

# **SPIRAL: Tuning DSP Transforms to Computing Platforms**

**José M. F. Moura**

**Carnegie Mellon**

## **Faculty**

- **Jeremy Johnson (Drexel)**
- **Robert Johnson (MathStar Inc.)**
- **David Padua (UIUC)**
- **Viktor Prasanna (USC)**
- **Markus Püschel (CMU)**
- **Manuela Veloso (CMU)**

## **Students**

- **Franz Franchetti (TU Vienna)**
- **Gavin Haentjens (CMU)**
- **Pinit Kumhom (Drexel)**
- **Neungsoo Park (USC)**
- **David Sepiashvili (CMU)**
- **Bryan Singer (CMU)**
- **Yevgen Voronenko (Drexel)**
- **Jianxin Xiong (UIUC)**

<http://www.ece.cmu.edu/~spiral>



# Sponsor

Carnegie Mellon



Work supported by DARPA (DSO), Applied & Computational Mathematics Program, OPAL, through grant managed by research grant DABT63-98-1-0004 administered by the Army Directorate of Contracting.

12/5/2002 IBM-Thomas T. J. Watson Res. Center



# Moore's Law and High(est) Performance Scientific Computing

(single processor, off-the-shelf)

- Moore's Law:**
- processor-memory bottleneck
  - short life cycles of computers
  - very complex architectures
    - vendor specific
    - special instructions (MMX, SSE, FMA, ...)
    - undocumented features

## Effects on software/algorithms:

- arithmetic cost model not accurate for predicting runtime (one cache miss = 10 floating point ops)
- better performance models hard to get
- best code is machine dependent (registers/caches size, structure)
- hand-tuned code becomes obsolete as fast as it is written
- compiler limitations
- full performance requires (in part) assembly coding



**Portable performance requires automation**



SPIRAL

# Automatic Performance Tuning: Research

Carnegie Mellon



## Linear Algebra:

- ATLAS (J. Dongarra et al.)
- LAPACK
- PhiPACK (J. Demmel et al.)

## Signal Processing:

- FFTW (M. Frigo and S. Johnson)
- SPIRAL

# SPIRAL

**Automates**

**Implementation**

- cuts development costs
- code less error-prone

**Optimization**

- systematic exploration of alternatives both at algorithmic and code level

**Platform-Adaptation**

- takes advantage of architecture specific features
- porting without loss of performance

**of DSP algorithms**

- are performance critical

**A library generator for highly optimized signal processing algorithms**

Carnegie Mellon

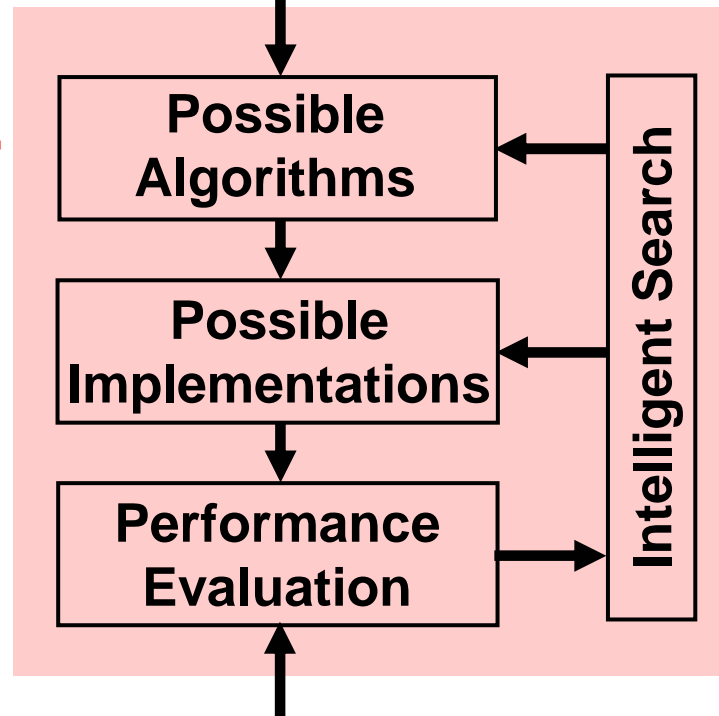


# SPIRAL Approach

given →

**DSP Transform**  
(DFT, DCT, Wavelets etc.)

SPIRAL Search Space



→ **adapted implementation**

given →

**Computing Platform**

(Pentium III, Pentium 4, Athlon,  
SUN, PowerPC, Alpha, ... )



# Organization

- **Mathematical Framework**

**Transforms, Rules, and Formulas**

- **Formula Generator**

**Transform → Algorithm**

- **SPL and SPL Compiler**

**Algorithm → Implementation**

- **Search Engine**

**How to find the best implementation**

- **SPIRAL system**

**Everything taken together**

- **Conclusions**

# 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)

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

- product of structured sparse matrices
- mathematical notation



# DSP Algorithms: Terminology

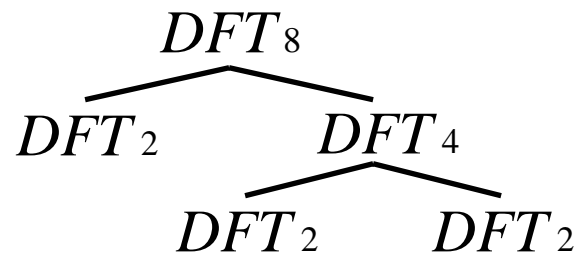


**Transform**  $DFT_n$  parameterized matrix

**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

# More Cooley-Tukey Rules

- DFT is symmetric  $\Rightarrow$  transpose the rule:

$$F_{RS} = L_S^{RS} (I_R \otimes F_S) T_S^{RS} (F_R \otimes I_S) \quad \text{CT rule transposed}$$

- Commuting tensor product factors

$$B \otimes A = L_n^{mn} (A \otimes B) L_m^{mn} \quad \text{A and B square size m and n}$$

- Commutation property  $\Rightarrow$  further variations

$$F_N = L_S^{RS} (I_S \otimes F_R) L_R^{RS} T_S^{RS} (I_R \otimes F_S) L_R^{RS}$$

$$F_N = (F_R \otimes I_S) T_S^{RS} L_S^{RS} (F_S \otimes I_R)$$

$$(F_2 \otimes I_2) x = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & -1 & 0 \\ 0 & 1 & 0 & -1 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad (I_2 \otimes F_2) = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & -1 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

- Different patterns for access, storage and flow of data

# DSP Transforms

discrete Fourier transform

$$DFT_n = [\exp(2kli\pi / n)]$$

Walsh-Hadamard transform

$$WHT_{2^k} = DFT_2 \otimes \dots \otimes DFT_2$$

discrete cosine and sine  
Transforms (16 types)

$$DCT^{(II)}_n = [\cos(k(l+1/2)\pi / n)]$$

$$DCT^{(IV)}_n = [\cos((k+1/2)(l+1/2)\pi / n)]$$

$$DST^{(I)}_n = [\sin(kl\pi / n)]$$

modified discrete  
cosine transform

$$MDCT_{n \times 2n} = [\cos((k+(n+1)/2)(l+1/2)\pi / n)]$$

two-dimensional transform

$$T \otimes T$$

discrete wavelet transform

$$DTWT_{2^n} = (DTWT_{2^{n-1}} \oplus I_{2^{n-1}}) \underbrace{L_{2^{n-1}}^{2^n} (I_{2^{n-1}} \otimes_{l-2} W)}_H$$

Others: filtering, Haar, Hartley, ...



# Rules = Breakdown Strategies

$$DCT_2^{(II)} \rightarrow \text{diag} \left( 1, 1 / \sqrt{2} \right) \cdot F_2$$

$$DCT_n^{(II)} \rightarrow P \cdot \left( DCT_{n/2}^{(II)} \oplus DCT_{n/2}^{(IV)} \right) \cdot \left( I_{n/2} \otimes F_2 \right)^Q$$

$$DCT_n^{(IV)} \rightarrow S \cdot DCT_n^{(II)} \cdot D$$

$$DCT_n^{(IV)} \rightarrow M_1 \cdots M_r$$

$$DFT_n \rightarrow \text{CosDFT}_n + j \cdot \text{SinDFT}_n$$

$$DFT_n \rightarrow B \cdot \left( DCT_{n/2}^{(I)} \oplus DST_{n/2}^{(I)} \right) \cdot C$$

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

$$\text{CosDFT}_n \rightarrow \cdots \text{CosDFT}_{n/2} \cdots DCT_{n/4}^{(II)} \cdots$$

$$\text{SinDFT}_n \rightarrow \cdots \text{SinDFT}_{n/2} \cdots DCT_{n/4}^{(II)} \cdots$$

$$WHT_{2^n} \rightarrow \prod_{i=1}^n \left( I_{2^{n_1+\dots+n_{i-1}}} \otimes WHT_{2^{n_i}} \otimes I_{2^{n_{i+1}+\dots+n_t}} \right)$$

$$MDCT_{n \times 2n} \rightarrow S \cdot DCT_n^{(IV)} \cdot P$$

$$DTWT_{2^n} = \left( DTWT_{2^{n-1}} \oplus I_{2^{n-1}} \right) \underbrace{L_{2^{n-1}}^{2^n} \left( I_{2^{n-1}} \otimes_{l-2} W \right)}_H$$

base case

recursive

translation

iterative

recursive

recursive

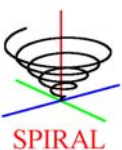
recursive

recursive

recursive

iterative/  
recursive

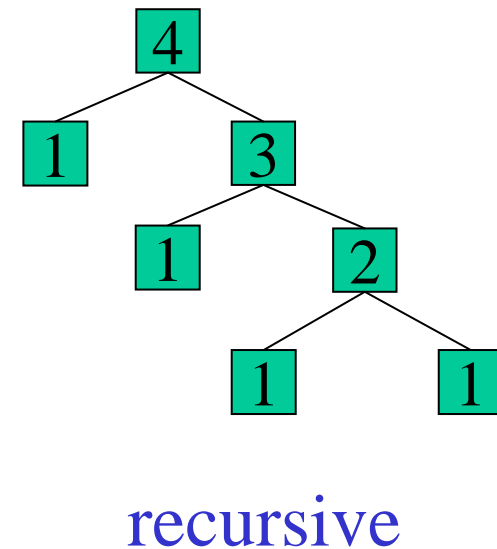
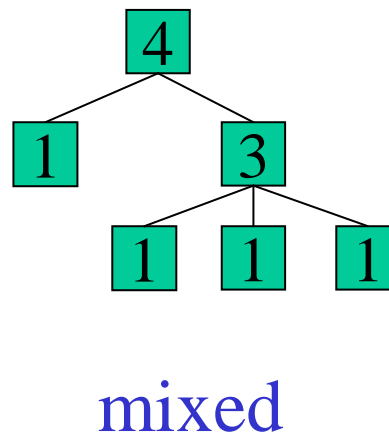
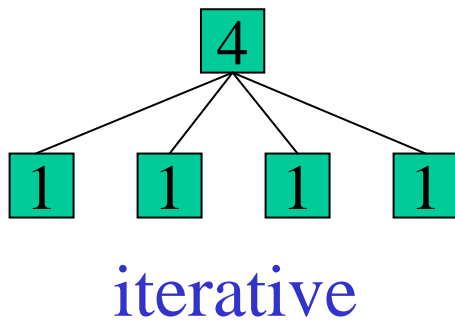
translation



