

Signalverarbeitungsalgorithmen von der Theorie zur Implementierung

Markus Püschel
Electrical and Computer Engineering
Carnegie Mellon University
(13. – 16. Dezember 2004, TU Wien)

Idee der Vorlesung

Diese Vorlesung: Signalverarbeitung
(Fouriertransformation, Filter)



Anwendung/Problem

Algorithmen

Implementierung

Computerarchitektur

Überblick über die Vorlesung

■ Teil I: Was ist Signalverarbeitung?

- Beispiele
- Was sind die Bausteine der Signalverarbeitung?
- Warum sind schnelle Implementierungen wichtig?

■ Teil II: Algorithmentheorie

- Lineare Transformationen
- Diskrete Fouriertransformation (DFT)
- FIR Filter

■ Teil III: Implementierung

- Probleme
- Algorithmus, Implementierung und Mikroarchitektur: Beispiel DFT
- Selbstadaptive Software: FFTW und SPIRAL

■ Abschließende Bemerkungen

Teil I: Was ist Signalverarbeitung?

Definitionen

■ Definition: Signalverarbeitung

- [Die Disziplin die sich befasst mit] der Darstellung, Transformation und Manipulation von Signalen und der Information die sie enthalten (Oppenheim, Schafer 1999)

■ Definition: Signal

- (In der Signalverarbeitung) Eine Funktion über einem Indexbereich

$$s : I \rightarrow \mathbb{K}, \quad i \mapsto s(i)$$

Typische Beispiele:

$$\mathbb{K} = \mathbb{R}, \mathbb{C}, GF(2)$$

(reelle, komplexe, Bit-signale)

$$I = \mathbb{R}, \mathbb{Z}, \{0, \dots, n - 1\}$$

(kontinuierliche, diskrete, endliche Signale)

Beispiele

■ Multimedia

- Sprache (1-D), Bild (2-D), Video (3-D)
- Qualitätsverbesserung, Kompression, Übertragung

■ Biometrie

■ Medizinische/Biobildverarbeitung

■ Computersehen

■ Kommunikation

Multimedia: Beispiel Bildkompression



Lena



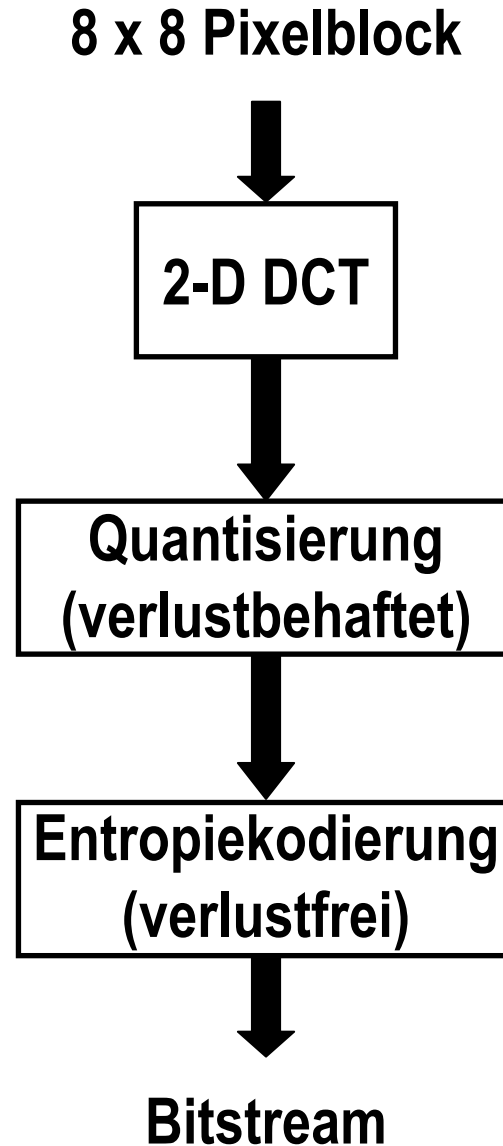
Pepper



Baboon

$512 \times 512 \times 3 \text{ bytes} = 768\text{KB}$
With JPEG, ~32KB

JPEG: Wie funktioniert's?



JPEG versus JPEG2000



original: 3MB



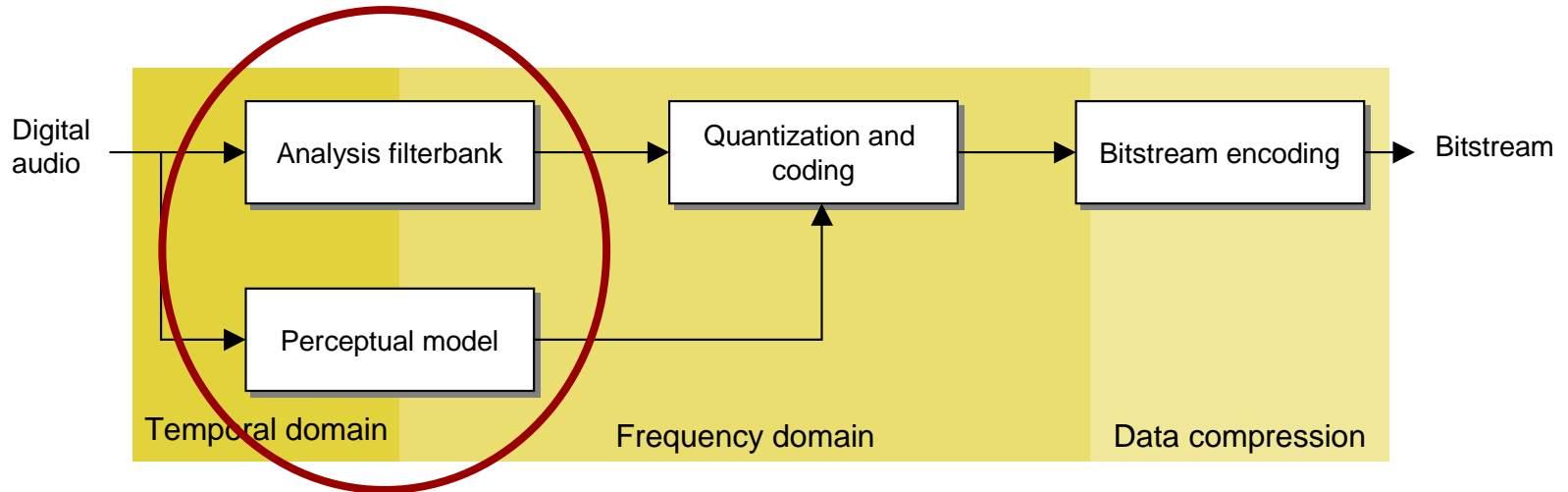
JPEG: 19KB (DCT basiert)



JPEG2000: 19KB (wavelet basiert)

Multimediakodierung

- MPEG-I bis MPEG-IV
- Enthält Standards für Audio, Bild und Video
- Beispiel: MPEG-II, layer III audio = MP3



DFT, MDCT, DCT

Beispiel: Biometrie



Facial Expression

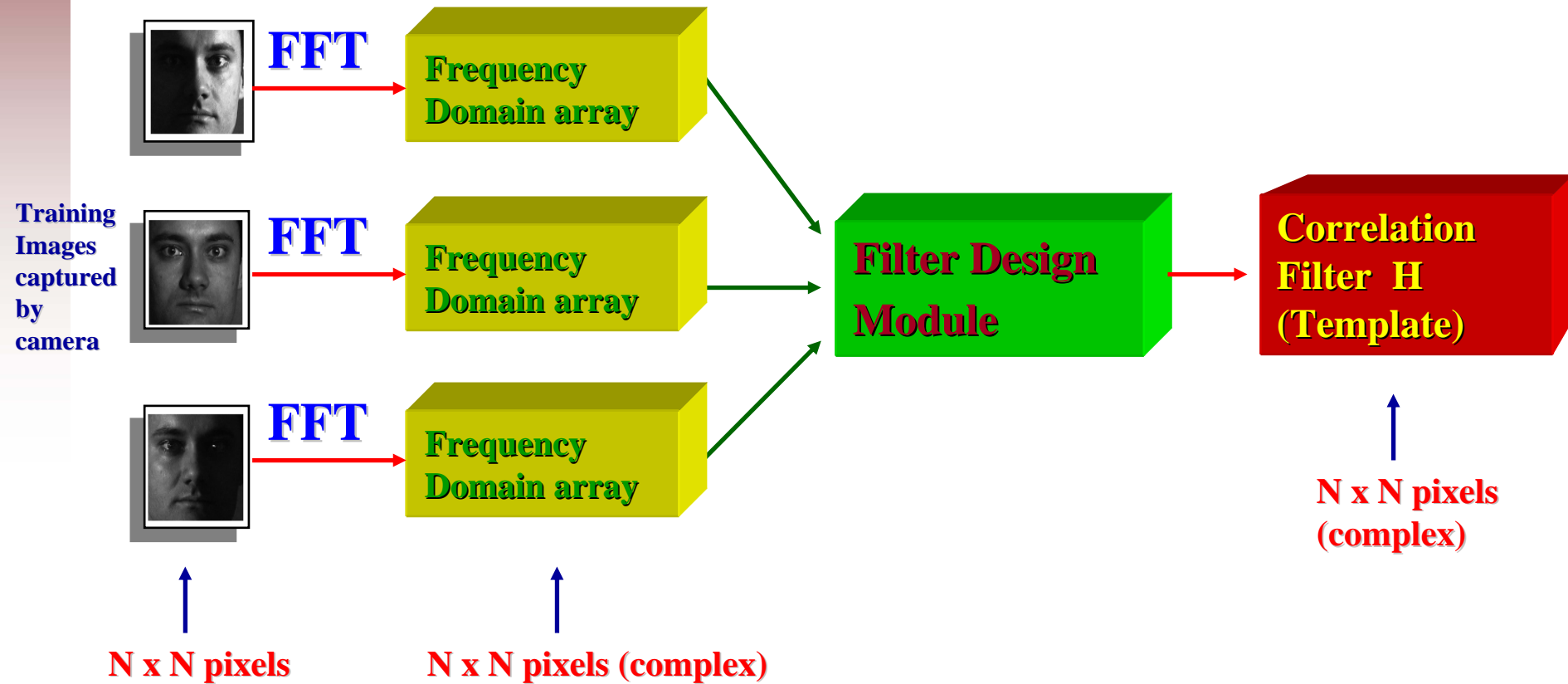


Fingerprints



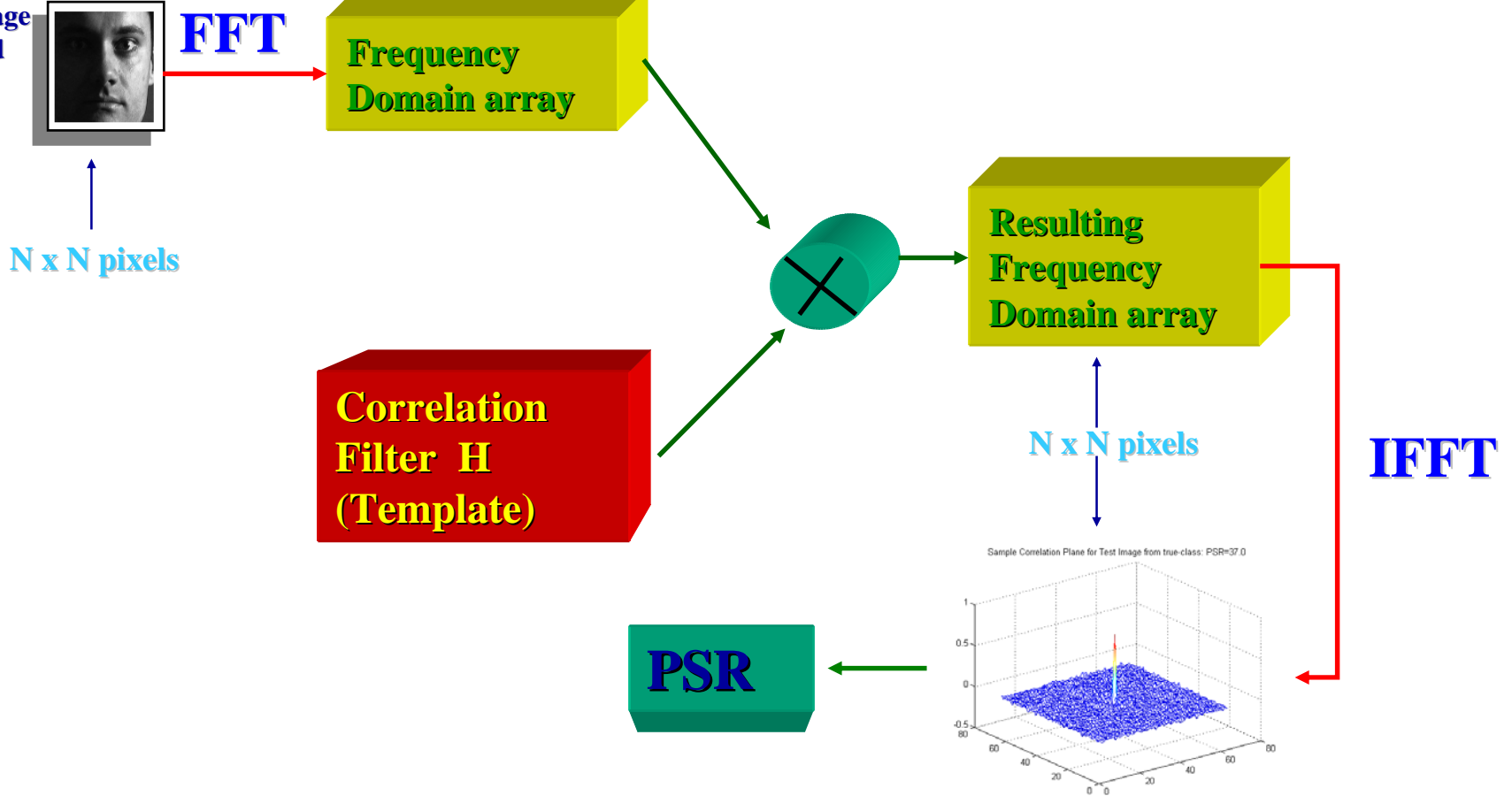
Illumination

Wie's Funktioniert: Registrierung

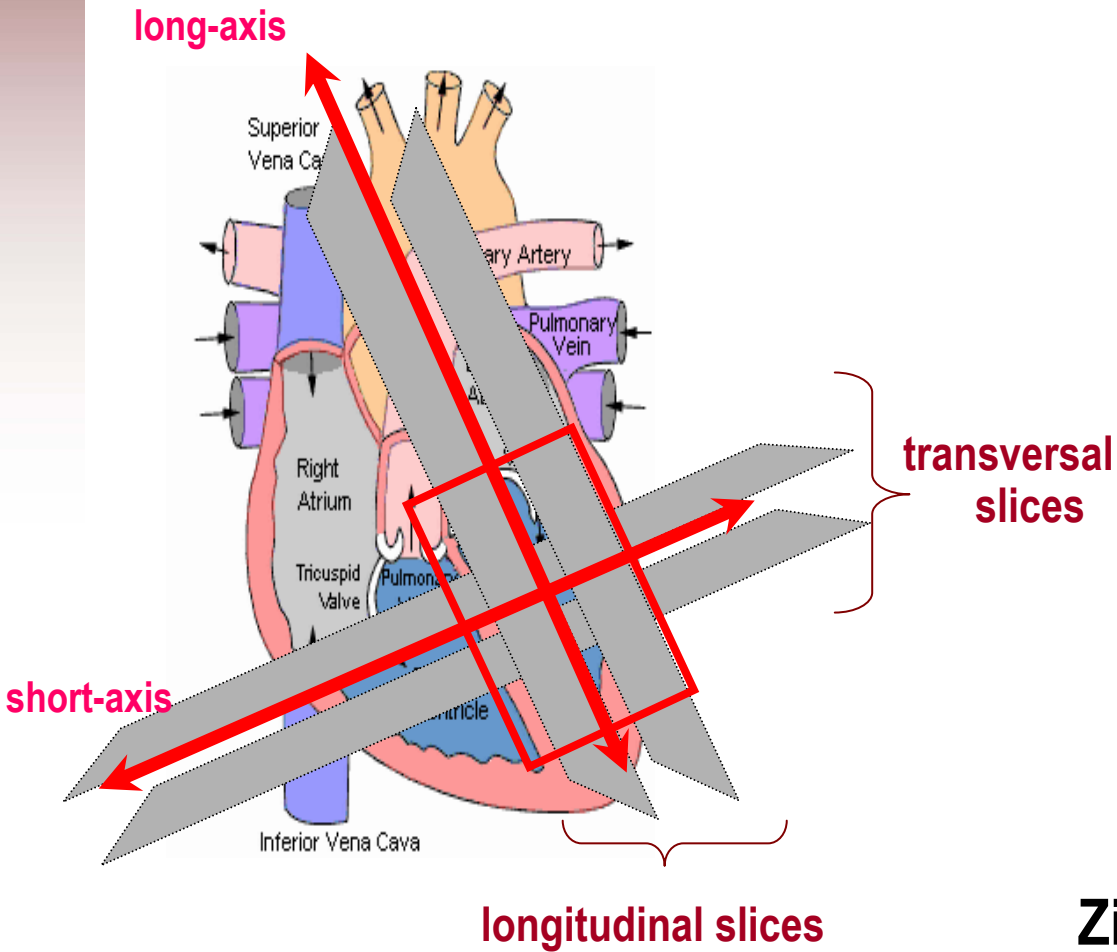
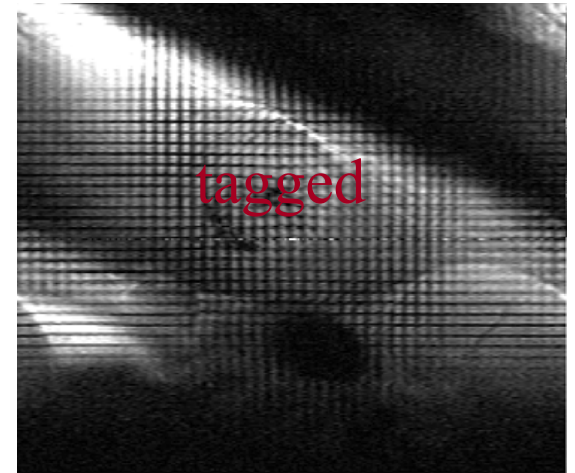


Wie's funktioniert: Identifikation

Test Image
captured
by
camera

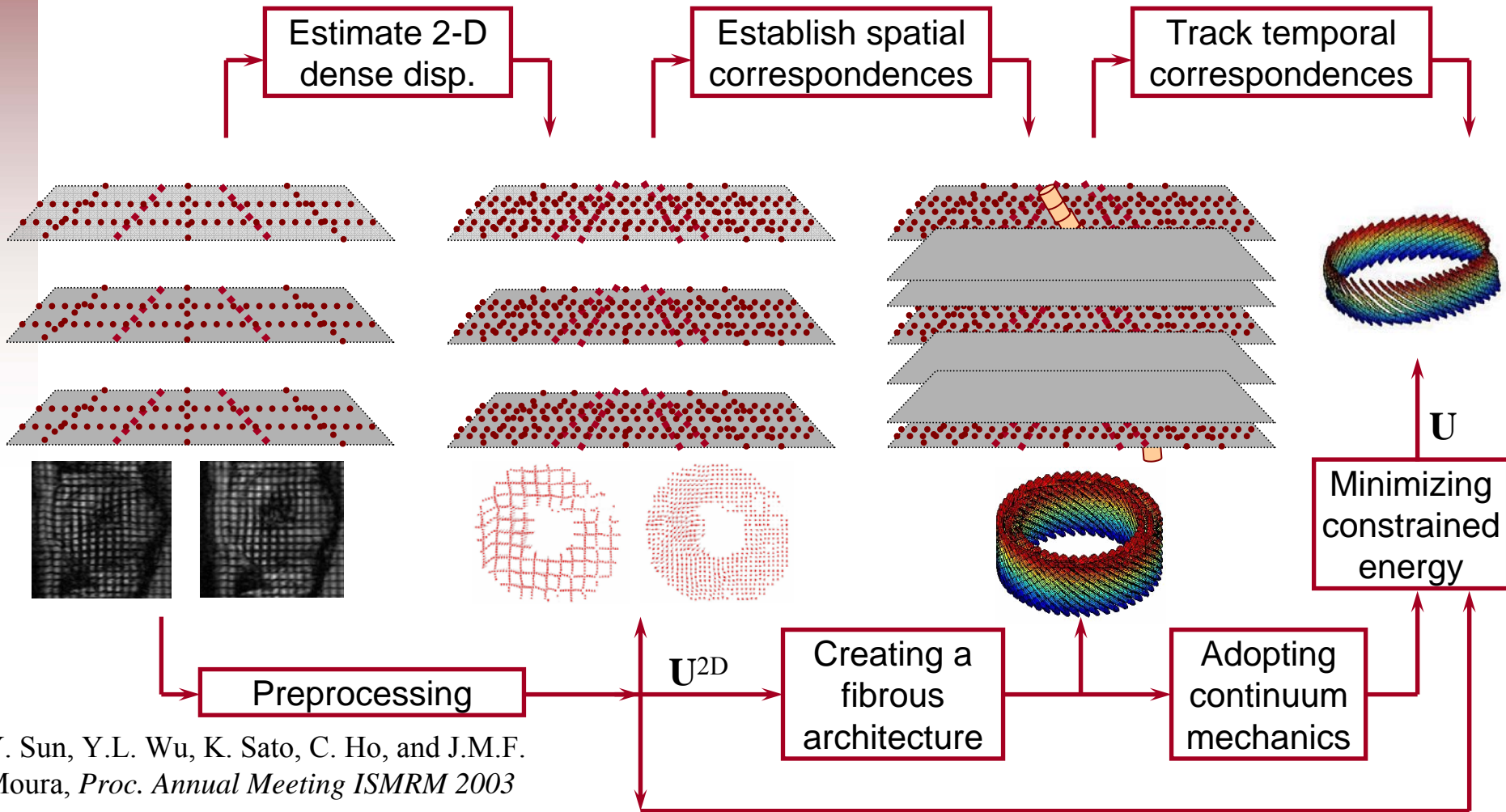


Beispiel: Cardiac MRI

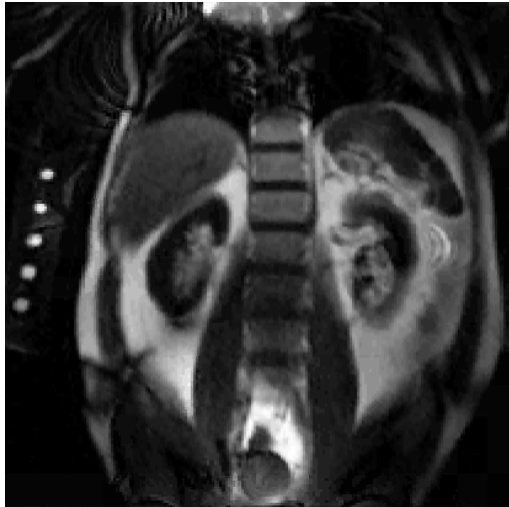


Ziel: 3-D-Movie von 2-D Daten

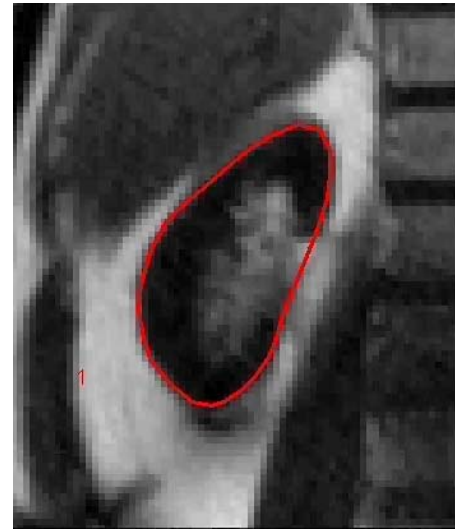
3-D Motion Estimation Procedure



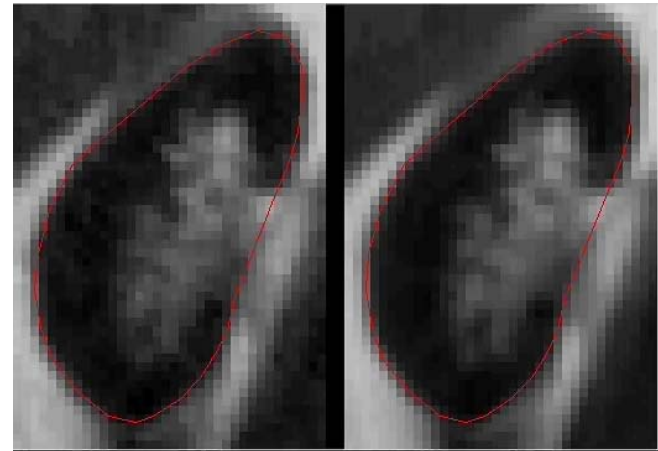
Beispiel: MRI



MRI



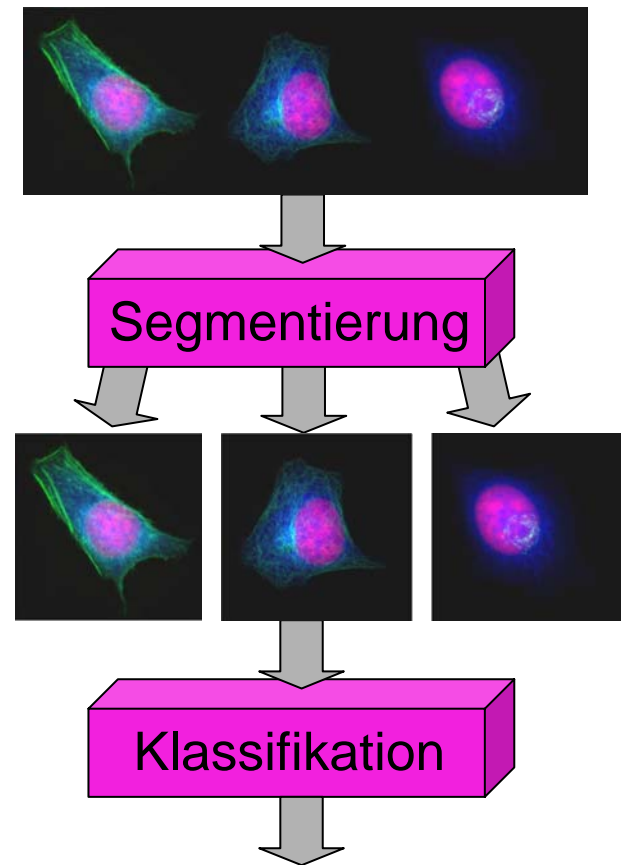
Tracking der Nieren



Herausrechnen der Bewegung

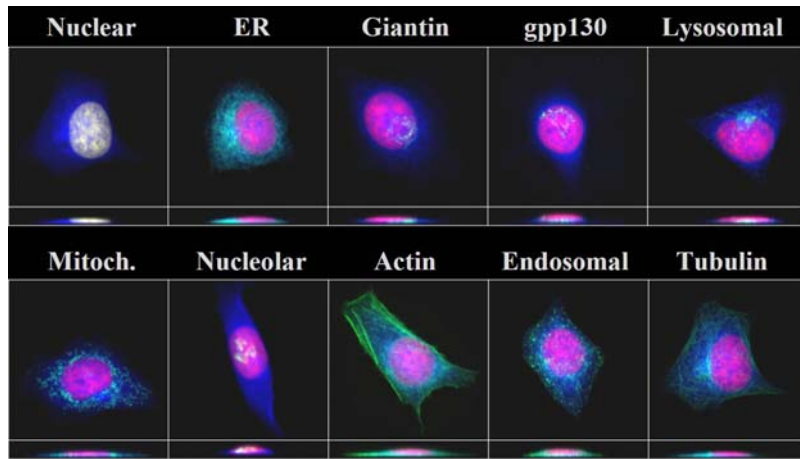
Beispiel: Biobildverarbeitung

- Ziel: Schnelle Identifikation von Proteinen mittels ihrer Verteilung in der Zelle
- Signalverarbeitung
 - Segmentierung
 - Klassifikation (Wavelets, Frames)



Das ist Tubulin!

