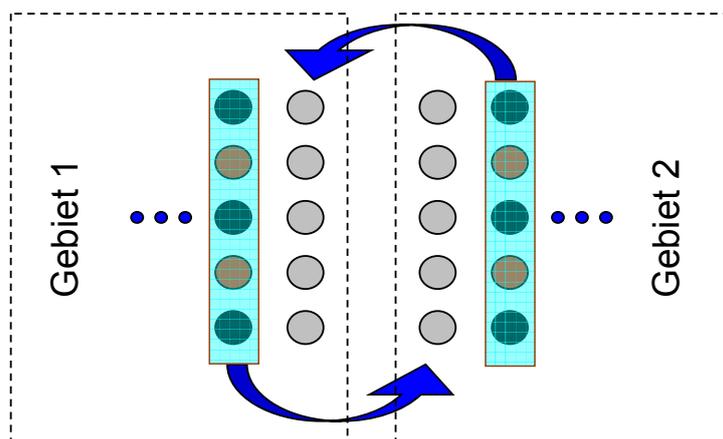
- Lösung der **Laplace-Gleichung**:

$$\Delta T = 0$$

- Numerik: Abb. auf NxN-Gitter
 - Algorithmus: "**Relaxation**"
T an einem Punkt = Mittelwert der 4 Nachbarn
 - **Iteration**: Update erst der roten, dann der schwarzen Punkte
- Parallelisierung mittels Gebietszerlegung
 - Was macht man mit den Randzellen eines Gebietes?



- Lösung des Randzellenproblems durch "Geisterzellen" (**ghost cells**)
 - jedes Gebiet hat eine zusätzliche Schicht Randzellen, wo es an ein anderes Gebiet grenzt
 - Nach jeder Iteration (red/black Sweep) werden die Geisterzellen mit dem Inhalt der Randzellen des Nachbarn **beschrieben**

