

# **Algorithms and Computation in Signal Processing**

**special topic course 18-799B  
spring 2005**

**24<sup>th</sup> and 25<sup>th</sup> Lecture Apr. 07 and 12, 2005**

Instructor: Markus Pueschel

TA: Srinivas Chellappa

# Research Projects

## ■ Presentations last week of April (26<sup>th</sup> and 28<sup>th</sup>)

- We distribute the dates in the lecture on the 12<sup>th</sup>
- Presentations 20 minutes + 5 minutes questions (~17-20 slides)

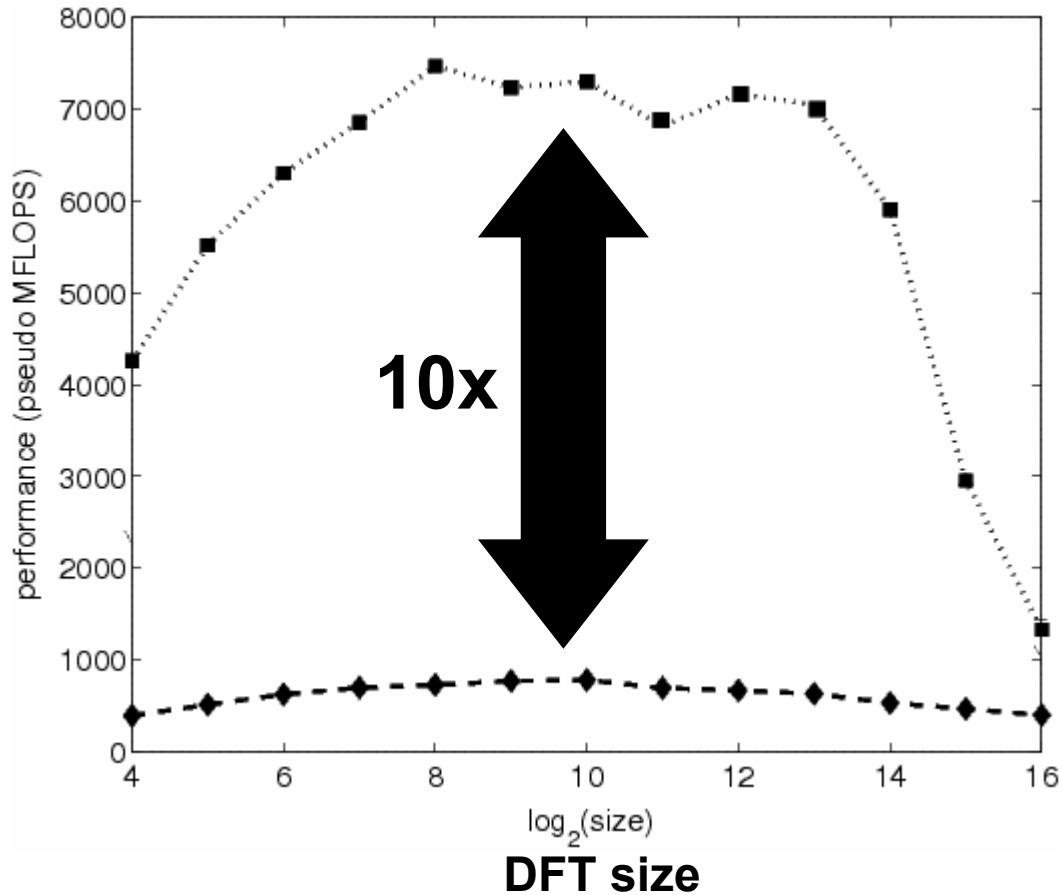
## ■ Research paper

- Due April 20<sup>th</sup>, the only thing that may be missing are some (but not all) experimental results
- You'll get feedback from me
- Final version with feedback incorporated due one week after your presentation

## ■ Remarks

- Follow guide to benchmarking!
- Try different sets of compiler flags to be sure
- Do a cost analysis

# The Problem Again



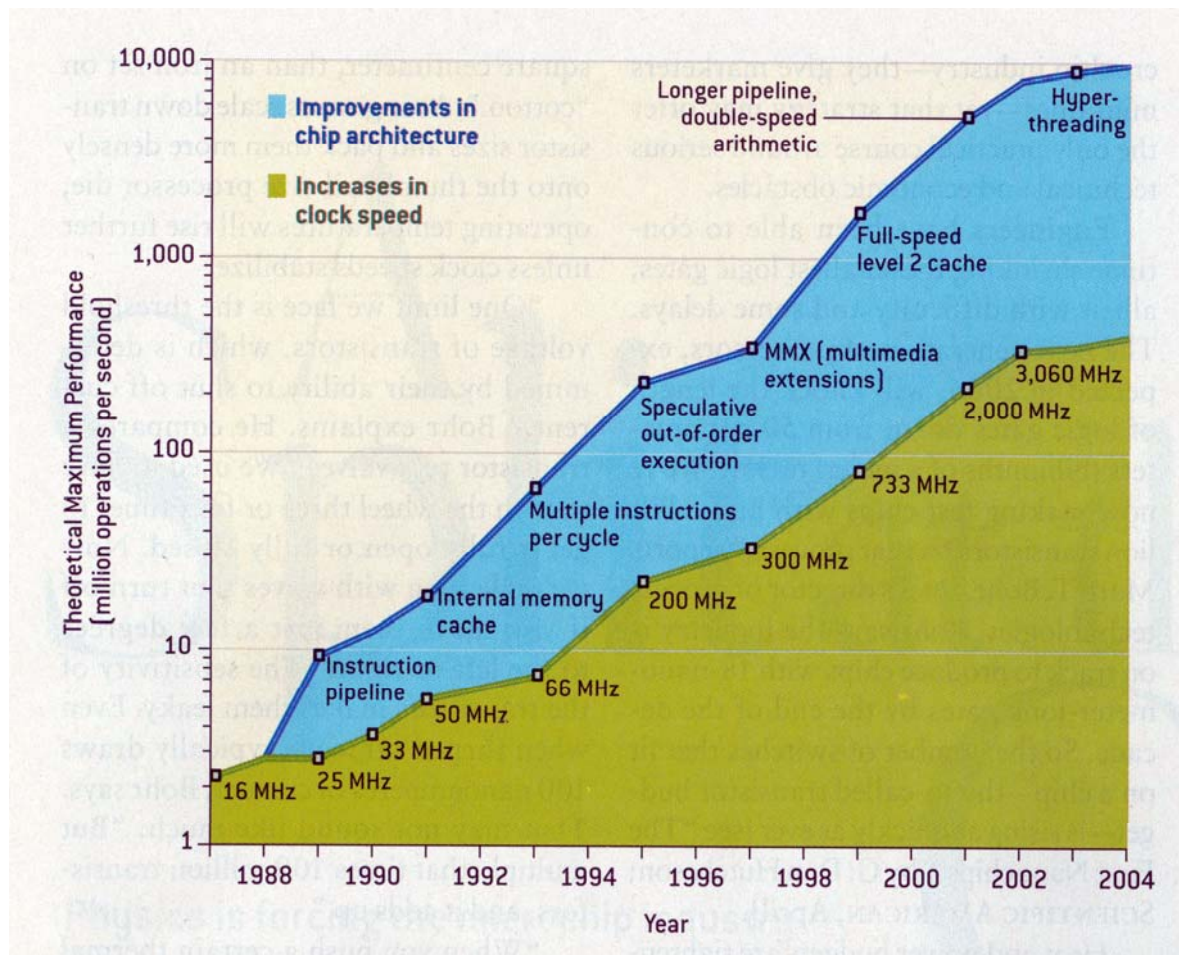
Intel vendor library  
(hand-optimized  
assembly code)  
**but also FFTW, SPIRAL  
generated code**

reasonable  
implementation  
(Numerical recipes.  
GNU scientific library)

***Writing fast numerical code is a tough problem***

# Moore's Law

- Moore's Law: exponential (x2 in ~18 months) increase number of transistors/chip



source: Scientific American, Nov 2004, p. 98

***But everything has its price ...***

# Moore's Law: Consequences

- **Computers are very complex**
  - multilevel memory hierarchy
  - special instruction sets beyond standard C programming model
  - undocumented hardware optimizations
- **Consequences:**
  - Runtime depends only roughly on the operations
  - Runtime behavior is hard to understand
  - Compiler development can hardly keep track
  - **The best code (and algorithm) is platform-dependent**
  - **It is very difficult to write really fast code**
- **Computers evolve fast**
  - Highly tuned code becomes obsolete almost as fast as it as written
- **It'll get rather worse: Multicoresystems**

# Solution #1: Brute Force

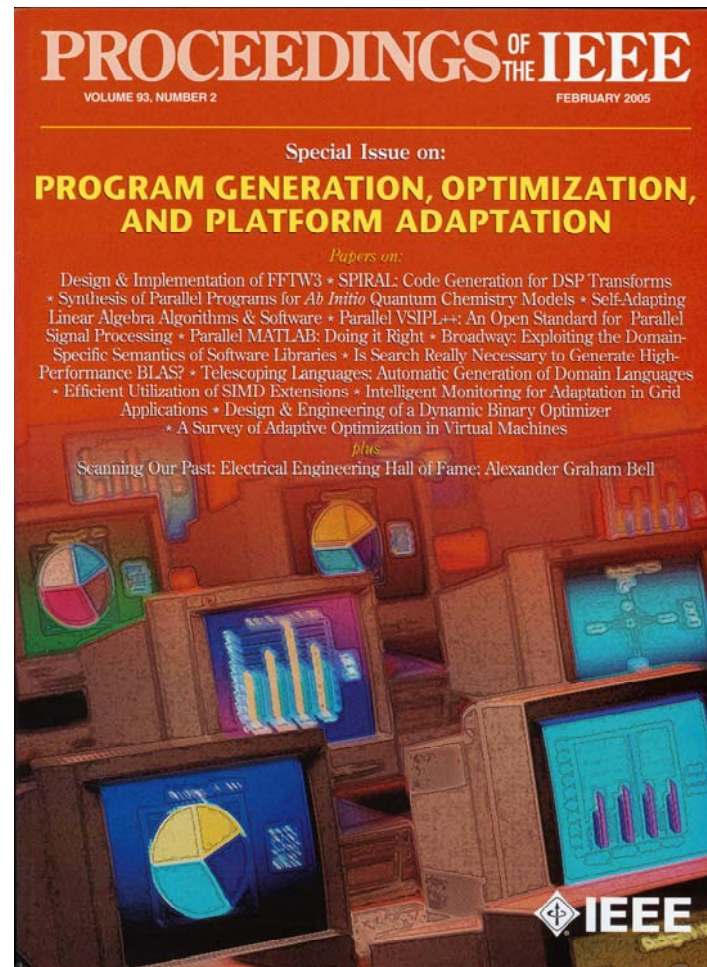
- Thousands of programmers hand-write and hand-tune (assembly) code for the same numerical problems and for every platform and whenever a new platform comes out?

*Hmm.....*

*(but it's current practice)*

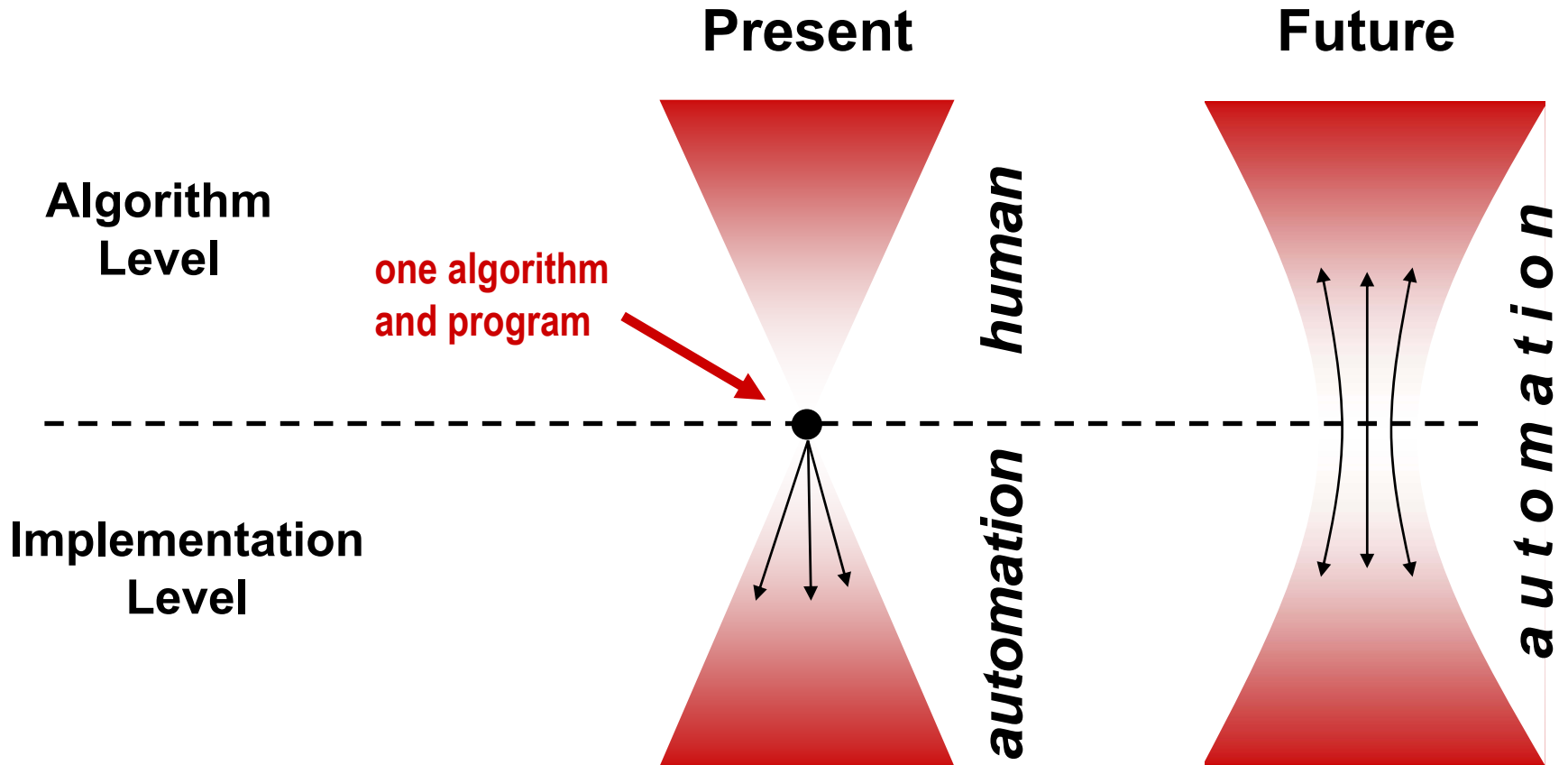
# Solution #2: New Approaches to Code Optimization and Code Creation

- **ATLAS**: Code generation/optimization for BLAS
- **SPARSITY/BeBop**: Code generation/optimization for sparse linear algebra routines
- **FFTW**: Self-adaptive DFT library + DFT kernel generator
- **SPIRAL**: Code generation/optimization for linear signal transforms



Proceedings of the IEEE special issue, Feb. 2005

# Common Philosophy



*a new breed of domain-aware approaches/tools  
push automation beyond what is currently possible  
applies for software and hardware design alike*



# SPIRAL [www.spiral.net](http://www.spiral.net)

## Sponsors:

DARPA

NSF ACR-0234293

NSF ITR/NGS-0325687

## and:

Cylab, CMU

Austrian Science Fund

Intel

ITRI, Taiwan

ENSCO, Inc.