Bilder



Beispiel: Computersehen



Suberbowl 2001 (Kanade et al.)

Beispiel: Kommunikation

- Ziel: Robustheit gegen Verluste

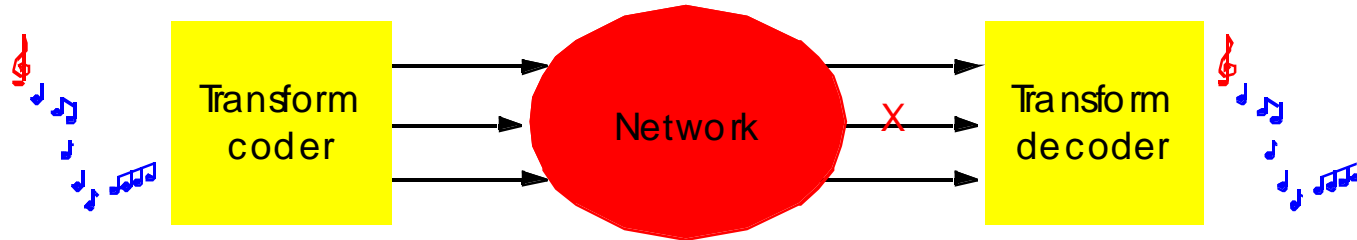


Foto-zur-Oma Problem

- Ziel: sende ein digitales Foto nach Italien
- Verfügbar: FedEx or regulärer Post "Kanal"
 - 📌 FedEx 99% verlässlich, Kosten \$39.99
 - 📌 Postal 80% verlässlich, Kosten \$3.40
- 1 Floppy pro Umschlag
- Foto braucht 2 Floppies (CDs sind noch nicht erfunden)

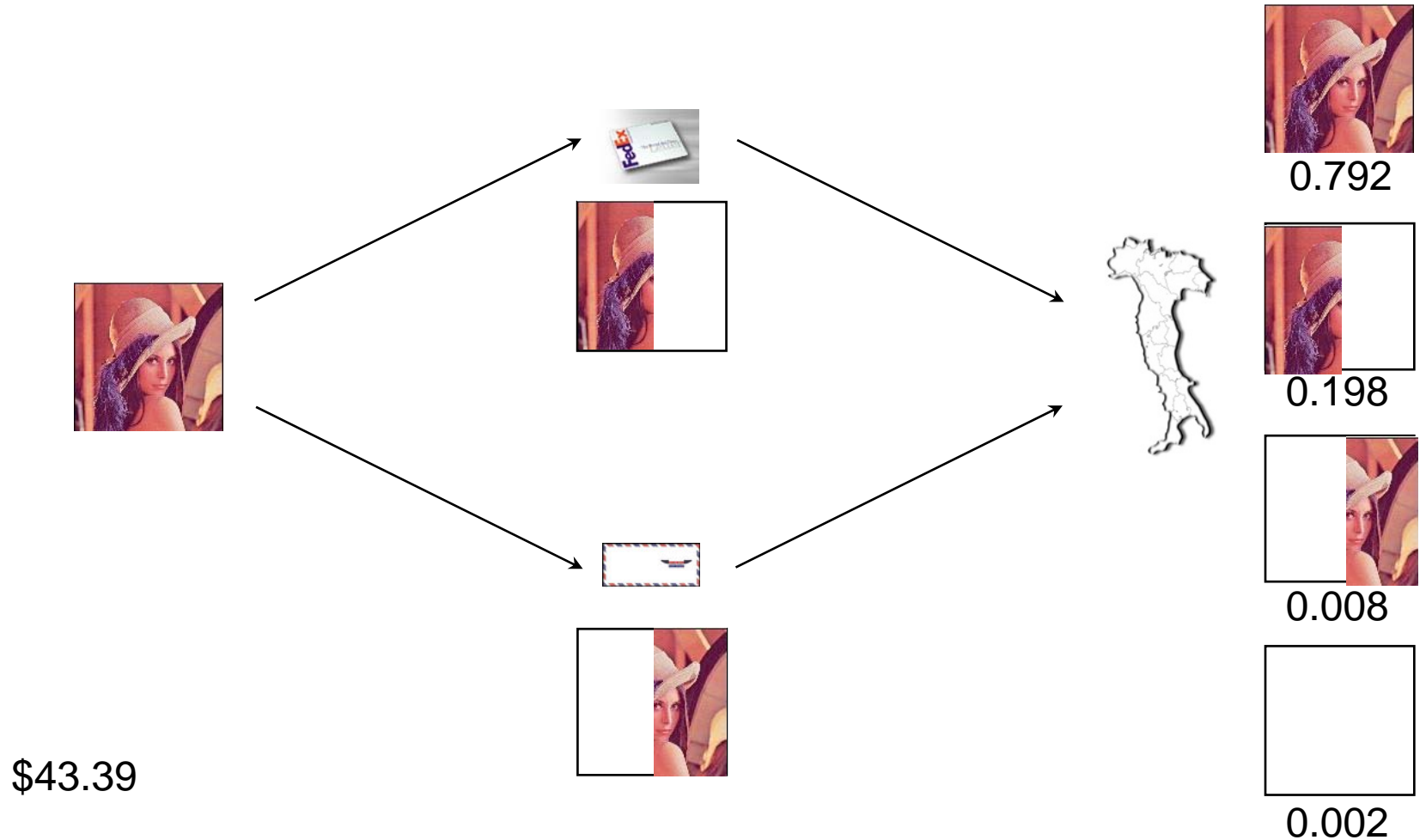


neue Freundin

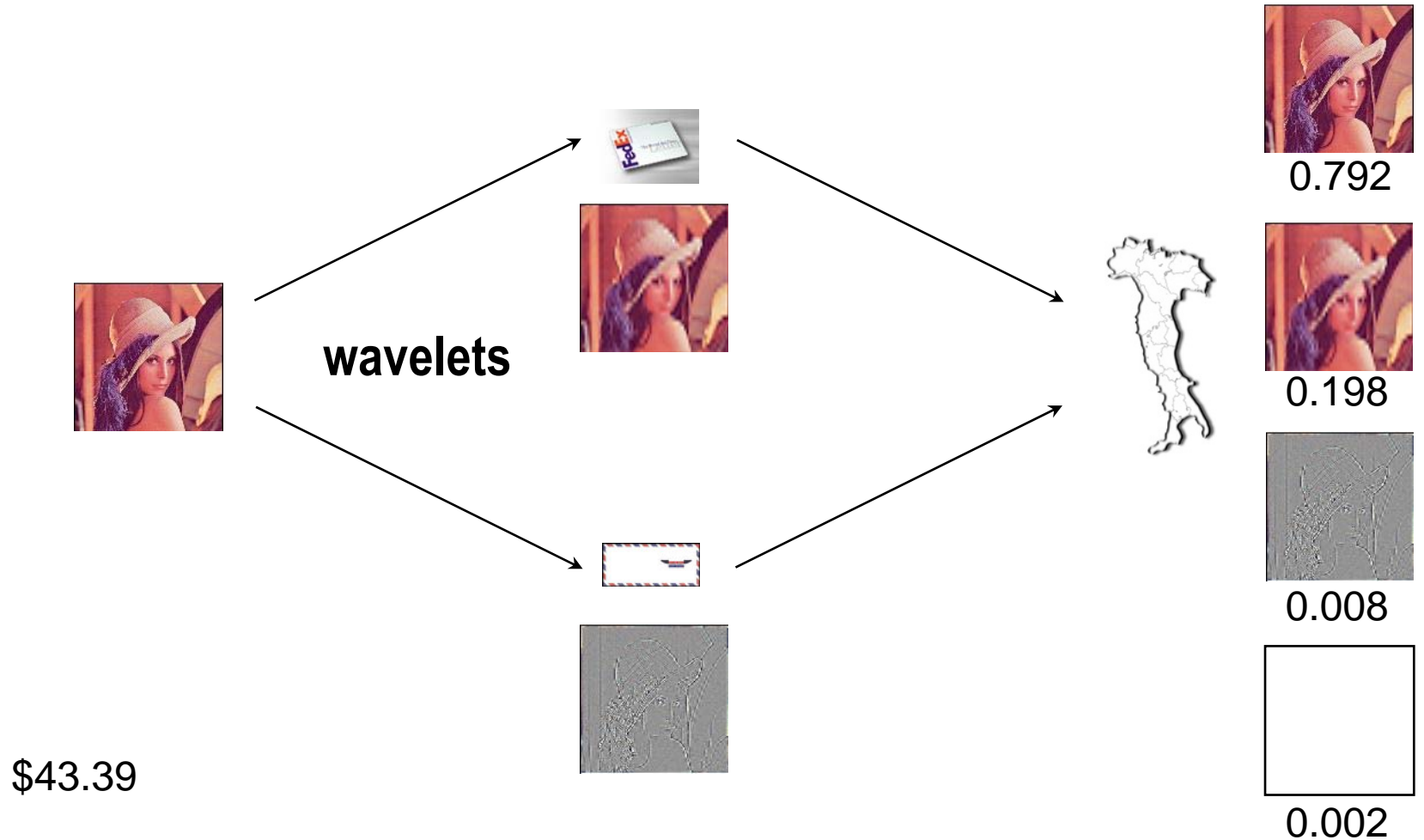


Oma wohnt in Italien

Heterogener Kanal (dumme Lösung)



Heterogener Kanal (kluge Lösung)



Zusammenfassung: Hauptbausteine der Signalverarbeitung

Filter:

FIR, IIR, Korrelation,
Filterbänke

Lineare Algebra:

Vektorsumme,
Matrix-Vektorprodukt,

...

Singulärwertzerlegung,
Matrixinversion,

...

Signaltransformationen:

DFT, DCT, wavelets, frames

Kodierung:

Huffman, arithmetisch,
Viterbi, LDPC

Implementierung

- **Unendliche Anforderungen an die Geschwindigkeit**
 - Sehr große Datensätze
 - Echtzeit
- **Vielfalt der Plattformen**
 - Hardware: ASIC, FPGA
 - Software
 - Ein- versus Mehrprozessorsystemen
 - Workstation versus eingebetteter Prozessor
 - Gleitpunkt versus Fixpunktarithmetik
 - Kombinierte Hardware/Software Plattformen
- **Problem: Implementierung schwierig, aufwendig, veraltet schnell**

In dieser Vorlesung: Einprozessor Workstations

Teil II: Algorithmentheorie

In dieser Vorlesung

- **Signaltransformationen**
(hauptsächlich diskrete Fouriertransformation)
- **FIR Filter**

Signaltransformationen

- Mathematisch: Matrix-Vektor Multiplikation

$$x \mapsto y = M \cdot x$$

Eingabevektor (Signal) \nearrow x
 Ausgabevektor (Signal) \nearrow y
 Transformation = Matrix \uparrow M

- Wird oft so geschrieben:

$$y_k = \sum_{l=0}^{n-1} m_{k,l} x_l, \quad M = [m_{k,l}]_{0 \leq k, l < n}$$

- Idee: Stelle Signal in einer anderen Basis dar für bessere Verarbeitung

Transformationen: Beispiele

$$\begin{aligned}
 \text{DFT}_n &= [e^{-2k\ell\pi i/n}]_{0 \leq k, \ell < n} \\
 \text{DCT-2}_n &= [\cos(k(2\ell + 1)\pi/2n)]_{0 \leq k, \ell < n}, \\
 \text{DCT-3}_n &= \text{DCT-2}_n^T \quad (\text{transpose}), \\
 \text{DCT-4}_n &= [\cos((2k + 1)(2\ell + 1)\pi/4n)]_{0 \leq k, \ell < n}, \\
 \text{IMDCT}_n &= [\cos((2k + 1)(2\ell + 1 + n)\pi/4n)]_{0 \leq k < 2n, 0 \leq \ell < n}, \\
 \text{RDFT}_n &= [r_{kl}]_{0 \leq k, \ell < n}, \quad r_{kl} = \begin{cases} \cos \frac{2\pi k\ell}{n}, & k \leq \lfloor \frac{n}{2} \rfloor \\ -\sin \frac{2\pi k\ell}{n}, & k > \lfloor \frac{n}{2} \rfloor \end{cases}, \\
 \text{WHT}_n &= \begin{bmatrix} \text{WHT}_{n/2} & \text{WHT}_{n/2} \\ \text{WHT}_{n/2} & -\text{WHT}_{n/2} \end{bmatrix}, \quad \text{WHT}_2 = \text{DFT}_2, \\
 \text{DHT} &= [\cos(2k\ell\pi/n) + \sin(2k\ell\pi/n)]_{0 \leq k, \ell < n}.
 \end{aligned}$$

Schnelle Algorithmen

- Transformation per Definition: $O(n^2)$
- Transformationen haben schnelle Algorithmen: $O(n \log(n))$
- Beispiel: FFT (fast Fourier transform)
- Diese existieren aufgrund von “Symmetrien” in den Matrizen

Die Geschichte schneller Algorithmen ...

- ...beginnt mit der FFT
- FFT ([Cooley-Tukey 1965](#)) wird oft bezeichnet als der Anfang der digitalen Signalverarbeitung
- **Geschichte:**
 - ~1805: FFT entdeckt von Gauss (nach [1])
(Fourier publizierte das Konzept Fourieranalyse in 1807!)
 - 1965: Neu entdeckt von Cooley-Tukey
 - 2002: James W. Cooley bekommt die IEEE Jack S. Kilby Signal Processing Medal
"For pioneering the Fast Fourier Transform (FFT) algorithm."

Wie wurden/werden Algorithmen Gefunden?

- Laaange die Matrix angucken bis man was sieht
- Beispiel: Cooley-Tukey FFT
- **Anderes Beispiel:** [G. Bi "Fast Algorithms for the Type-III DCT of Composite Sequence Lengths" IEEE Trans. SP 47\(7\) 1999](#)

Darstellung von Algorithmen

■ Darstellung von Algorithmen: 2 Schulen

- Folge von Summationen
- Strukturierte Matrixfaktorisierung: $M = M_1 \cdot M_2 \dots M_k$

■ Beispiel vorige Folie:

$$\text{DCT}_n = K_m^n \left(\bigoplus_{0 \leq i < k} \text{DCT}_m(r_i) \right) (\text{DCT}_k \otimes I_m) B_{n,k}$$

Beispiel C.-T. FFT: Was ist besser?

Das?:

$$\text{DFT}_n = L_{n_2}^n (I_{n_1} \otimes \text{DFT}_{n_2}) T_{n_1}^n (\text{DFT}_{n_1} \otimes I_{n_2})$$

Oder das?: $k = n_1 k_1 + k_2, j = n_2 j_1 + j_2$

$$y_{n_2 j_1 + j_2} = \sum_{k_1=0}^{n_1-1} \left(\omega_n^{j_2 k_1} \right) \left(\sum_{k_2=0}^{n_2-1} x_{n_1 k_2 + k_1} \omega_{n_2}^{j_2 k_2} \right) \omega_{n_1}^{j_1 k_1}$$

Rest an der Tafel

- Definition DFT
 - Interpretation mit dem Chinesischen Restesatz (CRS)
 - DFT Eigenschaften, Zirkulante, Toeplitzmatrizen
 - Herleitung der Cooley-Tukey FFT mit dem CRS
 - Andere DFT Algorithmen (FFTs)
 - Komplexität der DFT
-
- FIR Filter
 - FIR Filter Algorithmen

FFT Literatur

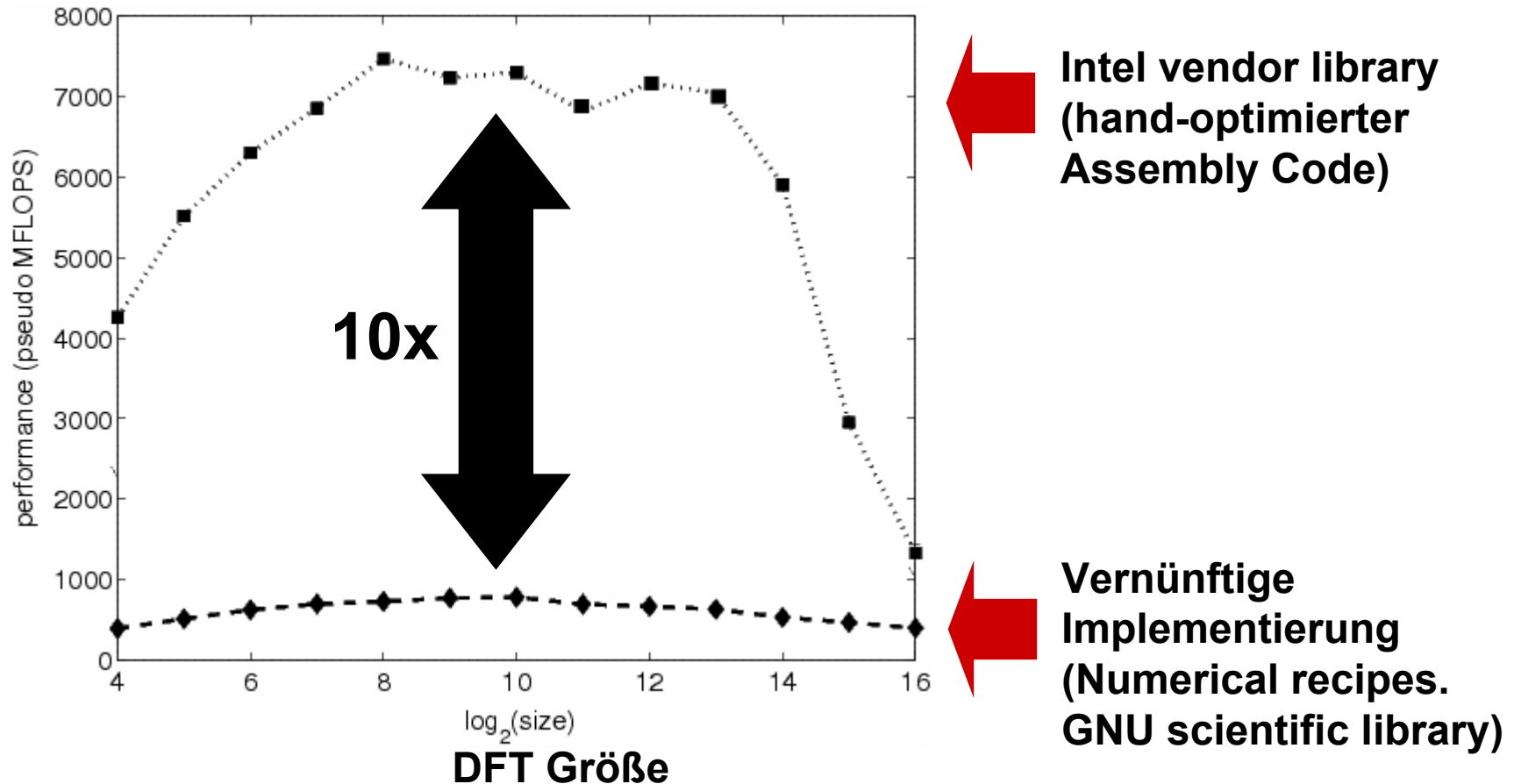
- **Nussbaumer (1982): Fast Fourier Transforms and Convolution Algorithms**
- **Van Loan (1992): Computational Frameworks for the Fast Fourier Transform**
- **Clausen/Baum (1993): Fast Fourier Transforms**
- **Tolimieri/An/Lu (1997): Algorithms for discrete Fourier transform and convolution**

Teil III: Implementierung

Implementierung: Überblick

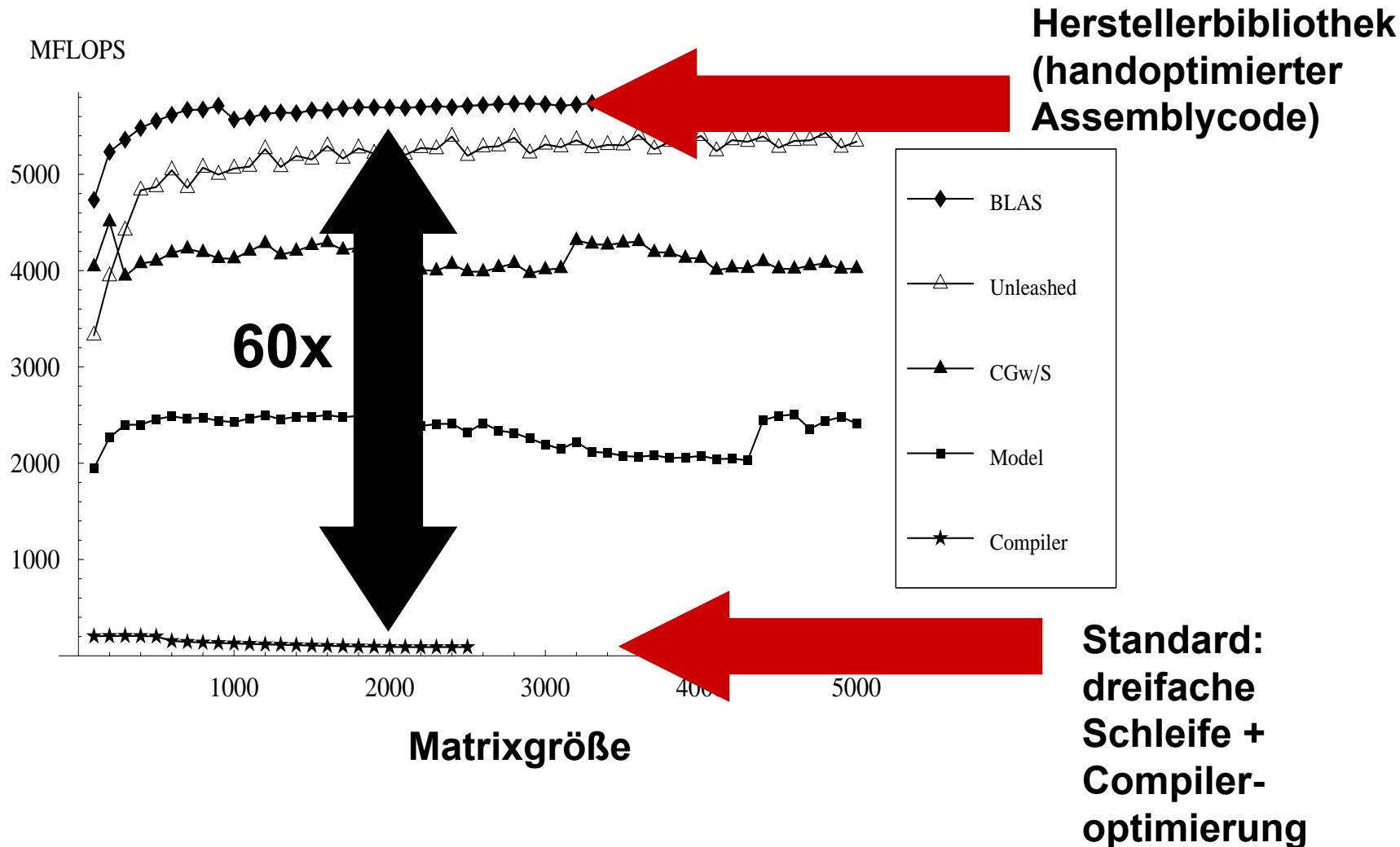
- **Das Problem: Moore's Law und Numerische Software**
- **Mikroarchitektur und Architektur**
- **Algorithmus, Implementierung und Mikroarchitektur:
Analyse der DFT und Herleitung von FFTW**
- **Softwaregenerierung und Experimente mit SPIRAL**
- **Zusammenfassung:
Wie schreibt man schnellen numerischen Code?**

Das Problem: Beispiel DFT auf Pentium 4



*Ok, aber die DFT ist ja auch kompliziert,
nehmen wir was einfacheres ...*

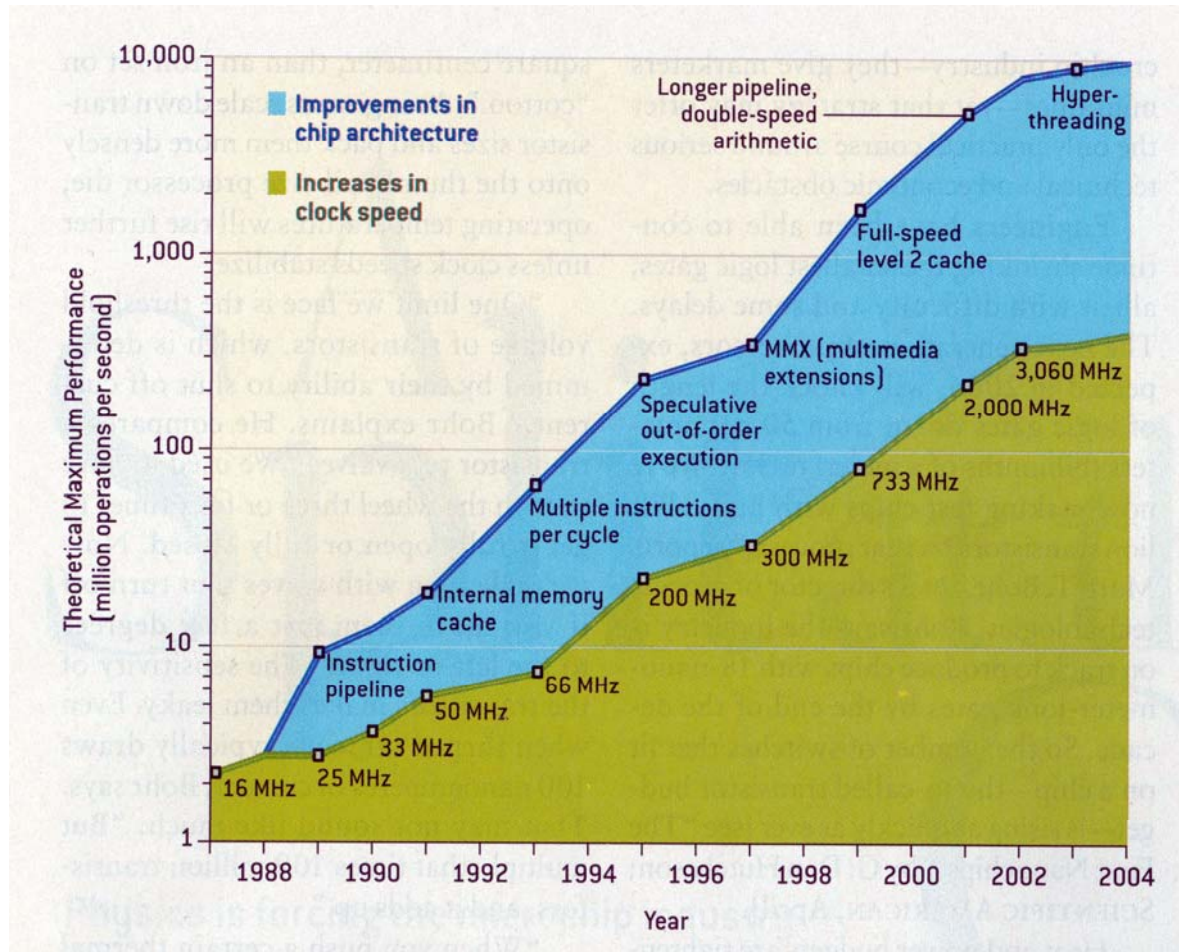
Das Problem: Matrix-matrix Multiplikation



Woran liegt das?