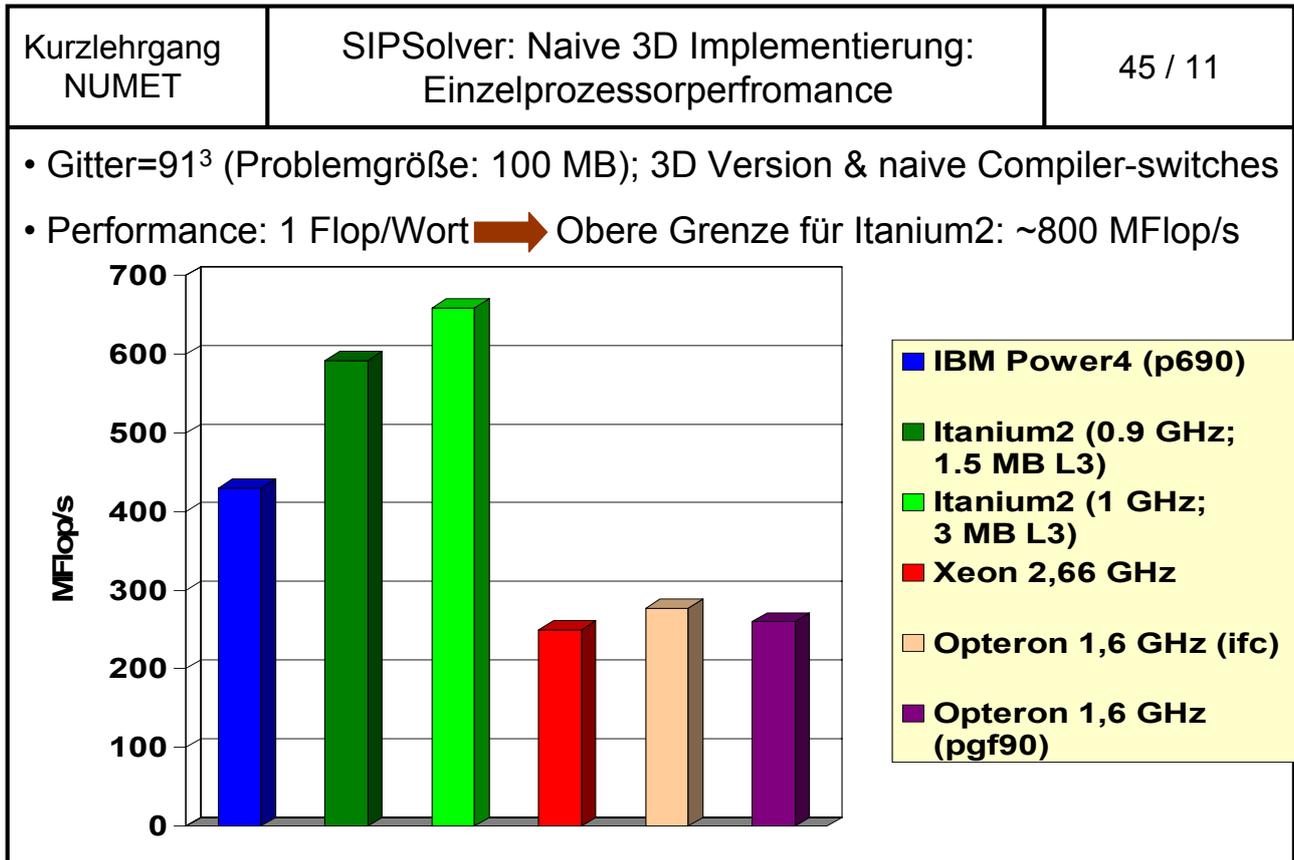


Kurzlehrgang NUMET	Gebietszerlegung und Kommunikation: Kostenmodelle	41 / 11
<ul style="list-style-type: none"> • Wirkt sich die Kommunikation auf die Skalierbarkeit aus? <ul style="list-style-type: none"> – Annahmen: keine Überlappung zwischen Kommunikation und Rechnung, nichtblockierendes Netzwerk, 3D-Problem • 2 Fälle: <ol style="list-style-type: none"> (1) strong scaling: Problemgröße konstant, N steigt → Latenz L wird wichtig <div style="display: flex; align-items: center; margin-top: 10px;"> <div style="flex: 1;"> <p>Kommunikationszeit = $k/N^{2/3}+L$ Skalierung sättigt langsamer und auf niedrigerem Niveau!</p> </div> <div style="border: 1px solid black; padding: 5px; margin-left: 10px;"> $S_p(N) = \frac{1}{s + \frac{1-s}{N} + kN^{-2/3} + L} \xrightarrow{N \gg 1} \frac{1}{s + L}$ </div> </div> (2) weak scaling: Problemgröße steigt proportional zu N <div style="display: flex; align-items: center; margin-top: 10px;"> <div style="flex: 1;"> <p>Arbeit pro CPU bleibt konstant, ebenso die Kommunikation k Performance skaliert langsamer, aber noch linear</p> </div> <div style="border: 1px solid black; padding: 5px; margin-left: 10px;"> $P(N) = \frac{s + pN}{s + p + k} = \frac{s + (1-s)N}{1 + k}$ </div> </div> 		

Kurzlehrgang NUMET		42 / 11
<p>CFD Applikationsperformance</p> <p>SIP Solver nach Stone</p>		

Kurzlehrgang NUMET	CFD kernel: Strong Implicit Solver	43 / 11
<ul style="list-style-type: none"> • <i>Strongly Implicit Procedure</i> (SIP) nach Stone wird in Finite-Volumen Paketen oftmals zur Lösung des linearen Gleichungssystems $Ax = b$ verwendet • Beispiele: <ul style="list-style-type: none"> – LESOCC, FASTEST, FLOWSI (LSTM, Erlangen) – STHAMAS3D (Kristalllabor, Erlangen) – CADiP (Theoret. Thermodynamik und Transportprozesse, Bayreuth) – ... • SIPSolver: 1) Unvollständige LU-Zerlegung von A 2) Serie von Vorwärts-/Rückwärts-Substitutionen und Residuenberechnungen • Ursprüngliches Testprogramm (M. Peric): ftp.springer.de:/pub/technik/peric • Optimization: HPC-Gruppe RRZE 		