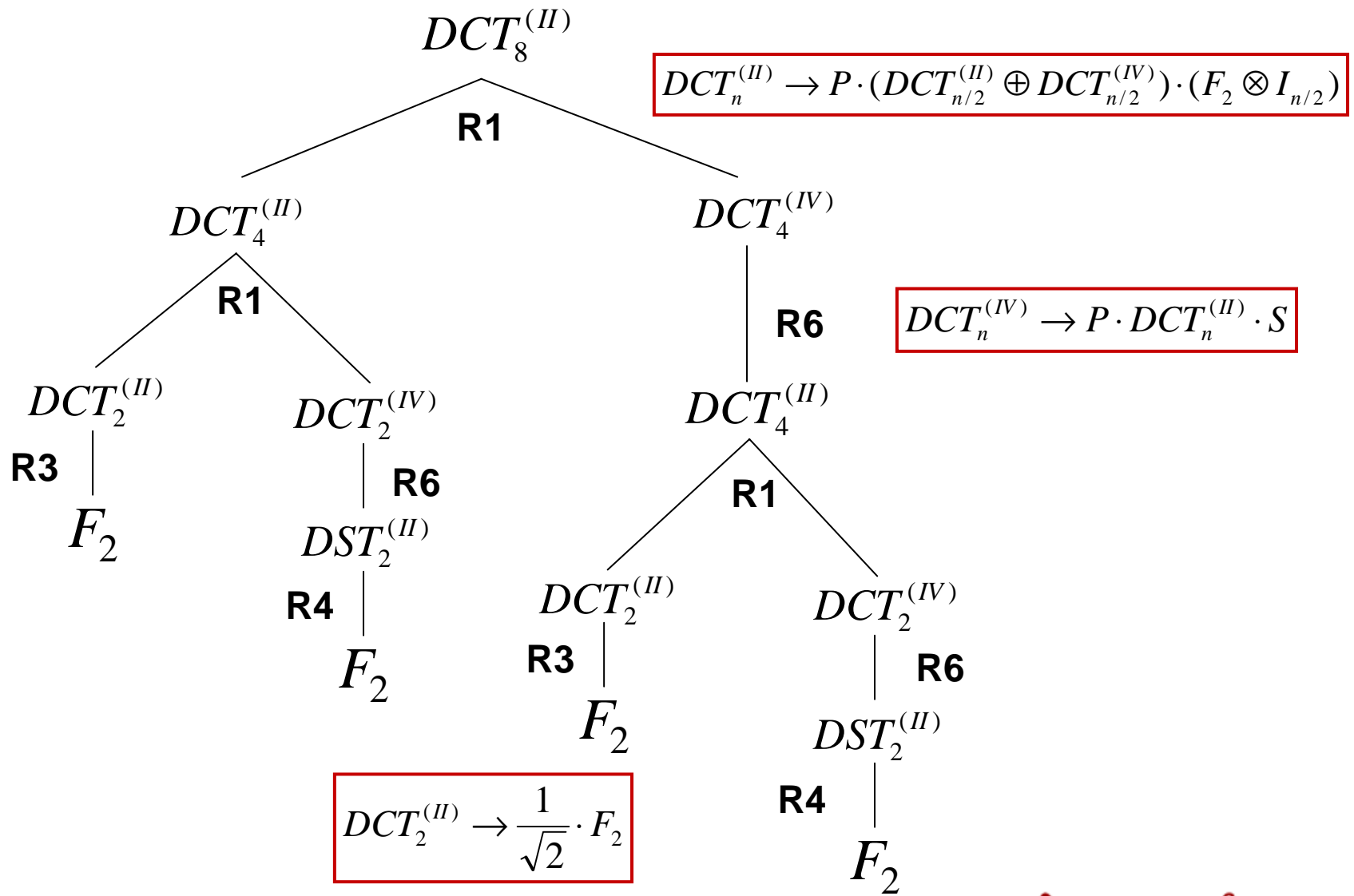
SPIRAL

# Partition Trees

Each WHT algorithm can be represented by a tree, where a node labeled by  $n$  corresponds to  $WHT_{2^n}$



# Algorithms = Rulertrees = Formulas



# Mathematical Framework for Filters and Discrete-Time Wavelet Transforms

- Overlapped direct sum of matrices  $A$  ( $m \times n$ ) and  $B$  ( $r \times s$ )

column overlap  $A \oplus_v B = \left[ \begin{array}{c} \boxed{A} \\ \underbrace{\boxed{B}}_v \end{array} \right] = (A \oplus B)(I_n \oplus_v I_s)$

row overlap  $A \oplus^v B = \left[ \begin{array}{c} \boxed{A} \quad \boxed{B} \\ \underbrace{\quad}_v \end{array} \right] = (I_m \oplus^v I_r)(A \oplus B)$

- Overlapped tensor product

column overlap  $I_k \otimes_v B = \underbrace{B \oplus_v B \oplus_v \dots \oplus_v B}_{k\text{-fold}}$

row overlap  $I_k \otimes^v B = \underbrace{B \oplus^v B \oplus^v \dots \oplus^v B}_{k\text{-fold}}$

# Properties of Overlapped Direct Sums and Tensor Products

- Column and row overlap direct sum

$$A \oplus^v B = \left[ \begin{array}{|c|} \hline A \\ \hline \underbrace{\quad}_v B \\ \hline \end{array} \right] = \left[ \begin{array}{|c|} \hline A^T \\ \hline \underbrace{\quad}_v B^T \\ \hline \end{array} \right]^T = (A^T \oplus_v B^T)^T$$

$$\mathbf{I}_n \oplus^v \mathbf{I}_s = (\mathbf{I}_n \oplus_v \mathbf{I}_s)^T$$

- Column and row overlapped tensor product

$$\mathbf{I}_k \otimes_v \mathbf{B}_{m \times n} = \left( \bigoplus_{l=1}^k \mathbf{B}_{m \times n} \right) \cdot \left( \bigoplus_{l=1}^k \mathbf{I}_n \right) = (\mathbf{I}_k \otimes \mathbf{B}_{m \times n}) \cdot (\mathbf{I}_k \otimes_v \mathbf{I}_n)$$

$$\begin{aligned} \mathbf{I}_k \otimes^v \mathbf{B}_{m \times n} &= \left( \bigoplus_{l=1}^k \mathbf{I}_m \right)^T \cdot \left( \bigoplus_{l=1}^k (\mathbf{B}_{m \times n})^T \right)^T = (\mathbf{I}_k \otimes_v \mathbf{I}_m)^T \cdot \bigoplus_{l=1}^k \mathbf{B}_{m \times n} \\ &= (\mathbf{I}_k \otimes^v \mathbf{I}_m) \cdot (\mathbf{I}_k \otimes \mathbf{B}_{m \times n}) \end{aligned}$$





# Filters

- Linear convolution in matrix form

$$C_n(\underline{h}) = \mathbf{I}_n \otimes^{l-1} [h_0 \quad h_1 \quad \cdots \quad h_{l-1}]^T$$

- Block convolutions

- Overlap-save rule

$$C_n(\underline{h}) = (\mathbf{I}_{n/b} \otimes^{l-1} \mathbf{I}_{b+l-1}) (\mathbf{I}_{n/b} \otimes C_b(\underline{h}))$$

- Overlap-add rule

$$C_n(\underline{h}) = (\mathbf{I}_{m/b} \otimes C_b^T(\underline{h})) (\mathbf{I}_{m/b} \otimes_{l-1} \mathbf{I}_{b+l-1}) E_{l-1, l-1} \quad m = n + l - 1$$

- Circular convolution rule for circulant  $K$

$$K(\underline{h}) = \text{DFT}_n^{-1} \cdot \text{diag}(\text{DFT}_n \cdot \underline{h}) \cdot \text{DFT}_n \quad \underline{h}\text{-first column of } K$$