Moore's Law

- Moore's Law: exponentielles (x2 in ~18 Monaten) Wachstum Anzahl Transistoren/Chip



source: Scientific American, Nov 2004, p. 98

Aber alles hat seinen Preis ...

Moore's Law: Auswirkungen

■ Computer sind komplex

- Mehrstufige Speicherhierarchie
- Spezialinstruktionen
- Undokumentierte Hardwareoptimierungen

■ Konsequenzen:

- Laufzeit hängt nur grob von der Anzahl Operationen ab
- Laufzeitverhalten schwer zu verstehen
- Compilerentwicklung kann kaum mithalten
- **Der schnellste Code (inkl. Algorithmus) ist computerabhängig**
- **Es ist sehr schwierig schnellen Code zu schreiben**

■ Computer entwickeln sich schnell

- Hochoptimierter Code veraltet fast so schnell (schneller?) wie man ihn schreibt

***Wie produziert man schnellen Code?
und mit erträglichem Aufwand?***

Implementierung: Überblick

- Das Problem: Moore's Law und Numerische Software
- **Mikroarchitektur und Architektur**
- Algorithmus, Implementierung und Mikroarchitektur:
Analyse der DFT und Herleitung von FFTW
- Softwaregenerierung und Experimente mit SPIRAL
- Zusammenfassung:
Wie schreibt man schnellen numerischen Code?

Definitionen

- **Architektur:** (auch ISA) Die Teile eines Prozessordesigns die man verstehen muss um Assemblycode zu schreiben, z.B. die Instruktionsspezifikation und die Register.
Die Teile des Designs die nur die Implementierung betreffen, z.B. Cachegröße und Taktfrequenz gehören nicht dazu.
Beispiel: ia32 (x86), ia64
- **Mikroarchitektur:** Implementierung der Architektur
Beispiel: [Pentium4 Mikroarchitekturen](#)

Speicherhierarchie, Cachestruktur

- Mach's besser an der Tafel

SIMD (Signal Instruction Multiple Data)

Vektorinstruktionen

■ Was ist das?

- Erweiterung der ISA. Datentypen und Instruktionen zum parallelen Rechnen mit kurzen (2-8) Vektoren von Integers oder Floats



■ Warum gibt's sie?

- **Nützlich:** viele Anwendungen (z.B. Multimedia) haben den nötigen feinkörnigen Parallismus – Code potentiell schneller
- **Machbar:** Chipdesigner haben genügend Transistoren zum spielen

Überblick Vektor ISAs

Vendor	Name	ν -way	Precision	Processor
Intel	SSE	4-way	single	Pentium III Pentium 4
Intel	SSE2	2-way	double	Pentium 4
Intel	SSE3	4-way 2-way	single double	Pentium 4
Intel	IPF	2-way	single	Itanium Itanium 2
AMD	3DNow!	2-way	single	K6
AMD	Enhanced 3DNow!	2-way	single	K7, Athlon XP Athlon MP
AMD	3DNow! Professional	4-way	single	Athlon XP Athlon MP
AMD	AMD64	2-way 4-way 2-way	single single double	Athlon 64 Opteron
Motorola	AltiVec	4-way	single	MPC 74xx G4
IBM	AltiVec	4-way	single	PowerPC 970 G5
IBM	Double FPU	2-way	double	PowerPC 440 FP2

Die Evolution von Intel Vektorinstruktionen

■ MMX (1996, Pentium)

- Nur Integers, 64-bit geteilt in 2 x 32 bis 8 x 8
- MMX Register = Floatregister
- Nicht mehr sehr wichtig mit modernen Grafikkarten

■ SSE (1999, Pentium III)

- Obermenge von MMX
- 4-weg Floatops, einfache Genauigkeit
- 8 neue 128 Bit Register
- 100+ Instruktionen

■ SSE2 (2001, Pentium 4)

- Obermenge von SSE
- "MMX" auf SSE Registern, 2 x 64
- 2-weg Floatops, doppelte Genauigkeit, gleiche Register wie 4-weg einfach

■ SSE3 (2004, Pentium 4 Prescott)

- Obermenge von SSE2
- Mehr 2-weg Vektorinstruktionen speziell für komplexen Code

Wie benutzt man Vektorinstruktionen?

- **Voraussetzung: feinkörniger Parallelismus**

- **Am einfachsten: Benutze existierende Bibliothek**

- **Selbst schreiben:**
 - Compilervektorisierung
 - Erweiterung von C (Intrinsics, compilerspezifisch)
 - Assembly schreiben (alles oder Teile inline)

Probleme

- **Verlangt richtig ausgerichtete Daten**
- **Datenumordnungen zerstören leicht die Laufzeit**
- **Verlangt oft eine Anpassung des Algorithmus**
- **Optimierungen/Implementierungen maschinenabhängig**
- **Verlangt gutes Verständnis der ISA + Mikroarchitektur**

Beispiel: 4x4 Matrix-Vektor Multiplikation

$$\begin{bmatrix} x & x & x & x \\ x & x & x & x \\ x & x & x & x \\ x & x & x & x \end{bmatrix} \cdot \begin{bmatrix} a \\ a \\ a \\ a \end{bmatrix}$$

4-weg Vektoren (ab SSE) mit Intrinsics:

```

__m128 BLAS_II_4x4 (__m128 *matrix, __m128 vector) {
    __m128 t[4];

    for (int i=0; i<4; i++)
        t[i] = _mm_mul_ps(matrix[i], vector);

    _MM_TRANSPOSE_PS(t[0], t[1], t[2], t[3]);

    return _mm_add_ps(_mm_add_ps(t[0], t[1]),
                      _mm_add_ps(t[2], t[3])); }

```

4 vmults

1 Transposition

3 vadds

Beispiel: Komplexe Multiplikation SSE3

**Komplexer C99 Code + Compilervektorisierung
funktioniert sehr gut**

**Komplexer Code hat natürlichen
2-weg Vektorparallelismus**

Und so sieht der Assemblycode aus

SSE3:

```

movapd   xmm0, XMMWORD PTR A
movddup  xmm2, QWORD PTR B
mulpd    xmm2, xmm0
movddup  xmm1, QWORD PTR B+8
shufpd   xmm0, xmm0, 1
mulpd    xmm1, xmm0
addsubpd xmm2, xmm1
movapd   XMMWORD PTR C, xmm2

```

SSE2:

```

movsd    xmm3, QWORD PTR A
movapd   xmm4, xmm3
movsd    xmm5, QWORD PTR A+8
movapd   xmm0, xmm5
movsd    xmm1, QWORD PTR B
mulsd    xmm4, xmm1
mulsd    xmm5, xmm1
movsd    xmm2, QWORD PTR B+8
mulsd    xmm0, xmm2
mulsd    xmm3, xmm2
subsd    xmm4, xmm0
movsd    QWORD PTR C, xmm4
addsd    xmm5, xmm3
movsd    QWORD PTR C, xmm5

```

In SSE2 ist Skalarcode besser

Wie schreibt man guten Vektorcode?

- **Nimm den “richtigen” Algorithmus und die “richtigen” Datenstrukturen**
 - Feinkörniger Parallelismus
 - Korrekte Ausrichtung im Speicher
 - Zusammenhängende Felder (arrays)
- **Verwende eine guten Compiler (z.B. vom Hersteller)**
- **Zuerst: probiere Compilervektorisierung**
 - Richtige Optionen + Pragmas (Information an den Compiler: z.B. Datenausrichtung, Schleifenunabhängigkeit)
 - Überprüfe Assemblycode und Laufzeit
- **Wenn nötig: schreibe selbst Vektorcode**
 - Zuerst die wichtigste (teuerste) Subroutine
 - Benutze Intrinsics, kein Inlineassembly
 - Wichtig: Verstehe ISA

Implementierung: Überblick

- Das Problem: Moore's Law und Numerische Software
- Mikroarchitektur und Architektur
- **Algorithmus, Implementierung und Mikroarchitektur:
Analyse der DFT und Herleitung von FFTW**
- Softwaregenerierung und Experimente mit SPIRAL
- Zusammenfassung:
Wie schreibt man schnellen numerischen Code?

Auch besser an der Tafel

FFTW Codeletbeispiel

- [No-twiddle-codelet, size 8](#)

FFTW Zusammenfassung

- **Eine** Bibliothek für FFTs beliebiger Größe, Dimension, reell und komplex
- Sehr schnell (Adaptierung + benutzt Vektorcode)
- Stabil, einfache Installation (configure – make)
- Alles in allem: tolle Software, sehr populär

Wie evaluiert man Code?

- **Vergleiche gegen die beste erhältliche Software**

- Stelle sicher daß die gleiche Zeitmessmethode verwendet wird

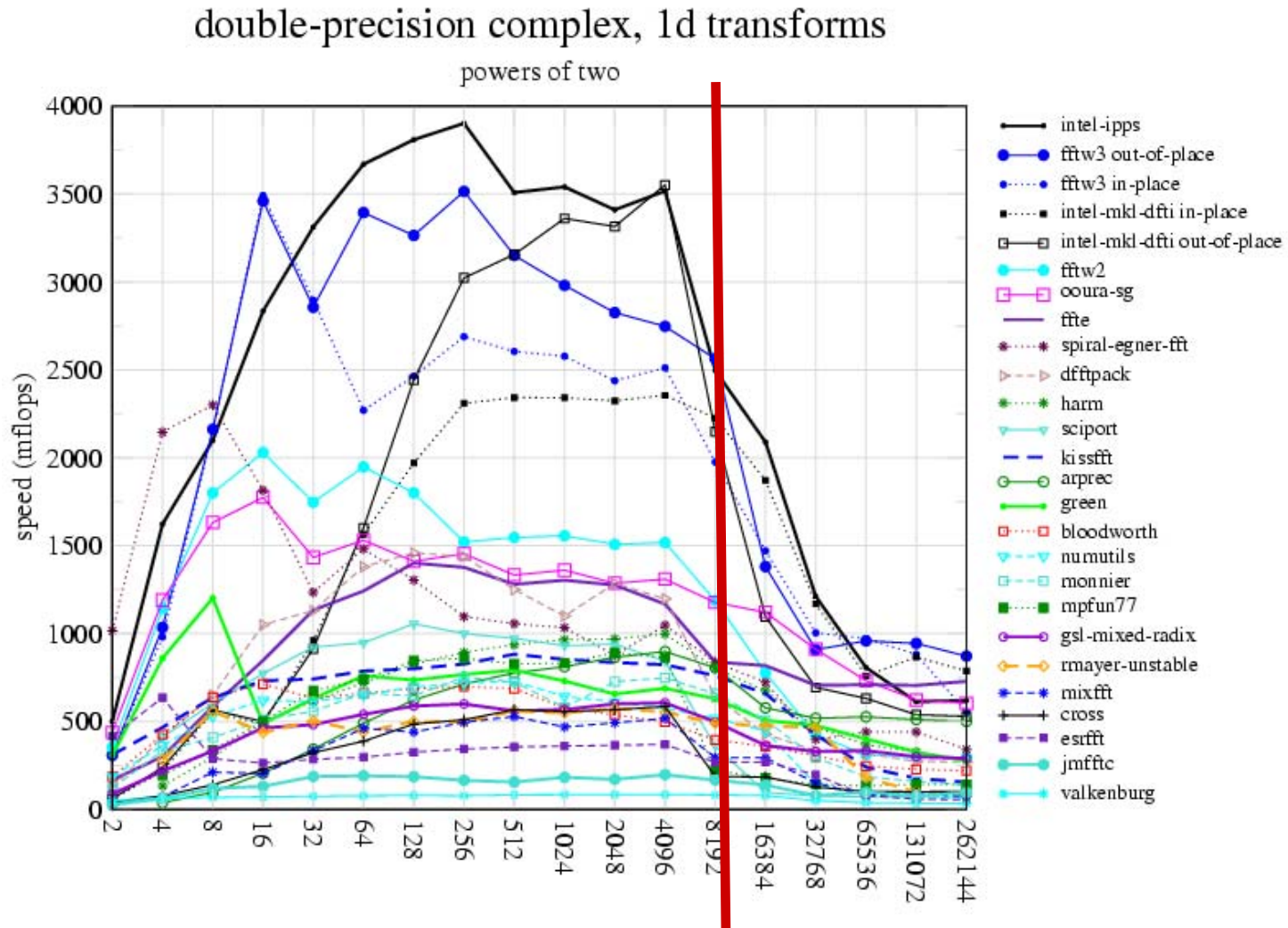
- **Berechne die Leistung (in MFLOPS)**

$$MFLOPS = \text{Anzahl Operationen} / (10^6 \cdot \text{Laufzeit [s]})$$

und vergleiche zur Maximalleistung der Maschine

- Erfordert die Berechnung/Messung der Anzahl der Operationen
- Maximalleistung aus dem Manual oder die richtigen Leute fragen
- Vorsicht: Code mit höherer Leistung kann langsamer sein

FFTW Benchmarks, Pentium 4

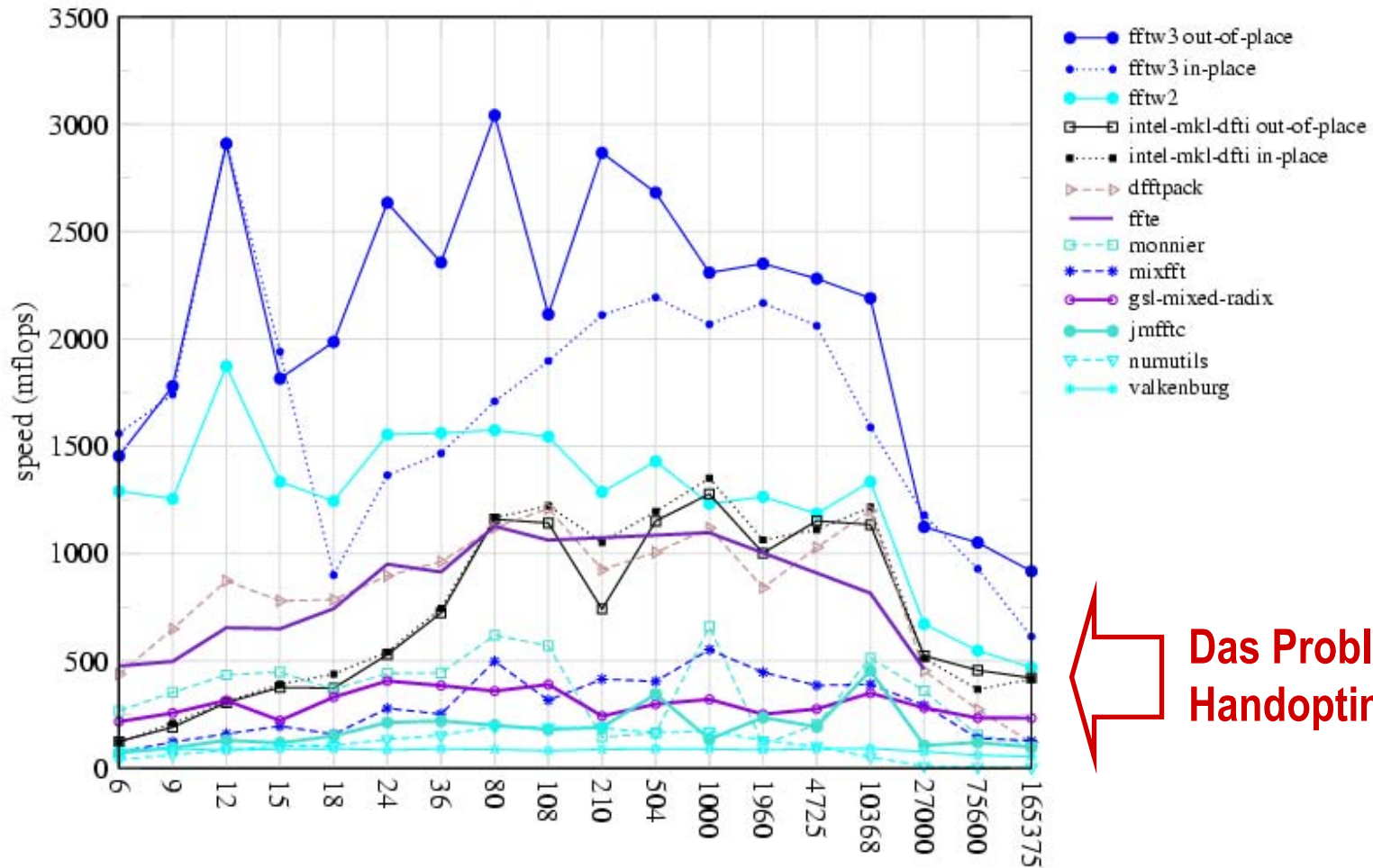


L2 Cacheknick

Quelle: www.fftw.org/speed

FFTW Benchmarks, Pentium 4

double-precision complex, 1d transforms
non-powers of two



Das Problem mit Handoptimierung

Andere FFTs: CT, GT, Rader, Bluestein

recht billig

$$\text{DFT}_{km} \rightarrow (\text{DFT}_k \otimes \text{I}_m) \text{T}_m^n (\text{I}_k \otimes \text{DFT}_m) \text{L}_k^n$$

teuer

$$\text{DFT}_{km} \rightarrow P_n (\text{DFT}_k \otimes \text{DFT}_m) Q_n, \quad \text{gcd}(k, m) = 1$$

$$\text{DFT}_p \rightarrow R_p^T (\text{I}_1 \oplus \text{DFT}_{p-1}) D_p (\text{I}_1 \oplus \text{DFT}_{p-1}) R_p$$

sehr teuer

$$\text{DFT}_n \rightarrow D_m \text{DFT}_m D'_m \text{DFT}_m D''_m, \quad m > 2n$$

Probleme:

- Verschmelzung der Permutationen
- erst kürzlich (teilweise) gelöst (Paper eingereicht)
- können nicht einfach in FFTW integriert werden

Implementierung: Überblick

- **Das Problem: Moore's Law und Numerische Software**
- **Mikroarchitektur und Architektur**
- **Algorithmus, Implementierung und Mikroarchitektur:
Analyse der DFT und Herleitung von FFTW**
- **Softwaregenerierung und Experimente mit SPIRAL**
- **Zusammenfassung:
Wie schreibt man schnellen numerischen Code?**

SPIRAL Information

Seit: Mitte 1998

Gefördert von:

DARPA

NSF ACR-0234293

NSF ITR/NGS-0325687

und:

Cylab, CMU

Austrian Science Fund

Intel

ITRI, Taiwan

ENSCO, Inc.

www.spiral.net

SPIRAL Team: Faculty

James C. Hoe (ECE, CMU)

Jeremy Johnson (CS, Drexel)

José M. F. Moura (ECE, CMU)

David Padua (CS, UIUC)

Markus Püschel (ECE, CMU)

Manuela Veloso (CS, CMU)

Robert W. Johnson (Quarry Comp. Inc.)

Studenten:

Bryan W. Singer (CS, CMU)

Jianxin Xiong (CS, UIUC)

Franz Franchetti (ECE, CMU, vorher TU Wien)

Aca Gacic (ECE, CMU)

Yevgen Voronenko (ECE, CMU)

Anthony Breitzman (CS, Drexel)

Kang Chen (CS, Drexel)

Pinit Kumhom (ECE, Drexel)

Adam Zelinski (ECE, CMU)

Peter Tummeltshammer (CS, TU Wien)

...

Kollaboration mit:

Christoph Überhuber (Math, TU Wien)

Spiral

- **Codegenerierung für Signaltransformationen (DFT, DCT, Filter, Wavelets, ...)**
- **Automatische Optimierung und Adaptierung von Algorithmus und Implementierung**
- **Verschiedene Codetypen (skalar, Vector, FMA, Fixpunkt, multiplikationsfrei, ...)**
- **Ermöglicht vielseitige Experimente um verschiedene Algorithmen zu testen**

SPIRAL: Ziele

- Ein flexibles, erweiterbares Codegenerierungssystem dass nicht veraltet (soweit möglich) für eine ganze Klasse von Algorithmen
- Suche nach Antworten auf die Fragen:
 - Bis zu welchem Grad kann man Handcodierung abschaffen?
 - Wie bringt man dem Computer Algorithmenwissen bei?

Codegenerierung und Adaptierung als Optimierungsproblem

T eine Signaltransformation

P die Zielplattform

$\mathcal{I} = \mathcal{I}(T, P)$ Menge aller Implementierungen von **T** auf **P**

$C = C(T, I, P)$ Laufzeit von der Implementierung **I** von **T** auf **P**

Die Implementierung von **T** adaptiert zu **P** ist:

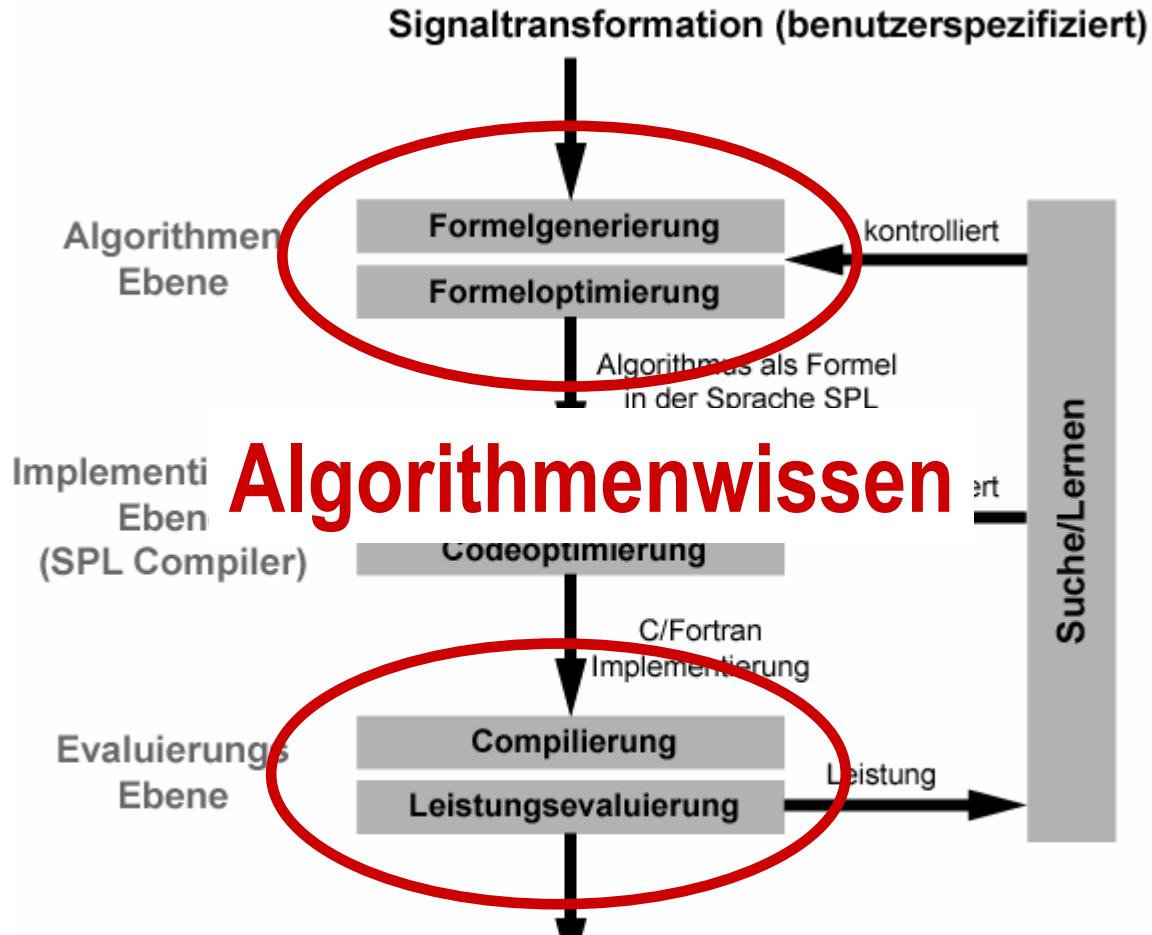
$$\hat{I} = \hat{I}(P) = \arg \min_{I \in \mathcal{I}(P)} C(T, P, I)$$

Probleme:

- Wie beschreibt und generiert man die Menge der Implementierungen?
- Wie minimiert man effizient C ?