Kurzlehrgang NUMET	SIPSolver: Datenabhängigkeiten und naive Implementierung	44 / 11
<p>Vorwärts-Substitution (naive 3D version)</p> <p>Datenabhängigkeit: $(i, j, k) \longleftarrow \{(i-1, j, k); (i, j-1, k); (i, j, k-1)\}$</p>		
<pre> do k = 2 , kMax do j = 2 , jMax do i = 2 , iMax RES(i,j,k) = {RES(i,j,k) \$ - LB(i,j,k)*RES(i,j,k-1) \$ - LW(i,j,k)*RES(i-1,j,k) \$ - LS(i,j,k)*RES(i,j-1,k) \$ }*LP(i,j,k) enddo enddo enddo </pre>		<p><u>Datenabhängigkeit</u></p> <ul style="list-style-type: none"> • verhindert Vektorisierung oder einfache Parallelisierung • erlaubt hohe örtliche Datenlokalität: Alle Elemente der Cachelines werden genutzt. • erlaubt reuse von Cachelines: RES(i, j, k), RES(i-1, j, k)
<p> Geeignete Implementierung für Cache-basierte Prozessoren!</p>		



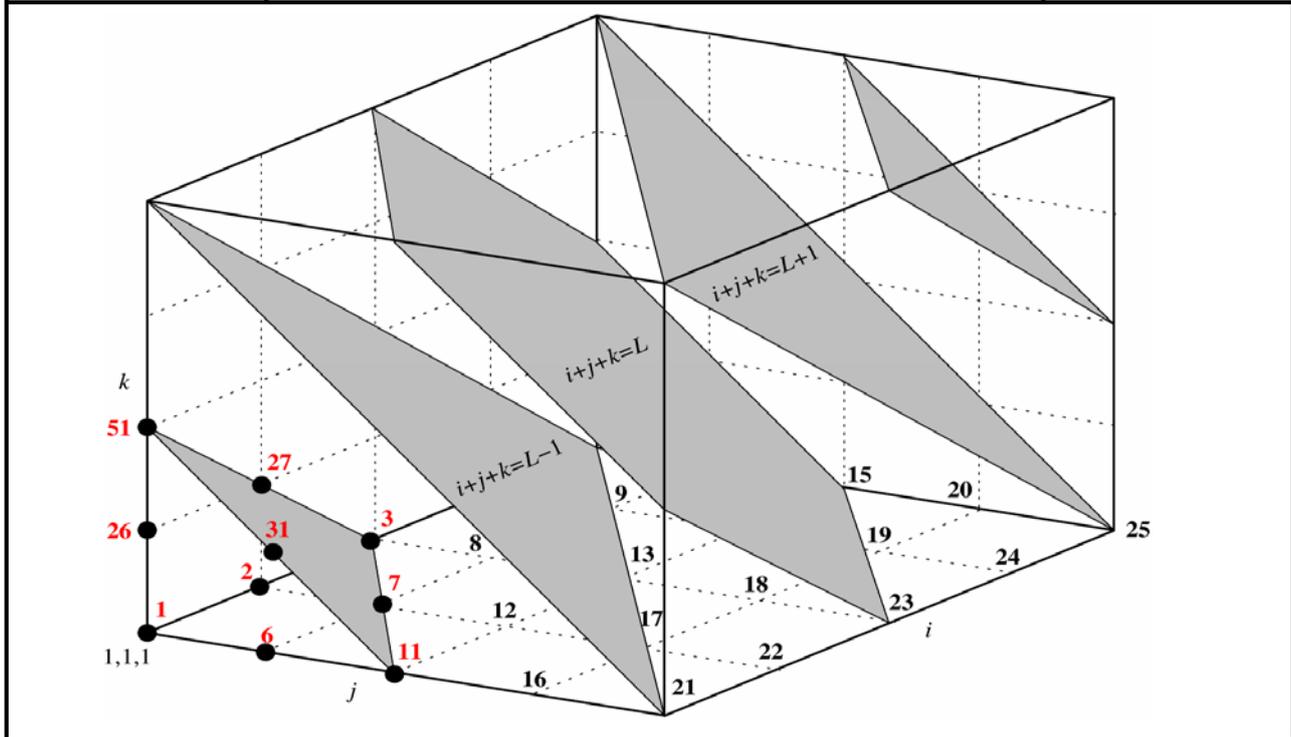
Kurzlehrgang NUMET	SIPsSolver: <i>hyperplane</i> Implementierung	46 / 11
-----------------------	--	---------