$n$ -length of  $\underline{h}$



# Discrete-Time Wavelet Transform

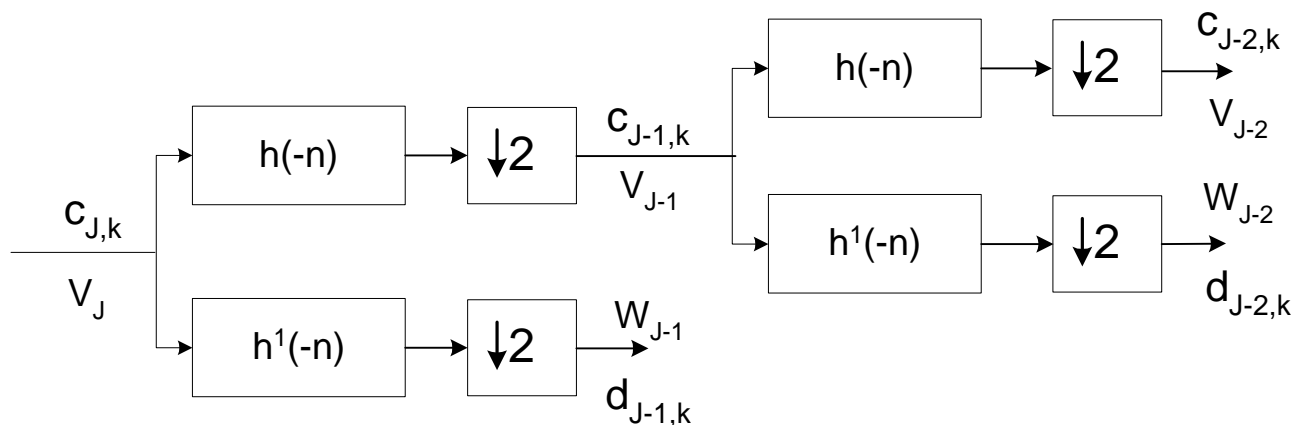
- DTWT scaling and wavelet expansion coefficients

$$c_{0,k} = \sum_n x_n h_J(2^J k - n) \quad \text{scaling coefficients}$$

$$d_{j,k} = \sum_n x_n h_{J-j}^1(2^{J-j} k - n) \quad j = 0, \dots, J-1 \quad \text{wavelet coefficients}$$

- Recursive algorithm (Mallat)

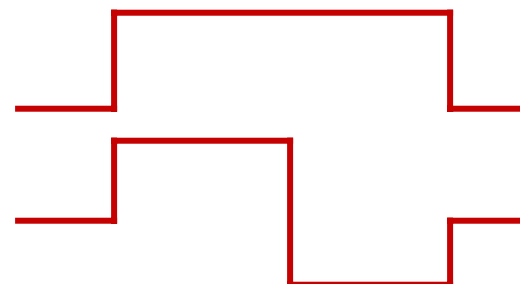
$$c_{j,k} = \sum_m h_{m-2k} c_{j+1,m} \quad d_{j,k} = \sum_m h_{m-2k}^1 c_{j+1,m}$$



# Haar Wavelets – Example

- Haar wavelets = square waves

$$h = \frac{1}{\sqrt{2}} [1, 1], \quad h^1 = \frac{1}{\sqrt{2}} [1, -1]$$



- First stage:  $V_2 \Rightarrow V_1 \oplus W_1$

$$H = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & -1 \end{bmatrix} = L_2^4 \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & -1 \end{bmatrix} = L_2^4 (I_2 \otimes F_2) \quad \begin{bmatrix} \underline{c}_1 \\ \underline{d}_1 \end{bmatrix} = H \underline{c}_2$$

- The process is repeated for the upper half of the output

$$HT_{2^n} = \left( HT_{2^{n-1}} \oplus I_{2^{n-1}} \right) \underbrace{L_{2^{n-1}}^2 (I_{2^{n-1}} \otimes F_2)}_H, \quad HT_2 = F_2$$

# Discrete-Time Wavelet Transform

- Discrete-Time Wavelet Transform (DTWT) rule

$$\text{DTWT}_{2^n} = \left( \text{DTWT}_{2^{n-1}} \oplus \text{I}_{2^{n-1}} \right) \underbrace{\text{L}_{2^{n-1}} \left( \text{I}_{2^{n-1}} \otimes_{l-2} W \right)}_H$$

- Scaling (lowpass) and wavelet (highpass) filter coefficients

$$W = \begin{bmatrix} h_0 & h_1 & \cdots & h_{l-1} \\ h'_0 & h'_1 & \cdots & h'_{l-1} \end{bmatrix}$$

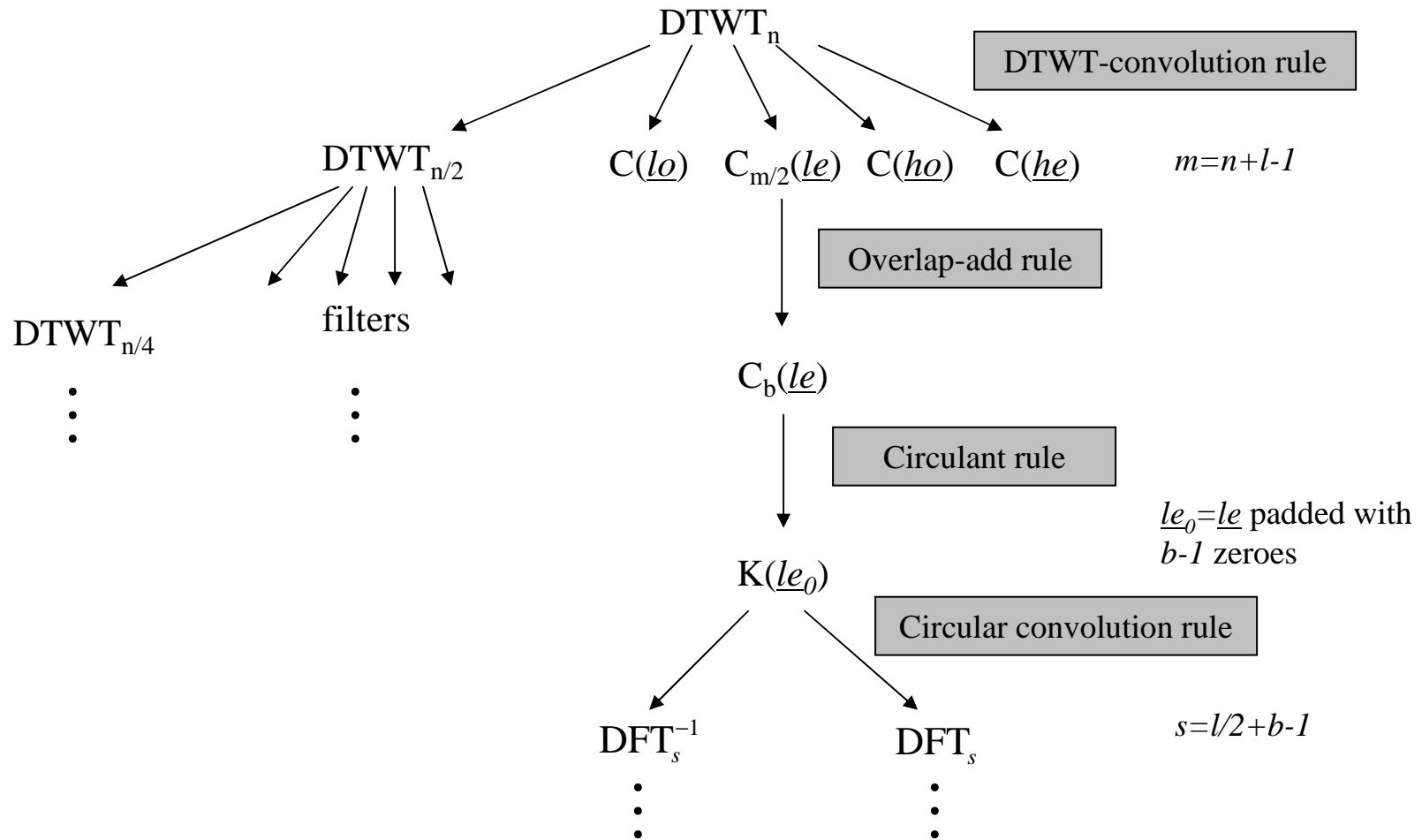
- DTWT - convolution rule

$$H = \left( \left( \begin{bmatrix} 1 & 1 \end{bmatrix} \otimes \text{I}_{m/2} \right) \cdot \left( C_{m/2}^T(\underline{lo}) \oplus C_{m/2}^T(\underline{le}) \right) \cdot \text{L}_2^n \oplus \right. \\ \left. \left( \begin{bmatrix} 1 & 1 \end{bmatrix} \otimes \text{I}_{m/2} \right) \cdot \left( C_{m/2}^T(\underline{ho}) \oplus C_{m/2}^T(\underline{he}) \right) \cdot \text{L}_2^n \right) \cdot \left( \begin{bmatrix} 1 \\ 1 \end{bmatrix} \otimes \text{I}_n \right)$$

$\underline{lo}$ -lowpass odd coeffs.,  $\underline{le}$ -lowpass even coeffs.

$\underline{ho}$ -highpass odd coeffs.,  $\underline{he}$ -highpass even coeffs.

# DTWT Ruletree – Example



# Formula for a DCT, size 16

$$\begin{aligned}
 & [(2, 16, 9, 5, 3)(4, 15, 8, 13, 7)(6, 14, 10, 12, 11), 16] \cdot \\
 & ((([(2, 8, 5, 3)(4, 7), 8] \cdot ((([(2, 4, 3), 4] \cdot (\text{diag}(1, \sqrt{\frac{1}{2}}) \cdot \text{DFT}_2) \oplus ([[(1, 2), 2] \cdot \text{R}_{\frac{13}{8}\pi})^{[(1,2),2]}] \cdot \\
 & (\mathbf{1}_2 \otimes \text{DFT}_2)^{[(2,4,3),4]}) \oplus (\text{diag}(\frac{1}{2\cos(\frac{1}{16}\pi)}, \frac{1}{2\cos(\frac{3}{16}\pi)}, \frac{1}{2\cos(\frac{5}{16}\pi)}, \frac{1}{2\cos(\frac{7}{16}\pi)}) \cdot (\mathbf{1}_2 \otimes \text{DFT}_2)^{[(2,4,3),4]}. \\
 & ((\text{DFT}_2 \cdot \text{diag}(1, \sqrt{\frac{1}{2}})) \oplus ([[(1, 2), 2] \cdot \text{R}_{\frac{13}{8}\pi})^{[(1,2),2]}] \cdot [(2, 3, 4), 4] \cdot \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 1 \end{bmatrix})^{[(1,4)(2,3),4]}). \\
 & (\mathbf{1}_4 \otimes \text{DFT}_2)^{[(2,8,5,3)(4,7),8]} \oplus ([[(2, 5, 4, 3, 7, 6, 8), 8] \cdot (\mathbf{1}_2 \otimes (\mathbf{1}_2 \oplus ([[(1, 2), 2] \cdot \text{R}_{\frac{7}{4}\pi})^{[(1,2),2]}])) \cdot \\
 & (\mathbf{1}_2 \otimes \text{DFT}_2 \otimes \mathbf{1}_2) \cdot (\mathbf{1}_4 \oplus ([[(1, 2), 2] \cdot \text{R}_{\frac{13}{8}\pi})^{[(1,2),2]}] \oplus ([[(1, 2), 2] \cdot \text{R}_{\frac{1}{8}\pi})^{[(1,2),2]}]) \cdot (\mathbf{1}_1 \otimes \text{DFT}_2 \otimes \\
 & \mathbf{1}_4) \cdot ([[(1, 2), 2] \cdot \text{R}_{\frac{49}{32}\pi})^{[(1,2),2]} \oplus ([[(1, 2), 2] \cdot \text{R}_{\frac{53}{32}\pi})^{[(1,2),2]}] \oplus ([[(1, 2), 2] \cdot \text{R}_{\frac{57}{32}\pi})^{[(1,2),2]}] \oplus ([[(1, 2), 2] \cdot \\
 & \text{R}_{\frac{61}{32}\pi})^{[(1,2),2]}]) \cdot [(2, 8)(4, 6), 8]^{[(1,8)(2,7)(3,6)(4,5),8]}). \\
 & (\mathbf{1}_8 \otimes \text{DFT}_2)^{[(2,16,9,5,3)(4,15,8,13,7)(6,14,10,12,11),16]}
 \end{aligned}$$