SPIRAL nützt die spezifische Struktur von Transformationsalgorithmen um die Optimierung automatisch durchzuführen

Spiral's Architektur

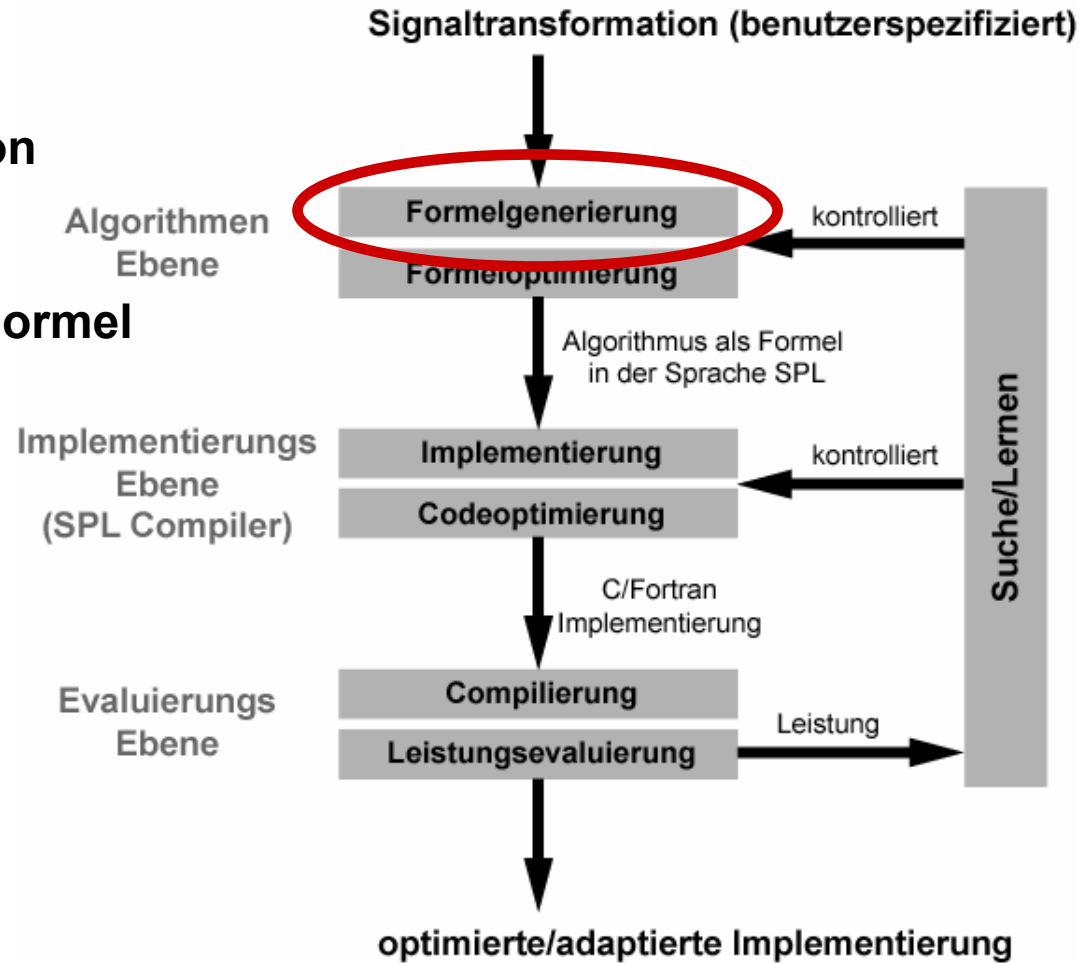


Information über den Computer

Von Transformation zum Algorithmus (Formel)

Eingabe:
Transformationsspezifikation

Ausgabe:
Schneller Algorithmus als Formel



Algorithmen: Beispiel DFT der Größe 4

Cooley/Tukey FFT (Größe 4):

$$\begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & i & -1 & -i \\ 1 & -1 & 1 & -1 \\ 1 & -i & -1 & i \end{bmatrix} = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & -1 & 0 \\ 0 & 1 & 0 & -1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & i \end{bmatrix} \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & -1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

Fouriertransformation

Diagonalmatrix (twiddles)

$$DFT_4 = \underbrace{(DFT_2 \otimes I_2)}_{\text{Kroneckerprodukt}} \cdot \underbrace{T_2^4}_{\text{Identität}} \cdot \underbrace{(I_2 \otimes DFT_2)}_{\text{Identität}} \cdot \underbrace{L_2^4}_{\text{Permutation}}$$

- Mathematische Darstellung zeigt die Struktur des Algorithmus:
SPL = Signal Processing Language
- Enthält die notwendige Information um Code zu generieren

SPL: Definition (BNF)

$\langle \text{spl} \rangle ::= \langle \text{generic} \rangle \mid \langle \text{symbol} \rangle \mid \langle \text{transform} \rangle \mid$	
$\langle \text{spl} \rangle \dots \langle \text{spl} \rangle \mid$	(product)
$\langle \text{spl} \rangle \oplus \dots \oplus \langle \text{spl} \rangle \mid$	(direct sum)
$\langle \text{spl} \rangle \otimes \dots \otimes \langle \text{spl} \rangle \mid$	(tensor product)
$\mathbf{I}_n \otimes_k \langle \text{spl} \rangle \mid \mathbf{I}_n \otimes^k \langle \text{spl} \rangle \mid$	(overlapped tensor product)
$\langle \text{spl} \rangle \mid$	(conversion to real)
\dots	
$\langle \text{generic} \rangle ::= \text{diag}(a_0, \dots, a_{n-1}) \mid \dots$	
$\langle \text{symbol} \rangle ::= \mathbf{I}_n \mid \mathbf{J}_n \mid \mathbf{L}_k^n \mid \mathbf{R}_\alpha \mid \mathbf{F}_2 \mid \dots$	
$\langle \text{transform} \rangle ::= \text{DFT}_n \mid \text{WHT}_n \mid \text{DCT-2}_n \mid \text{Filt}_n(h[z]) \mid \dots$	

Einige Definitionen:

$$A \oplus B = \begin{bmatrix} A & \\ & B \end{bmatrix}$$

$$A \otimes B = [a_{k,l} B], \quad \text{where } A = [a_{k,l}]$$

$$\mathbf{F}_2 = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$$

$$\mathbf{R}_\alpha = \begin{bmatrix} \cos \alpha & \sin \alpha \\ -\sin \alpha & \cos \alpha \end{bmatrix}$$

$$I_S \otimes_k A = \begin{bmatrix} \boxed{A} & & & \\ & \boxed{A} & & \\ & & \dots & \\ & & & \boxed{A} \end{bmatrix}$$

$$I_S \otimes^k A = \begin{bmatrix} \boxed{A} & & & \\ & \boxed{A} & & \\ & & \boxed{A} & \\ & & & \dots \\ & & & & \boxed{A} \end{bmatrix}$$

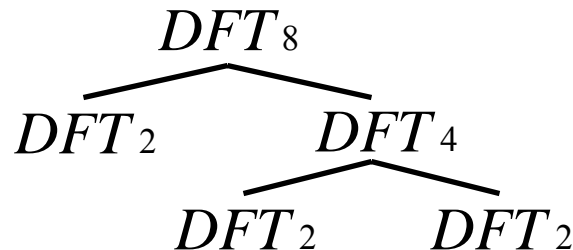
Algorithmen: Spiral's Terminologie

Transformation DFT_n parametrisierte Matrix

Regel $DFT_{nm} \rightarrow (DFT_n \otimes I_m) \cdot D \cdot (I_n \otimes DFT_m) \cdot P$

- Zerlegungsregel (teile-und-herrsche)
- Produkt dünnbesetzter Matrizen

Regelbaum



- rekursive Anwendung von Regeln
- definiert einen Algorithmus
- effiziente Darstellung
- einfache Manipulierung

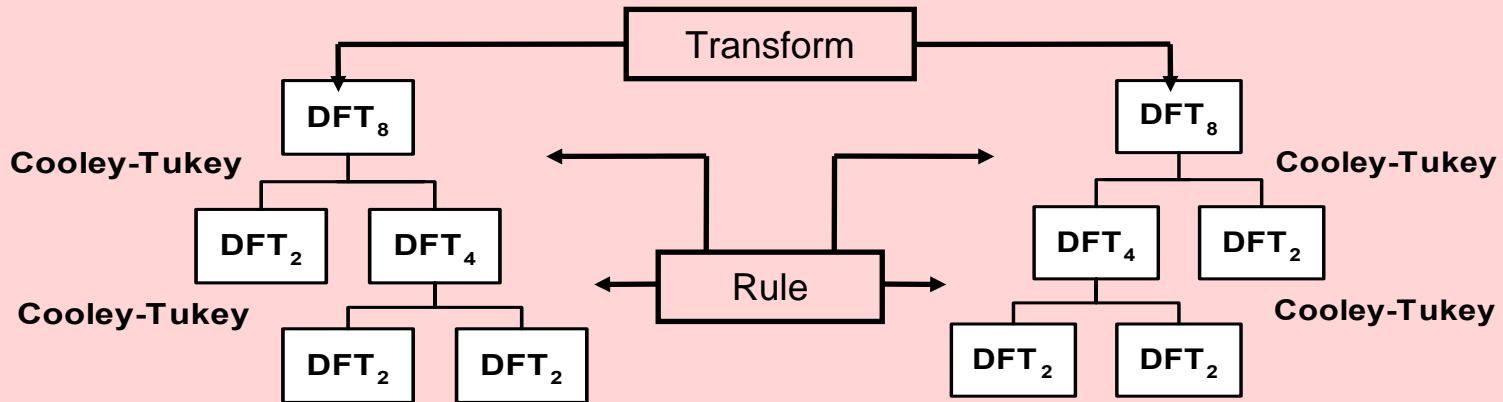
Formel

$$DFT_8 = (F_2 \otimes I_4) \cdot D \cdot (I_2 \otimes (I_2 \otimes F_2 \cdots)) \cdot P$$

- wenig Symbole
- definiert einen Algorithmus
- kann in Code übersetzt werden

Regelbäume = Formeln = Algorithmen

Regelbäume



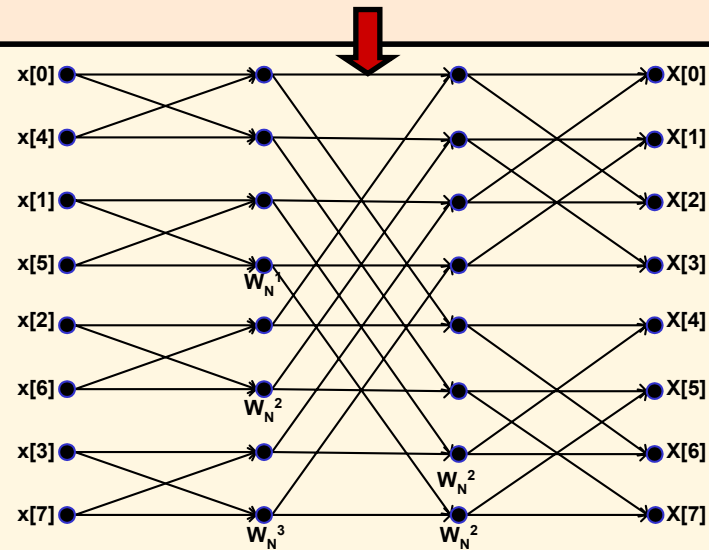
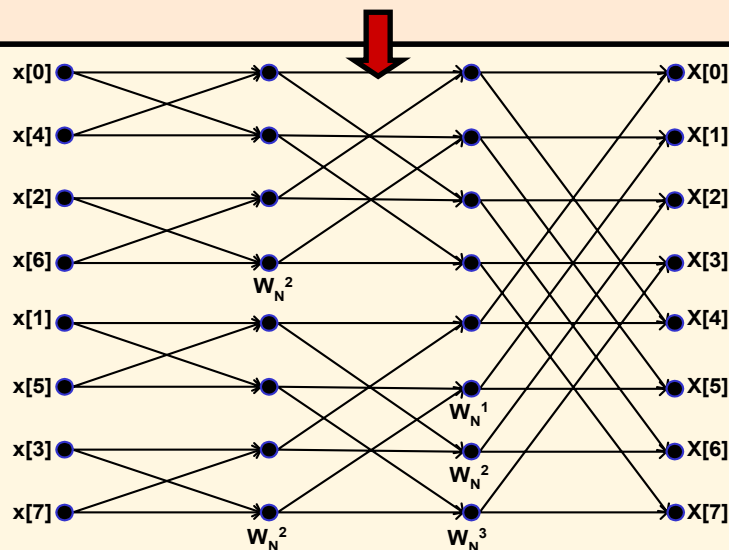
Formeln

$$\text{DFT}_8 = (\text{DFT}_2 \otimes I_4) \cdot T_2^8 \cdot (I_2 \otimes \text{DFT}_4) \cdot L_4^8$$

$$\text{DFT}_8 = (\text{DFT}_4 \otimes I_2) \cdot T_4^8 \cdot (I_4 \otimes \text{DFT}_2) \cdot L_2^8$$

$$(\text{DFT}_2 \otimes I_2) \cdot T_2^4 \cdot (I_2 \otimes \text{DFT}_2) \cdot L_2^4$$

Algorithmen



Höhere Abstraktionsebene

Einige Transformationen

$$\text{DCT-2}_n = \left[\cos(k(2l + 1)\pi/2n) \right]_{0 \leq k, l < n},$$

$$\text{DCT-3}_n = \text{DCT-2}_n^T \quad (\text{transpose}),$$

$$\text{DCT-4}_n = \left[\cos((2k + 1)(2l + 1)\pi/4n) \right]_{0 \leq k, l < n},$$

$$\text{IMDCT}_n = \left[\cos((2k + 1)(2l + 1 + n)\pi/4n) \right]_{0 \leq k < 2n, 0 \leq l < n},$$

$$\text{RDFT}_n = [r_{kl}]_{0 \leq k, l < n}, \quad r_{kl} = \begin{cases} \cos \frac{2\pi kl}{n}, & k \leq \lfloor \frac{n}{2} \rfloor \\ -\sin \frac{2\pi kl}{n}, & k > \lfloor \frac{n}{2} \rfloor \end{cases},$$

$$\text{WHT}_n = \begin{bmatrix} \text{WHT}_{n/2} & \text{WHT}_{n/2} \\ \text{WHT}_{n/2} & -\text{WHT}_{n/2} \end{bmatrix}, \quad \text{WHT}_2 = \text{DFT}_2,$$

$$\text{DHT} = \left[\cos(2kl\pi/n) + \sin(2kl\pi/n) \right]_{0 \leq k, l < n}.$$

Spiral “kennt” zur Zeit 36 Transformationen

Filter

$$h(z) = h_l z^{-l} + \dots + h_0 z + \dots + h_{-r} z^r$$

FIR Filtertransformation

$$\mathbf{Filt}_n(h(z)) = \begin{bmatrix} h_l & \dots & h_{-r} & & & \\ & h_l & \dots & h_{-r} & & \\ & & \ddots & & \ddots & \\ & & & h_l & \dots & h_{-r} \end{bmatrix}$$

Faltungstransformation

$$\mathbf{Conv}_n(h(z)) = \begin{bmatrix} h_{-r} & & & & & \\ & \ddots & & & & \\ & & h_l & \dots & h_{-r} & \\ & & & \ddots & & \ddots \\ & & & & h_l & \dots & h_{-r} \\ & & & & & \ddots & \vdots \\ & & & & & & h_l \end{bmatrix}$$

Triviale Regel:

$$\mathbf{Filt}_n(h(z)) \rightarrow \mathbf{I}_s \otimes_{l+r} \mathbf{Filt}_1(h(z))$$

Zirkulante

$$\mathbf{C}_n(h(z)) = [h_{(i-j) \bmod n}]_{n \times n}$$

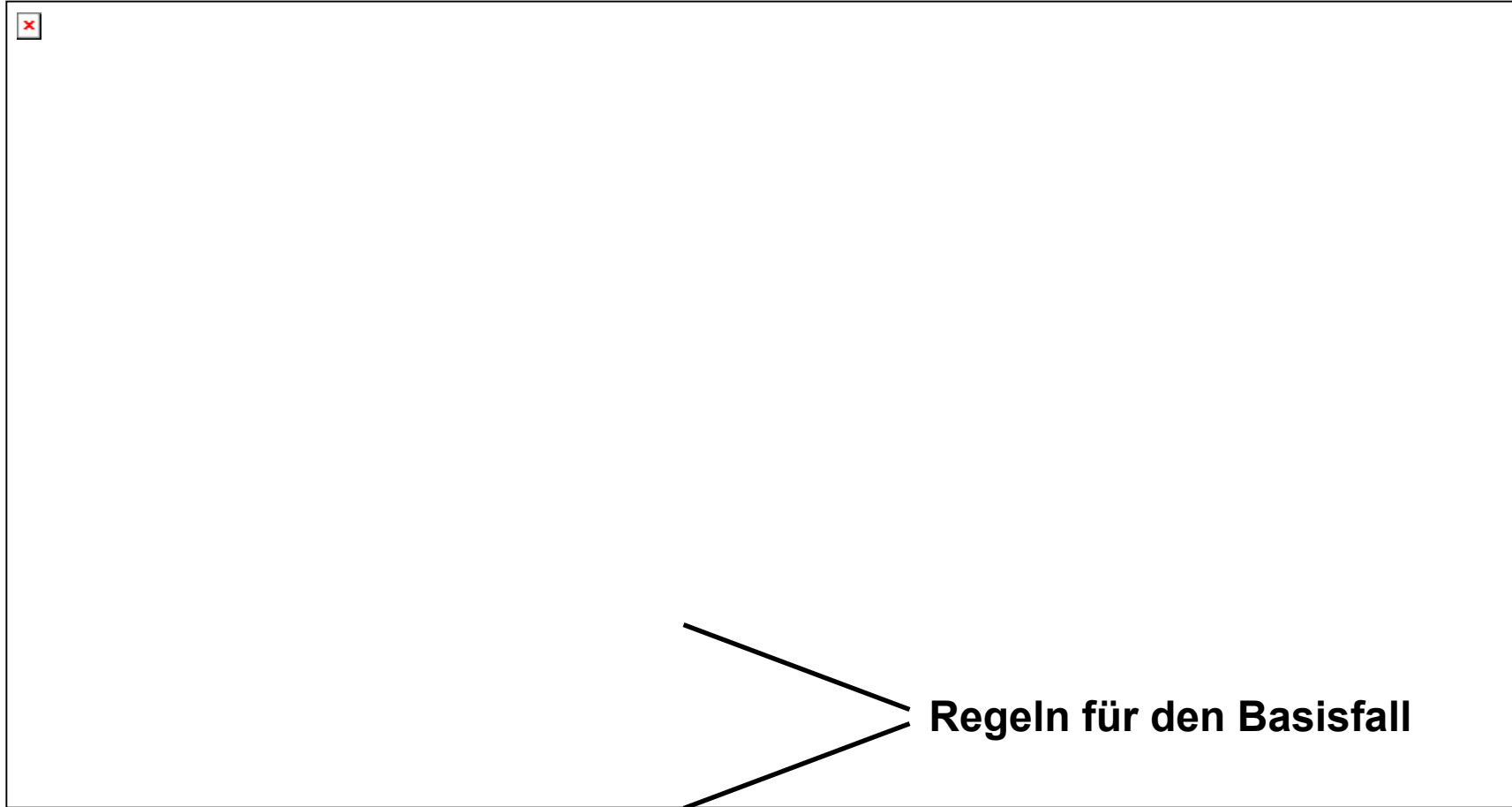
$$\mathbf{C}_n(h(z)) = \mathbf{Filt}_n^{\text{per}}(h(z))$$

Toeplitz

$$\mathbf{T}_n(h(z)) = [h_{i-j}]_{n \times n}$$

$$\mathbf{T}_n(h(z)) = \mathbf{Filt}_n^{\text{zero}}(h(z))$$

Einige Zerlegungsregeln



Spiral "kennt" 100+ Regeln

Regeln für Filter

“Blocken” Regel

$$\text{Filt}_n(h(z)) \rightarrow \mathbf{I}_{\lfloor \frac{n}{b} \rfloor} \otimes_{l+r} \left(\left[\begin{array}{c} | \\ | \\ \vdots \\ | \end{array} \right]_{i=0}^{\lfloor \frac{l+r}{b} \rfloor} \mathbf{T}_b(h(z)z^{l-ib}) \oplus^k \mathbf{T}_k(h(z)z^{l-\lfloor \frac{l+r}{b} \rfloor b-k}) \right)$$

Karatsubaregel

$$\text{Filt}_n(h(z)) \rightarrow \mathcal{L}_{\frac{n}{2}}^n \text{Filt}_{\frac{n}{2}} \left(\begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \end{bmatrix} \right) \cdot \text{Filt}_{\frac{n}{2}} \left(\begin{bmatrix} h_0(z) \\ h_1(z) \\ h_0(z) + h_1(z) \end{bmatrix} \right) \cdot \text{Filt}_{\frac{n}{2} + \frac{r+l-1}{2}} \left(\begin{bmatrix} 1 & -1 \\ z & -1 \\ 0 & 1 \end{bmatrix} \right) \cdot \mathcal{L}_2^{n+r+l}$$

“Einbetten in Zirkulante” Regel

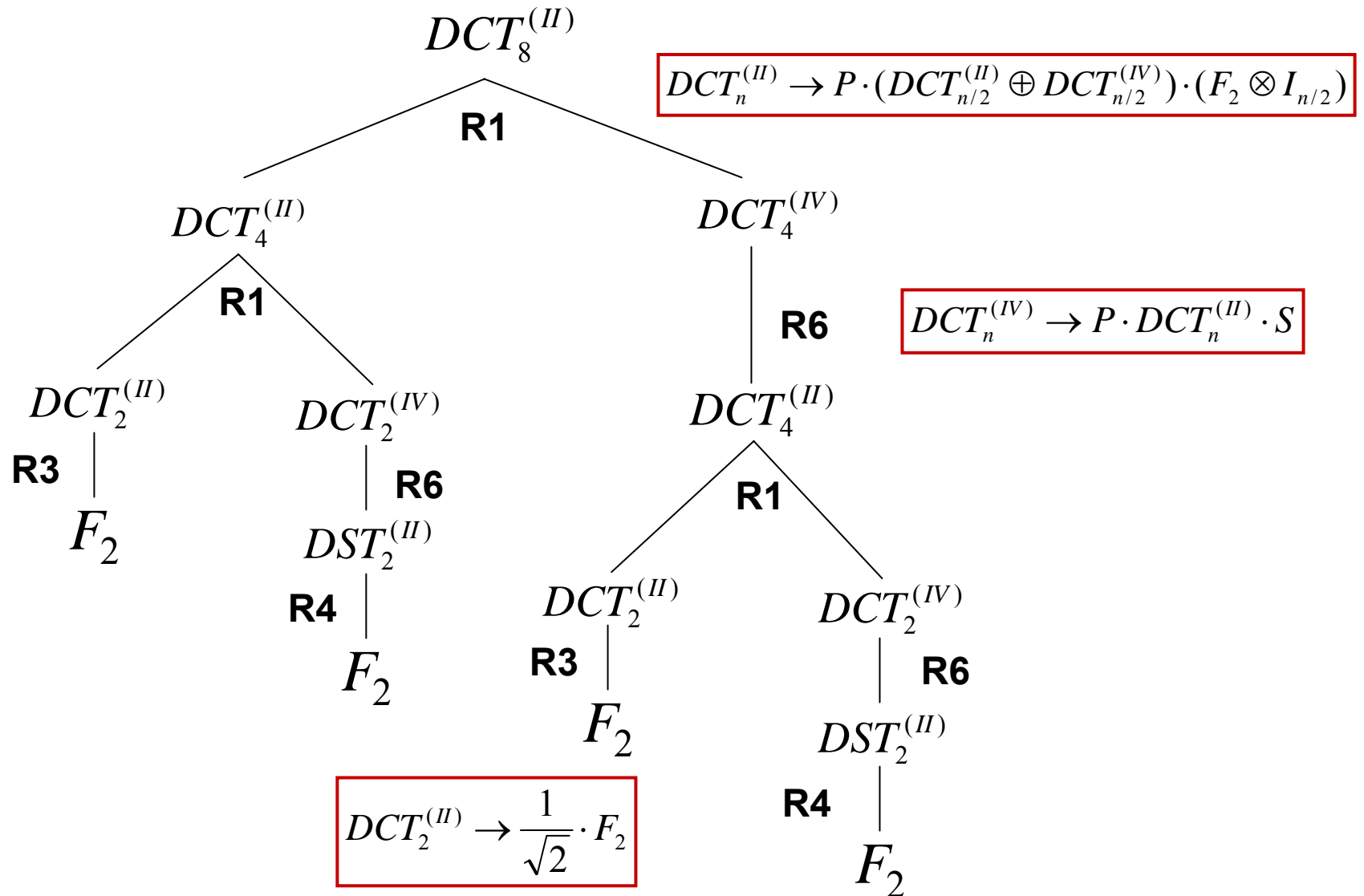
$$\text{Filt}_n(h(z)) \rightarrow \mathbf{R}_{n,l,r}^{\text{zero}} \cdot \mathbf{C}_{n+l+r}(h(z))$$

“Zirkulante über DFT” Regel

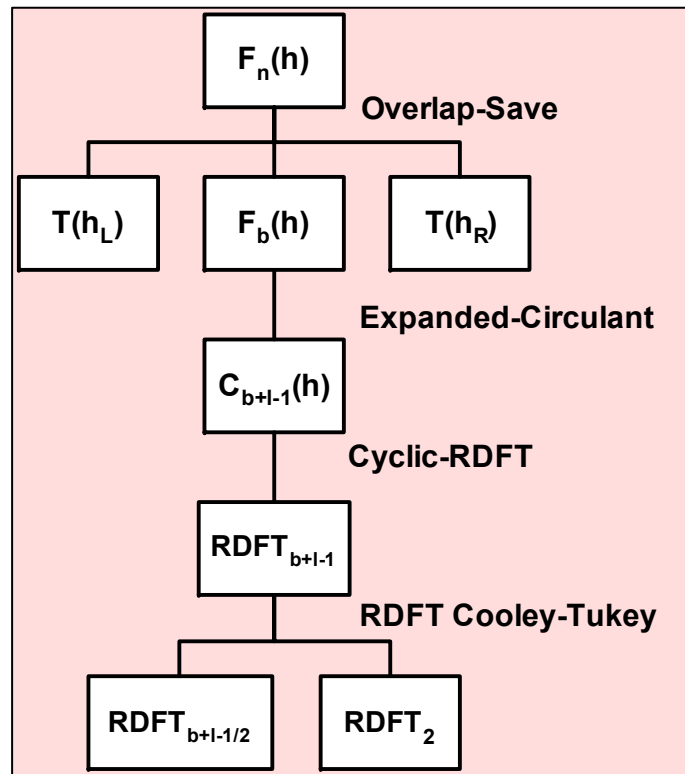
$$\mathbf{C}_n(h(z)) \rightarrow \mathbf{RDFT}_n^{-1} \cdot X(\hat{\mathbf{h}}) \cdot \mathbf{RDFT}_n,$$

$$\hat{\mathbf{h}} = \mathbf{RDFT}_n \cdot \mathbf{h}$$

Regelbaumbeispiel (I)