~40 Publications

## Overview paper:

Markus Püschel, José M. F. Moura, Jeremy Johnson, David Padua, Manuela Veloso, Bryan Singer, Jianxin Xiong, Franz Franchetti, Aca Gacic, Yevgen Voronenko, Kang Chen, Robert W. Johnson, and Nick Rizzolo,  
**SPIRAL: Code Generation for DSP Transforms**,  
 Proceedings of the IEEE

## Team:

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.)

Christoph Überhuber (Math, TU Wien)

## Students/PostDocs:

Bryan W. Singer (CS, CMU)

Jianxin Xiong (CS, UIUC)

Srinivas Chellappa (ECE, CMU)

Franz Franchetti (ECE, CMU, before TU Vienna)

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 Vienna)

...

# Spiral

- Code generation from scratch for linear digital signal processing (DSP) transforms (DFT, DCT, DWT, filters, ....)
- Automatic optimization and platform-tuning at the algorithmic level and the code level
- Different code types supported (scalar, vector, FMA, fixed-point, multiplierless, ...)

**Goal:** A flexible, extensible code generation framework that can survive time (to whatever extent possible) for an entire domain of algorithms

**Research question:** To what extent is it possible to abolish handcoding and handoptimization?

# Code Generation and Tuning as Optimization Problem

$T$  a DSP transform to be implemented

$P$  the target platform

$\mathcal{I} = \mathcal{I}(T, P)$  set of possible implementations of  $T$  on  $P$

$C = C(T, I, P)$  cost of implementation  $I$  of  $T$  on  $P$

The implementation of  $T$  that is tuned to  $P$  is given by:

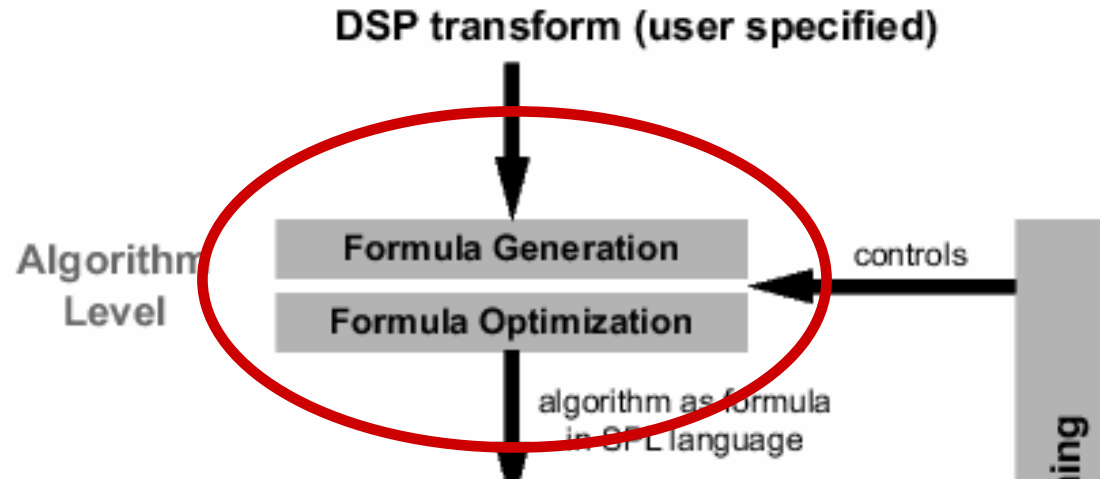
$$\hat{I} = \hat{I}(P) = \arg \min_{I \in \mathcal{I}(P)} C(T, P, I)$$

## Problems:

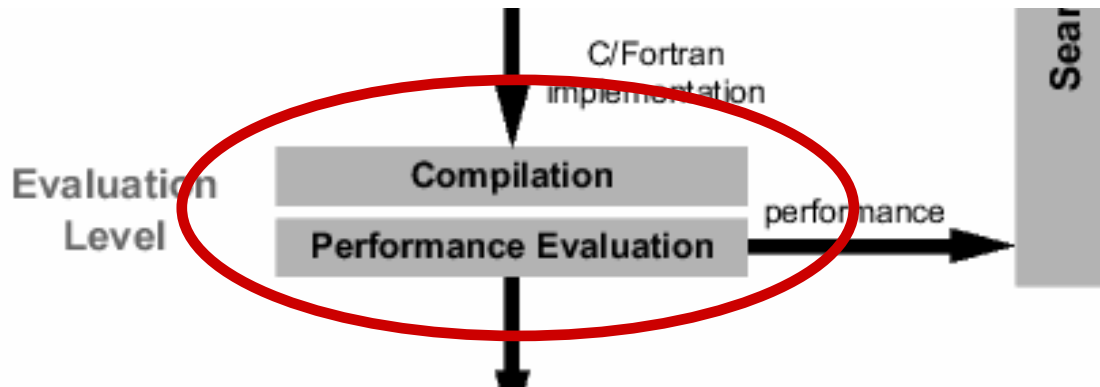
- How to characterize and generate the set of implementations?
- How to efficiently minimize  $C$ ?

**Spiral exploits the domain-specific structure to implement a solver for this optimization problem**

# Spiral's architecture



**Domain knowledge:  
Generating algorithms & manipulating algorithms**



**Architecture knowledge:  
by evaluating runtime**

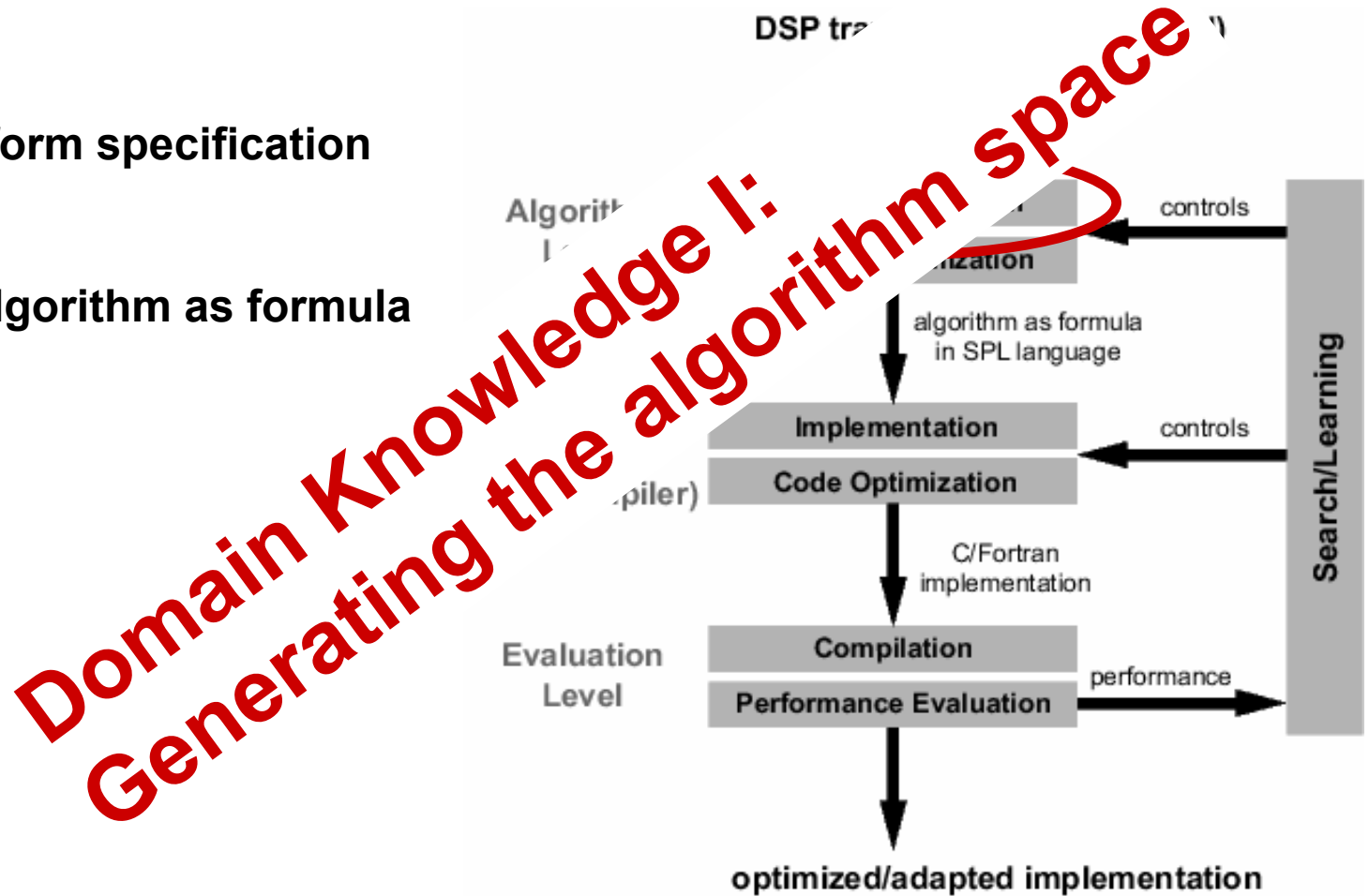
# From Transform to Algorithm (Formula)

Input:

Transform specification

Output:

Fast algorithm as formula



# DSP Algorithms: Example 4-point DFT

Cooley/Tukey FFT (size 4):

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

Fourier transform

Diagonal matrix (twiddles)

$$DFT_4 = (DFT_2 \otimes I_2) \cdot T_2^4 \cdot (I_2 \otimes DFT_2) \cdot L_2^4$$

Kronecker product

Identity

Permutation

- mathematical notation exhibits structure: **SPL (signal processing language)**
- Suitable for **computer representation**
- contains **all information** to generate code

# SPL: Definition (BNF)

- Description language for linear DSP algorithms
- 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)
	$\dots$	
$\langle \text{generic} \rangle ::=$	$\text{diag}(a_0, \dots, a_{n-1}) \mid \dots$	
$\langle \text{symbol} \rangle ::=$	$\mathbf{I}_n \mid \mathbf{J}_n \mid \mathbf{L}_k^n \mid \mathbf{R}_\alpha \mid \mathbf{F}_2 \mid \dots$	
$\langle \text{transform} \rangle ::=$	$\mathbf{DFT}_n \mid \mathbf{WHT}_n \mid \mathbf{DCT-2}_n \mid \mathbf{Filt}_n(h[z]) \mid \dots$	

## Some Definitions:

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

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

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

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

$$\mathbf{I}_n \otimes A = \begin{bmatrix} A & & & \\ & A & & \\ & & \dots & \\ & & & A \end{bmatrix}$$

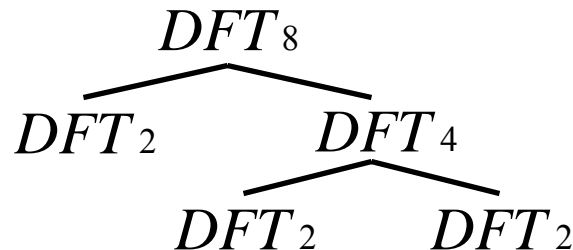
# DSP Algorithms: Spiral Terminology

**Transform**  $DFT_n$  parameterized matrix

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

- a breakdown strategy
- product of sparse matrices

**Ruletree**



- recursive application of rules
- uniquely defines an algorithm
- efficient representation
- easy manipulation

**Formula**

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

- few constructs and primitives
- uniquely defines an algorithm
- can be translated into code



# Some Transforms

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

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

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

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

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

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

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

**Spiral currently contains 36 transforms**

# Some Breakdown Rules

$$\text{DFT}_n \rightarrow (\text{DFT}_k \otimes \mathbf{I}_m) \mathbf{T}_m^n (\mathbf{I}_k \otimes \text{DFT}_m) \mathbf{L}_k^n, \quad n = km$$

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

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

$$\begin{aligned} \text{DCT-3}_n \rightarrow & (\mathbf{I}_m \oplus \mathbf{J}_m) \mathbf{L}_m^n (\text{DCT-3}_m(1/4) \oplus \text{DCT-3}_m(3/4)) \\ & \cdot (\mathbf{F}_2 \otimes \mathbf{I}_m) \begin{bmatrix} \mathbf{I}_m & 0 \oplus -\mathbf{J}_{m-1} \\ \frac{1}{\sqrt{2}}(\mathbf{I}_1 \oplus 2\mathbf{I}_m) \end{bmatrix}, \quad n = 2m \end{aligned}$$

$$\text{DCT-4}_n \rightarrow S_n \text{DCT-2}_n \text{diag}_{0 \leq k < n} (1 / (2 \cos((2k + 1)\pi / 4n)))$$

$$\text{IMDCT}_{2m} \rightarrow (\mathbf{J}_m \oplus \mathbf{I}_m \oplus \mathbf{I}_m \oplus \mathbf{J}_m) \left( \left( \begin{bmatrix} 1 \\ -1 \end{bmatrix} \otimes \mathbf{I}_m \right) \oplus \left( \begin{bmatrix} -1 \\ -1 \end{bmatrix} \otimes \mathbf{I}_m \right) \right) \mathbf{J}_{2m} \text{DCT-4}_{2m}$$

$$\text{WHT}_{2^k} \rightarrow \prod_{i=1}^t (\mathbf{I}_{2^{k_1 + \dots + k_{i-1}}} \otimes \text{WHT}_{2^{k_i}} \otimes \mathbf{I}_{2^{k_{i+1} + \dots + k_t}}), \quad k = k_1 + \dots + k_t$$

$$\text{DFT}_2 \rightarrow \mathbf{F}_2$$

$$\text{DCT-2}_2 \rightarrow \text{diag}(1, 1/\sqrt{2}) \mathbf{F}_2$$

$$\text{DCT-4}_2 \rightarrow \mathbf{J}_2 \mathbf{R}_{13\pi/8}$$

**Base case rules**

**Spiral contains 100+ rules**

# Some Breakdown Rules for Filters

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

$$\mathbf{Filt}_n(h(z)) \rightarrow \mathbf{L}_{\frac{n}{2}} \mathbf{Filt}_{\frac{n}{2}} \left( \left[ \begin{array}{ccc} 1 & 0 & 1 \\ 0 & 1 & 1 \end{array} \right] \right).$$

$$\mathbf{Filt}_{\frac{n}{2}} \left( \left[ \begin{array}{c} h_0(z) \\ h_1(z) \\ h_0(z) + h_1(z) \end{array} \right] \right) \cdot \mathbf{Filt}_{\frac{n}{2} + \frac{r+l-1}{2}} \left( \left[ \begin{array}{cc} 1 & -1 \\ z & -1 \\ 0 & 1 \end{array} \right] \right) \cdot \mathbf{L}_2^{n+r+l}$$

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

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

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

# Formula- (Algorithm) generation

**Transform:**

DCT-2<sub>4</sub>



**Ruletree:**

DCT-2<sub>4</sub>

DCT-4<sub>2</sub>

DCT-2<sub>2</sub>

(many possibilities)



**Formula:**

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



**Remaining task**

**(fast)**

**C Code:**

```
void sub(double *y, double *x) {
    double f0, f1, f2, f3, f4, f7, f8, f10, f11;
    f0 = x[0] - x[3];
    f1 = x[0] + x[3];
    f2 = x[1] - x[2];
    f3 = x[1] + x[2];
    f4 = f1 - f3;
    y[0] = f1 + f3;
    y[2] = 0.7071067811865476 * f4;
    f7 = 0.9238795325112867 * f0;
    f8 = 0.3826834323650898 * f2;
    y[1] = f7 + f8;
    f10 = 0.3826834323650898 * f0;
    f11 = (-0.9238795325112867) * f2;
    y[3] = f10 + f11;
}
```

# Set of Algorithms

## ■ Given a transform:

- Apply breakdown rules recursively until all occurring transforms are expanded
- Choice of rules at each step yields (usually) exponentially large algorithms space:
  - about equal in operations count
  - differ in data flow