# Helpful Concept

## DSP Transforms

## Formal Languages

DSP transform (of size)	↔	Non-terminal symbol (with attribute)
Rule	↔	Rule (production)
Formula/Algorithm	↔	Element in Language (only terminals)

# Mathematical Framework: Summary

- fast algorithms represented as **ruletrees** (easy generation/manipulation) and as **formulas** (can be translated into code)
- formulas built from **few** constructs and primitives
- **many** different algorithms/formulas generated from **few** rules (combinatorial explosion)
- these algorithms are (essentially) **equal in arithmetic cost**, but **differ in data flow**





# Organization

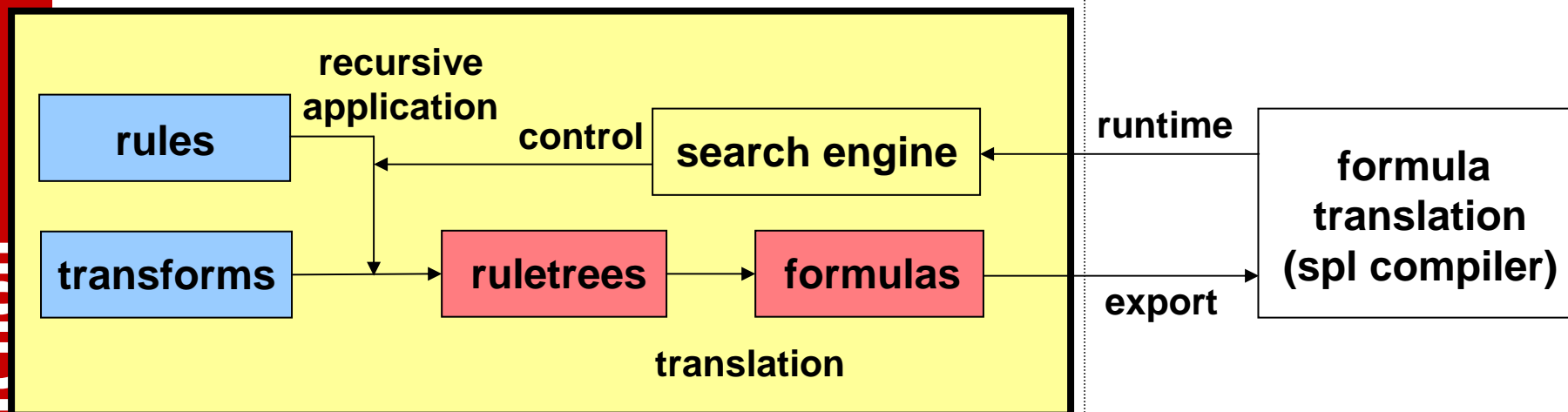


- **Mathematical Framework**  
Transforms, Rules, and Formulas
- **Formula Generator**  
Transform  $\rightarrow$  Algorithm
- **SPL and SPL Compiler**  
Algorithm  $\rightarrow$  Implementation
- **Search Engine**  
How to find the best implementation
- **SPIRAL system**  
Everything taken together
- **Conclusions**

# Formula Generation

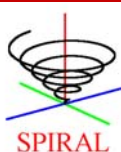
 data base (**extensible!**)  
 data type

## Formula Generator



- written in GAP/AREP (computer algebra system)
- all computation/manipulation is **symbolic**
- **exact** arithmetic
- **easy extensible** rule and transform data base
- **verification** of rules and formulas

**cut here for other optimization problems**



SPIRAL

# Formula Generator: Summary

## Concept:

- easy extensible
- formulas as ruletrees (efficient representation, easy manipulation)
- interface to spl
- allows implementation of search methods

## Implementation:

- efficient implementation moves bottleneck to spl compiler
- fully symbolic derivation/manipulation of formulas
- verification
- tools to analyze structure, arithmetic cost of formulas
- interfaces with search engine



# Organization

- **Mathematical Framework**  
Transforms, Rules, and Formulas
- **Formula Generator**  
Transform  $\rightarrow$  Algorithm
- **SPL and SPL Compiler**  
Algorithm  $\rightarrow$  Implementation
- **Search Engine**  
How to find the best implementation
- **SPIRAL system**  
Everything taken together
- **Conclusions**

# Formulas in SPL

••••

```
( compose
  ( diagonal ( 2*cos(1/16*pi) 2*cos(3/16*pi) 2*cos(5/16*pi) 2*cos(7/16*pi) ) )
  ( permutation ( 1 3 4 2 ) )
  ( tensor
    ( I 2 )
    ( F 2 )
  )
  ( permutation ( 1 4 2 3 ) )
  ( direct_sum
    ( compose
      ( F 2 )
      ( diagonal ( 1 sqrt(1/2) ) )
    )
    ( compose
      ( matrix
        ( 1 1 0 )
        ( 0 (-1) 1 )
      )
      ( diagonal ( cos(13/8*pi)-sin(13/8*pi) sin(13/8*pi) cos(13/8*pi)+sin(13/8*pi) ) )
      ( matrix
        ( 1 0 )
        ( 1 1 )
        ( 0 1 )
      )
    )
  )
  ( permutation ( 2 1 ) )
)
```

••••

# SPL Syntax (Subset)



- matrix operations:
  - (compose formula formula ...)
  - (tensor formula formula ...)
  - (direct\_sum formula formula ...)
- direct matrix description:
  - (matrix (a11 a12 ...) (a21 a22 ...) ...)
  - (diagonal (d1 d2 ...))
  - (permutation (p1 p2 ...))
- parameterized matrices:
  - (I n)
  - (F n)
- scalars:
  - 1.5, 2/7, cos(..), w(3), pi, 1.2e-04
- definition of new symbols: ← allows extension of SPL
  - (define name formula)
  - (template formula (i-code-list))
- directives for code generation
  - #codetype real/complex
  - #unroll on/off ← controls loop unrolling

# SPL Compiler, 4-point FFT



```
(compose (tensor (F 2) (I 2)) (T 4 2)
(tensor (I 2) (F 2)) (L 4 2))
```

#codetype

complex

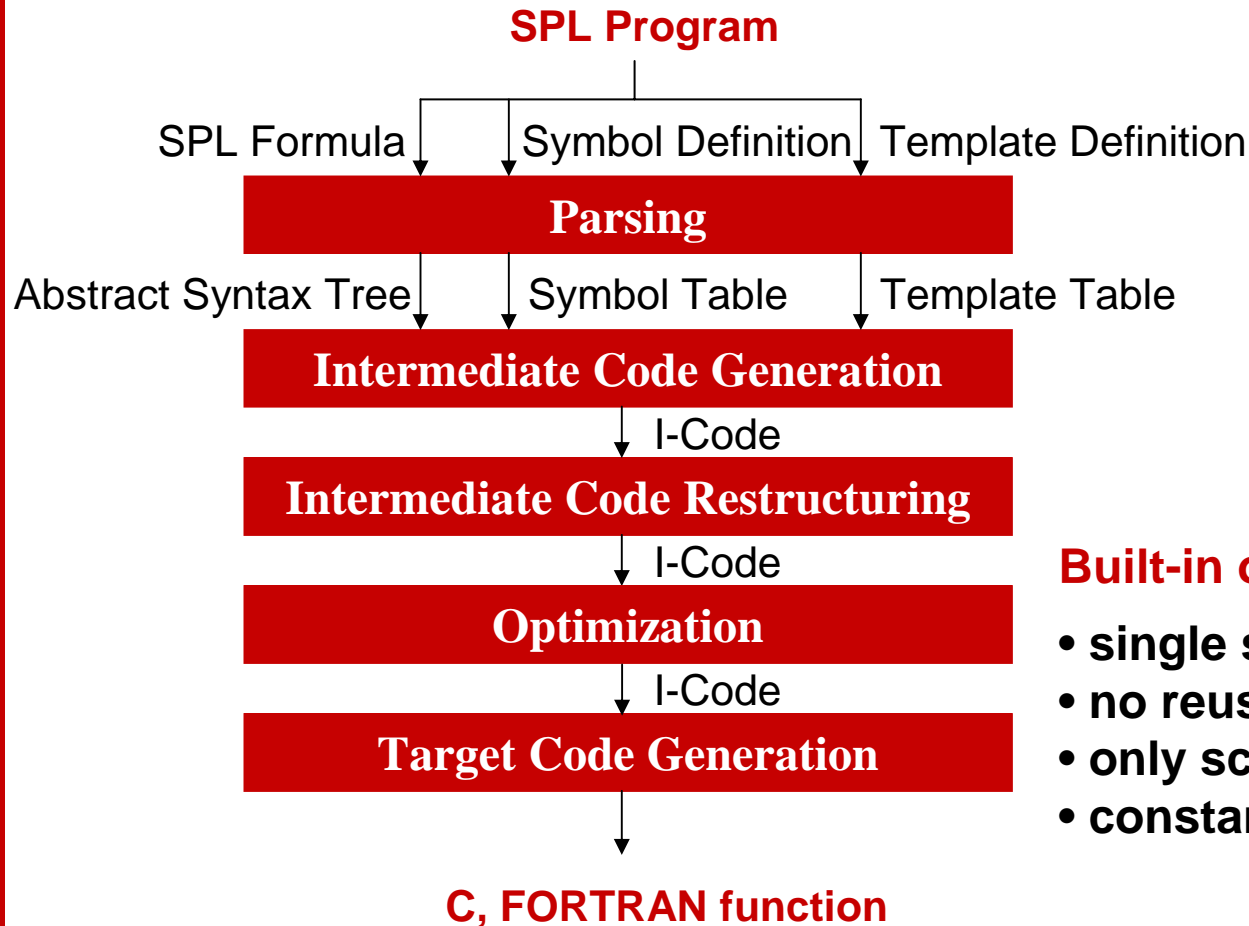
real

```
f0 = x(1) + x(3)
f1 = x(1) - x(3)
f2 = x(2) + x(4)
f3 = x(2) - x(4)
f4 = (0.00d0,-1.00d0)*f(3)
y(1) = f0 + f2
y(2) = f0 - f2
y(3) = f1 + f4
y(4) = f1 - f4
```

```
r0 = x(1) + x(5)
r1 = x(1) - x(5)
r2 = x(2) + x(6)
r3 = x(2) - x(6)
r4 = x(3) + x(7)
r5 = x(3) - x(7)
r6 = x(4) + x(8)
r7 = x(4) - x(8)
y(1) = r0 + r4
y(2) = r1 + r5
y(3) = r0 - r4
y(4) = r1 - r5
y(5) = r2 + r7
y(6) = r3 - r6
y(7) = r2 - r7
y(8) = r3 + r6
```