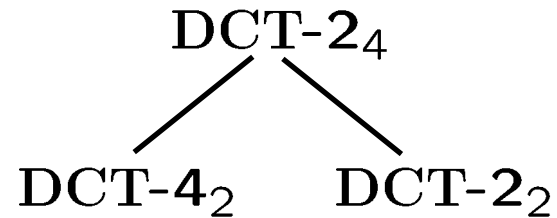


Regelbaumbeispiel (II)



Formelbeispiel

Regelbaum:



Formel:

$$\text{DCT-2}_4 = L_2^4(\text{diag}(1, 1/\sqrt{2}) F_2 \oplus J_2 R_{13\pi/8})(F_2 \otimes I_2)(I_2 \oplus J_2)$$

Algorithmen- (Formel) generierung

■ Gegeben eine Transformation

- Wende Zerlegungsregeln an bis alle Transformationen expandiert sind
- Die Auswahl der Regeln produziert eine exponentielle Menge von Algorithmen
 - ungefähr gleich bzgl. arithmetischer Kosten (nicht für Filter)
 - unterschiedlicher Datenfluss

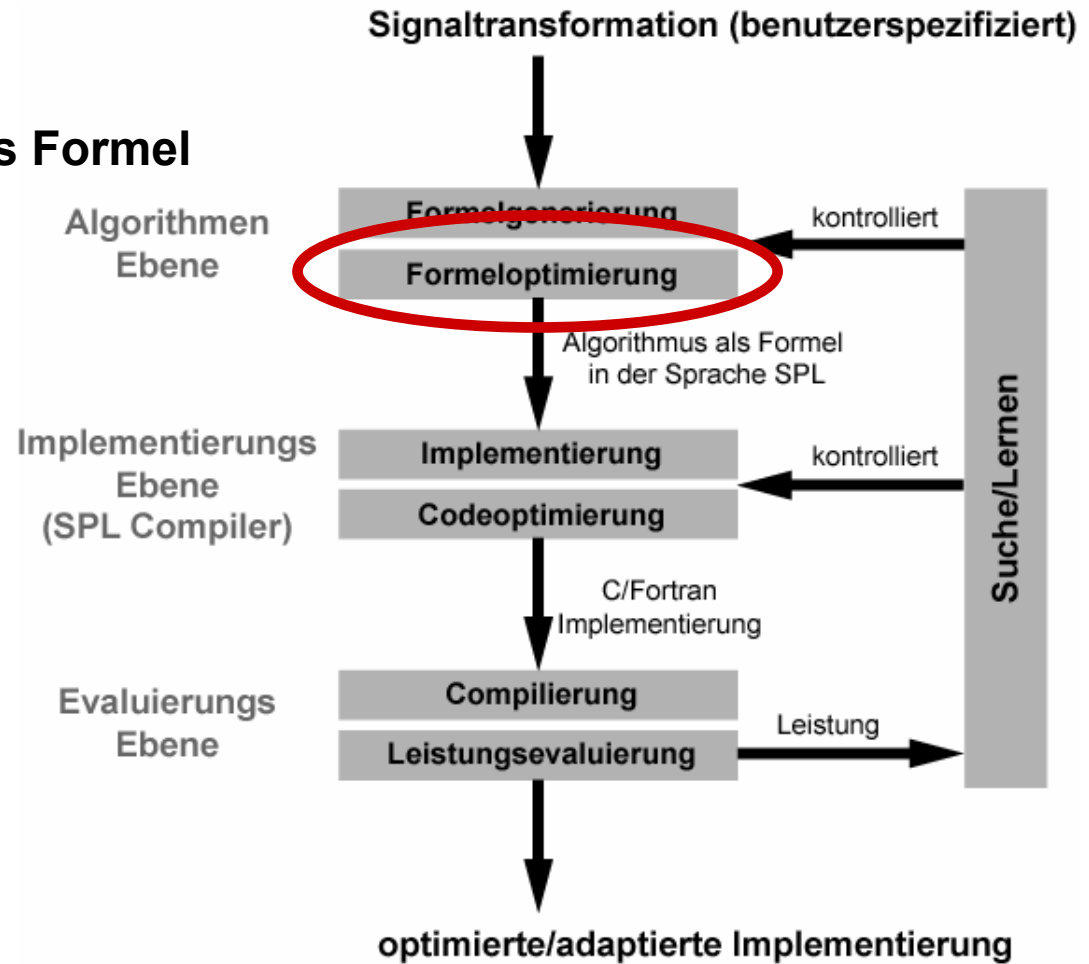
exponentiell

k	# DFTs, Größe 2^k	# DCT IV, Größe 2^k
1	1	1
2	6	10
3	40	126
4	296	31242
5	27744	1924443362
6	162570361280	7343815121631354242
7	$\sim 1.01 \cdot 10^{27}$	$\sim 1.07 \cdot 10^{38}$
8	$\sim 2.31 \cdot 10^{61}$	$\sim 2.30 \cdot 10^{76}$
9	$\sim 2.86 \cdot 10^{133}$	$\sim 1.06 \cdot 10^{153}$

Algorithmen- (Formel) optimierung

Eingabe:
Schneller Algorithmus als Formel

Ausgabe:
Optimierte Formel



Motivation: Schleifenverschmelzung

∞
L4
 $(I_4 \otimes F_2)$

direkte
Codegener.

```
void I4xF2_L84(double *y, double *x) {
    double t[8];
    for (int i=0; i<8; i++)
        t[i==7 ? 7 : (i*4)%7] = x[i];
    for (int i=0; i<4; i++){
        y[2*i] = t[2*i] + t[2*i+1];
        y[2*i+1] = t[2*i] - t[2*i+1];
    }
}
```

kein Compiler kann das

Lösung:
Formel-
manipulierung

```
void I4xF2_L84(double *y, double *x) {
    for (int j=0; j<4; j++){
        y[2*j] = x[j] + x[j+4];
        y[2*j+1] = x[j] - x[j+4];
    }
}
```

$$\sum_{j=0}^3 S(j)_{4 \otimes v_2} F_2 G_{v_2 \otimes (j)_4}$$

Optimierung auf der Formelebene

■ Hauptziele:

- Verschmelzung iterative Schritte (Schleifen)
- Strukturmanipulierung für Vektorinstruktionen

■ Überwindet Compilerlimitierungen

■ Möglich durch Σ -SPL und mathematische Regeln

■ Implementierung durch mehrfache Ebenen von Rewritingsystemen

■ Gibt SPIRAL “Algorithmenwissen”

Σ -SPL: Generating Functions

- **Symbolische Funktionen**

Identität $v_n : \begin{cases} \{0, \dots, n-1\} \rightarrow \{0, \dots, n-1\} \\ i \mapsto i \end{cases}$

Basis Funktion $(j)_n : \begin{cases} \{0\} \rightarrow \{0, \dots, n-1\} \\ i \mapsto j \end{cases}$

Stride Funktion $\ell_m^{mn} : \begin{cases} \{0, \dots, mn-1\} \rightarrow \{0, \dots, mn-1\} \\ i \mapsto \begin{cases} (im) \bmod (mn-1) & \text{if } 0 \leq i < mn-1 \\ mn-1 & \text{if } i = mn-1 \end{cases} \end{cases}$

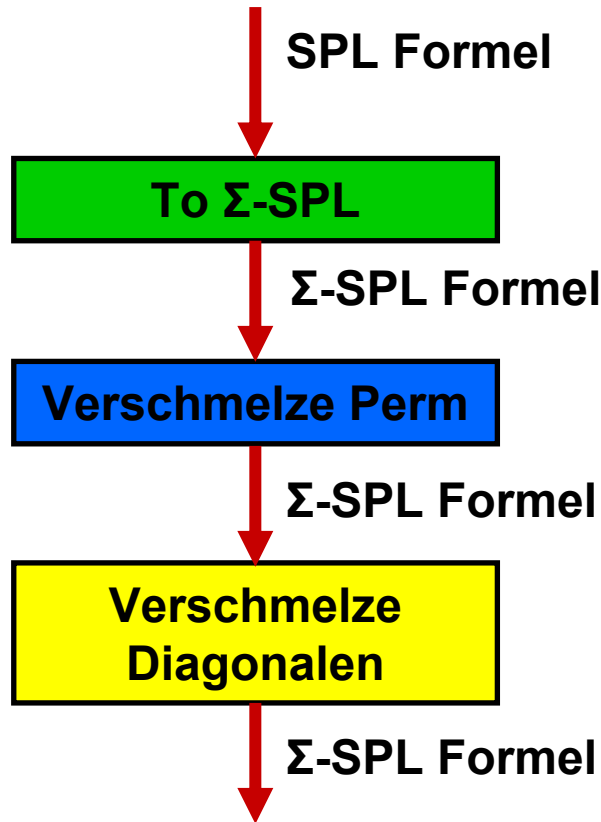
Twiddle Funktion $t_n^{mn} : \begin{cases} \{0, \dots, mn-1\} \rightarrow \mathbb{C} \\ i \mapsto \omega_{mn}^{(i \bmod n) \lfloor \frac{i}{n} \rfloor} \end{cases}$

- **Funktionsoperatoren**

$$(j)_m \otimes v_n : \begin{cases} \{0, \dots, n-1\} \rightarrow \{0, \dots, mn-1\} \\ i \mapsto jn + i \end{cases}$$

$$v_n \otimes (j)_m : \begin{cases} \{0, \dots, n-1\} \rightarrow \{0, \dots, mn-1\} \\ i \mapsto j + im \end{cases}$$

Beispiel: Schleifenoptimierung



$$T_n^{mn} (I_m \otimes C_n) L_m^{mn}$$

$$\text{diag} \left(t_n^{mn} \right) \left(\sum_{j=0}^{n-1} S_{(j)_m \otimes v_n} C_n G_{(j)_m \otimes v_n} \right) \text{perm} \left(\ell_m^{mn} \right)$$

$$\text{diag} \left(t_n^{mn} \right) \left(\sum_{j=0}^{n-1} S_{(j)_m \otimes v_n} C_n G_{v_n \otimes (j)_m} \right)$$

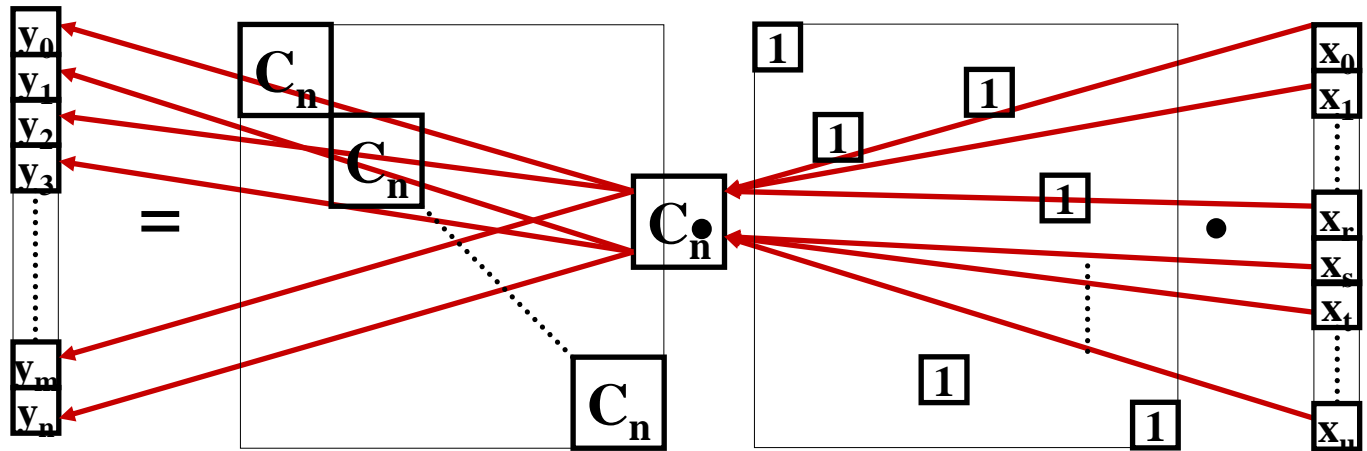
$$\sum_{j=0}^{n-1} S_{(j)_m \otimes v_n} \text{diag} \left(t_n^{mn} \circ \left((j)_m \otimes v_n \right) \right) C_n G_{v_n \otimes (j)_m}$$

Regeln: $G_r \text{perm}(\pi) = G_{\pi \circ r}$, $\ell_m^{mn} \circ \left((j)_m \otimes v_n \right) = v_n \otimes (j)_m$

$$\text{diag}(f) S_w = S_w \text{diag}(f \circ w)$$

Visualisierung

$$(\mathbf{I}_m \otimes \mathbf{C}_n) \mathbf{L}_m^{mn} = \sum_{j=0}^{n-1} \mathbf{S}_{(j)_m \otimes \mathbf{I}_n} \mathbf{C}_n \mathbf{G}_{\mathbf{I}_n \otimes (j)_m}$$



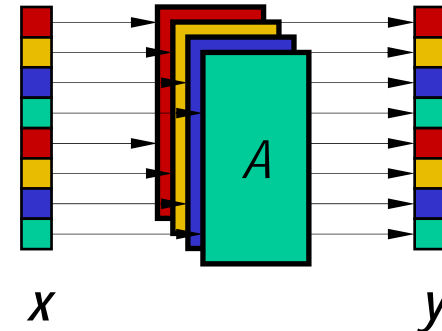
Iteration 0

Vektorcodegenerierung von SPL Formeln

Natürlich vektorisierbar

$$A \otimes I_4$$

Vektorlänge \nearrow



Verallgemeinerung (auch komplett vektorisierbar):

$$\prod_{i=1}^k P_i D_i (A_i \otimes I_v) E_i Q_i$$

P_i, Q_i

D_i, E_i

A_i

v

Permutationen

Diagonalmatrizen

beliebige Formel

Vektorlänge

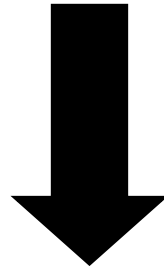
Vektorisierung in zwei Schritten:

1. Formelmanipulierung
2. Codegenerierung (Vektorcode + C code)

Beispiel DFT

Standard FFT

$$(\text{DFT}_k \otimes \text{I}_m) \text{T}_m^n (\text{I}_k \otimes \text{DFT}_m) \text{L}_k^n$$



Formelmanipulierung

$$\left(\text{I}_{\frac{mn}{\nu}} \otimes \text{L}_{\nu}^{2\nu} \right) \left(\overline{\text{DFT}_m \otimes \text{I}_{\frac{n}{\nu}} \otimes \text{I}_{\nu}} \right) \overline{\text{T}}_n^{/mn}$$

$$\left(\text{I}_{\frac{m}{\nu}} \otimes \left(\text{L}_{\nu}^{2n} \otimes \text{I}_{\nu} \right) \left(\text{I}_{\frac{2n}{\nu}} \otimes \text{L}_{\nu}^{\nu^2} \right) \left(\overline{\text{DFT}_n \otimes \text{I}_{\nu}} \right) \right) \left(\text{L}_{\frac{m}{\nu}}^{\frac{mn}{\nu}} \otimes \text{L}_{\frac{2}{2}}^{2\nu} \right)$$

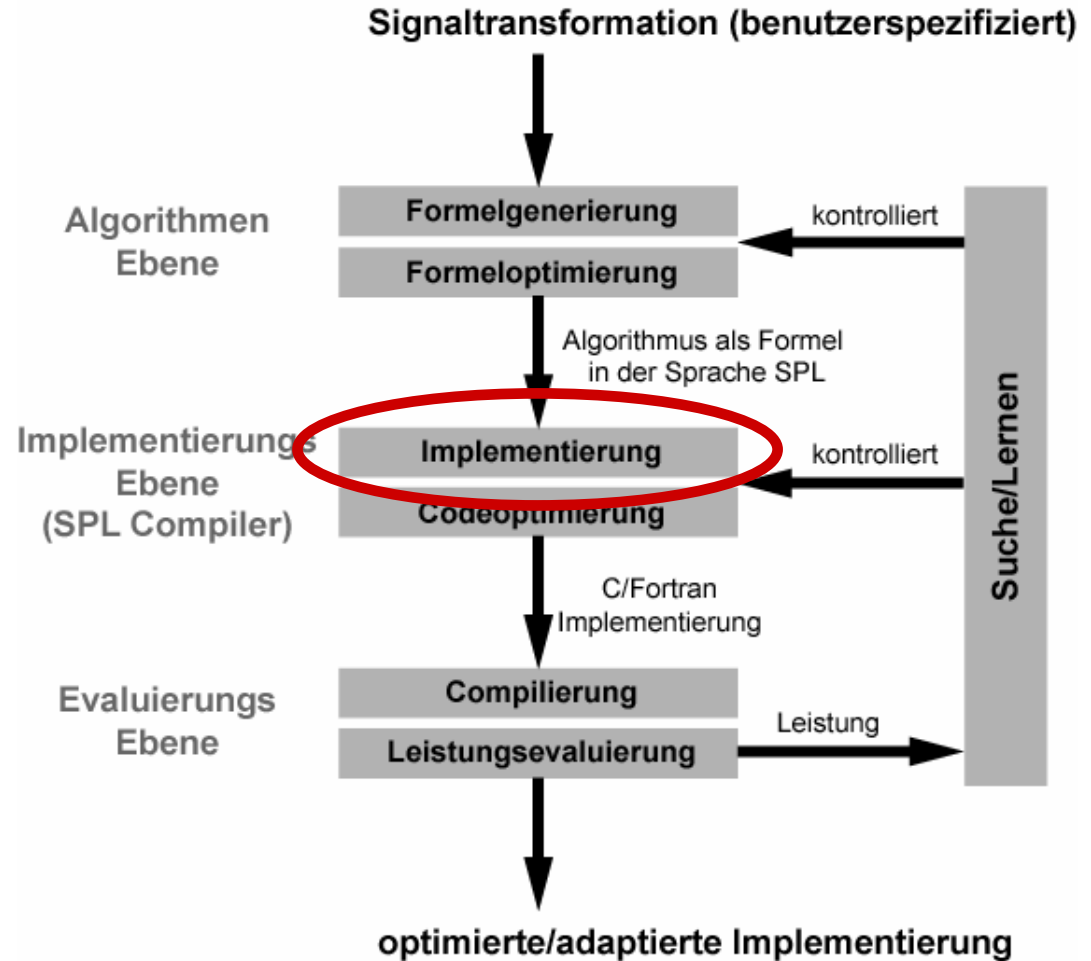
Vektor FFT für ν -weg Vektorinstruktionen

Optimierter Algorithmus (Formel) zu Code

Eingabe:
Optimierte Formel

Ausgabe:
Code

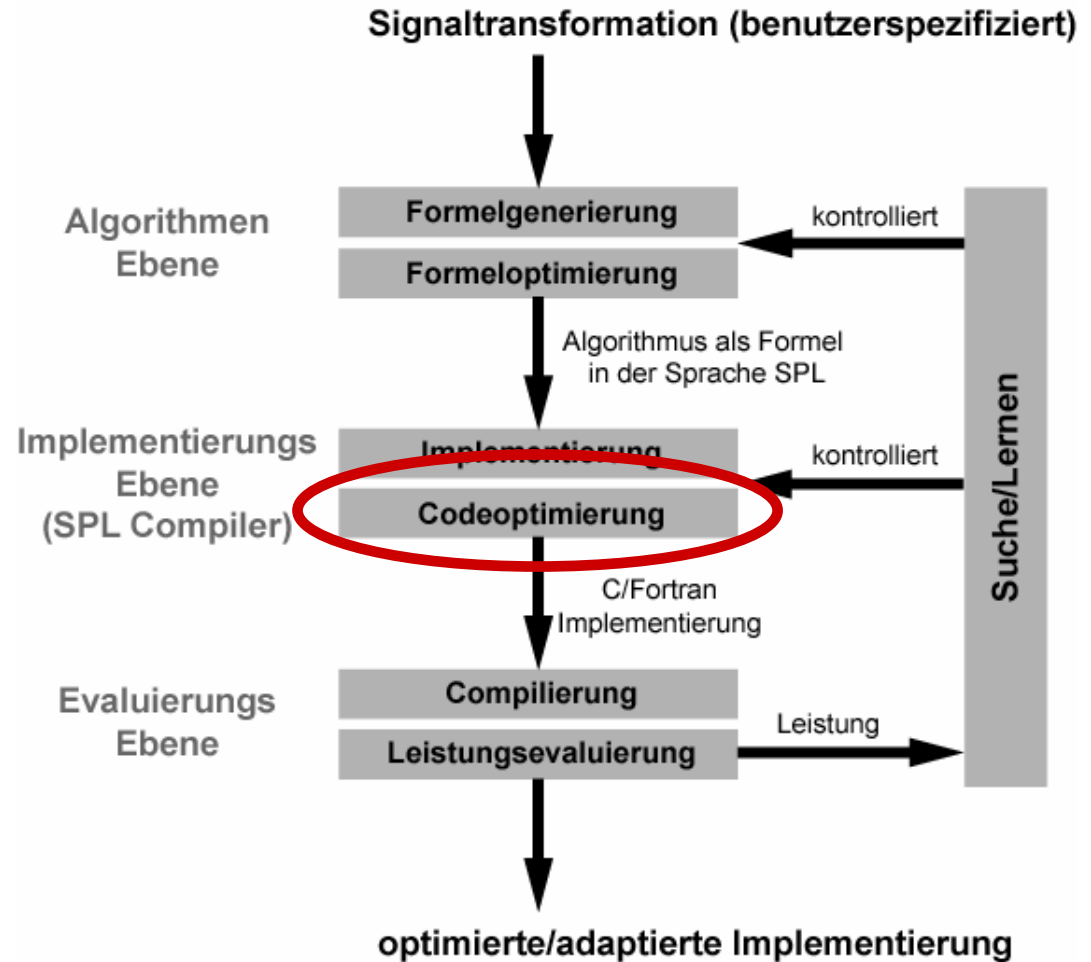
Einfach



Codeoptimierung

Eingabe:
Code

Ausgabe:
Optimierter Code



Optimierungen auf der Codeebene

- Vorbereitung von Konstanten
- Schleifenaufrollen (kontrolliert vom Suchmodul)
- Inlining der Konstanten
- SSA code, Array Skalarisierung, CSE
- Codescheduling für Lokalität (optional)
- Umwandlung in FMA Code (optional)
- Umwandlung in Fixpunktcode (optional)

- Am Ende: Unparsing nach C (or Fortran)

Codeevaluierung

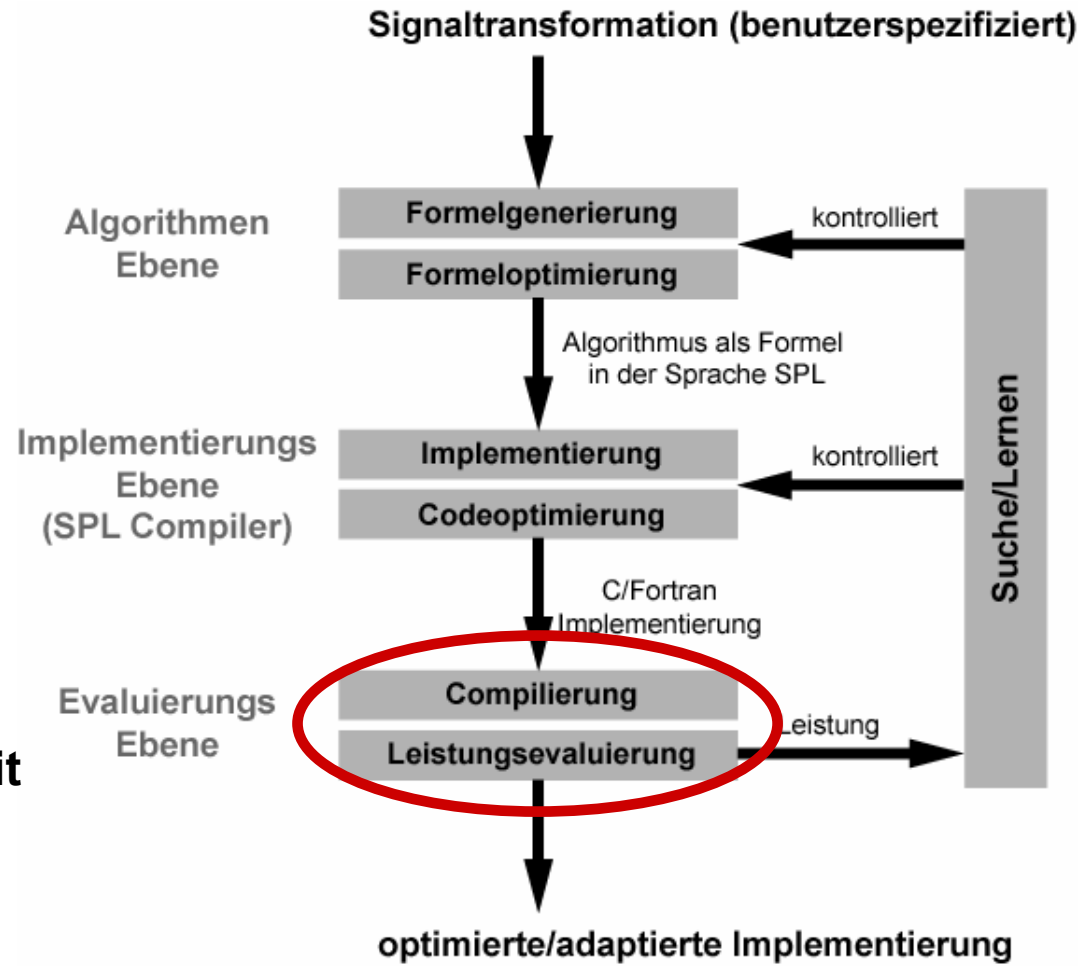
Eingabe:
Optimierter C code

Ausgabe:
Leistungszahl

Klar

Beispiele:

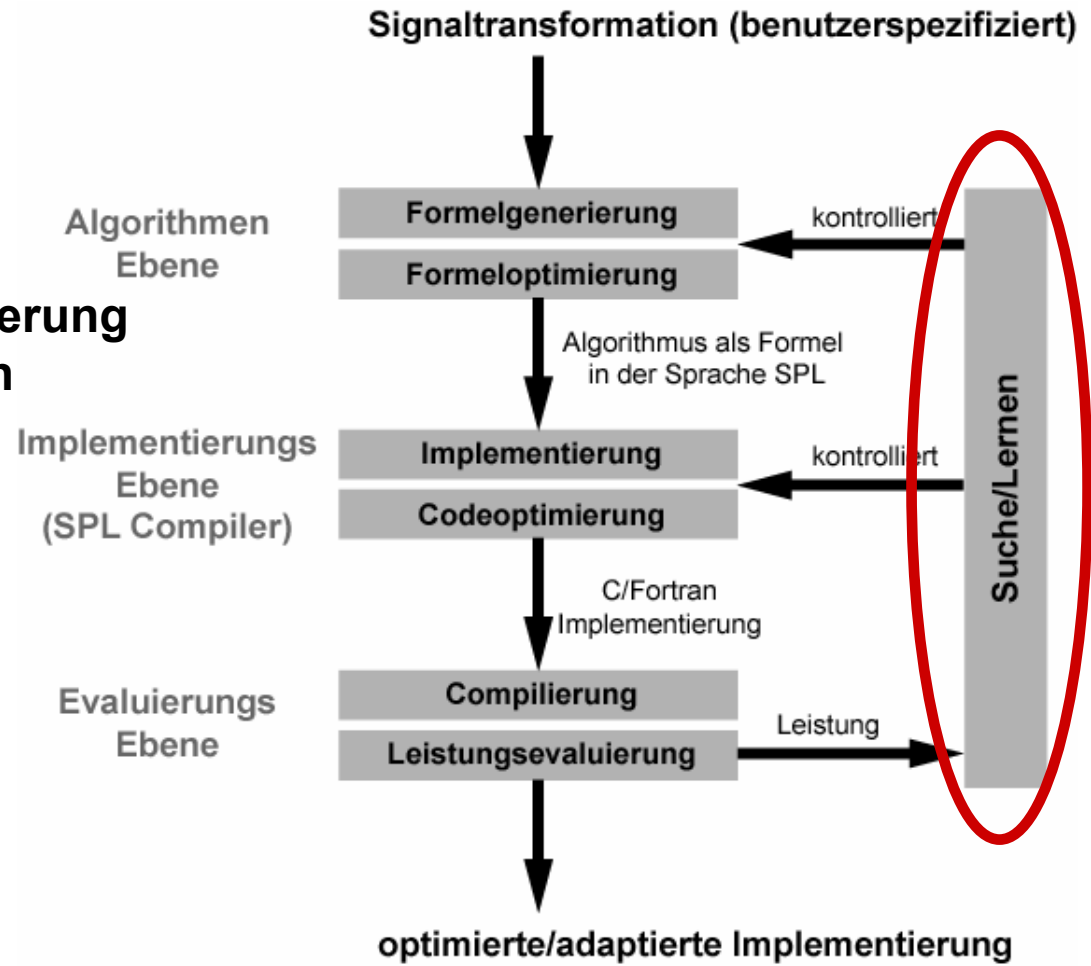
- Laufzeit
- Numerische Genauigkeit
- Arithmetische Kosten



Search (Learning) for the Best

Eingabe:
Leistungszahl

Ausgabe:
Kontrolle der Formelgenerierung
Kontrolle der Codeoptionen



Suche

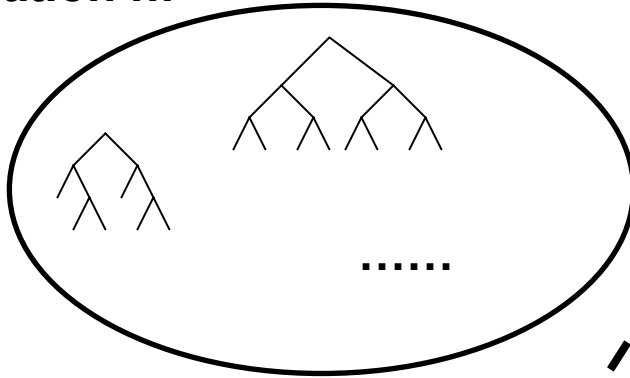
- **Suche über:**
 - Algorithmische Optionen (Auswahl der Zerlegungsregeln)
 - Codeoptionen (Grad des Schleifenaufrollens)

- **Suche arbeitet mit der Regelbaumdarstellung eines Algorithmus**
 - transformationsunabhängig
 - effizient

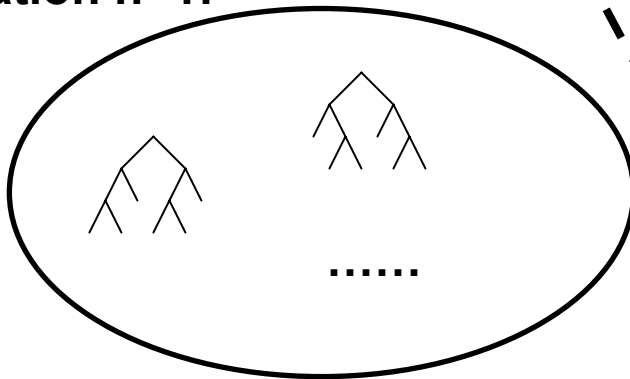
- **Suchmethoden**
 - Vollständige Suche
 - Dynamisches Programmieren (DP)
 - Zufällige Suche
 - Hill Climbing
 - STEER (ein evolutionärer Algorithmus)

STEER

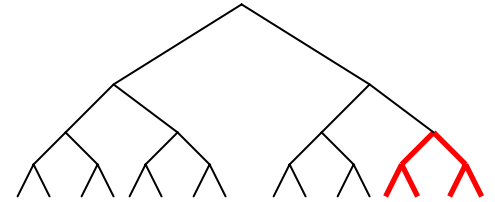
Population n:



Population n+1:

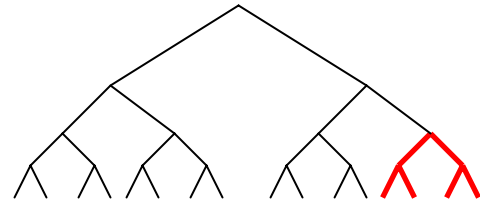


Mutation



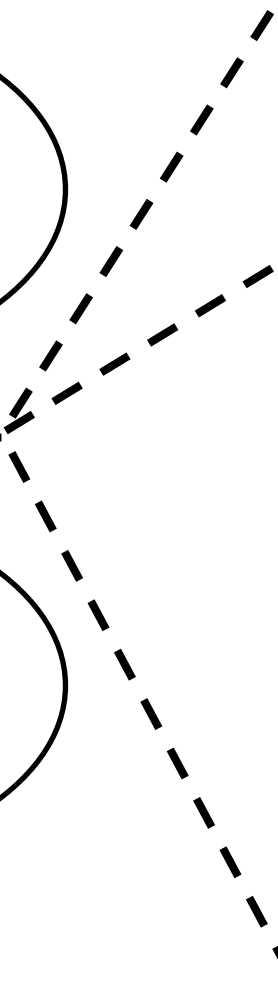
expandiere
anders

Kreuzung



tausche
Expansionen

Überleben des Fittesten



(Maschinen-) Lernen

■ Vorgehensweise:

- Generiere eine Menge von (z.B. 1000) Algorithmen und deren Laufzeiten (für eine Transformation und Größe); repräsentiere den Algorithmus durch Merkmale
- Von diesen Daten (Paare Merkmale/Laufzeiten), lerne eine Menge von Algorithmusdesignregeln
- Von dieser Menge (dem Wissen), generiere schnelle Algorithmen (Theorie der Markov Entscheidungsprozesse)

■ Evaluierung:

- Getestet für die WHT und DFT
- Von den generierten Daten für eine Größe (2^{15}) war es möglich, die besten Algorithmen für mehrere Größen (2^{12} - 2^{18}) zu erzeugen

Bryan Singer and Manuela Veloso

Learning to Construct Fast Signal Processing Implementations

Journal of Machine Learning Research, 2002, Vol. 3, pp. 887-919

Experimente mit SPIRAL

(kurze Demo ?)

Wie evaluiert man Code?

■ Vergleiche gegen die beste erhältliche Software

- Stelle sicher daß die gleiche Zeitmessmethode verwendet wird

■ Berechne die Leistung (in MFLOPS)

$$MFLOPS = \text{Anzahl Operationen} / (10^6 \cdot \text{Laufzeit [s]})$$

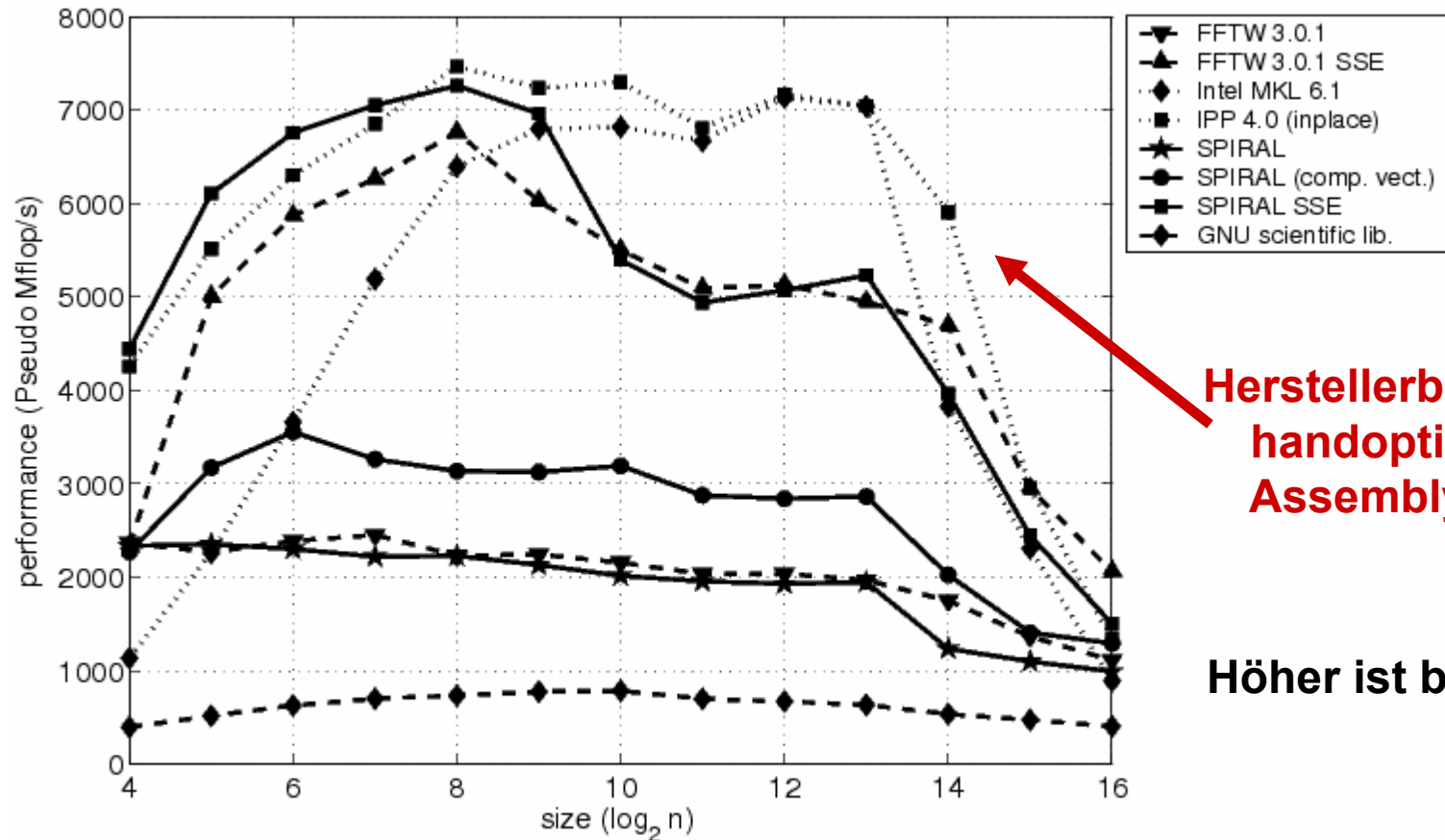
und vergleiche mit Maximalleistung der Maschine

- Erfordert die Berechnung/Messung der Anzahl der Operationen
- Maximalleistung aus dem Manual oder die richtigen Leute fragen
- Vorsicht: Code mit höherer Leistung kann langsamer sein

Benchmarks

Benchmark: DFT, 2er-Potenzen

P4, 3.2 GHz,
icc 8.0



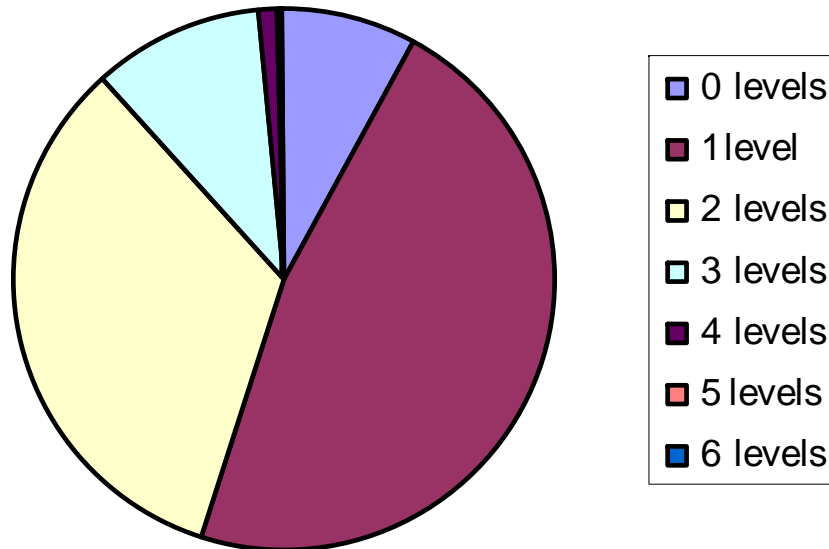
**Herstellerbibliothek:
handoptimierter
Assemblycode?**

Höher ist besser

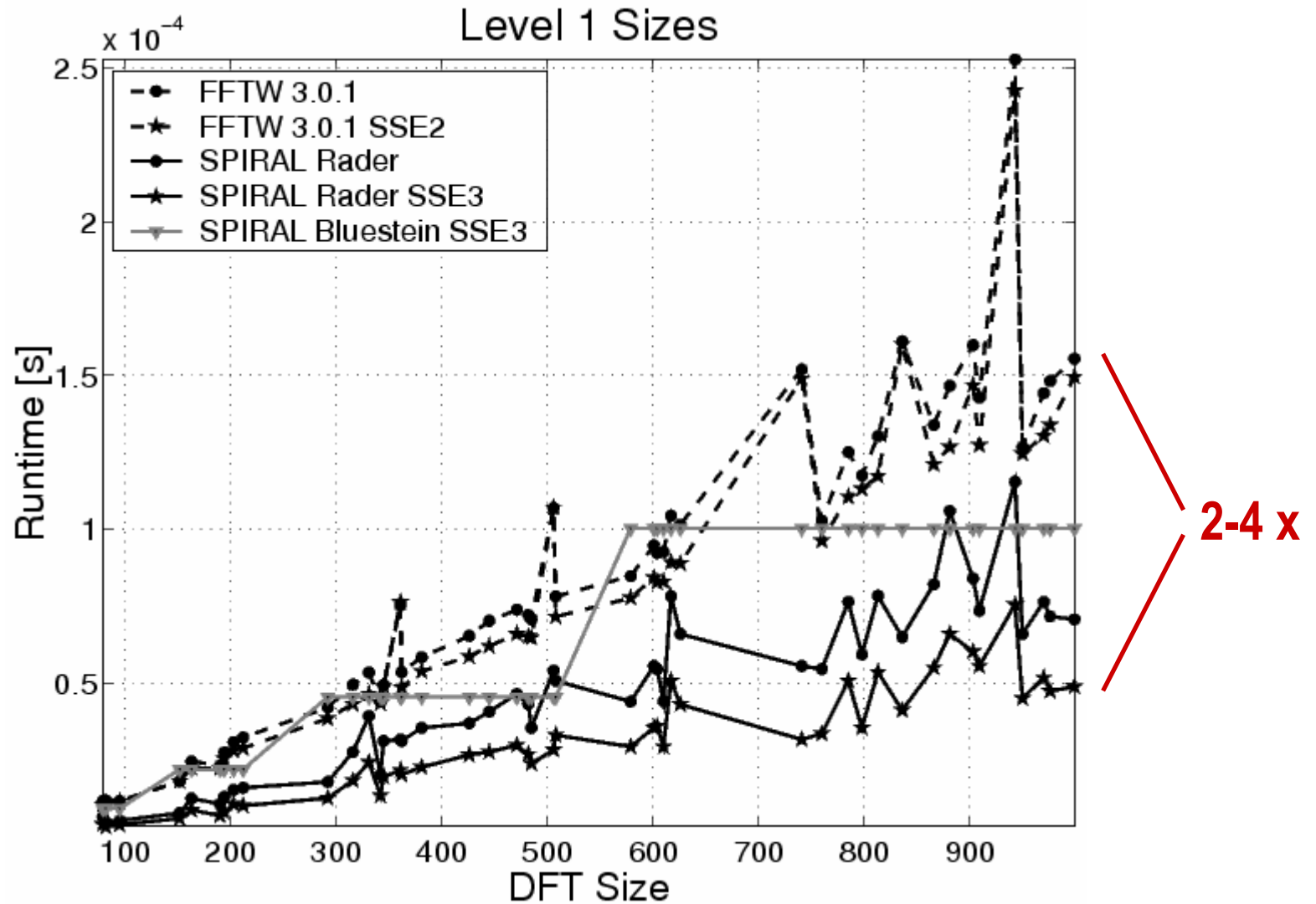
Einfache Genauigkeit

Benchmark: DFT, Andere Größen

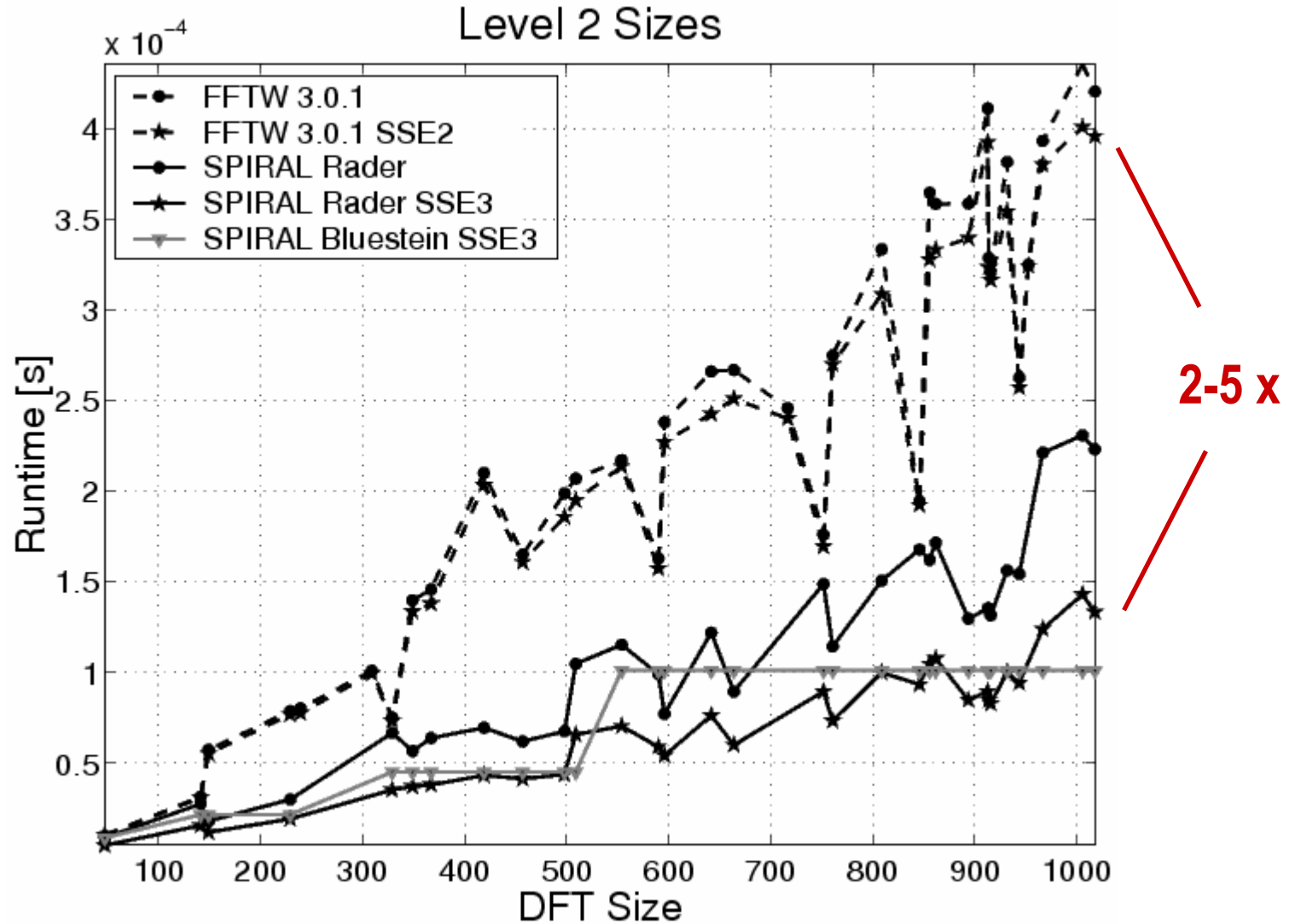
- Teile Größen ein nach Anzahl von nötigen Raderschritten
- $n < 8192$