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

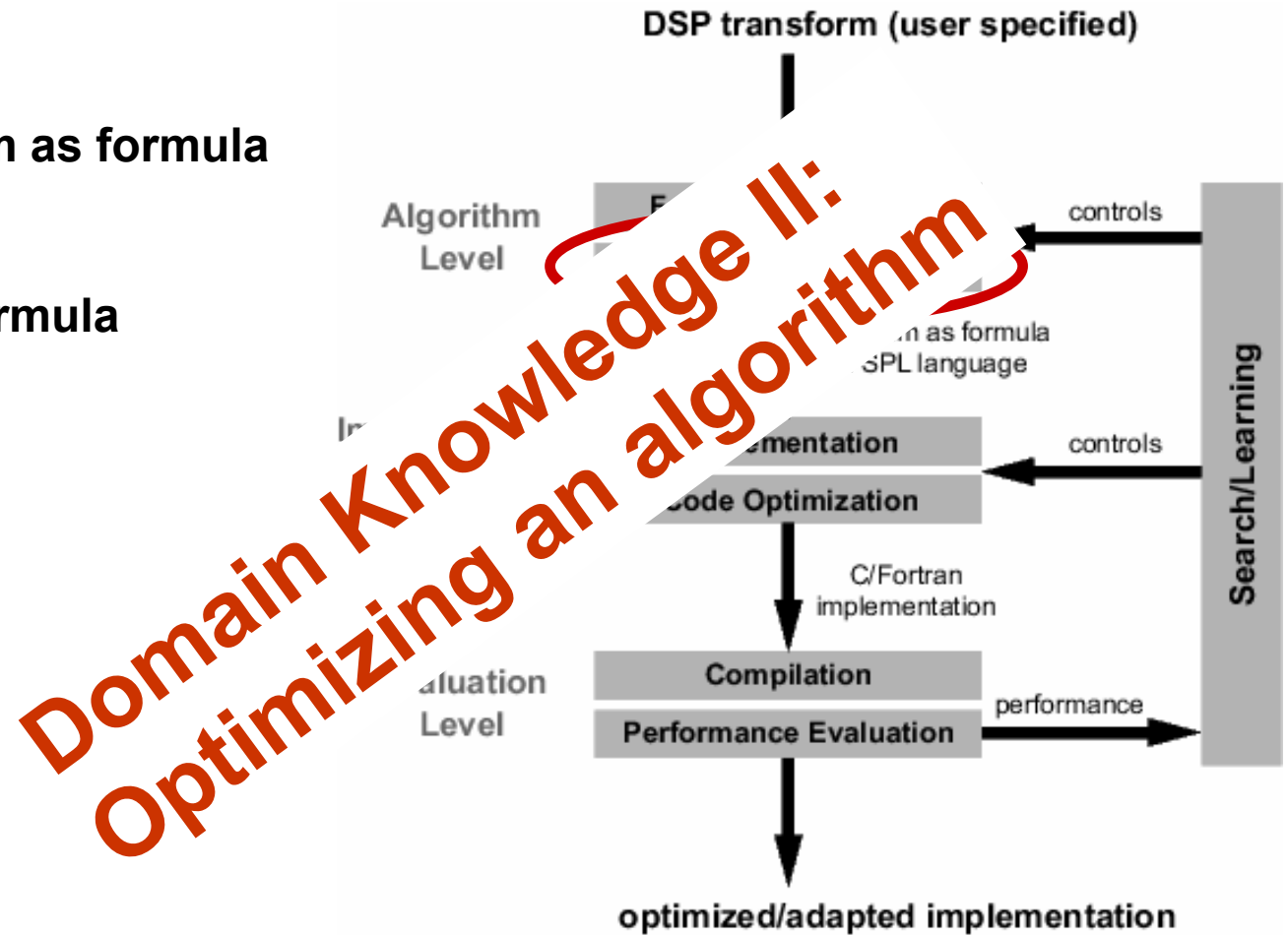
# From Algorithm (Formula) to Optimized Algorithm

Input:

Fast algorithm as formula

Output:

Optimized formula



# Motivation: Loop Fusion

$8$   
 $L_4$   
 $(I_4 \otimes F_2)$

direct  
code generation

```
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]; } }
```

no compiler does that

**Solution:  $\Sigma$ -SPL and  
Formula manipulation**

```
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]; } }
```

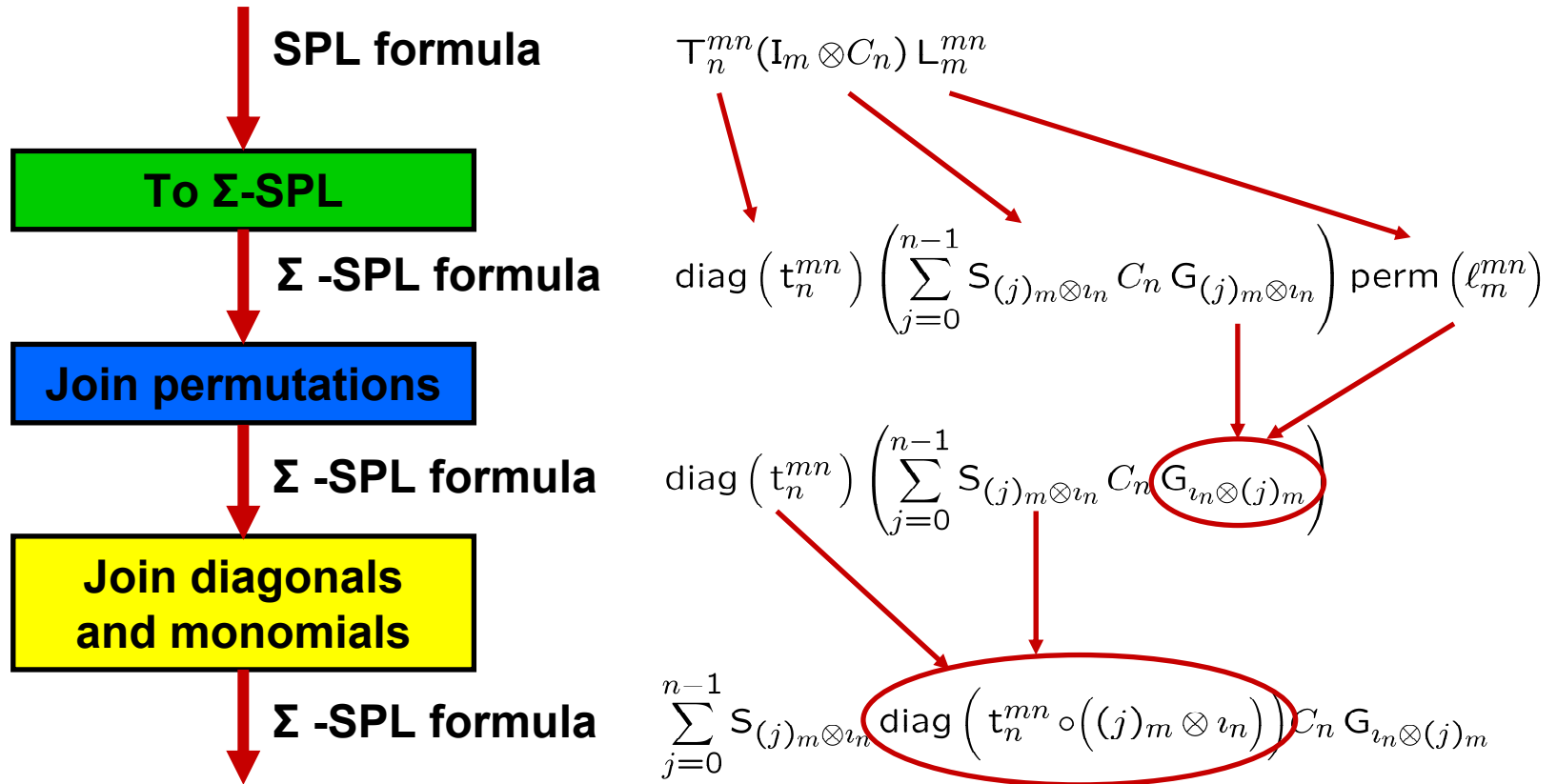
$\sum_{j=0}^3 S_{(j)_4 \otimes I_2} F_2 G_{I_2 \otimes (j)_4}$

# Formula Level Optimization

- **Main goals:**
  - Fusing iterative steps (fusing loops), e.g., permutations with loops
  - Improving structure (data flow) for SIMD instructions
  
- **Overcomes compiler limitations**
  
- **Formula manipulation through mathematical rules**
  
- **Implemented using multiple levels of rewriting systems**
  
- **Puts math knowledge into the system**



# Structure of Loop Optimization



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

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

# Loop Fusion Beyond Cooley-Tukey

Main DFT recursion (breakdown rules):

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

rather cheap



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

expensive

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

very expensive

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

How to fuse permutations from different combinations of rules?

# Example

- Consider the DFT formula fragment

$$(I_p \otimes (I_1 \oplus (I_r \otimes \text{DFT}_s) L_r^{rs}) W_p) V_{p,q}$$

Cooley-Tukey

Rader

Good-Thomas

- In  $\Sigma$ -SPL:

$$\sum_{j_1=0}^{p-1} \left( S_{((j_1)_p \otimes \iota_q) \circ (0)_+^{1 \rightarrow q} \circ \iota_1} G_{v^{p,q} \circ ((j_1)_p \otimes \iota_q) \circ \bar{w}_{1,g}^q \circ (0)_+^{1 \rightarrow q}}$$

$$+ \sum_{j_0=0}^{r-1} S_{((j_1)_p \otimes \iota_q) \circ (1)_+^{q-1 \rightarrow q} \circ ((j_0)_r \otimes \iota_s)} \text{DFT}_s$$

**Complicated  
array access**

$$G_{v^{p,q} \circ ((j_1)_p \otimes \iota_q) \circ \bar{w}_{1,g}^q \circ (1)_+^{q-1 \rightarrow q} \circ \ell_r^{rs} \circ ((j_0)_r \otimes \iota_s)} \Bigg) \cdot$$

# Example (cont'd)

- After index function simplification:

$$\sum_{\substack{j_1=0 \\ b_1=qj_1}}^{p-1} \left( S_{h_{0,q}^{p \rightarrow pq} \circ (j_1)_p} G_{\hbar_{0,q}^{p \rightarrow pq} \circ (j_1)_p} + \right. \\ \left. \sum_{\substack{j_0=0 \\ \phi_1=g^{j_0}}}^{r-1} S_{h_{qj_1+sj_0+1,1}^{s \rightarrow pq}} \text{DFT}_s G_{\hbar_{b_1,p}^{q \rightarrow pq} \circ w_{\phi_1,g^s}^{s \rightarrow q}} \right)$$

**Simplified  
array access**

# Example (cont'd)

## ■ Generated Code

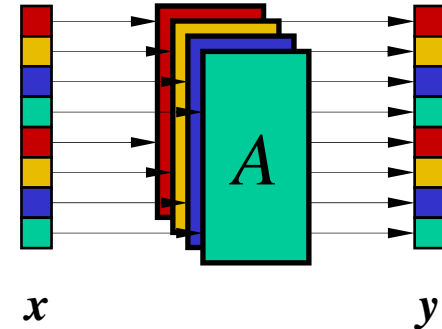
```
// Input: _Complex double x[28], output: y[28]
int p1, b1;
for(int j1 = 0; j1 <= 3; j1++) {
    y[7*j1] = x[(7*j1%28)];
    p1 = 1; b1 = 7*j1;
    for(int j0 = 0; j0 <= 2; j0++) {
        y[b1 + 2*j0 + 1] =
            x[(b1 + 4*p1)%28] + x[(b1 + 24*p1)%28];
        y[b1 + 2*j0 + 2] =
            x[(b1 + 4*p1)%28] - x[(b1 + 24*p1)%28];
        p1 = (p1*3%7);
    }
}
```

# Vector code generation from SPL formulas

Naturally vectorizable construct

$$A \otimes I_4$$

vector length



(Current) generic construct **completely vectorizable**:

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

$P_i$	$Q_i$	permutations
$D_i$	$E_i$	diagonals
$A_i$		arbitrary formulas
$v$		SIMD vector length

Vectorization in two steps:

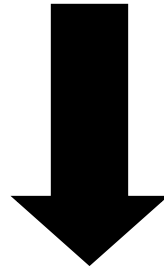
1. Formula manipulation using manipulation rules
2. Code generation (vector code + C code)

**Formula manipulation overcomes compiler limitations**

# Example DFT

## Standard FFT

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



**Formula manipulation**

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

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

**Vector FFT for  $\nu$ -way vector instructions**

# Implementation of Formula Generation and Manipulation

- Implementation using a computer algebra system (GAP)
- SPL/ $\Sigma$ -SPL implemented as recursive data types
- Exact representation of  $\sin()$ ,  $\cos()$ , etc.
- Symbolic computation enables exact verification of rules



# From Optimized Algorithm (Formula) to Code

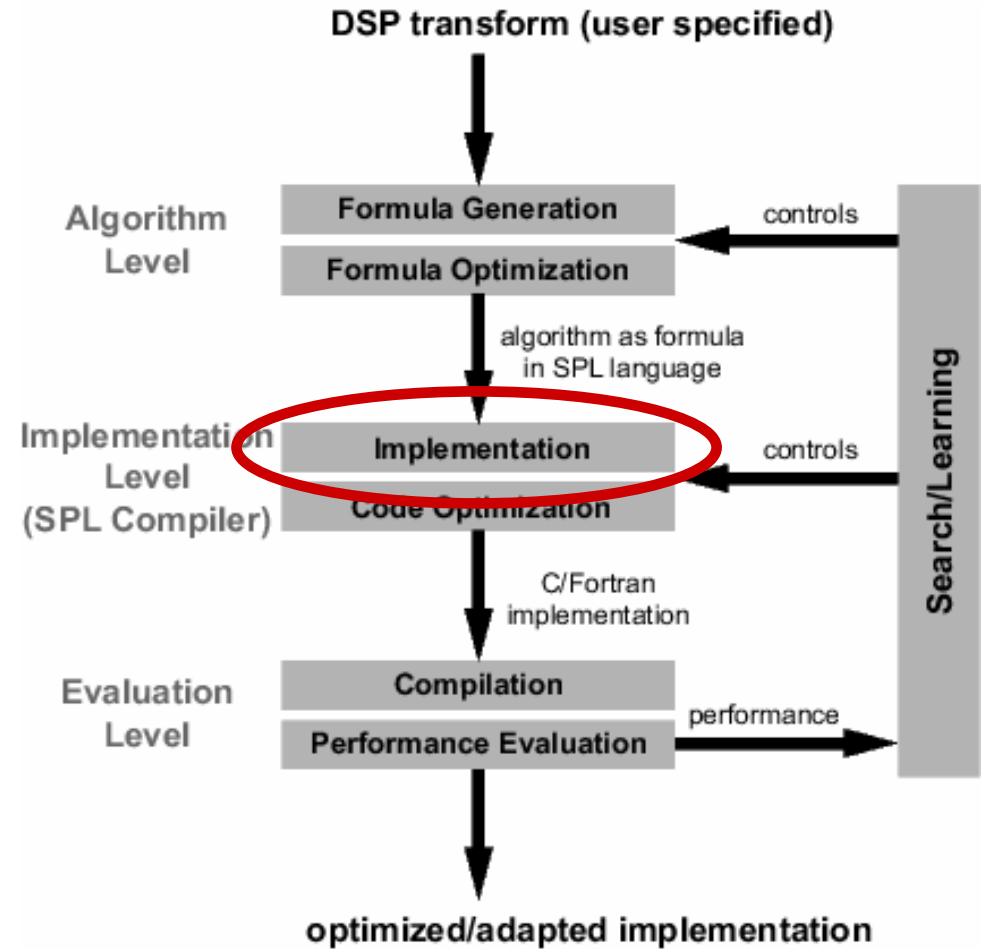
Input:

Optimized formula

Output:

Intermediate Code

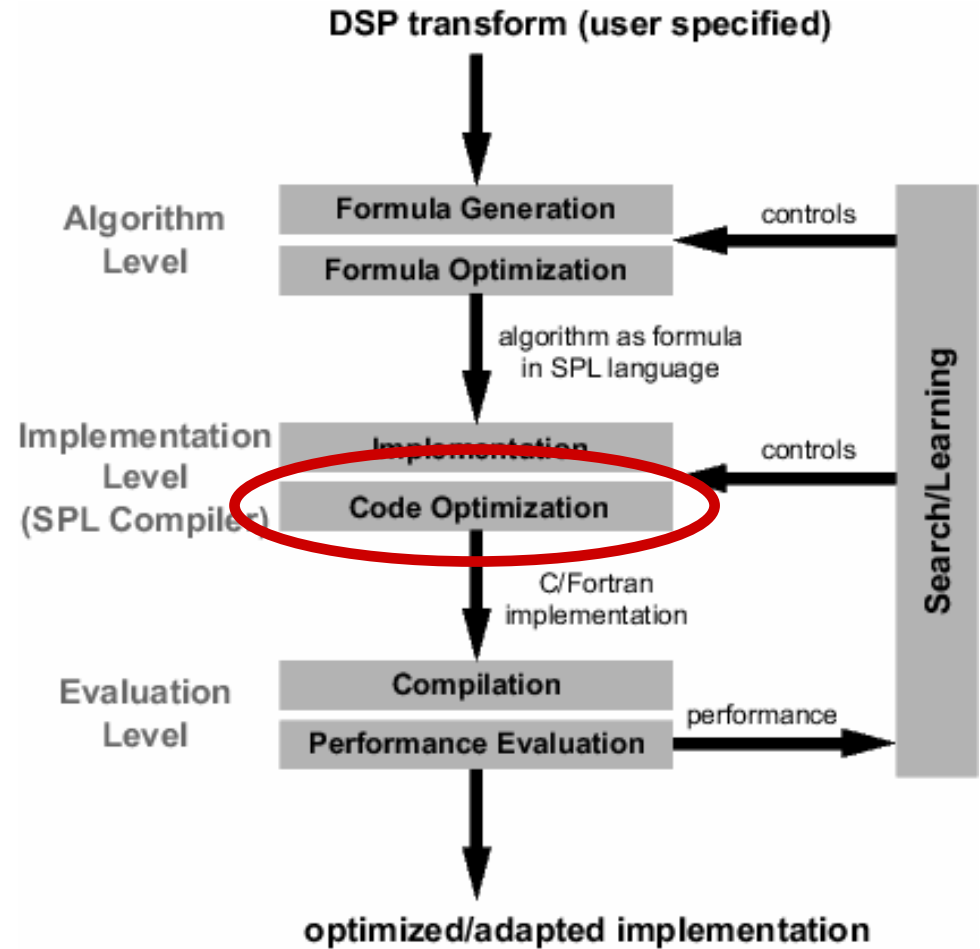
**Straightforward**



# From Code to Optimized Code

**Input:**  
Intermediate Code

**Output:**  
Optimized C code



# Code Level Optimizations

- **Precomputation of constants**
- **Loop unrolling (controlled by search module)**
- **Constant inlining**
- **SSA code, scalar replacement, algebraic simplifications, CSE**
- **Code reordering for locality (optional)**
- **Conversion to FMA code (optional)**
- **Conversion to fixed point code (optional)**
- **Conversion to multiplierless code (optional)**
  
- **Finally: Unparsing to C (or Fortran)**

# Conversion to FMA code

- FMA (fused multiply-add) or MAC (multiply accumulate) instructions:  $y = \pm ax \pm b$
- Extension of the instruction set + specialized execution units
- As fast as a single add or multiply
- Conversion of linear algorithms to FMA code: **blackboard**
- Paper: Yevgen Voronenko and Markus Püschel  
[Automatic Generation of Implementations for DSP Transforms on Fused Multiply-Add Architectures](#)  
Proc. (ICASSP) 2004

# Evaluating Code

Input:

Optimized C code

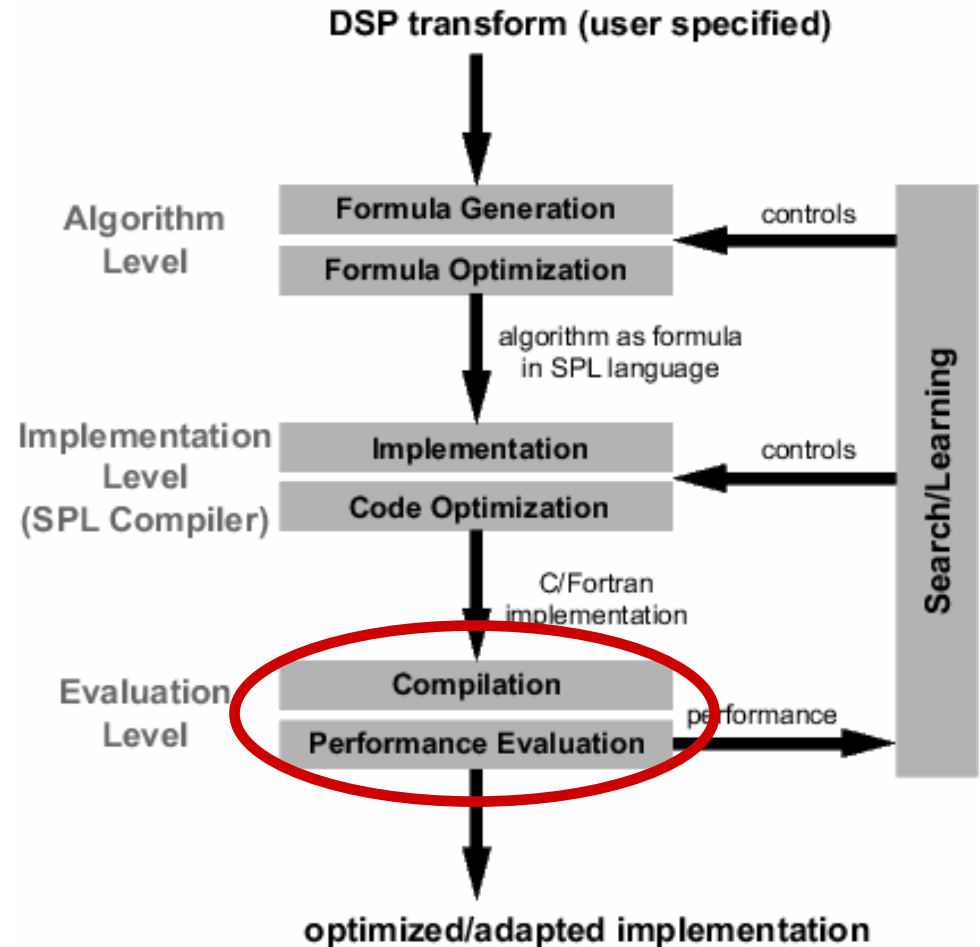
Output:

Performance Number

**Straightforward**

Examples:

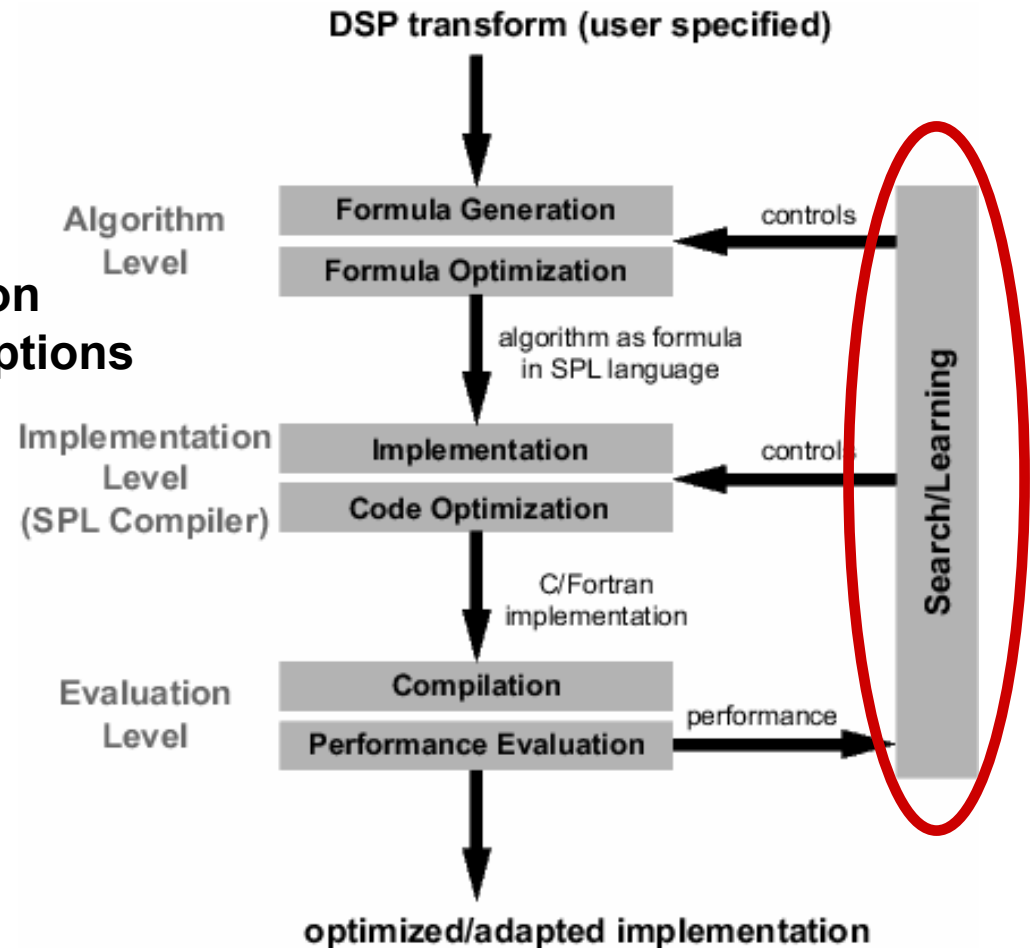
- runtime
- accuracy
- operations count



# Search (Learning) for the Best

**Input:**  
Performance Number

**Output:**  
Controls Formula Generation  
Controls Implementation Options

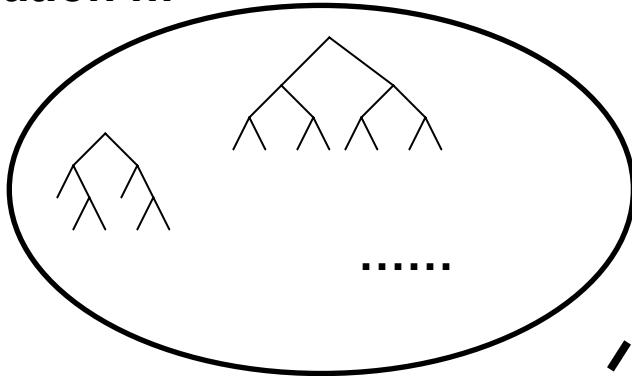


# Search Methods

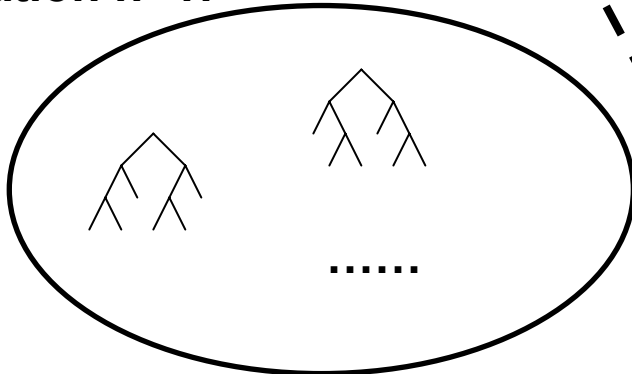
- **Search over:**
  - Algorithmic degrees of freedom  
(choice of breakdown rules)
  - Implementation degrees of freedom  
(degree of unrolling)
  
- **Operates with the ruletree representation of an algorithm**
  - transform independent
  - efficient
  
- **Search Methods**
  - Exhaustive Search
  - Dynamic Programming (DP)
  - Random Search
  - Hill Climbing
  - STEER (an evolutionary algorithm)

# STEER: Evolutionary Search

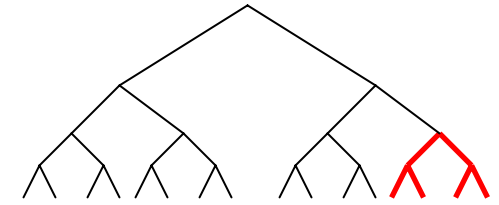
Population n:



Population n+1:

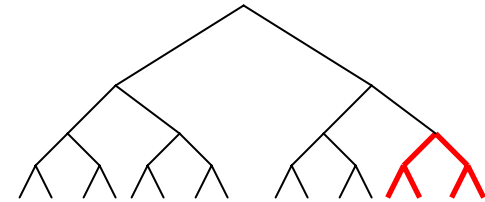


**Mutation**

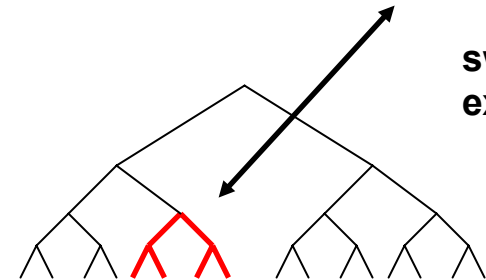


expand differently

**Cross-Breeding**



swap expansions



**Survival of Fittest**



# Learning

## ■ Procedure:

- Generate a set of (1000 say) algorithms and their runtimes (one transform, one size); represent algorithms by features
- From this data (pairs of features and runtimes), learn a set of algorithm design rules
- From this set, generate best algorithms (theory of Markov decision processes)

## ■ Evaluation:

- Tested for WHT and DFT
- From data generated for one size ( $2^{15}$ ) could construct best algorithms across sizes ( $2^{12}$ - $2^{18}$ )