fast algorithm  
as  
formula  
as  
SPL program

# SPL Compiler: Summary



## Built-in optimizations:

- single static assignment code
- no reuse of temporary vars
- only scalar temporary vars
- constants precomputed

Extensible through templates



# Templates

```
( template
  ( F n ) [ n >= 1 ]
  ( do i=0,n-1
    y(i)=0
    do j=0,n-1
      y(i)=y(i)+W(n,i*j)*x(j)
    end
  end ) )
```

Pattern

Condition

I-code



SPIRAL

# Code Generation and Template Matching

(**F 2**) matches pattern (**F n**) and assigns 2 to **n**.

Because  $n=2$  satisfies the condition  $n \geq 1$ , the following i-code is generated from the template:

```
do i = 0,1
  y(i) = 0
  do j = 0,1
    y(i) = y(i)+W(2,i*j)*x(j)
  end
end
```

```
Y(0)=x(0)+x(1)
y(1)=x(0)-x(1)
```

Unrolling & Optimization



# SIMD Short Vector Extensions



- Extension to instruction set architecture
- Available on most current architectures (SSE on Pentium, AltiVec on Motorola G4)
- Originally for multimedia (like MMX for integers)
- Requires fine grain parallelism
- **Large potential speed-up**

## Problems:

- SIMD instructions are architecture specific
- No common API (usually assembly hand coding)
- Performance **very sensitive** to memory access
- Automatic vectorization **very limited**

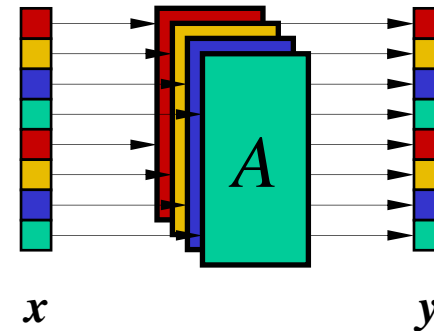


# Vector Code Generation from SPL Formulas

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)





# SPL Compiler: Vector Code Generation

$$DFT_{16} = (DFT_4 \otimes I_4) \cdot T_4^{16} \cdot (I_4 \otimes DFT_4) \cdot L_4^{16}$$



**Symbolic vectorization  
(automatic formula manipulation)**

$$\overline{DFT}_{16} = \left( (I_4 \otimes L_4^8) \cdot (\overline{DFT}_4 \otimes I_4) \cdot \overline{T}_4^{16} \right) \cdot \left( (I_4 \otimes L_2^8) (L_4^{16} \otimes I_2) (I_4 \otimes L_4^8) \cdot (\overline{DFT}_4 \otimes I_4) \cdot (I_4 \otimes L_2^8) \right)$$



**Mapping to C code + vector API**

```
LOAD_VECT(x10, x + 0);
LOAD_VECT(x14, x + 16);
f0 = SIMD_SUB(x10, x14);
LOAD_VECT(x11, x + 4);
LOAD_VECT(x15, x + 20);
f1 = SIMD_SUB(x11, x15);
...
y17 = SIMD_SUB(f1, f4);
STORE_L_8_4(y16, y17, y + 24);
y12 = SIMD_SUB(f0, f5);
y13 = SIMD_ADD(f1, f4);
STORE_L_8_4(y12, y13, y + 8);
```

**SSE  
SSE2  
AltiVec**

...

# Organization



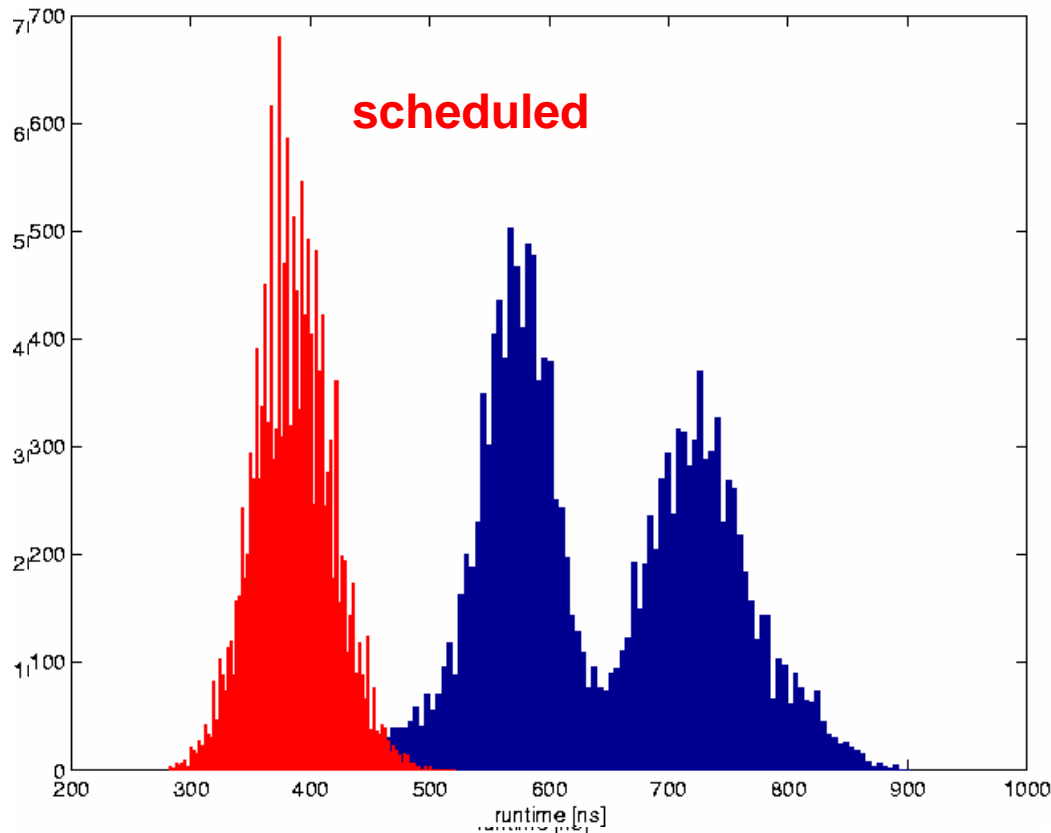
- **Mathematical Framework**  
Transforms, Rules, and Formulas
- **Formula Generator**  
Transform  $\rightarrow$  Algorithm
- **SPL and SPL Compiler**  
Algorithm  $\rightarrow$  Implementation
- **Search Engine**  
How to find the best implementation
- **SPIRAL system**  
Everything taken together
- **Conclusions**