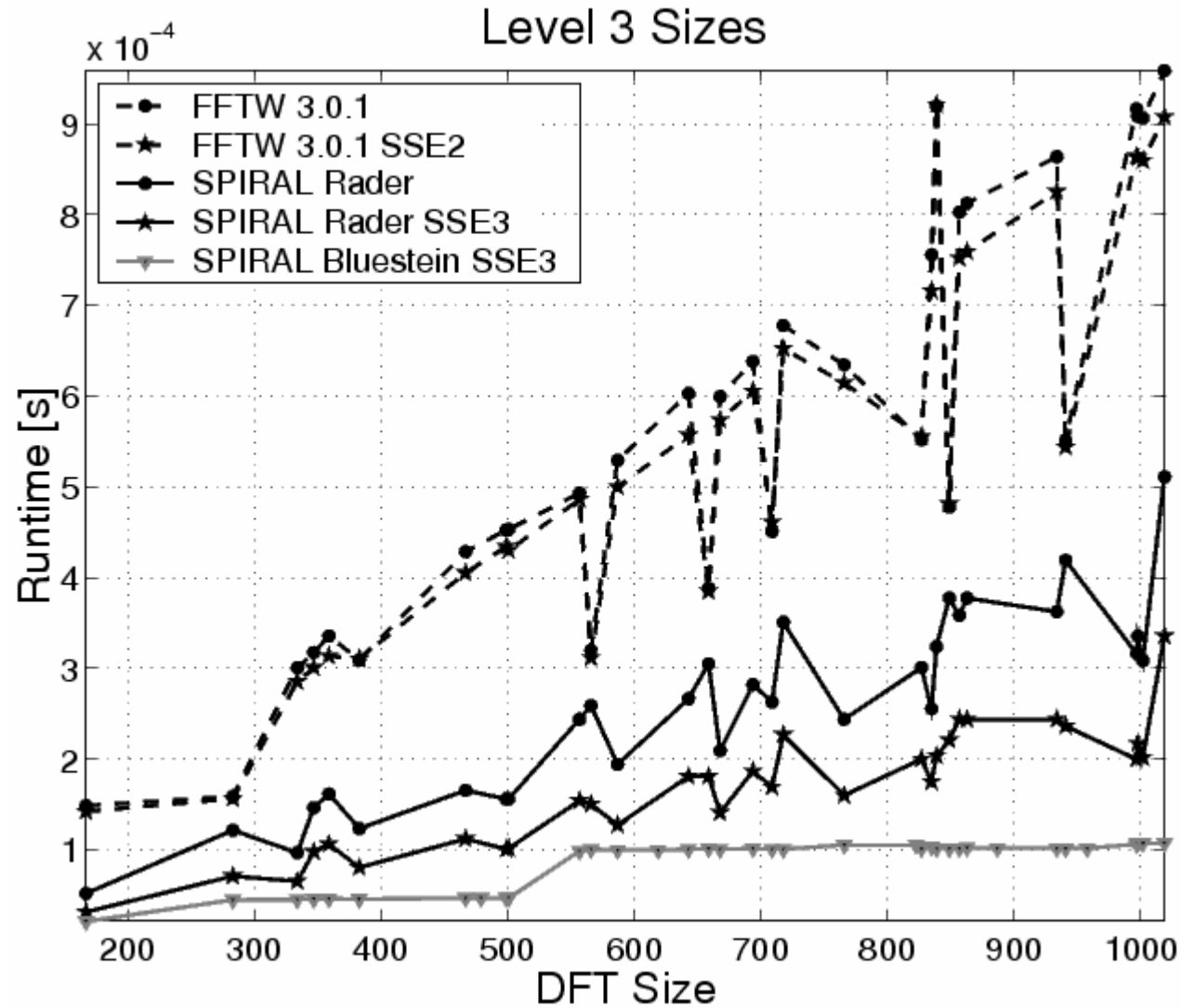
Benchmark: DFT, Level 1 Größen



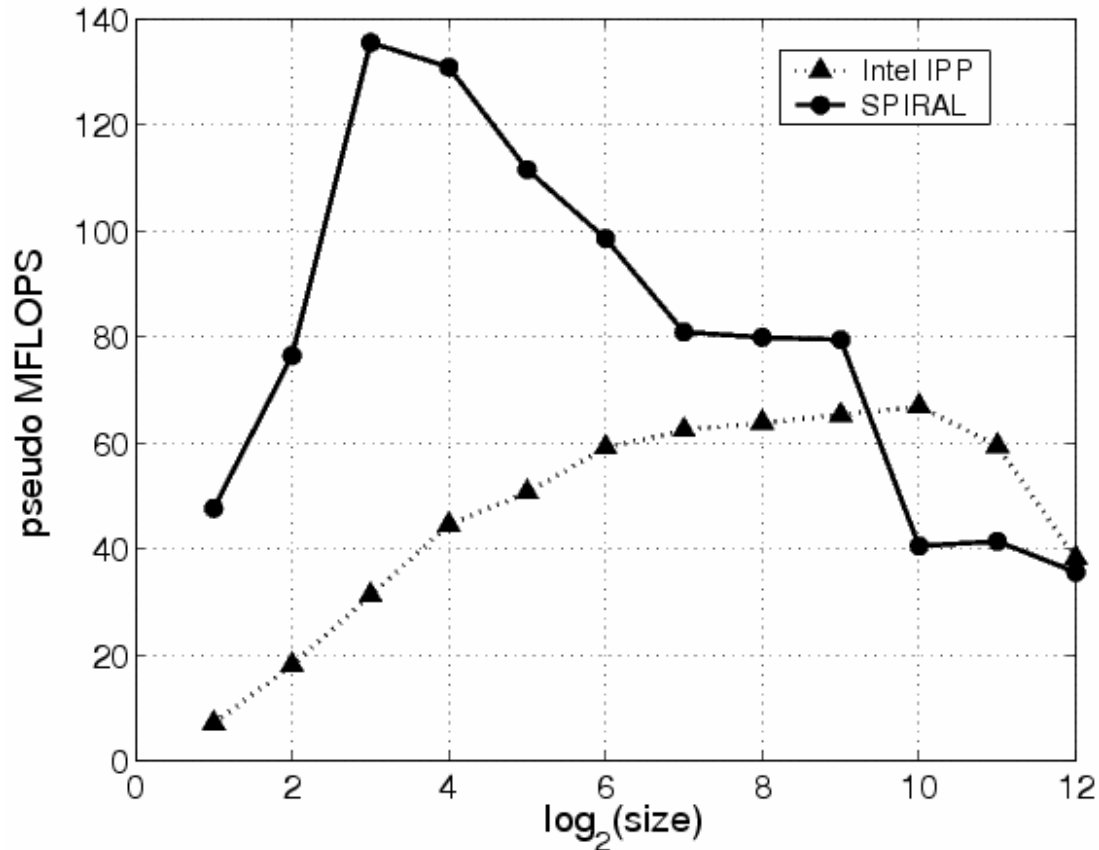
Benchmark: DFT, Level 2 Größen



Benchmark: DFT, Level 3 Größen



Benchmark: Fixpunkt DFT, IPAQ



IPAQ
Xscale arch.
400 MHz
hat nur Fixpunkt

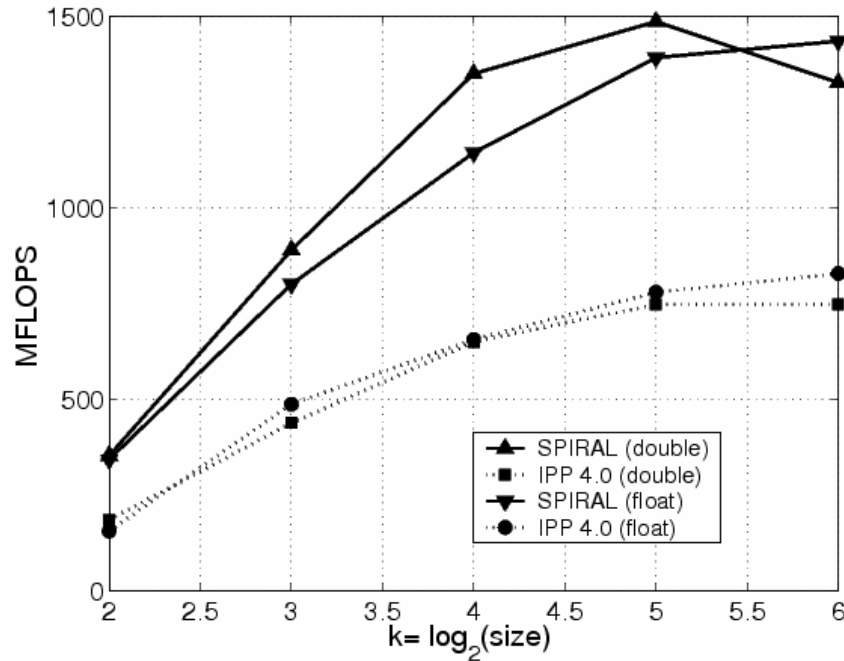
Höher ist besser

Intel hat sich weniger bemüht?

Benchmark: DCT

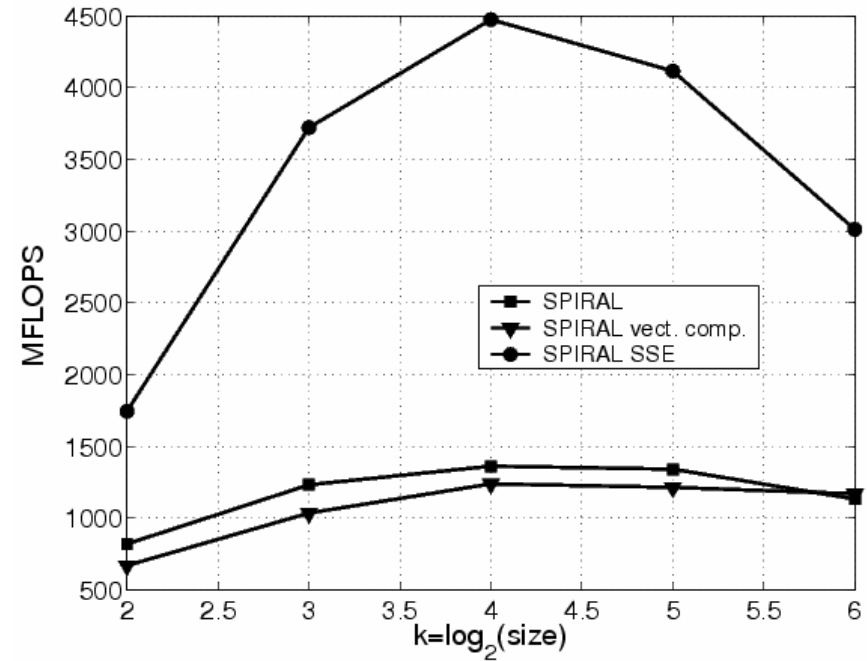
P4, 3.2 GHz,
icc 8.0

1-D DCT



Scalar code

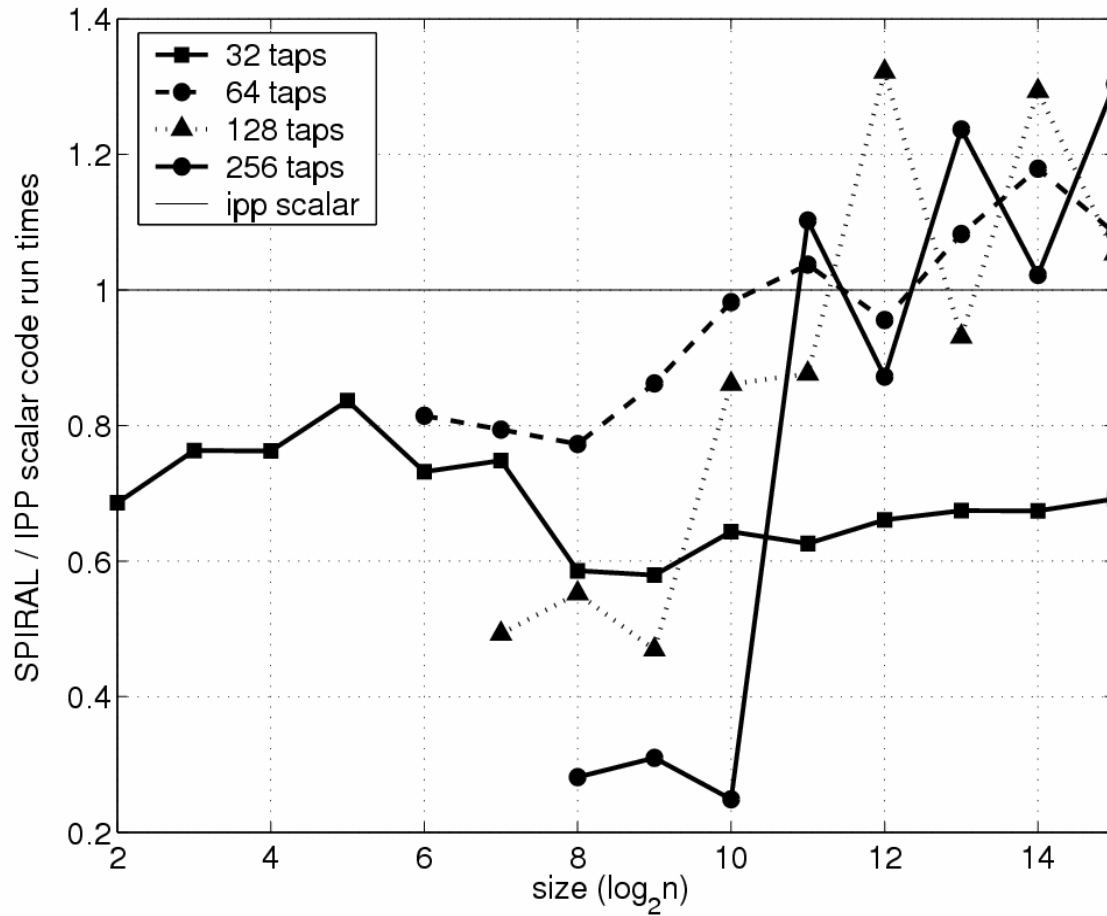
2-D DCT



Scalar vs. SSE code

- Faktor 2 gegenüber Herstellerbibliothek
- Ein weiterer Faktor 3 mit Vektorinstruktionen

Benchmark: Filter

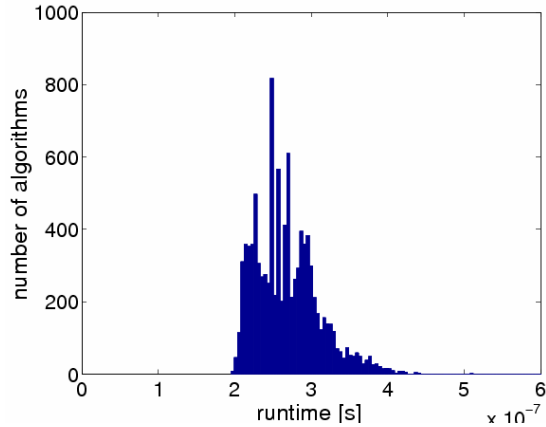


Lehrreiche Experimente

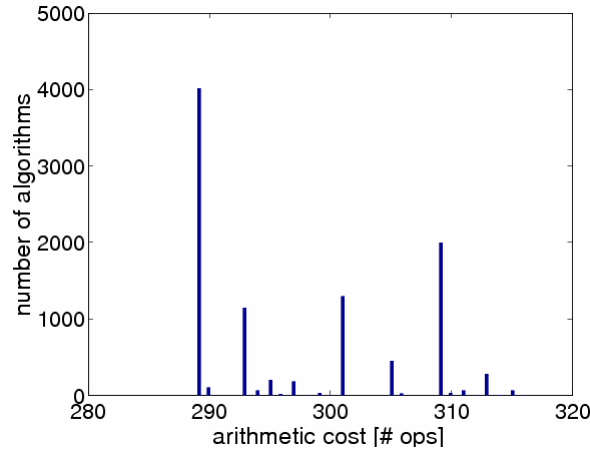
Laufzeitverteilung: DCT, size 32

Histogramme, 10,000 Algorithmen

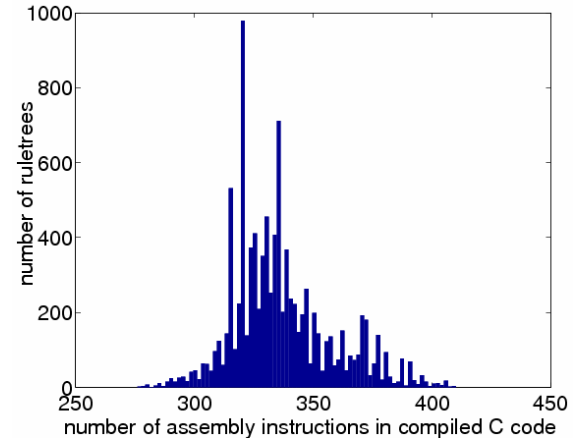
P4, 3.2 GHz,
gcc



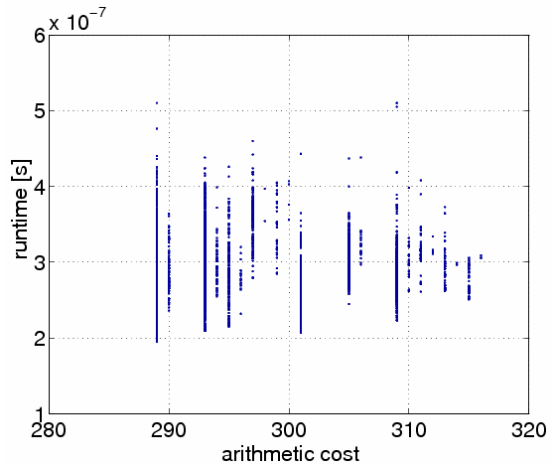
Laufzeit: x2



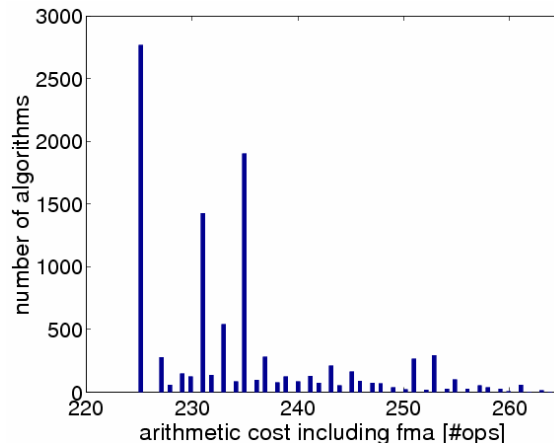
#Ops: x1.08



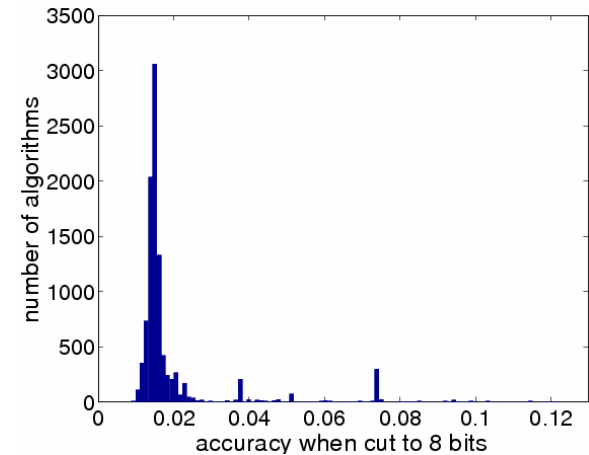
#Assembly Instr.: x1.5



**#Ops und Laufzeit:
keine Korrelation**



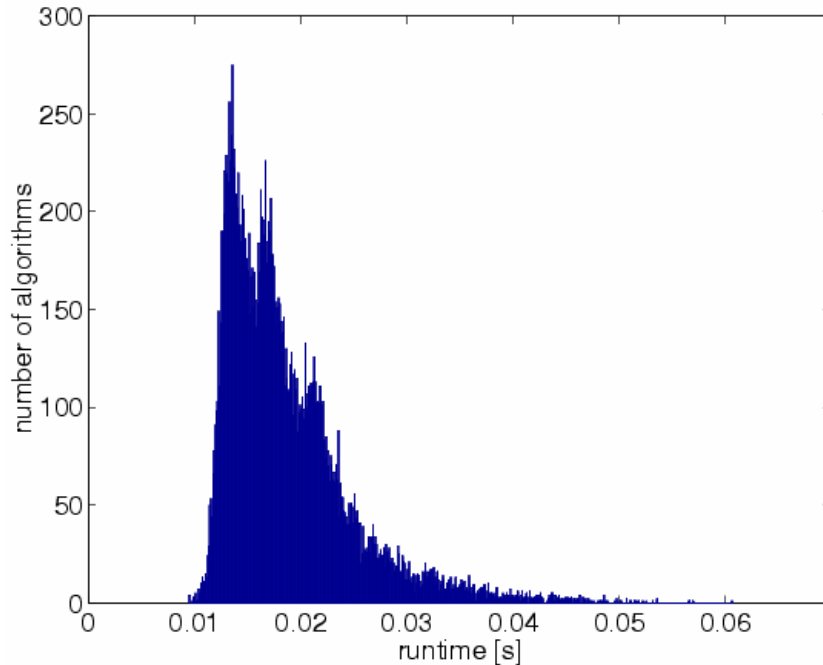
#FMA Ops: x1.2



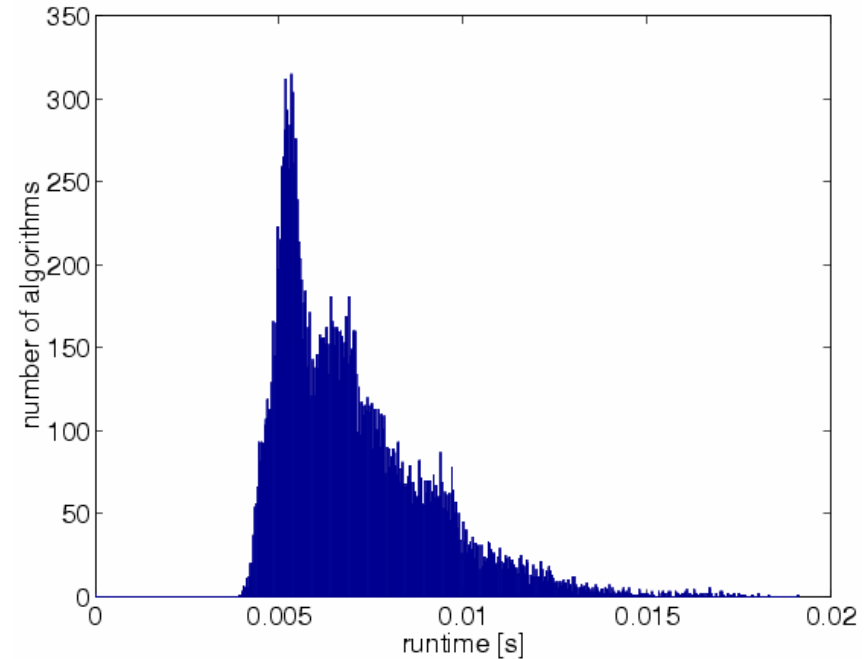
**Numerische Genauigkeit:
x10, die meisten x2**

Laufzeitverteilung: DFT 2^{16} Histogramme, 20,000 Algorithms

P4, 3.2 GHz,
icc 8.0



Generierter Skalarcode

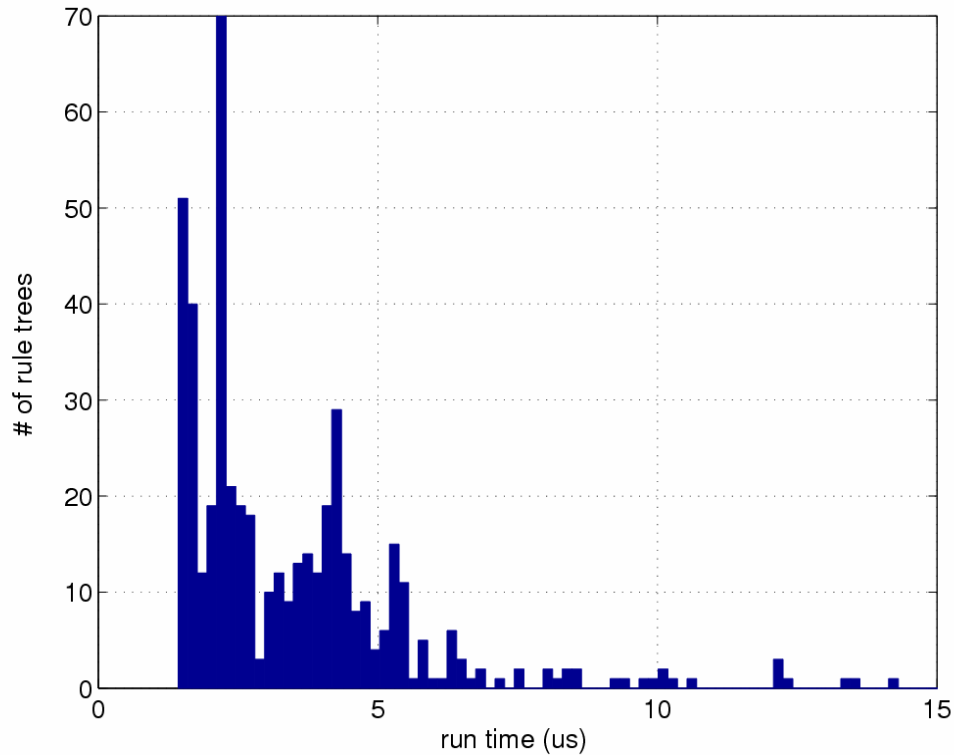


**Generierter SSE Code
(4-weg Vektor, einfache Genauigkeit)**

- **Alle Algorithmen werden schneller**
- **x 2.5 Verbesserung**

Laufzeitverteilung: Filter(128, 16)

Pentium 4 – 3.2



Filter: Zeitbereichmethoden

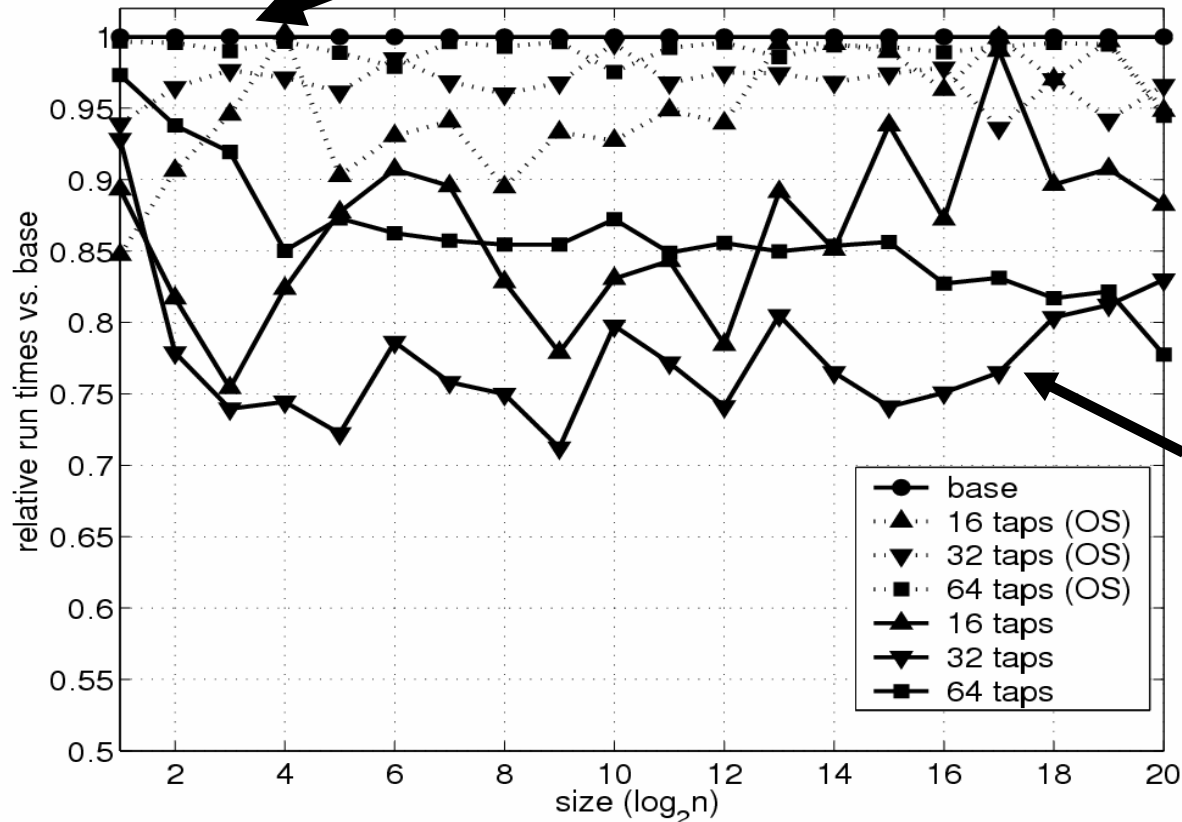
Xeon-1.7

$$\text{Filt}_n(h(z)) \rightarrow I_n \otimes_{k-1} (h_0, \dots, h_{k-1})$$

```
for i = 0..n-1
```