Bryan Singer and Manuela Veloso

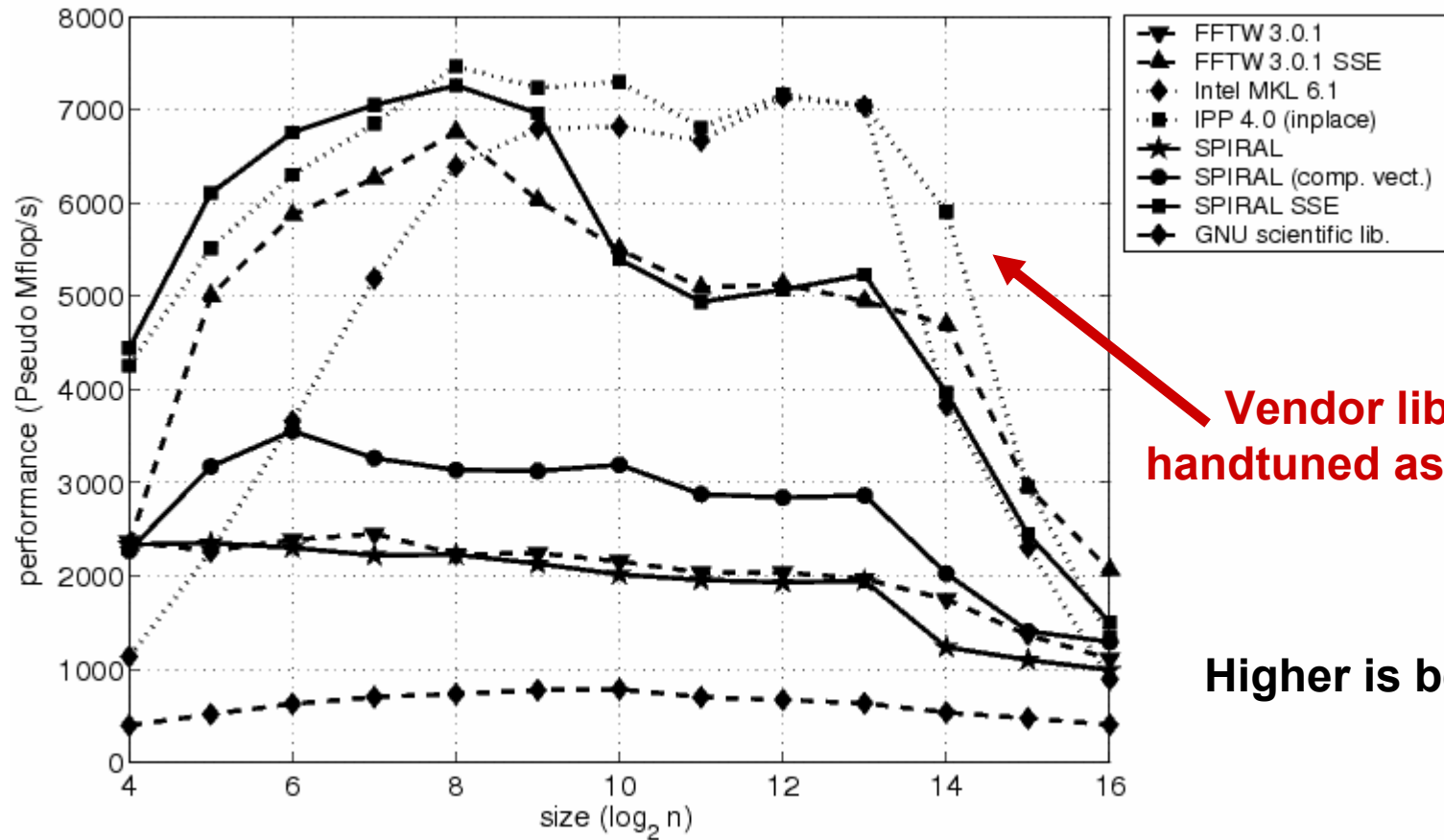
**Learning to Construct Fast Signal Processing Implementations**

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

# Benchmarks

# Benchmark: DFT, 2-powers

P4, 3.2 GHz,  
icc 8.0



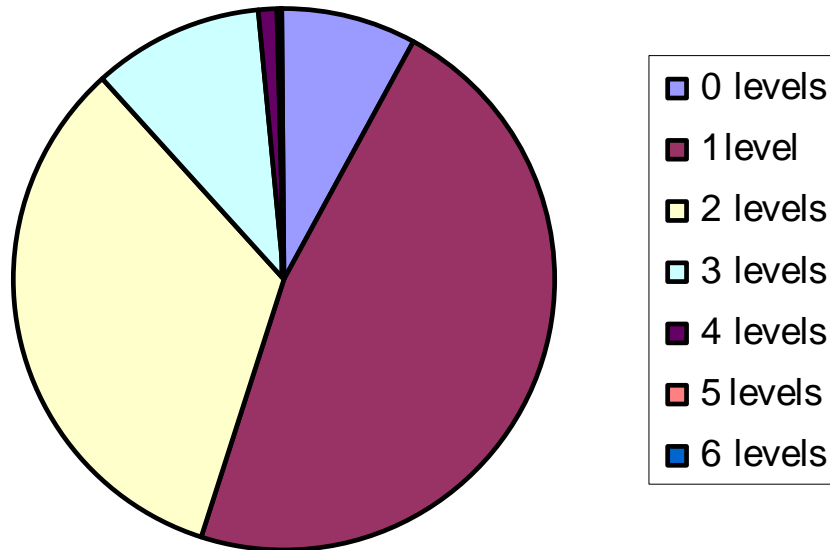
**Vendor library:  
handtuned assembly?**

**Higher is better**

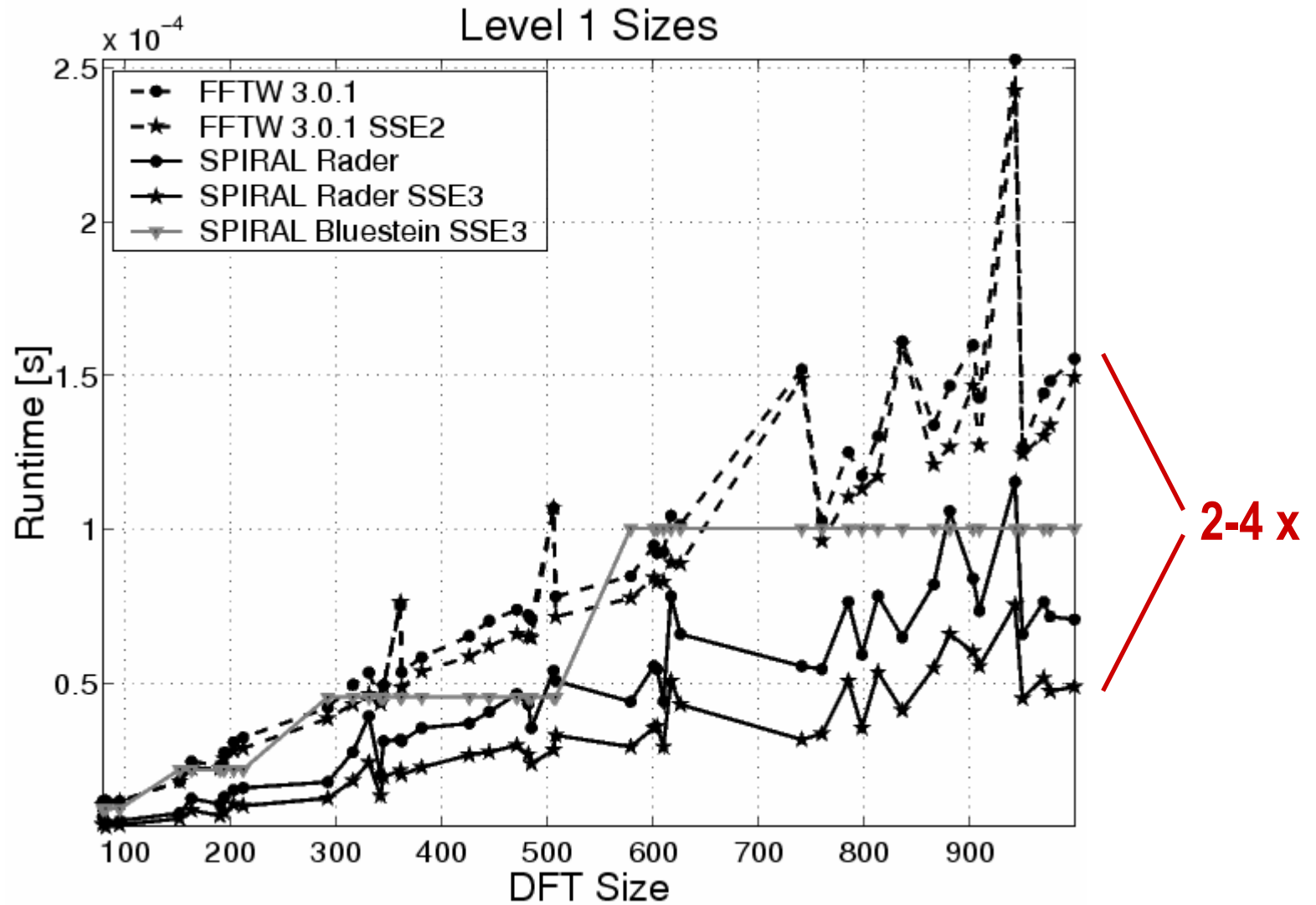
**Single precision**

# Benchmark: DFT, Other Sizes

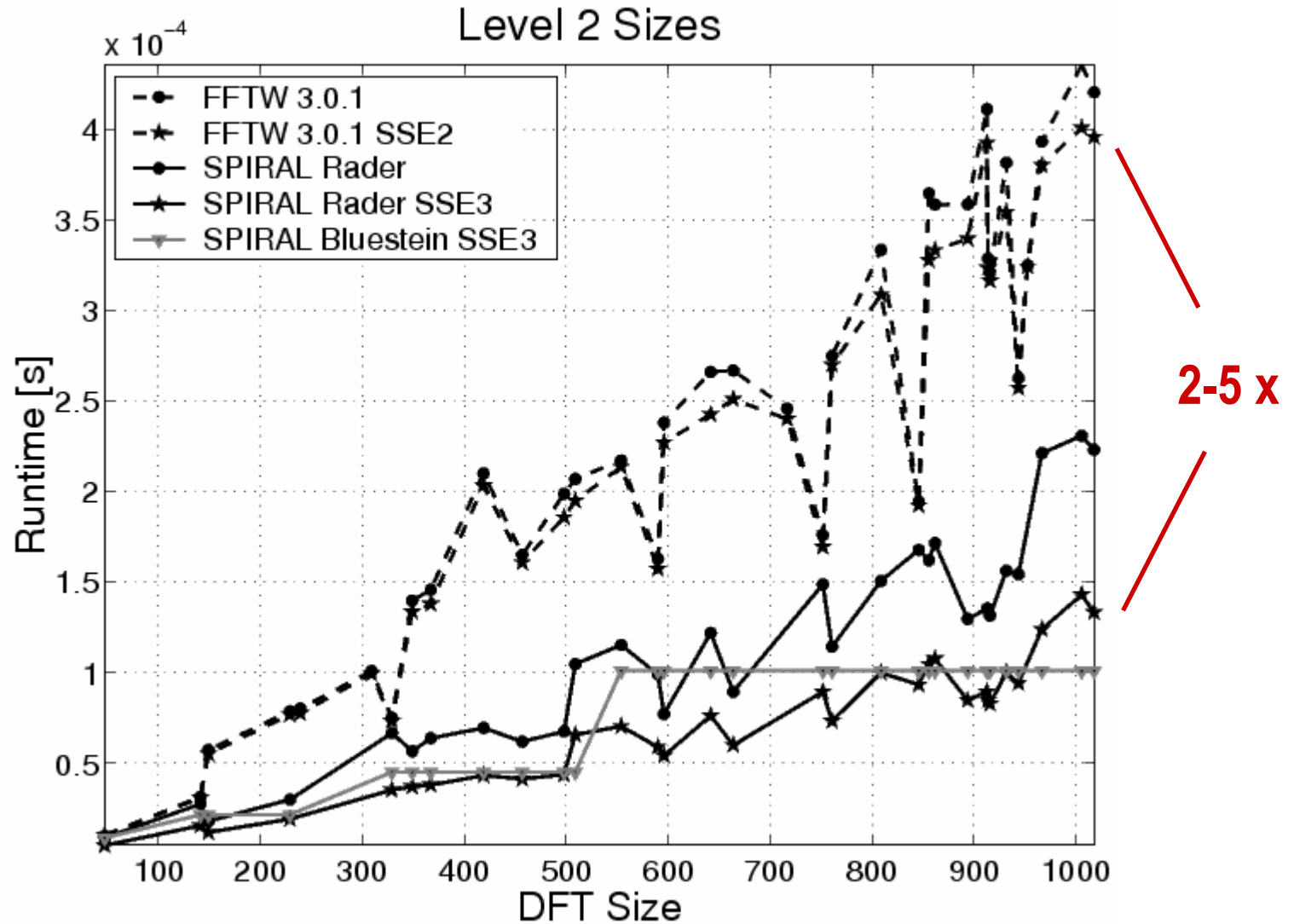
- Divide sizes into levels by number of necessary Rader steps
- $n < 8192$