Elimination der Datenabhängigkeit $(i, j, k) \leftarrow \{(i-1, j, k); \dots\}$ durch Umgruppierung:

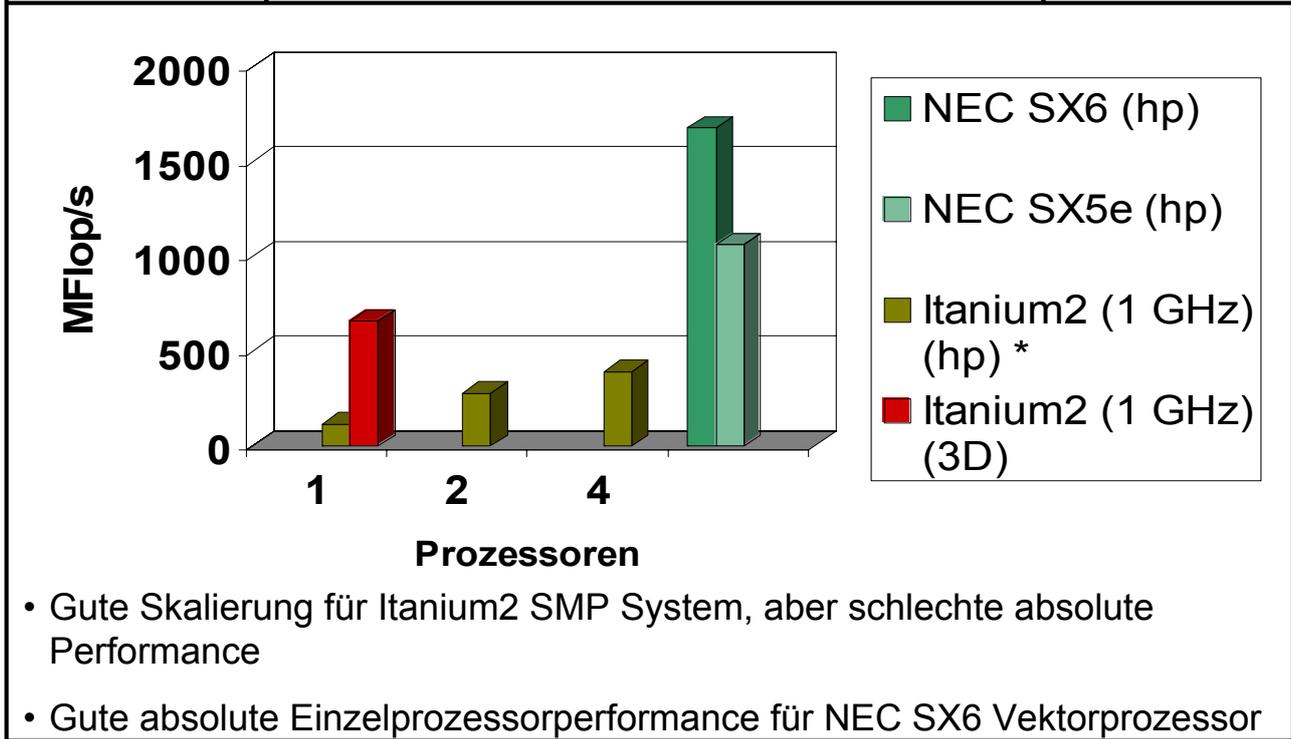
<p><u>Innerste Schleife über eine <i>hyperplane</i>:</u></p> <p>Alle Punkte $i+j+k=1$ (mit festem 1)</p> <ul style="list-style-type: none"> • liegen in der <i>hyperplane</i> 1, • sind voneinander unabhängig • greifen nur auf Datenpunkte in der <i>hyperplane</i> 1-1 (Vorwärts-substitution) zu. 	<pre> do l=1,hyperplanes n=ICL(l) do m=n+1,n+LM(l) ijk=IJKV(m) RES(ijk) = (RES(ijk) - \$ LB(ijk)*RES(ijk-ijMax) - \$ LW(ijk)*RES(ijk-1) - \$ LS(ijk)*RES(ijk-iMax)) \$ *LP(ijk) enddo enddo </pre>
<ul style="list-style-type: none"> • Einfache Vektorisierung/ Parallelis. der inneren Schleife • Verlust der Datenlokalität 	