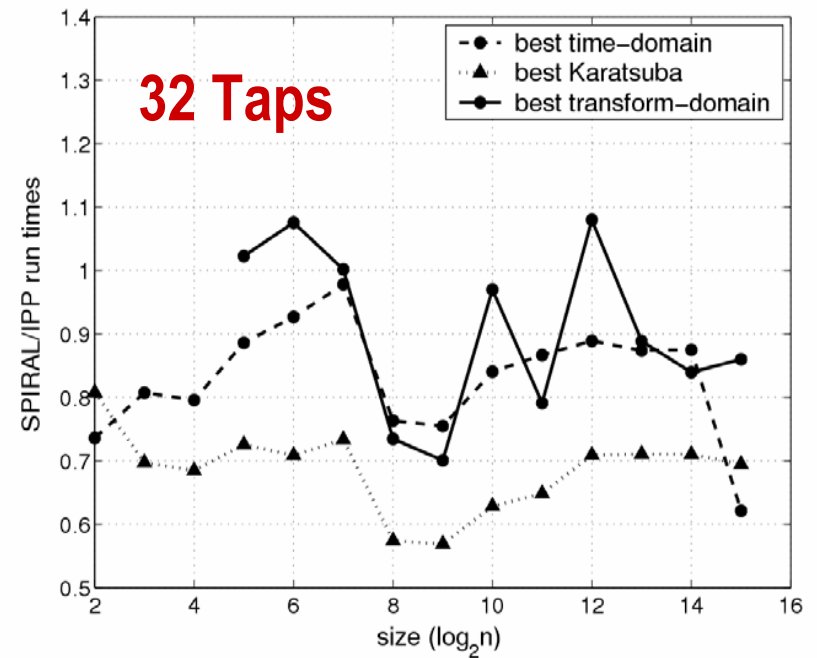
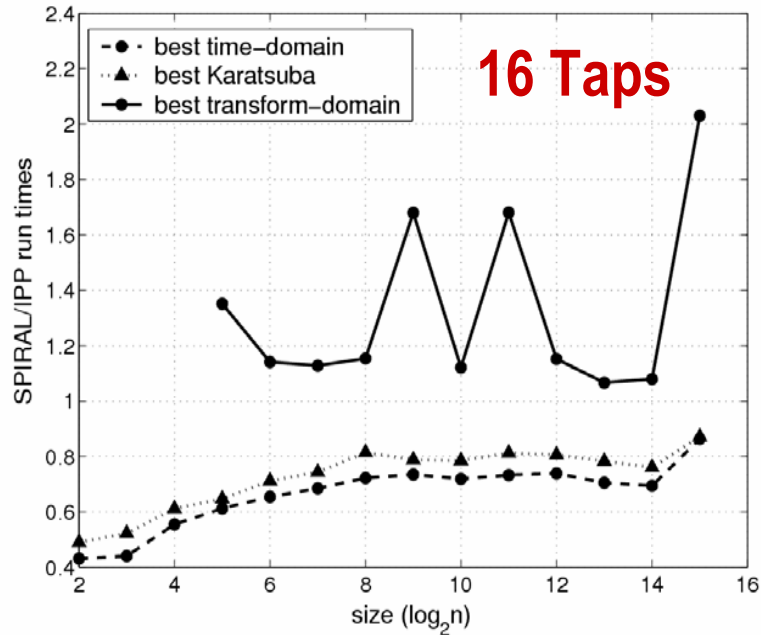
```
  y[i] = h[0]*x[0+i]+h[1]*x[1+i]+...+h[n-1]*x[n-1+i]
```

```
end
```



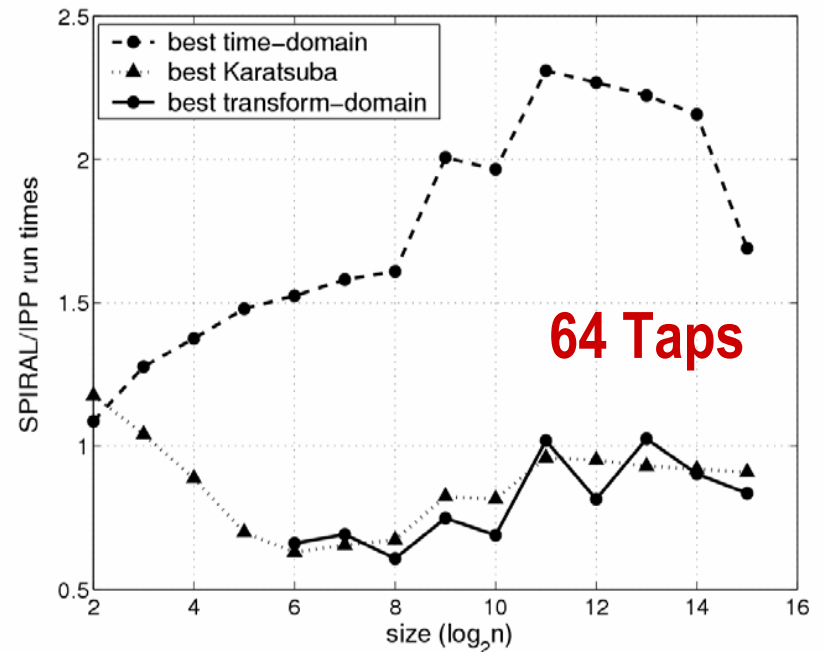
best blocking strategy

Filter: Alle Methoden

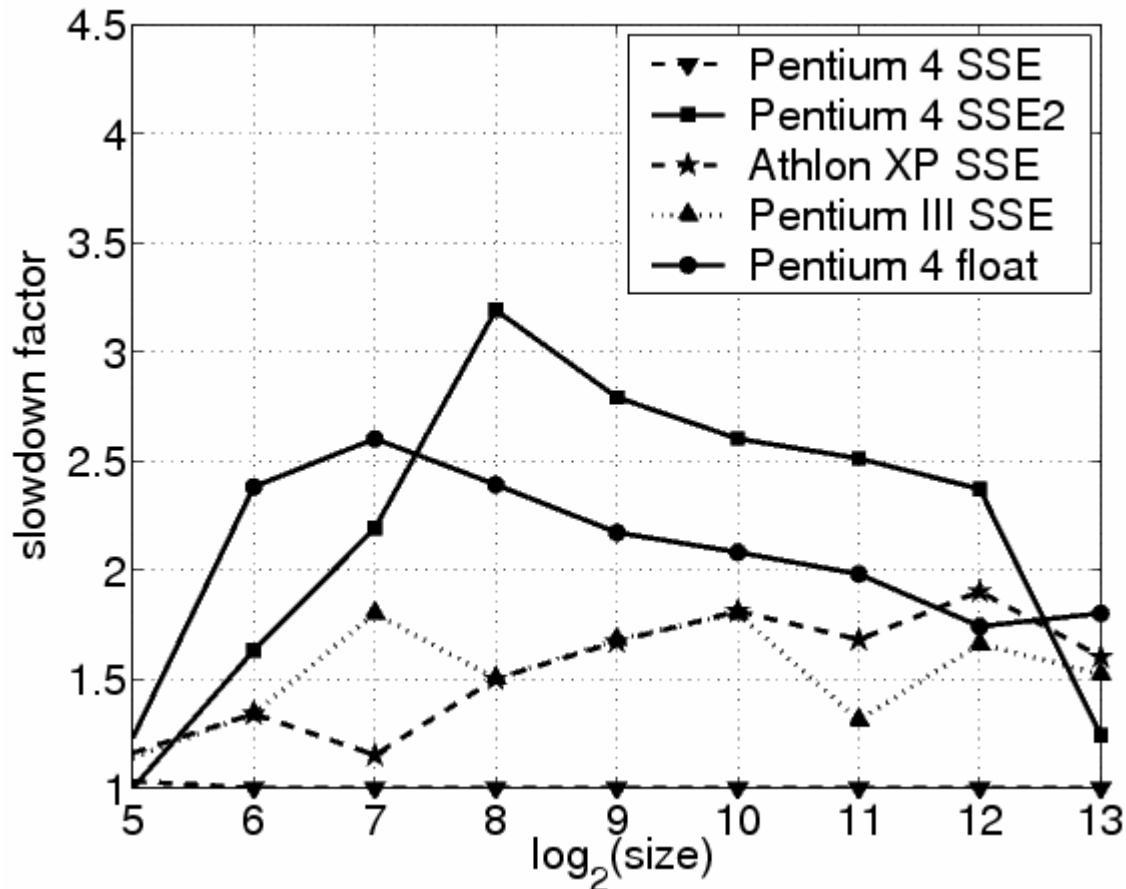


Athlon XP 1.73

- 16: Zeitbereich gewinnt
- 32: Karatsuba gewinnt
- 64: Karatsuba/DFT ~gleich



Plattformabhängigkeit: DFT



50% Verlust durch Portierung von PIII auf P4

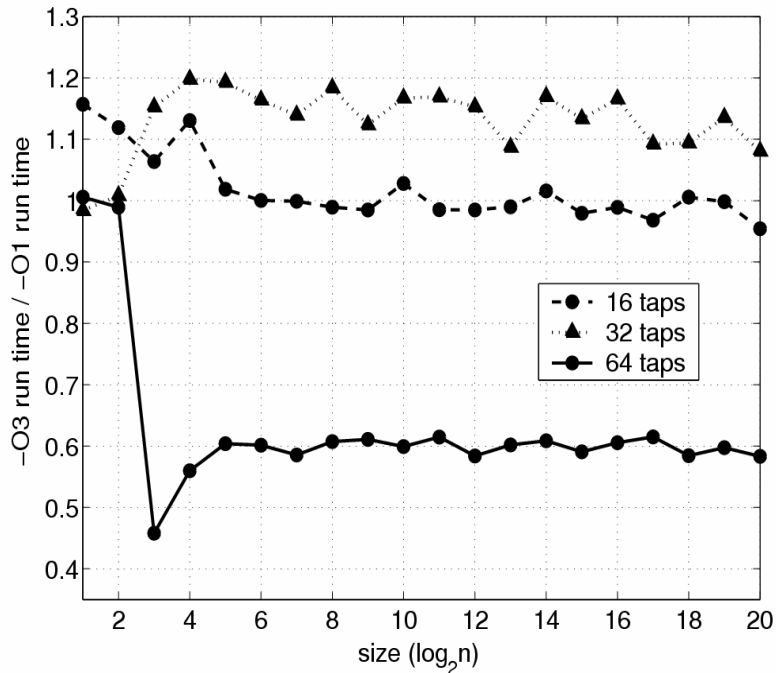
Plattformabhängigkeit: Filter

	16-tap	32-tap	64-tap	128-tap
Pentium 4 3.0GHz Northwood	Blocking	Karatsuba	RDFT	RDFT
Pentium 4 3.6GHz Prescott	Blocking	Karatsuba	Karatsuba	RDFT
Macintosh	Karatsuba	Karatsuba	RDFT	RDFT
Xeon 1.7 GHz	Blocking	Blocking	Blocking	RDFT
Athlon 1.73GHz	Karatsuba/ Blocking	Karatsuba	Karatsuba/ RDFT	RDFT

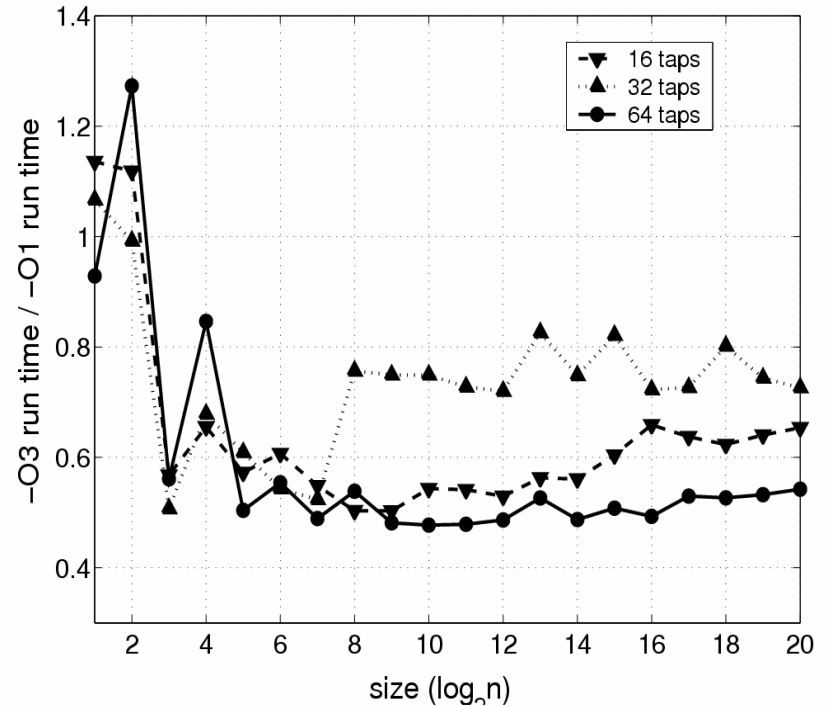
Compileroptionen: Filter

Macintosh - GNU C 3.3 (Apple)

Blocking/nesting



+ Karatsuba

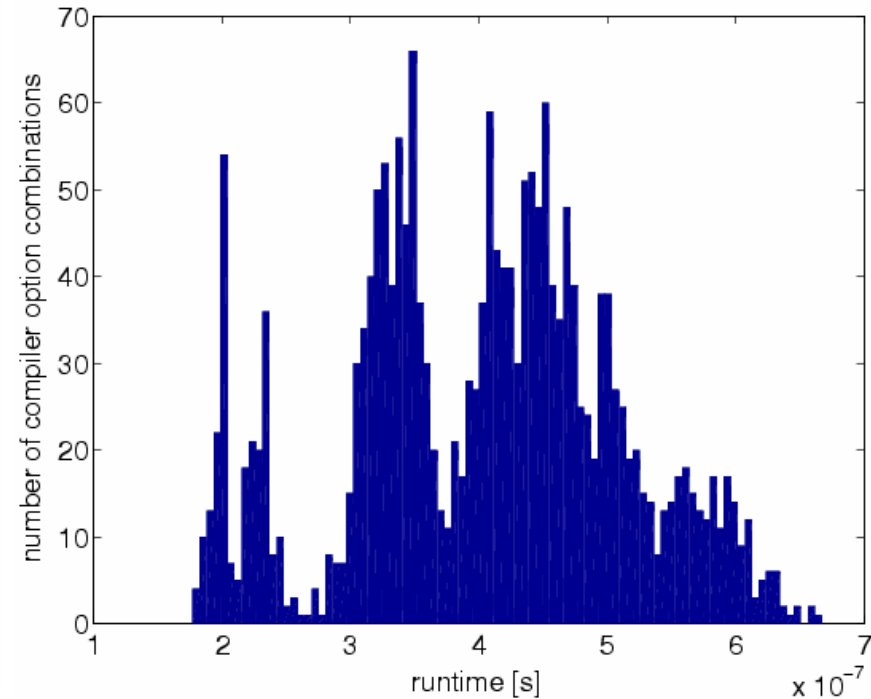
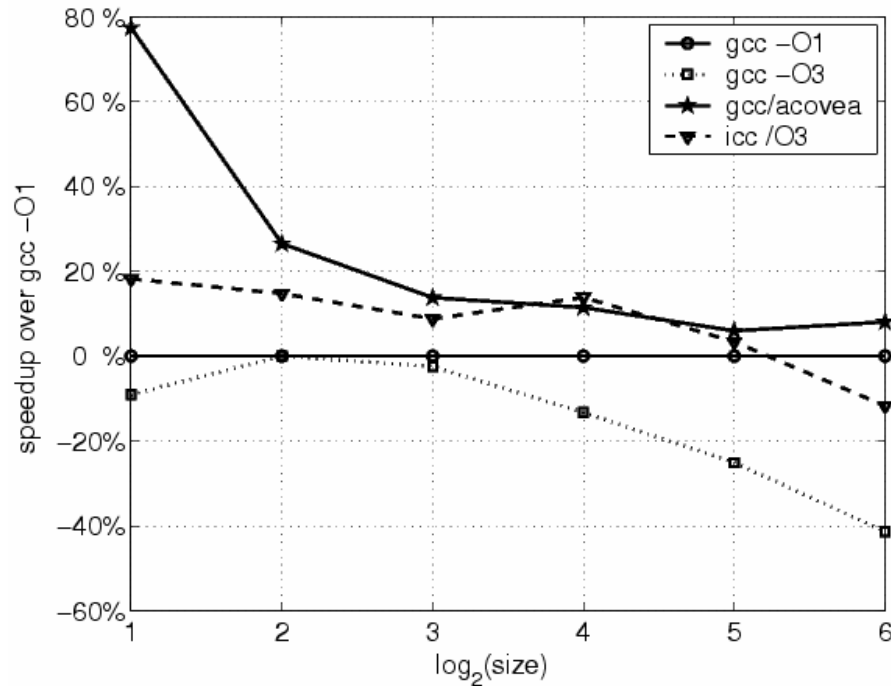


```
gcc {-O1/-O3} -fomit-frame-pointer -std=c99 -fast -mcpu=7450
```

Compileroptionen DCT

P4, 3.2 GHz,
gcc

ACOVEA: Evolutionäre Suche nach Compileroptionen (gcc hat ~500)



**Laufzeithistogramm
zufällige Compileroptionen
inkl. -O1 -march=pentium4**

10% Verbesserung des besten, Spiral-generierten Codes

SPIRAL: Designprinzipien

- Codegenerierung und –adaptierung als Optimierungsproblem auf der Menge aller Algorithmen und Implementierungen

SPIRAL nützt die Struktur der Algorithmen zur Lösung

- Deklarative Sprache (SPL + Regeln) zur Computerrepräsentation des Algorithmenwissens

*Ermöglicht Algorithmengenerierung and -optimierung
(bringt dem System die Mathematik bei)*

- Compiler um Mathe in Code zu übersetzen

- Suche und Lernen in der Menge der Alternativen zur Optimierung

*Schließt den feedback loop
gibt dem System die "Intelligenz"*

- Flexible, erweiterbar

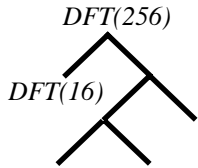
Jeder Schritt in der Codegenerierung ist explizit in SPIRAL

SPIRAL: Abstraktionsebenen

$DFT(256)$

Transformation

- Problemspezifikation



Regelbaum

- natürlicher erster Schritt der Algorithmengenerierung
- speichereffizient
- einfache Manipulation in der Suche

$\dots (I_4 \oplus J_4) D_8 (F_2 \otimes I_4) \dots$

SPL

- Strukturbeschreibung des Algorithmus
- zeigt Parallelismus, Vektorisierbarkeit
- Manipulation für Vektorcode

$\dots \sum_{j=0}^{n-1} S_{(j)m \otimes n} C_n G_{(j)m \otimes n}$

Σ -SPL

- macht Schleifen explizit
- explizite Indexfunktionen
- Schleifenoptimierung

```
...
for (i=0; i<8; i++) {
  t[2i+4]=x[i-4]+x[i+4];
  ... }

```

C Code

- Standard DAG/Compileroptimierungen

Philosophie: Optimiere auf der richtigen Abstraktionsebene

Implementierung: Überblick

- **Das Problem: Moore's Law und Numerische Software**
- **Mikroarchitektur und Architektur**
- **Algorithmus, Implementierung und Mikroarchitektur:
Analyse der DFT und Herleitung von FFTW**
- **Softwaregenerierung und Experimente mit SPIRAL**
- **Zusammenfassung:
Wie schreibt man schnellen numerischen Code?**

Schneller Numerischer Code: Richtlinien (I)

■ Vermeide offensichtliche Fehler

- Kenne die verfügbaren Algorithmen: Verwende gute Algorithmen
- Vorberechnung von Konstanten (z.B., Sinus, Kosinus)
- Gib dem Compiler eine Chance: Schreibe einfachen Code
- Vermeide komplizierte Daten/Kontroll/Infrastrukturen
- Verstehe wo die Laufzeit bleibt: Codeprofiling
- Verwende gute Compileroptionen, probiere Alternativen

■ Datenlokalität für Speicherhierarchie (Caches)

- Rekursiv ist besser als iterativ
- Kenne die Cachegrößen!

■ Blockoptimierung (Block=Codebaustein der auf kleinem Datenbereich operiert)

- Ohne Schleifen und Kontrollstrukturen (aufgerollter Code)
- Dann wichtig: ordne Instruktionen für Lokalität + evtl. andere Optimierungen
- Eventuell: Prüfe den Assemblycode

Schneller Numerischer Code: Richtlinien (II)

■ Mögliche Suche über Alternativen

- akzeptiere dass man nicht alles durch Wissen entscheiden kann
- Such über eine relevante Teilmenge von Algorithmen und/oder Implementierungsweisen

■ Verwende Vektorinstruktionen

- versuch es zuerst mit dem Compiler (konsultiere Manual für Codestil, Pragmas, Alignment, etc.)
- dann verwende Intrinsics, kein Inlineassembly
- umfangreiche Datenumordnungen zerstören die Laufzeit
- prüfe den Assemblycode

■ Nach jeder dieser Optimierungen, überprüfe ob Du noch den richtigen Algorithmus verwendest

Abschließende Bemerkungen

Wichtigste Erkenntnisse (Algorithmen)

- DFT und die Cooley-Tukey FFT können mit dem Chinesischen Restesatz erklärt und hergeleitet werden

$$\text{DFT}_n \leftrightarrow \mathbb{C}[x]/(x^n - 1)$$

$$\text{FFT} \leftrightarrow x^n - 1 = (x^k)^m - 1 \text{ dekomponiert für } n = km$$

- Es gibt sehr viele verschiedenen Algorithmen für die gleiche Transformation (oder Filter)
 - ~ gleiche Anzahl an Operationen
 - unterschiedliche Struktur/Eigenschaften = unterschiedliche Laufzeit und unterschiedliche Eignung für Cachestrukturen, Speizleinstruktionen, etc.

Wichtigste Erkenntnisse (Implementierung)

■ Schneller Code ist schwierig zu schreiben

- man verliert schnell einen Faktor 10-100
- man braucht gutes Verständnis von Algorithmen, Software/Compilern, und der Mikroarchitektur
- Computer entwickeln sich sehr schnell

■ Operationen zählen gibt wenig Information über die Laufzeit

■ In der Zukunft wird es eher schlimmer

- Das Ende von Moore's Law rückt näher
- Mögliche Zukunft: Multicoresysteme.
Siehe Scientific American Artikel (Nov. 2004) "A Split at the Core";
Untertitel: "[...]That is bad news for software companies."

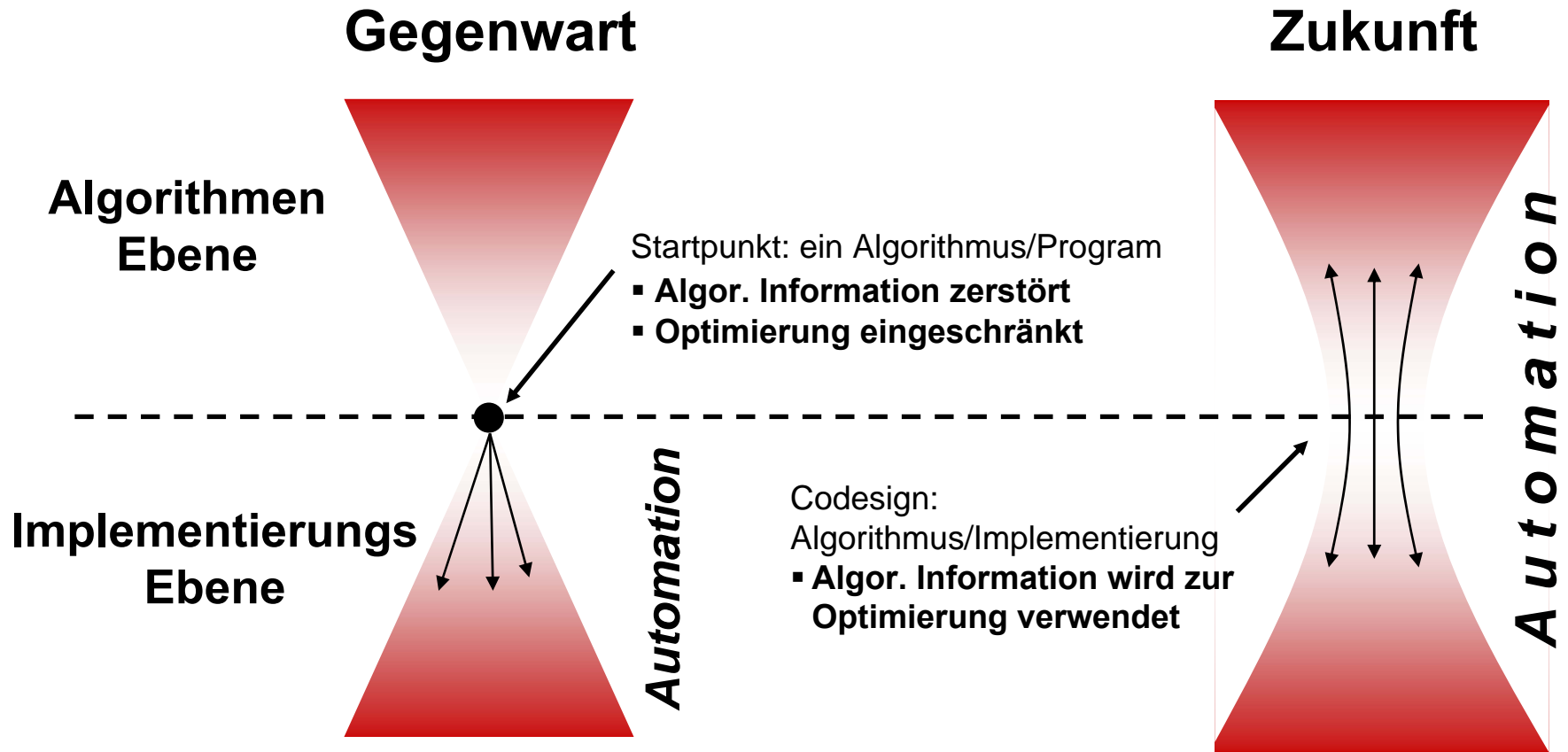
■ Schneller Code mit wenig Aufwand ist aktuelles Forschungsproblem

Überblick aktuelle Forschung

- **Linear Algebra:**
 - LAPACK, ATLAS (Whaley, Dongarra et al., U. Tennessee)
 - BeBOP (Vuduc, Yelick, U. Berkeley)
- **Tensorberechnungen (Quantumchemie):** Sadayappan, Baumgartner et al., Ohio State
- **Finite Elemente Methoden: Fenics (U. Chicago)**
- **Suchen:** Padua, UIUC
- **Signalverarbeitung:**
 - FFTW
 - SPIRAL
 - VSIPL (Kepner, Lebak et al., MIT Lincoln Labs)
- **Neue Compilermethoden (problemspezifisch):**
 - Modelbasiertes ATLAS (Yotov, Pingali, Cornell U.)
 - Broadway (Lin, Guyer, U. Texas Austin)
 - SIMD Optimierungen (Franchetti, Ueberhuber et al., Univ. Techn. Vienna)
 - Teleskopsprachen (Kennedy et al., Rice)
- **Siehe auch: Proceedings of the IEEE Special Issue “Program Generation, Optimization, and Adaptation” (Feb. 05, link auf Spiral Webseite)**

Wir haben gerade erst angefangen

Philosophie



eine neue Generation von problemspezifische Methoden/Werkzeugen

Automation auf allen Abstraktionsebenen

1954: der Fortran Compiler wurde erfunden

2004: heute

**2054: wenn Menschen immer noch code für FFTs etc.
selbst schreiben müssen,
vielleicht sogar in Assembly,
dann haben wir was falsch gemacht**

Das war's, danke für's kommen