# Why Search?

Carnegie Mellon



Tov problem



DCT IV- size  $2^4$

~31000 formulas

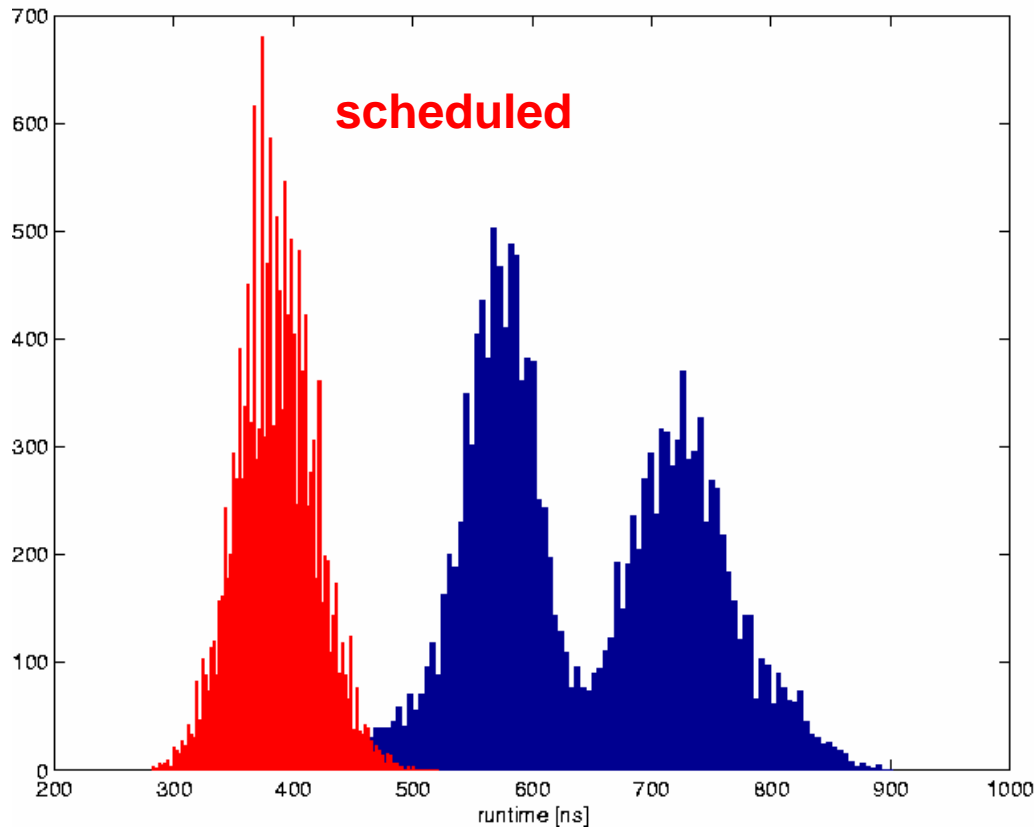
**Search in algorithm space in SPIRAL:**

**Exhaustive & Random Search, DP, Hill climbing, Genetic Algorithms**

**Beyond search: design of optimal tree**

# Why Search?

Carnegie Mellon



**Toy problem:**

**DCT, type IV, size 16**

**~31000 formulas**

- **maaaany different formulas**
- **large spread in runtimes, even for modest size**
- **precisely equal arithmetic cost**
- **best formula is platform-dependent**



# Number of Formulas/Algorithms

k	# DFT, size $2^k$	# DCT-IV, size $2^k$
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	$\sim 2.86 \cdot 10^{133}$	$\sim 1.06 \cdot 10^{153}$



- differ in data flow not in arithmetic cost
- exponential search space



# Search Methods Available in SPIRAL

- Exhaustive Search
- Dynamic Programming (DP)
- Random Search
- Hill Climbing
- STEER (similar to a genetic algorithm)

	Possible Sizes	Formulas Timed	Results
Exhaust	Very small	All	Best
DP	All	10s-100s	(very) good
Random	All	User decided	fair/good
Hill Climbing	All	100s-1000s	Good
STEER	All	100s-1000s	(very) good

## Search over

- algorithm space and
- implementation options (degree of unrolling)

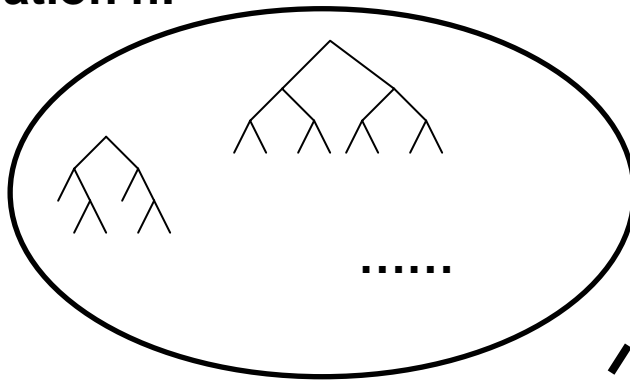


# STEER

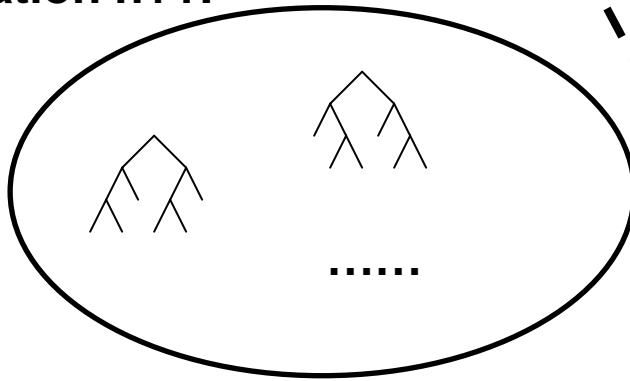
Carnegie Mellon



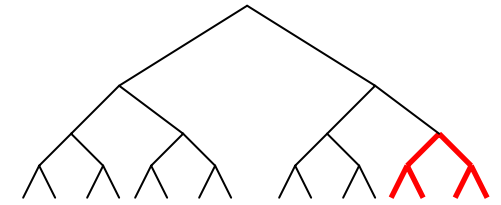
Population n:



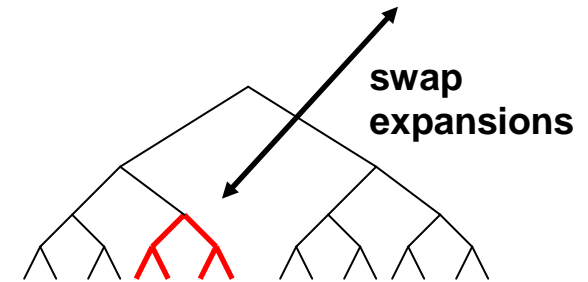
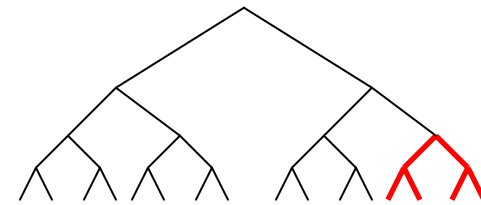
Population n+1:



Mutation



Cross-Breeding



Survival of Fittest

# Learning to Generate Fast Algorithms

- **Learns from given dataset** (formulas+runtimes) how to design a fast algorithm (breakdown strategy)
- Learns from a transform of **one size**, generates the best algorithm for **many sizes**
- Tested for DFT and WHT

## Fast Formula Generation Results

Size	Number of Formulas Generated	Generated Included the Fastest Known	Top N Fastest Known Formulas in Generated
$2^{12}$	101	yes	77
$2^{13}$	86	yes	4
$2^{14}$	101	yes	70
$2^{15}$	86	yes	11
$2^{16}$	101	yes	68
$2^{17}$	86	yes	15
$2^{18}$	101	yes	25
$2^{19}$	86	yes	16
$2^{20}$	101	yes	16