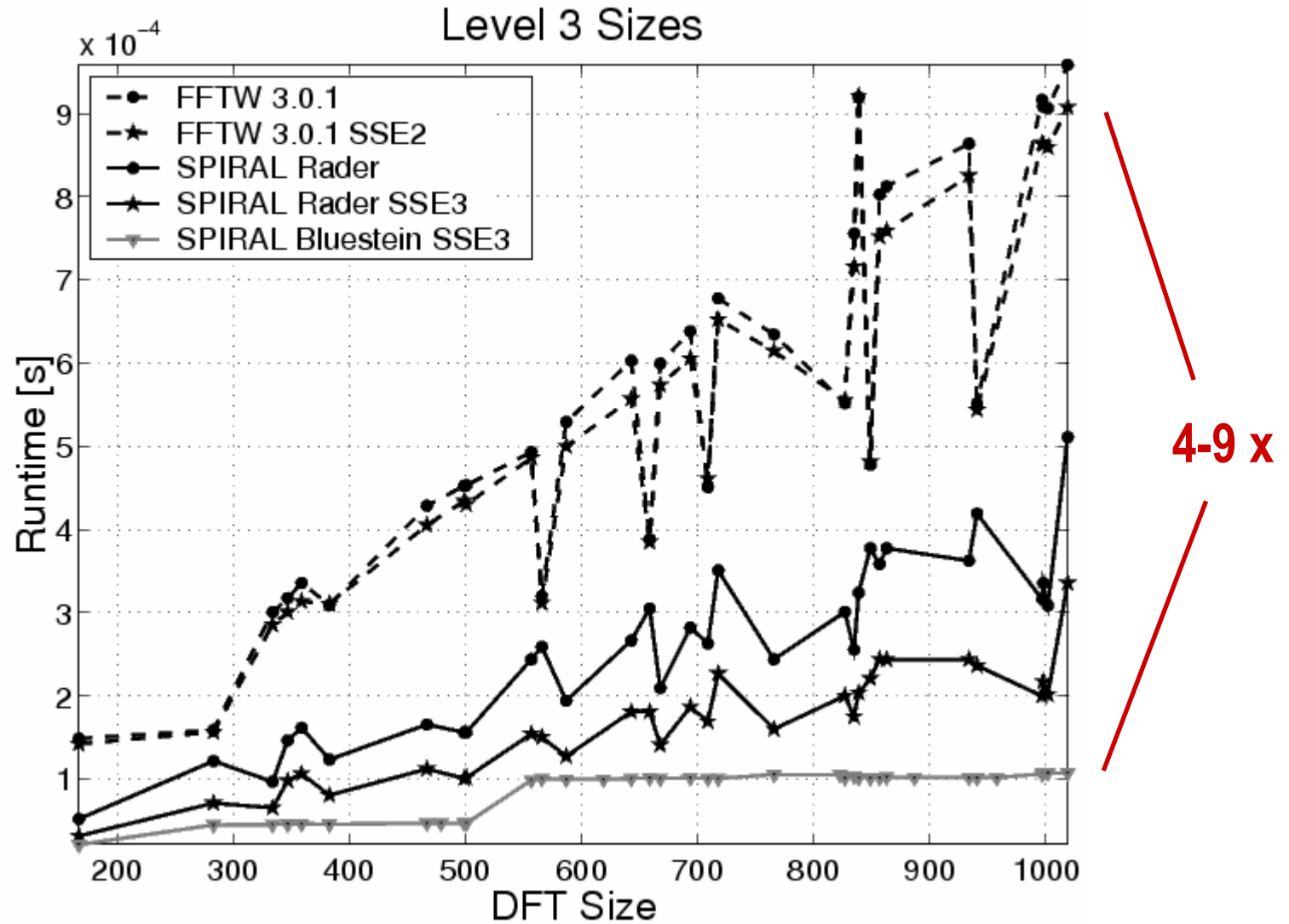
# Benchmark: DFT, Level 1 Sizes



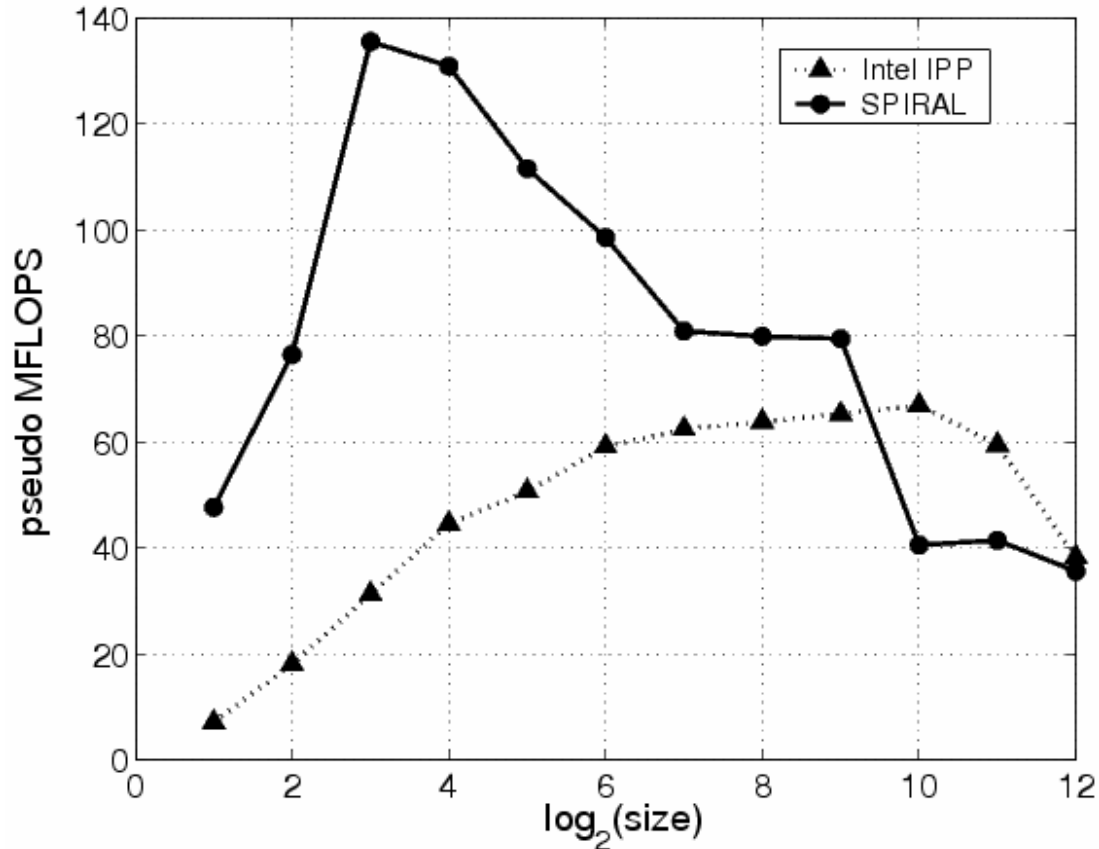
# Benchmark: DFT, Level 2 Sizes



# Benchmark: DFT, Level 3 Sizes



# Benchmark: Fixed Point DFT, IPAQ



**IPAQ**  
**Xscale arch.**  
**400 MHz**  
**Has only fixed point**

**Higher is better**

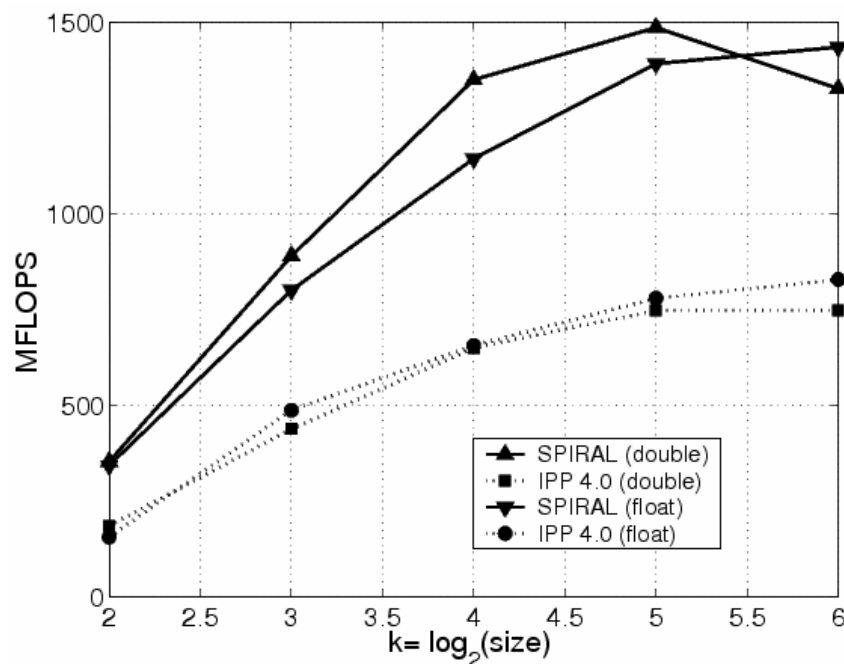
**Intel spent less effort?**



# Benchmark: DCT

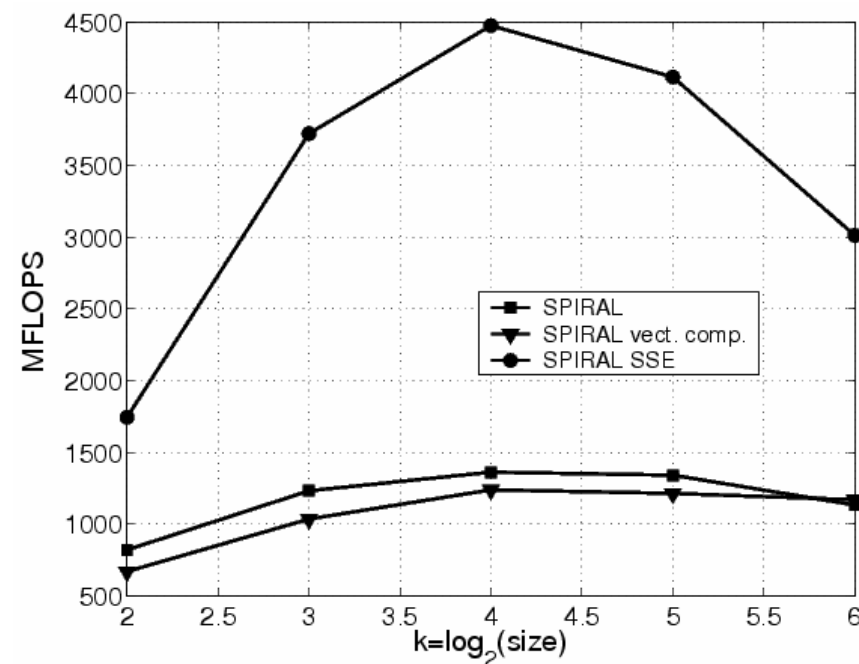
P4, 3.2 GHz,  
icc 8.0

## 1-D DCT



Scalar code

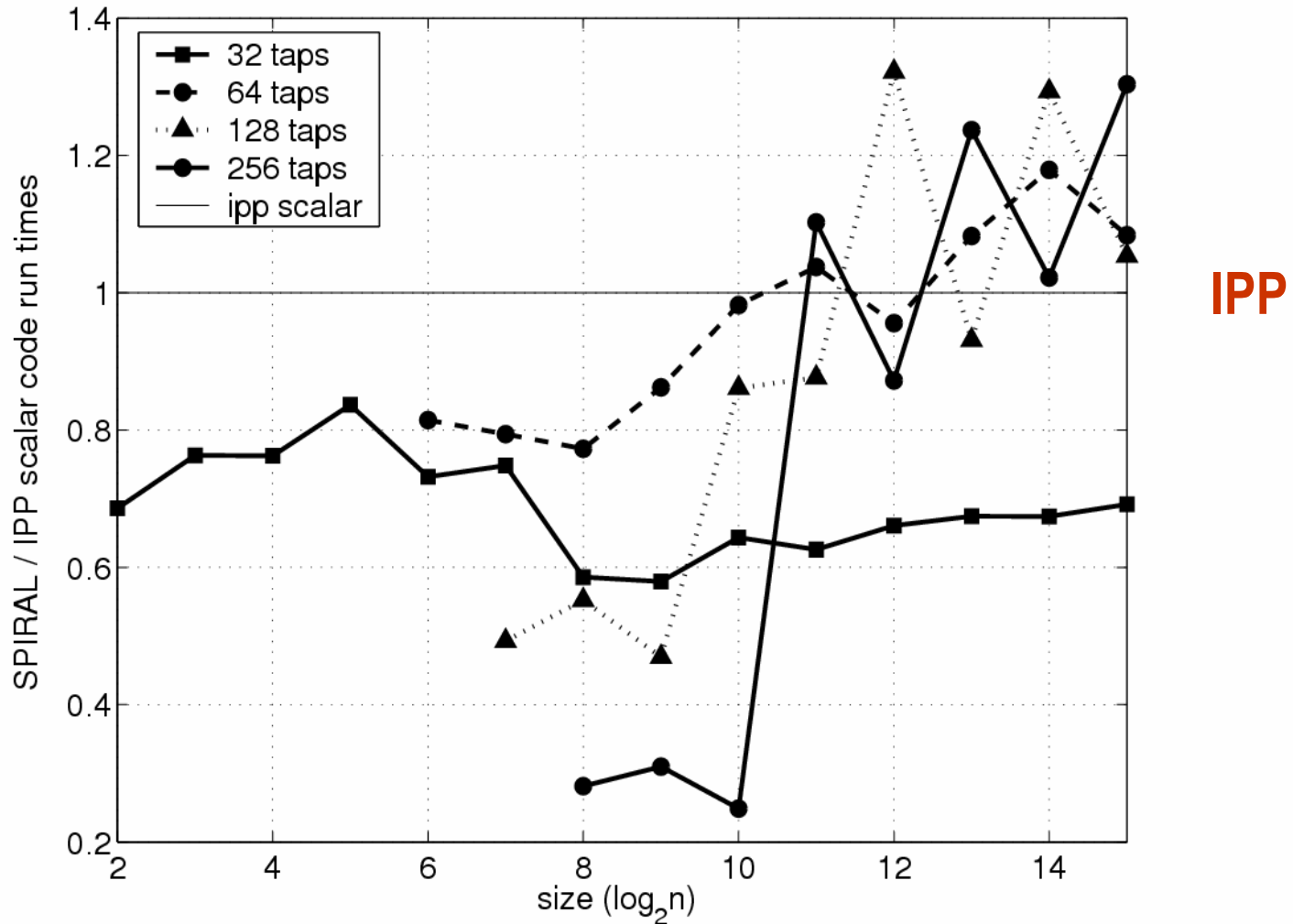
## 2-D DCT



Scalar vs. SSE code

- This is not the latest IPP
- Spiral gains a factor of 2 to vendor library
- Another factor of 3 with 2D and vector instructions