→ Geeignete Implementierung für **Vektorprozessoren**

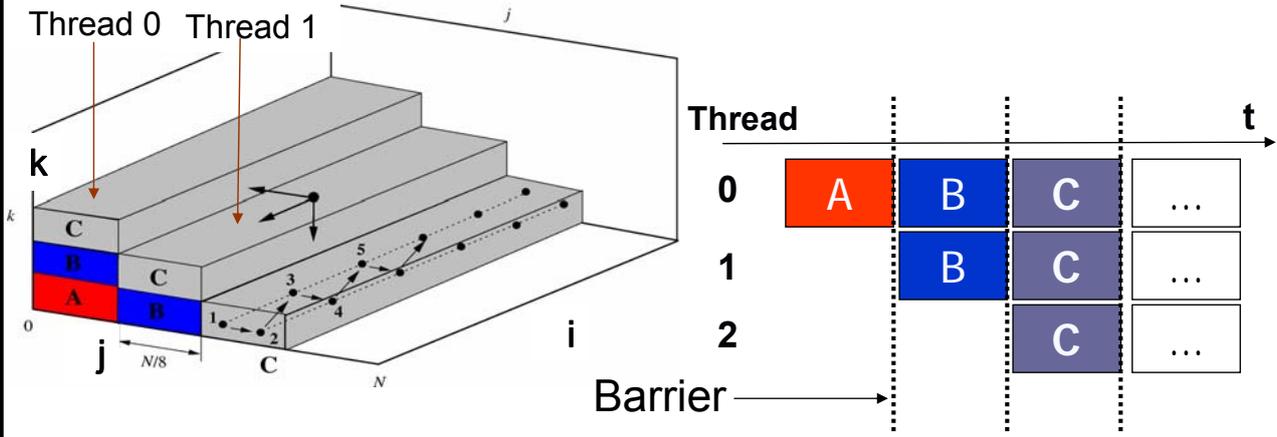
Kurzlehrgang NUMET	SIPsolver: <i>hyperplane</i> Implementierung	47 / 11
-----------------------	---	---------



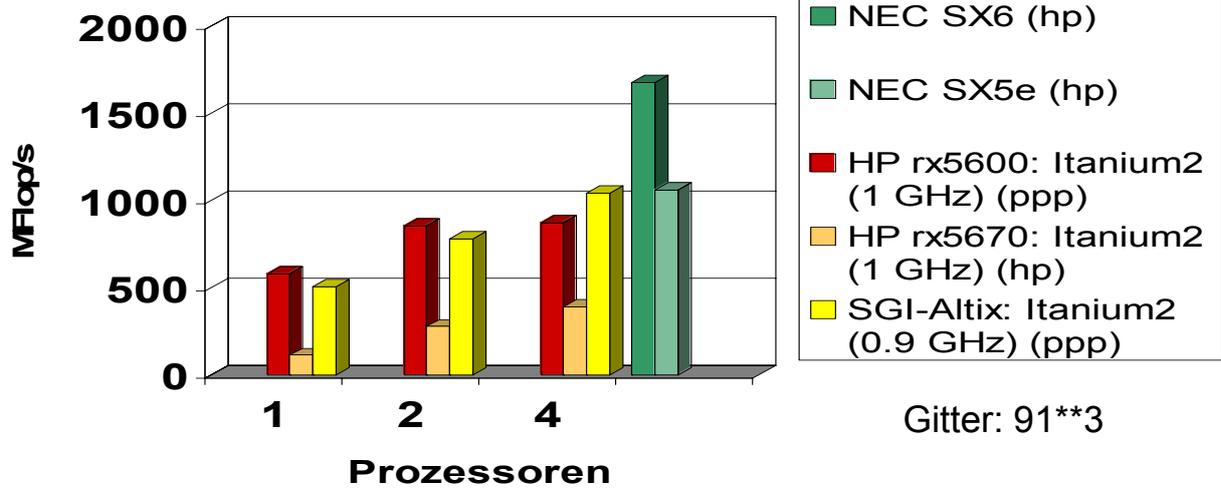
Kurzlehrgang NUMET	SIPsolver: <i>hyperplane</i> Leistungszahlen	48 / 11
-----------------------	--	---------