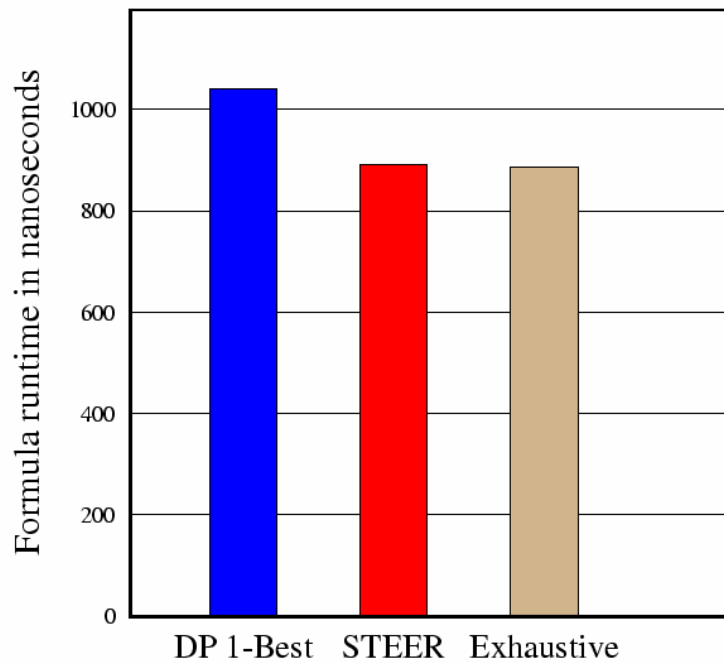
# Some Experimental Results

# DCT Type IV Size 16

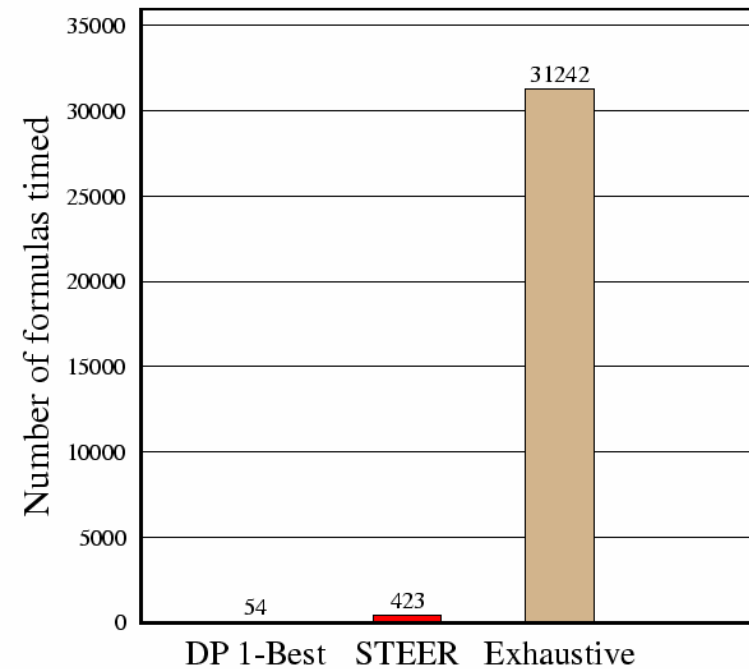
Carnegie Mellon



### Fastest Found Formulas



### Number of Formulas Timed

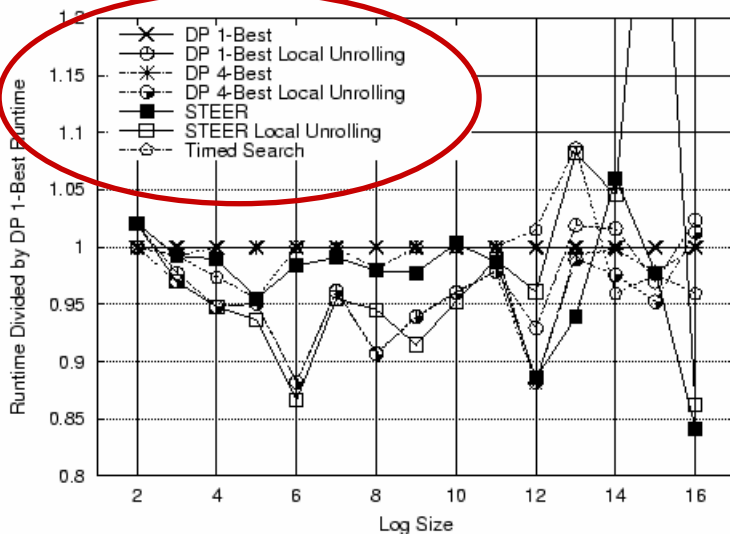


# Experimental Results

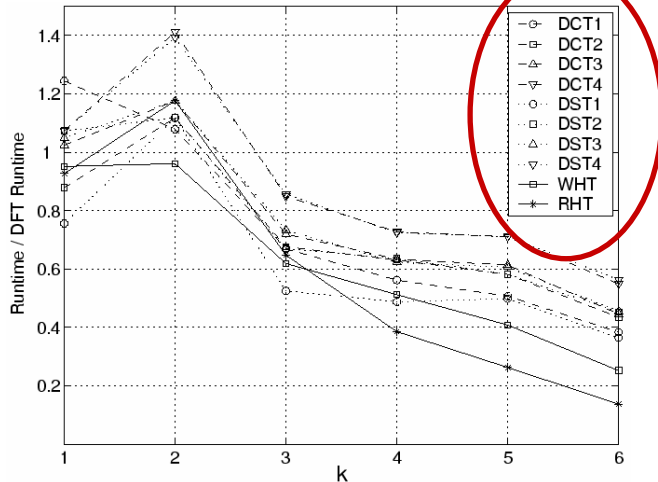
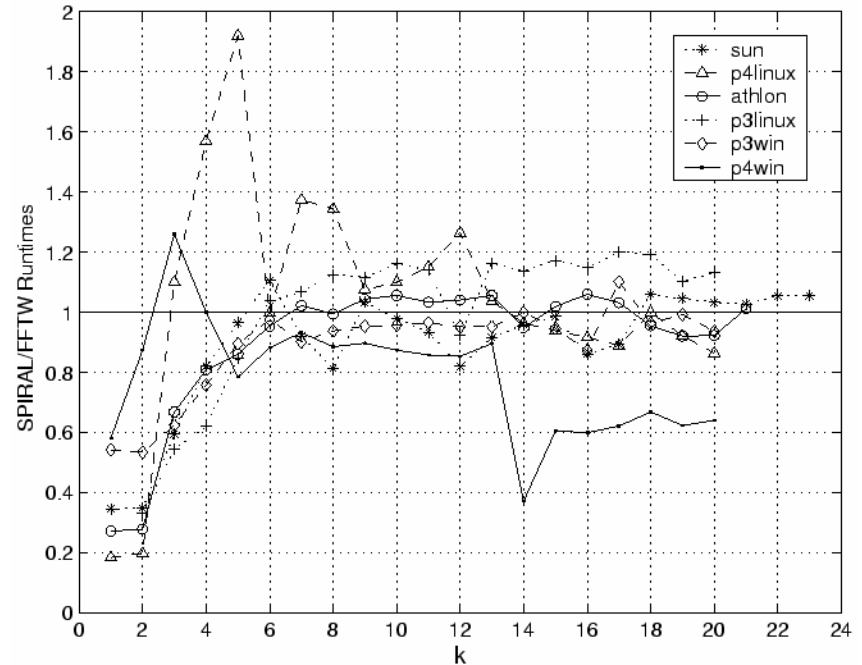
Carnegie Mellon



search methods  
(applicable to all transforms)



high performance code  
(compared with FFTW)