# Benchmark: Filter (Relative to IPP)



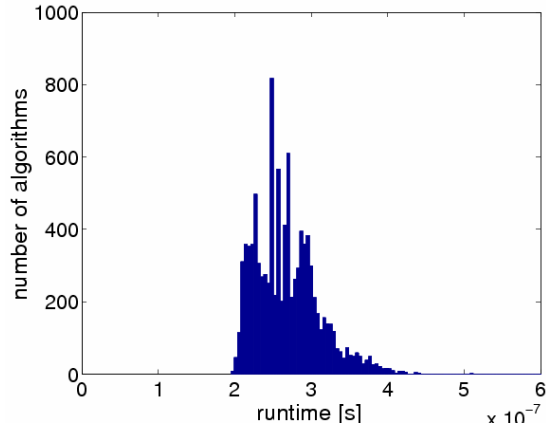
Lower = better

# Instructive Experiments

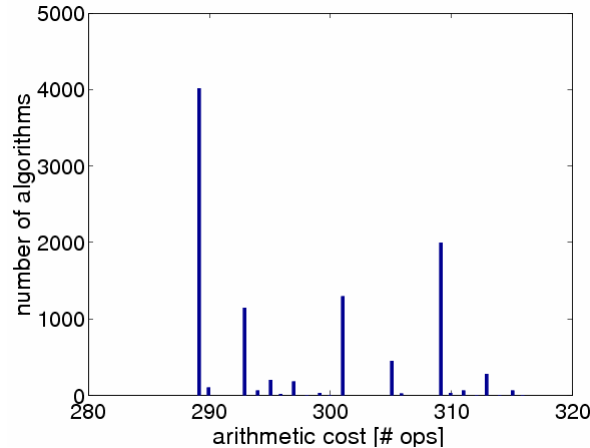
# Performance Spread: DCT, size 32

## Histograms, 10,000 algorithms

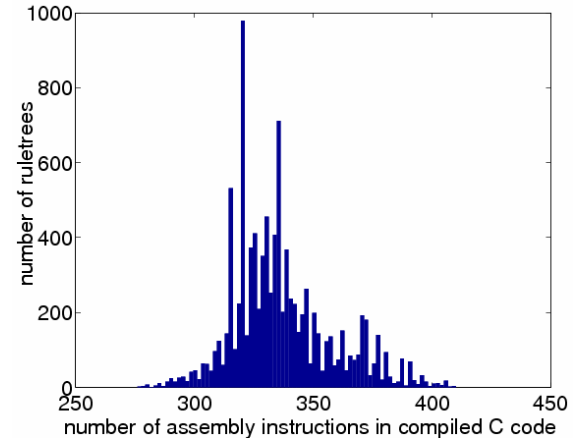
P4, 3.2 GHz,  
gcc



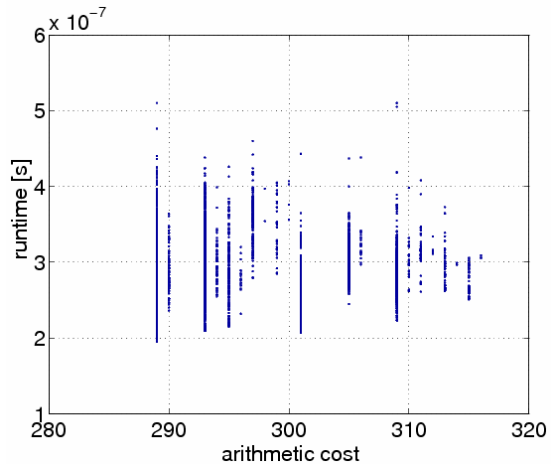
**runtime: x2**



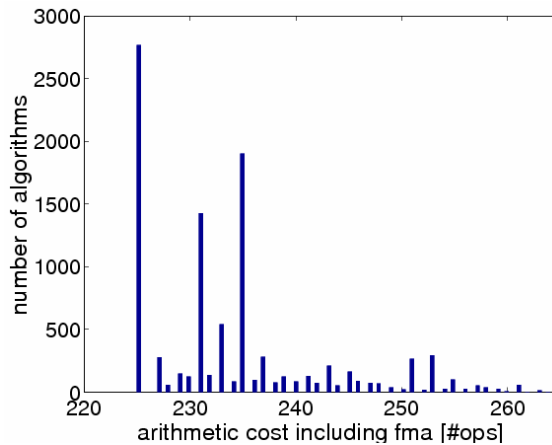
**#ops: x1.08**



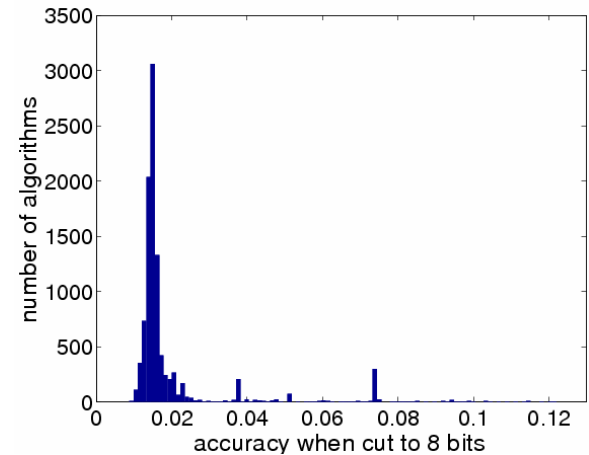
**#assembly instr: x1.5**



**#ops vs. runtime:  
no correlation**



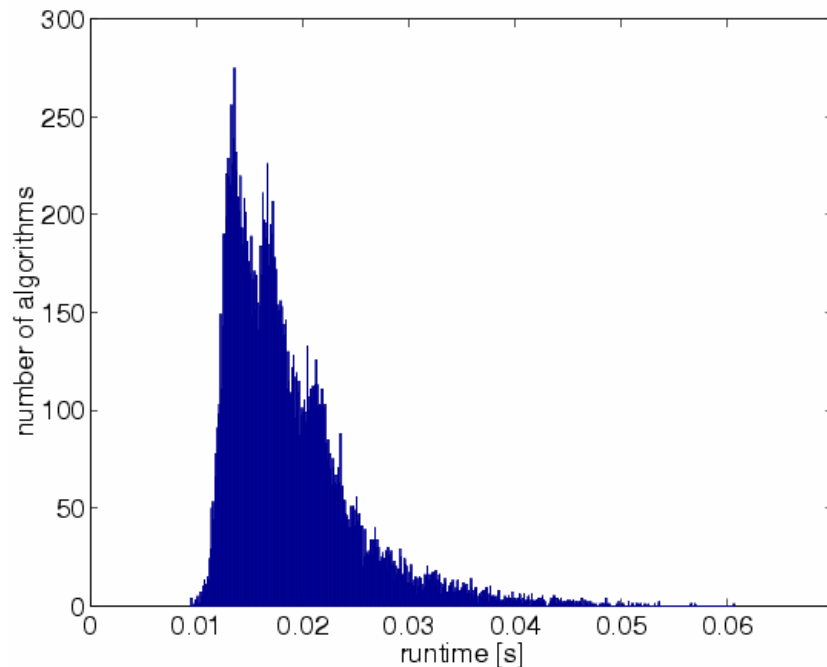
**#fma ops: x1.2**



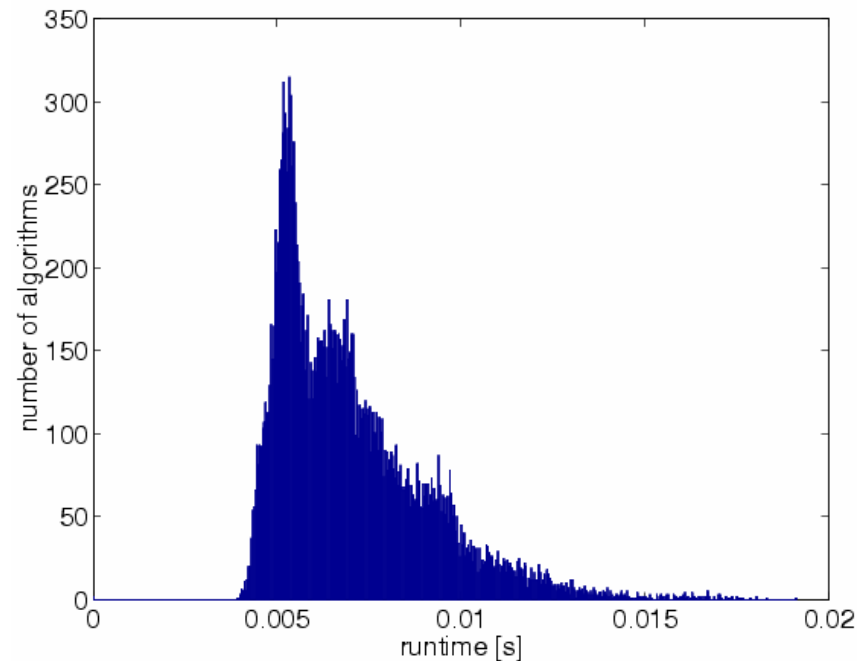
**accuracy: x10, most x2**

# Performance Spread: DFT 2<sup>16</sup> Histograms, 20,000 Algorithms

P4, 3.2 GHz,  
icc 8.0



**Generated scalar code**

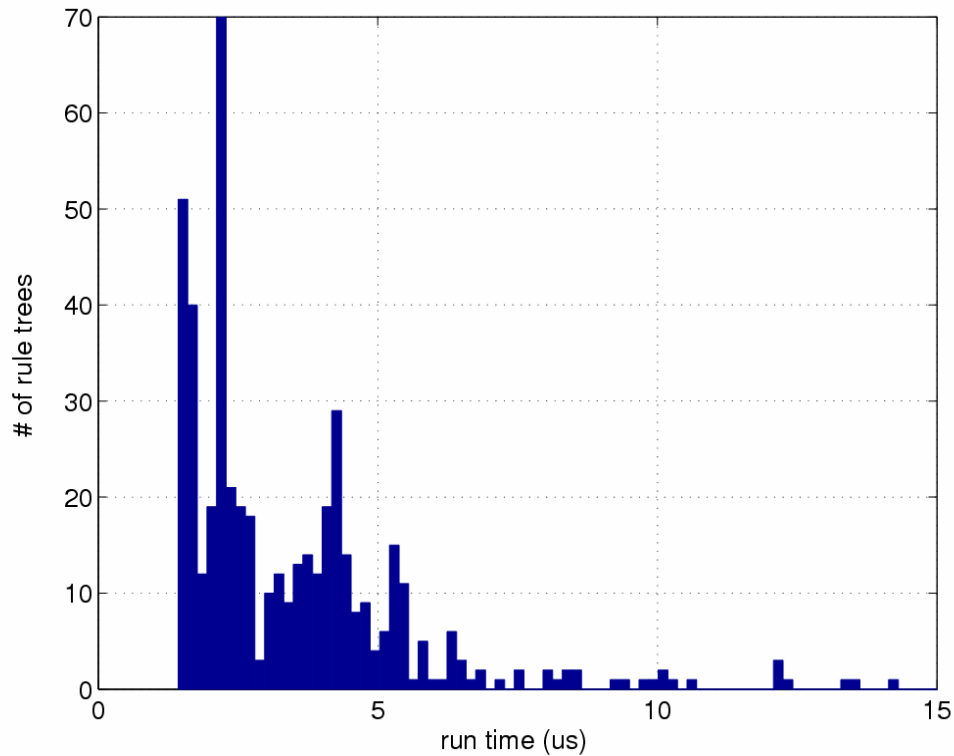


**Generated SSE code  
(4-way vector single precision)**

- **Generality of vectorization (all algorithms improve)**
- **Efficiency of vectorization (x 2.5 gain)**

# Performance Spread: Filter(128, 16)

Pentium 4 – 3.2



# Filter: Time Domain Methods

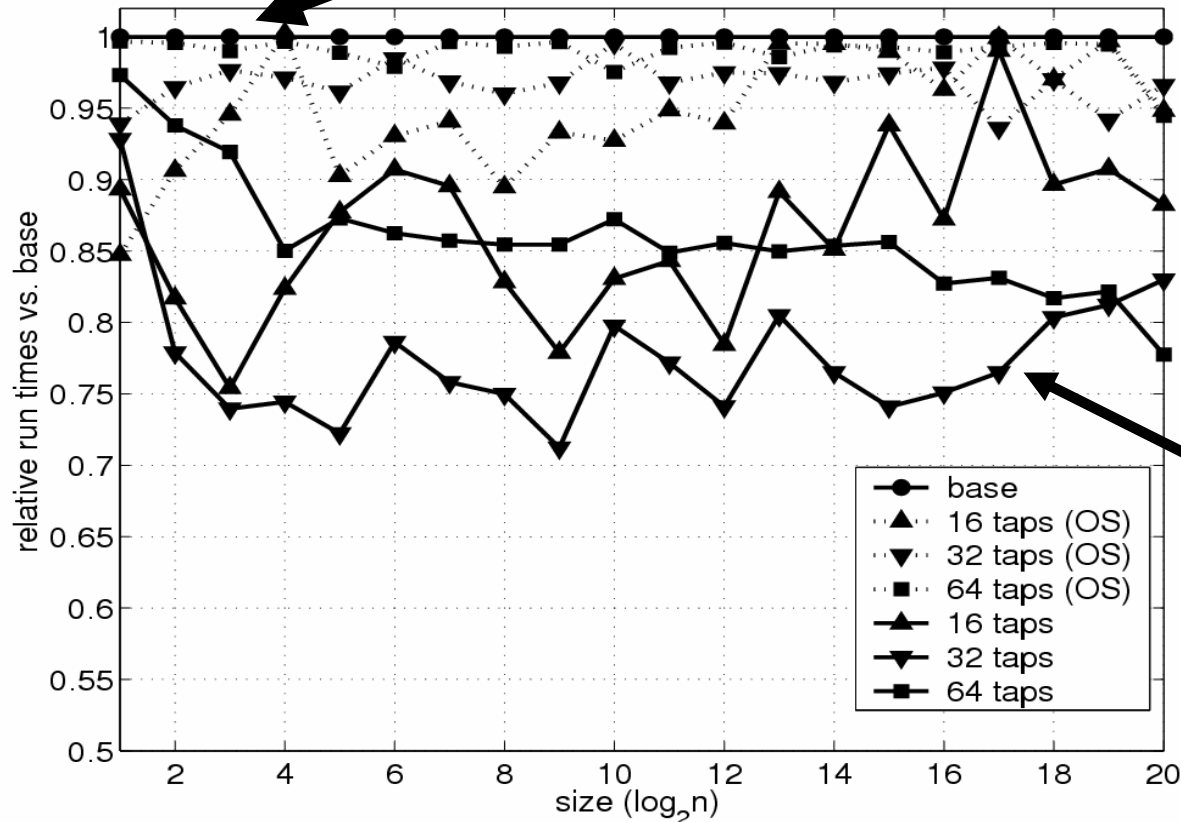
$$\text{Filt}_n(h(z)) \rightarrow I_n \otimes_{k-1} (h_0, \dots, h_{k-1})$$

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

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

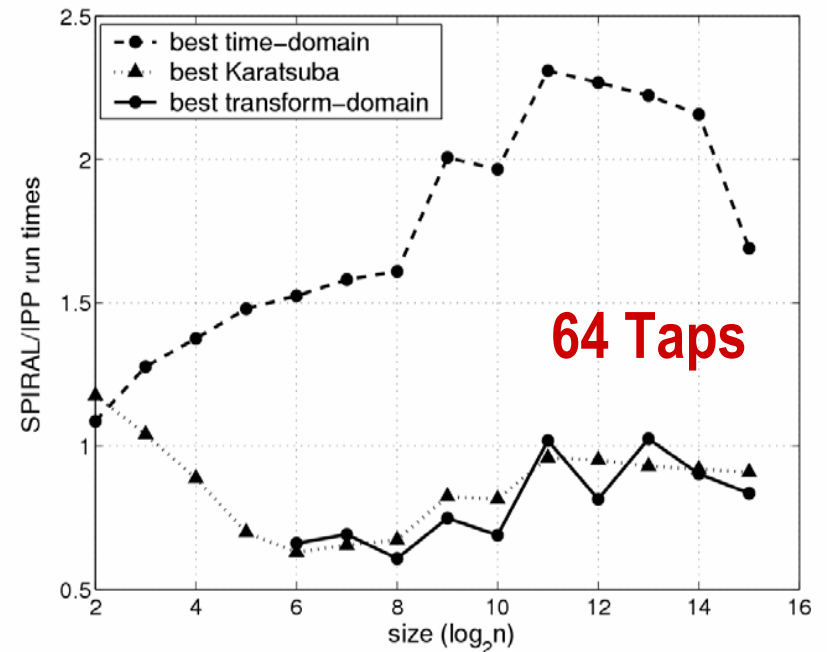
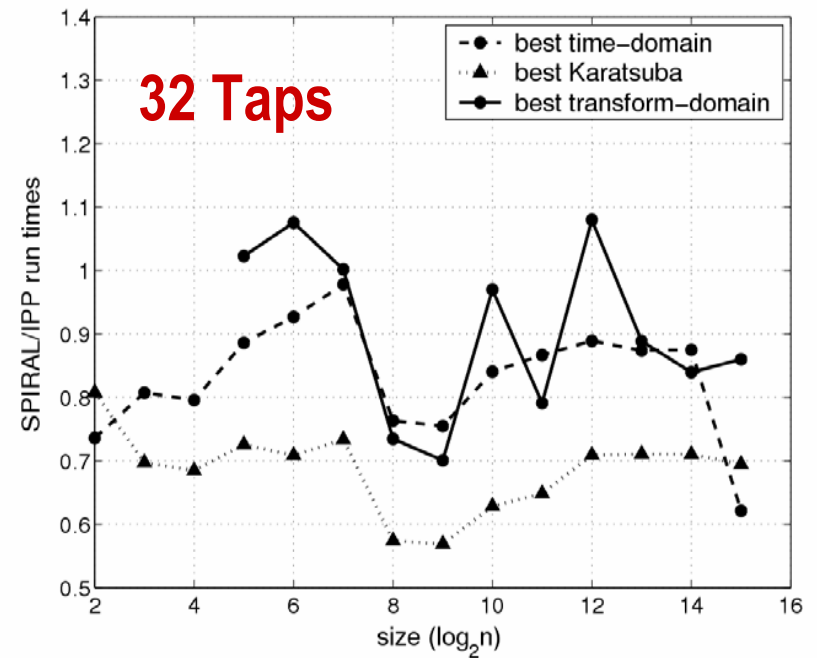
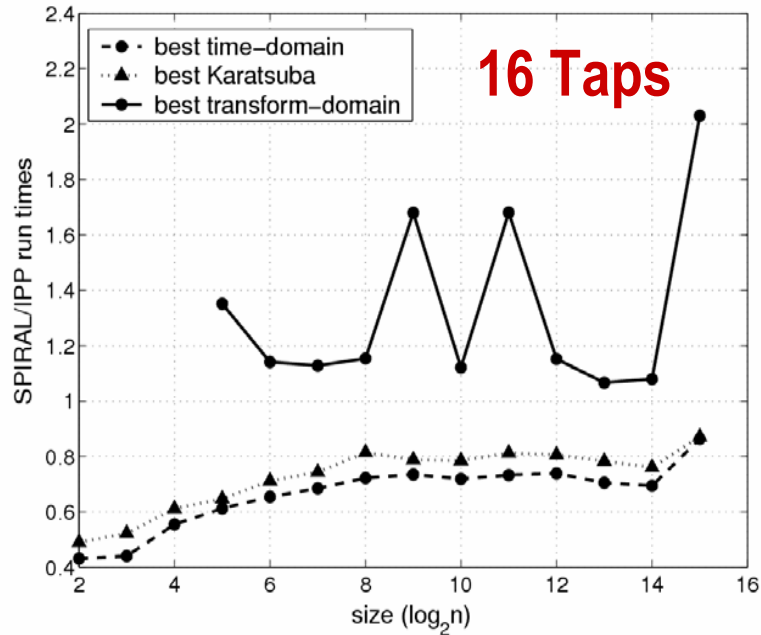
```
end
```