Kurzlehrgang NUMET	SIP Solver: <i>pipeline parallel processing</i> Idee und Implementierung	49 / 11
<p>Ziel: Parallelisierung (Shared-Memory) und Erhalt der Datenlokalität der 3D Implementierung</p> <p>Vorgehensweise:</p> <ul style="list-style-type: none"> • 3D Implementierung $(i, j, k) \leftarrow \{(i-1, j, k); \dots\}$ • Parallelisierung in j Richtung • Aufheben der Datenabhängigkeiten durch Synchronisieren nach jedem Schritt in k - Richtung <pre data-bbox="735 264 1422 913"> \$omp parallel private(...) do l = 2, kMax+numThreads-2 threadID=OMP_GET_THREAD_NUM() k = 1 - threadID if((k.ge.2).and.(k.le.kMaxM)) then do j = jS(threadID), jE(threadID) do i = 2, iMax RES(i,j,k)={RES(i,j,k-1)- LB(i,j,k)*RES(i,j,k-1) ... } enddo enddo endif \$omp barrier enddo \$omp end parallel </pre>		

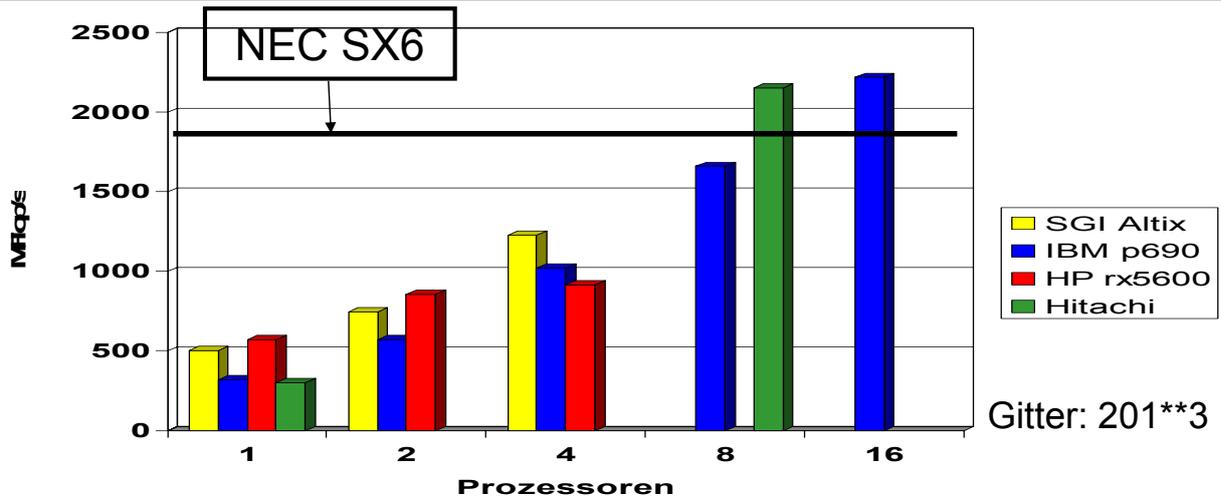
Kurzlehrgang NUMET	SIP Solver: <i>pipeline parallel processing</i> Schematische Darstellung	50 / 11
 <p>The 3D diagram shows a stack of layers (A, B, C) along the k-axis. The j-axis is divided into segments of size $N/8$. Two threads, Thread 0 and Thread 1, are shown processing different parts of the j-axis. The Gantt chart shows the execution timeline for threads 0, 1, and 2. Thread 0 starts with block A, followed by B and C. Thread 1 starts with B, followed by C. Thread 2 starts with C. Vertical dashed lines indicate barriers between threads.</p> <ul style="list-style-type: none"> • Parallelisierung sinnvoll, falls $k_{\text{Max}}, j_{\text{Max}} \gg \#\text{Threads}$ • Automatische Parallelisierung der 3D Version durch Hitachi-Compiler! • Beste Version für große SMP Knoten mit Cache-basierten Prozessoren 		

Kurzlehrgang NUMET	SIPSolver: SMP Skalierbarkeit (Pipeline Parallel Processing)	51 / 11
-----------------------	---	---------



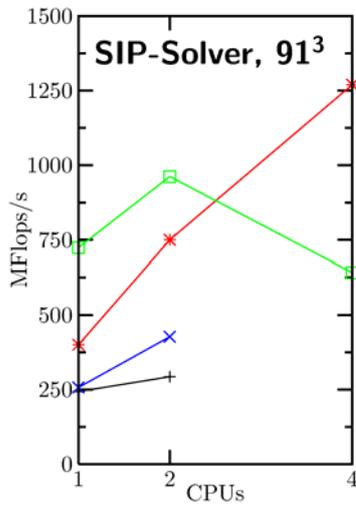
- HP rx5670: 4-fach Itanium2-System mit nur einem Speicherkanal für 4 Prozessoren (Speed-Up(4)=1,5)
- SGI Altix: 32-fach Itanium2-System mit einem Speicherkanal für **jeweils 2** Prozessoren (Speed-Up(4)=2,0)

Kurzlehrgang NUMET	SIPSolver: SMP Performance vs. Vektorprozessor	52 / 11
-----------------------	---	---------

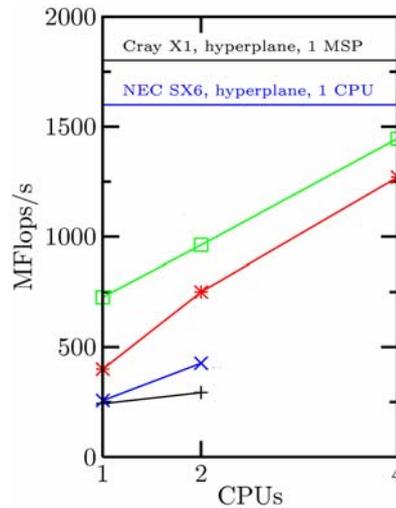


- IBM p690: 32-fach Power4-System mit einem Speicherkanal für **jeweils 2** Prozessoren (Speed-Up(16)=7)
- Hitachi SR8000: 8-fach „Power3“-System mit einem Speicherkanal pro Prozessor (Speed-Up(8)=7.2)

• Performancevergleich



Falsches Placement



Korrektes Placement

Altix

- Konkrete Maßnahme: korrekte Parallelisierung der Initialisierungsschleife

!\$omp parallel do private(i,j)



```
do k=1,kMax
  do j=1,jMax
    do i=1,iMax
      T(i,j,k)=0.
    enddo
  enddo
enddo
```

!\$omp parallel do private(i)



```
do k=1,kMax
  do j=1,jMax
    do i=1,iMax
      T(i,j,k)=0.
    enddo
  enddo
enddo
```

Kurzlehrgang NUMET	Zusammenfassung	55 / 11
<p>Ergebnisse für CFD-Applikationen:</p> <ul style="list-style-type: none">• Speicherbandbreite ist die entscheidende Einflussgröße• Effiziente Nutzung schneller, Vermeidung langsamer Datenpfade ist Pflicht!• Parallelisierung muss immer auch parallele Effizienz berücksichtigen<ul style="list-style-type: none">– Kommunikationsaufwand beachten– strong/weak scaling• Parallele Skalierbarkeit der Problemstellung ist Voraussetzung, um SMP Systeme/Cluster in Konkurrenz zu Vektorprozessoren zu setzen• Bei ccNUMA-Systemen ist auf korrektes Placement zu achten!		

Kurzlehrgang NUMET		56 / 11
<p style="text-align: center;">Vielen Dank!</p>		

HPC-Dienste und Rechner:

- RRZE: <http://www.rrze.uni-erlangen.de/hpc/>
- LRZ: <http://www.lrz-muenchen.de/services/compute/hlr/>
- HLR Stuttgart: <http://www.hlrs.de/>
- NIC Jülich: <http://www.fz-juelich.de/nic/>

HPC allgemein:

- KONWIHR: <http://konwihr.in.tum.de/>
- TOP500: <http://www.top500.org/>
- Benchmarks: <http://www.spec.org/>
- LSTM: <http://www.lstm.uni-erlangen.de/>