Xeon-1.7



best blocking strategy

# Filter: All Methods

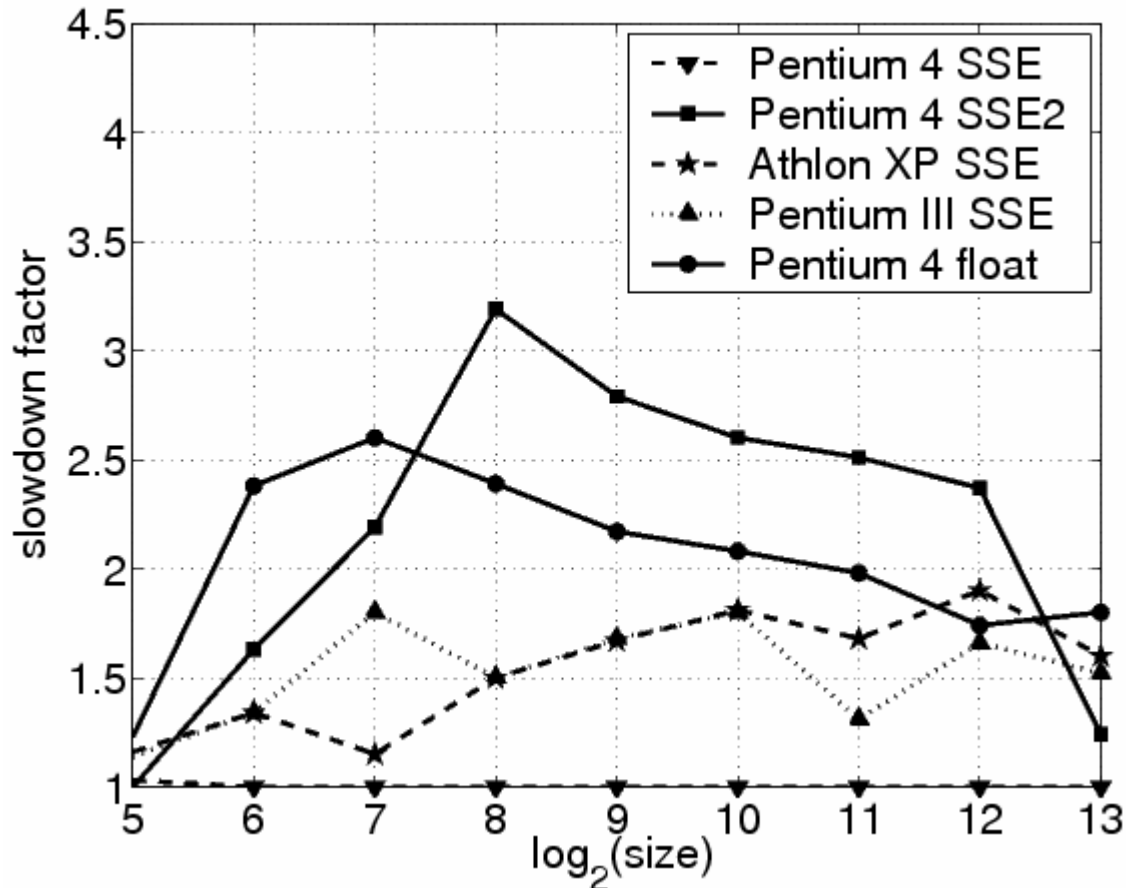


## Athlon XP 1.73

- 16: Time domain wins
- 32: Karatsuba wins
- 64: Karatsuba/DFT ~equal



# Platform Dependency: DFT



**50% Loss by porting from PIII to P4**

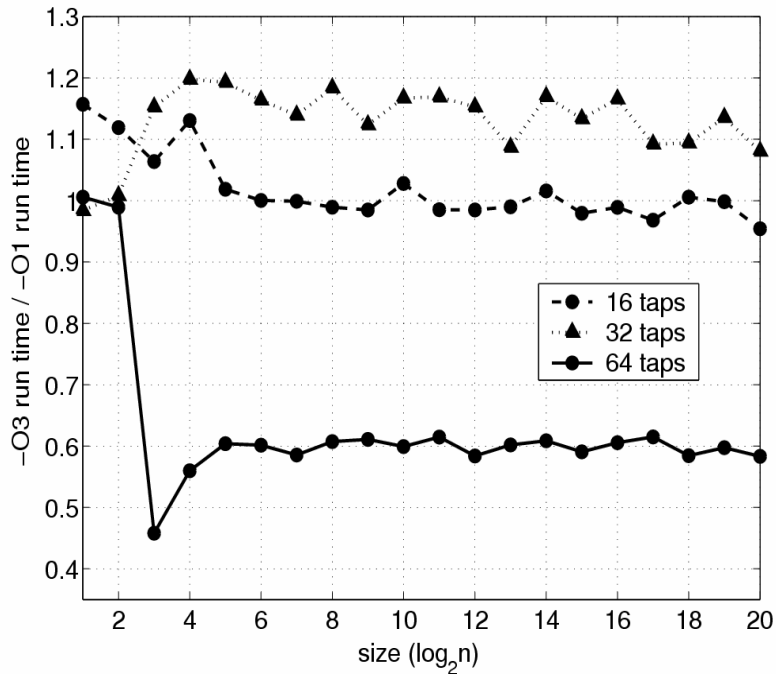
# Platform Dependency: Filter

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

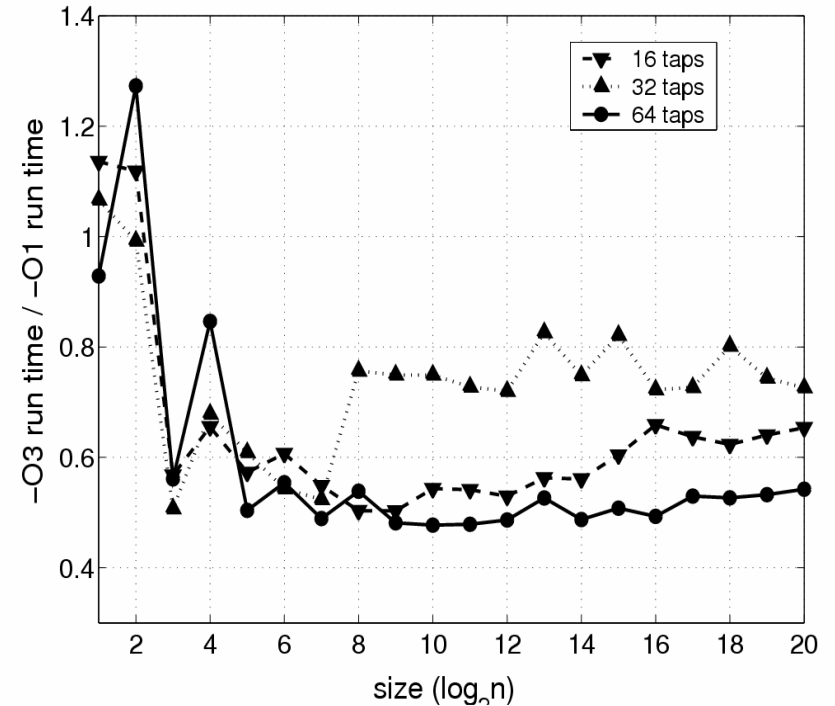
# Compileroptions: Filter

Macintosh - GNU C 3.3 (Apple)

## Blocking/nesting



## + Karatsuba

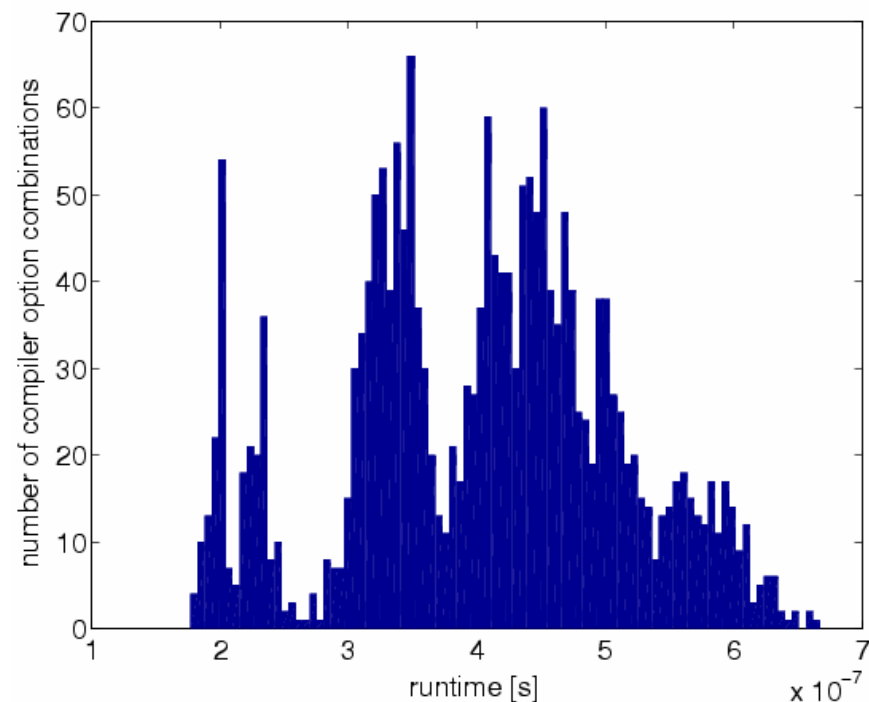
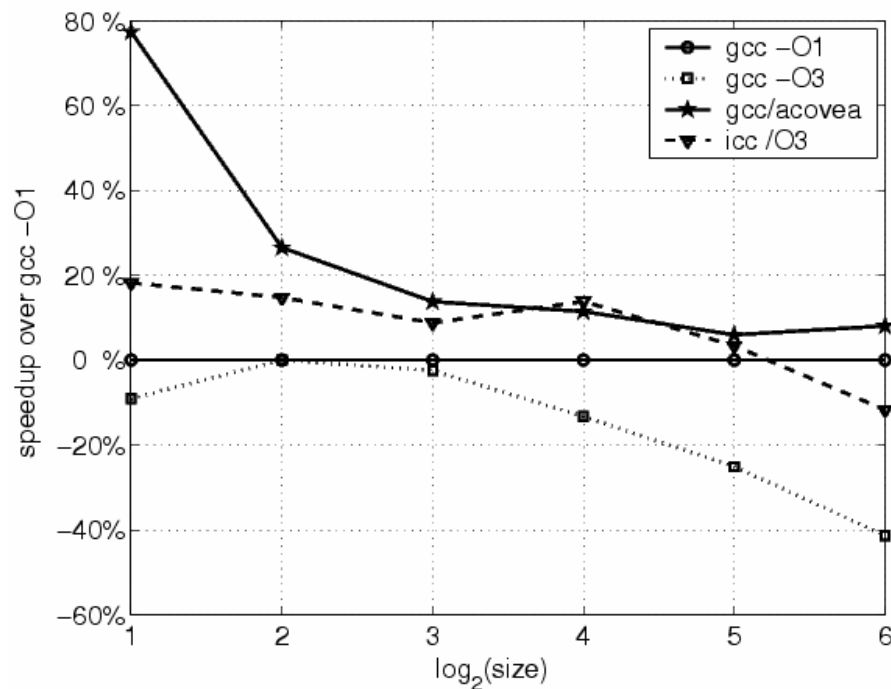


```
gcc {-O1/-O3} -fomit-frame-pointer -std=c99 -fast -mcpu=7450
```

# Compileroptions DCT

P4, 3.2 GHz,  
gcc

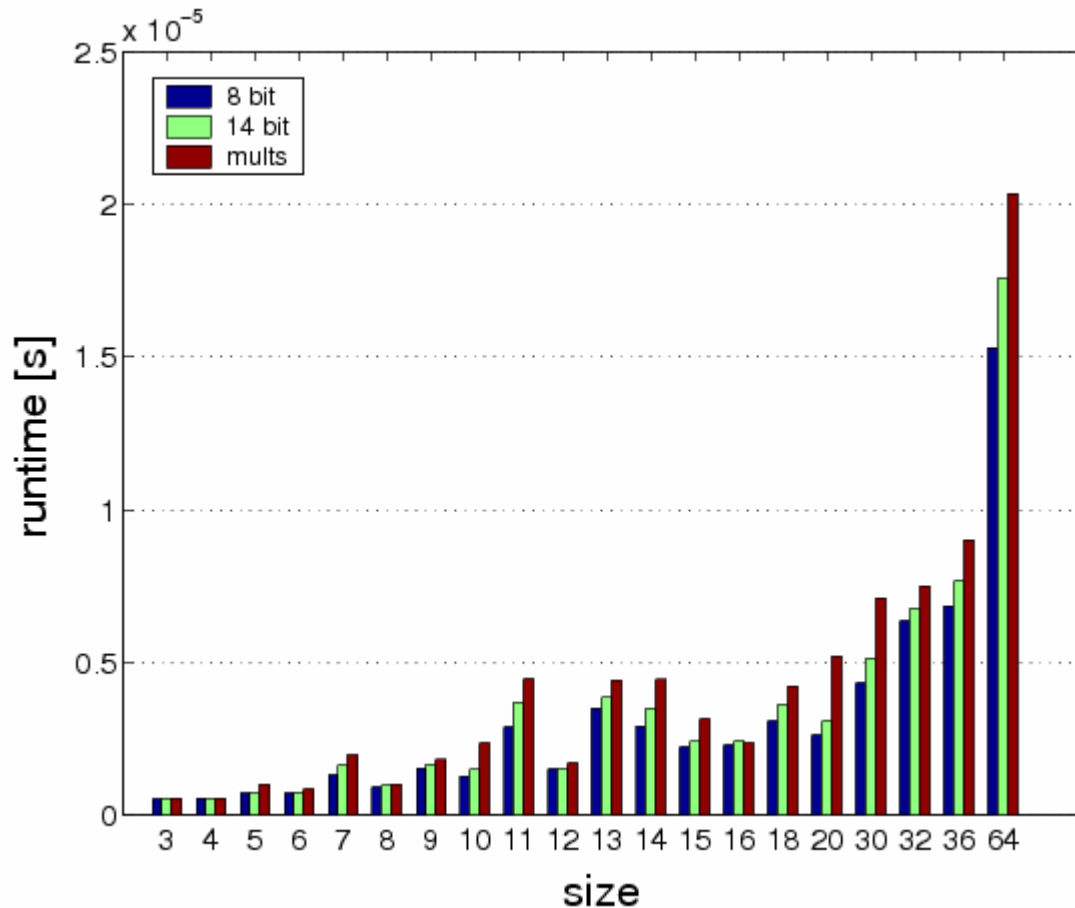
**ACOVEA: Evolutionary search for best compiler flags (gcc has ~500)**



**Runtime histogram**  
**Random compiler flags**  
**incl. -O1 -march=pentium4**

**10% improvement of best Spiral generated code**

# Multiplierless DFT, IPAQ



**IPAQ**  
**Xscale arch.**  
**400 MHz**  
**Fixed-point only**

- fixed-point constant multiplications replaced by adds and shifts
- trade-off runtime and precision

# Summary

- Code generation and tuning as optimization problem over the algorithm and implementation space

*Exploit the structure of the domain to solve it*

- Declarative framework for computer representation of the domain-knowledge

*Enables algorithm generation and optimization*

*(teaches the system the math; does not become obsolete?)*

- Compiler to translate math description into code
- Search and learning to explore implementation space

*Closes the feedback loop  
gives the system "intelligence"*

- Extensible, versatile

*Every step in the code generation is exposed*