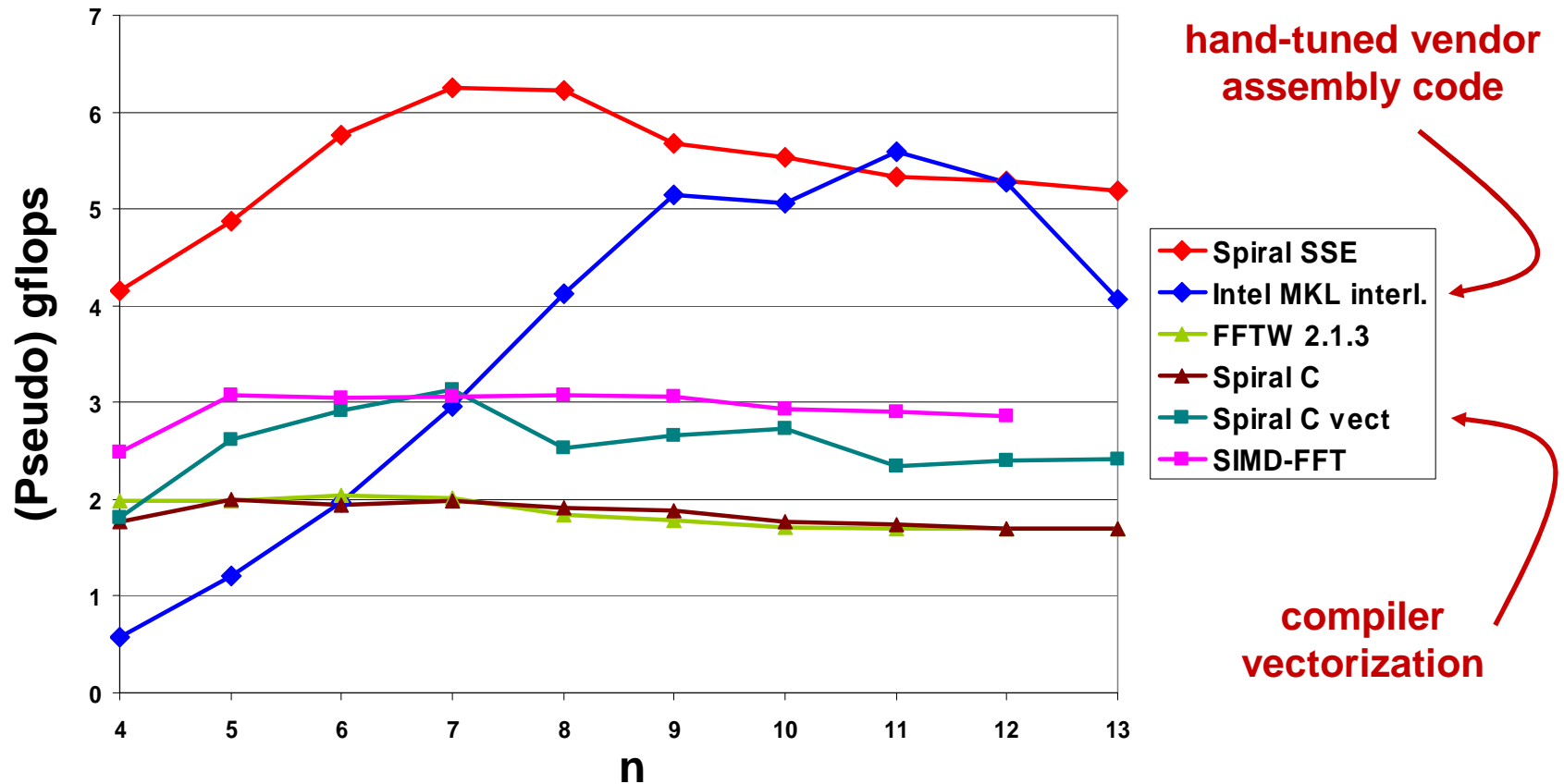
different transforms



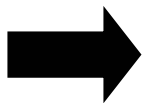
# Some Experimental Results



# Generated DFT Code: Pentium 4, SSE



DFT  $2^n$  single precision, Pentium 4, 2.53 GHz, using Intel C compiler 6.0



**speedups (to C code) up to factor of 3.1**

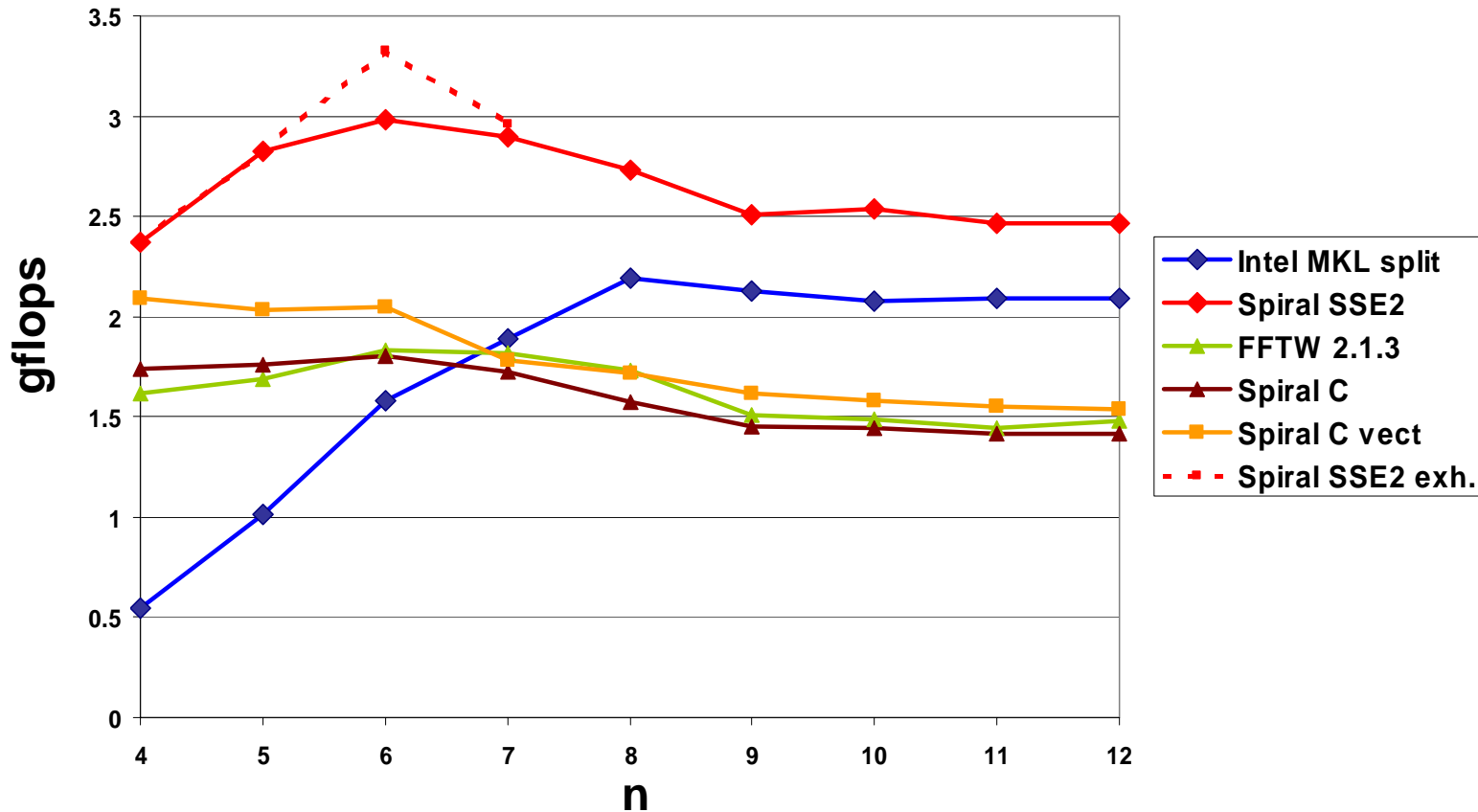
12/5/2002 IBM-Thomas T. J. Watson Res. Center



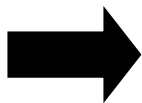
SPIRAL

# Generated DFT Code: Pentium 4, SSE2

Carnegie Mellon



DFT  $2^n$  *double* precision, Pentium 4, 2.53 GHz, using Intel C compiler 6.0



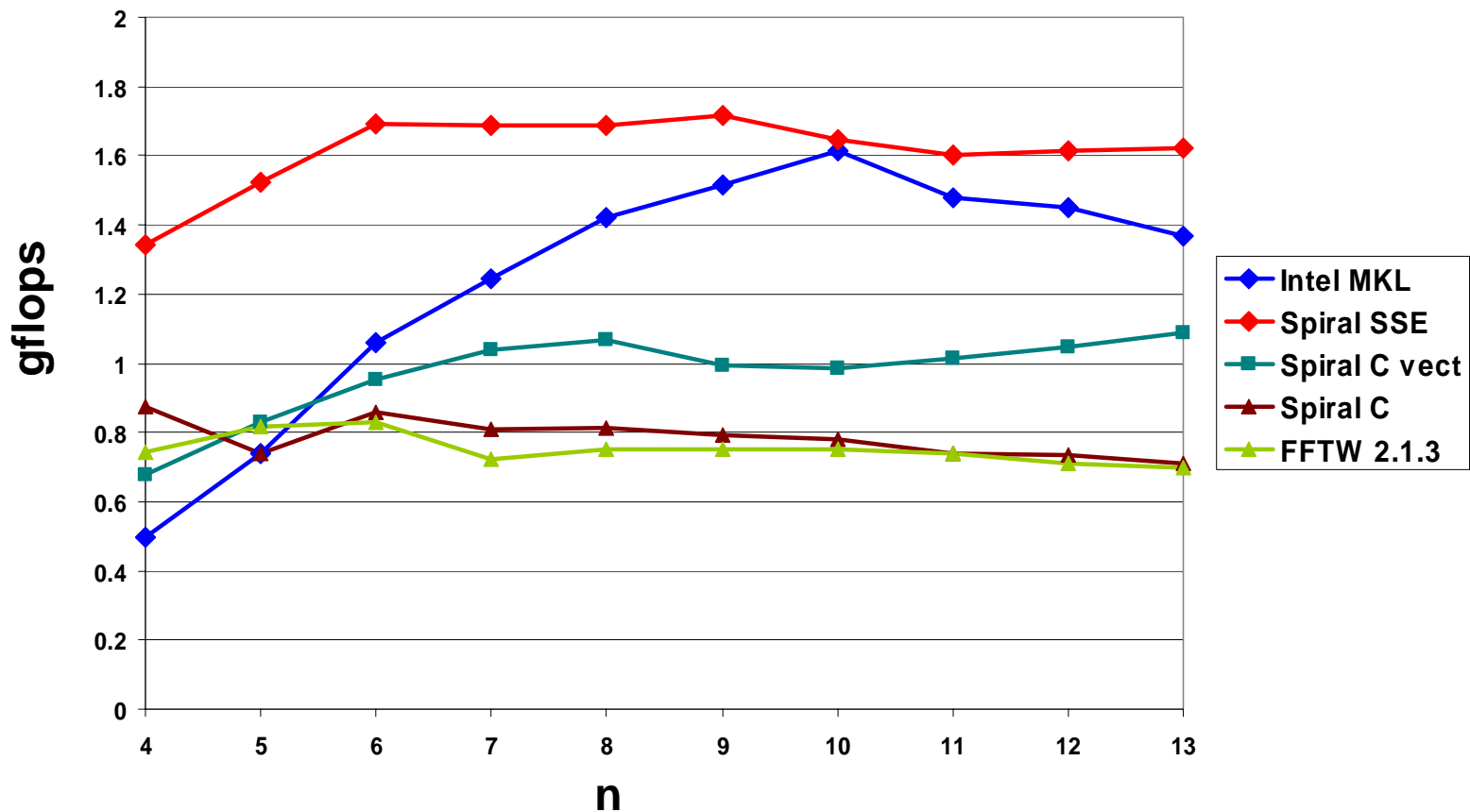
**speedups (to C code) up to factor of 1.8**



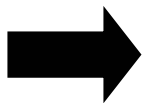
12/5/2002 IBM-Thomas T. J. Watson Res. Center



# Generated DFT Code: Pentium III, SSE



DFT  $2^n$  single precision, Pentium III, 1 GHz, using Intel C compiler 6.0



**speedups (to C code) up to factor of 2.1**

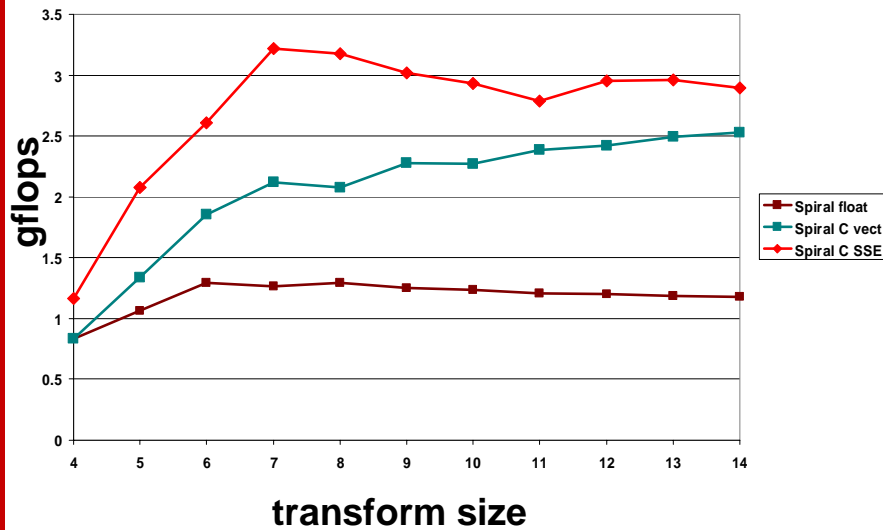


# Other transforms



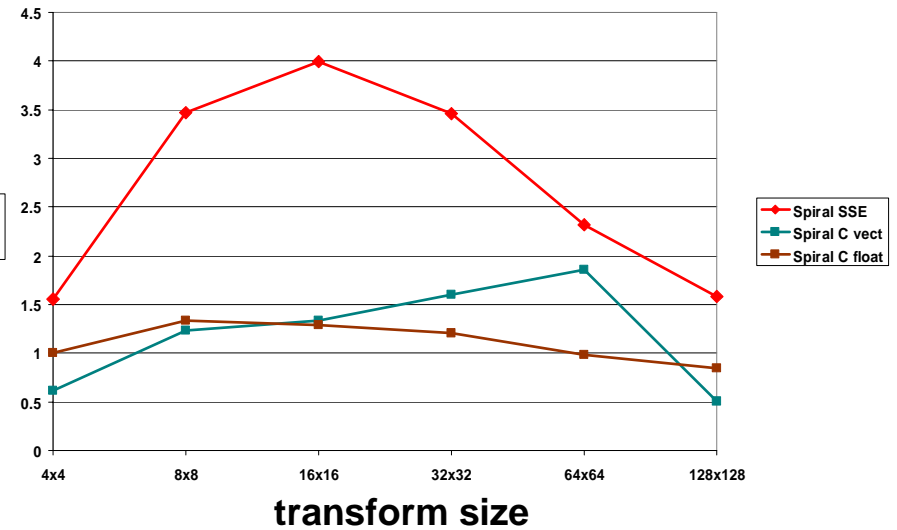
## WHT $2^n$

Pentium 4, 2.53 GHz, SSE

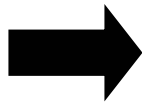


## 2-dim DCT $2^n \times 2^n$

Pentium 4, 2.53 GHz, SSE



- WHT has only additions
- very simple transform

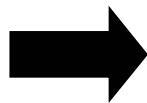


**speedups (to C code) up to factor of 3**

# Best DFT Trees, size $2^{10} = 1024$



	Pentium 4 float	Pentium 4 double	Pentium III float	AthlonXP float
scalar				
C vect				
SIMD				

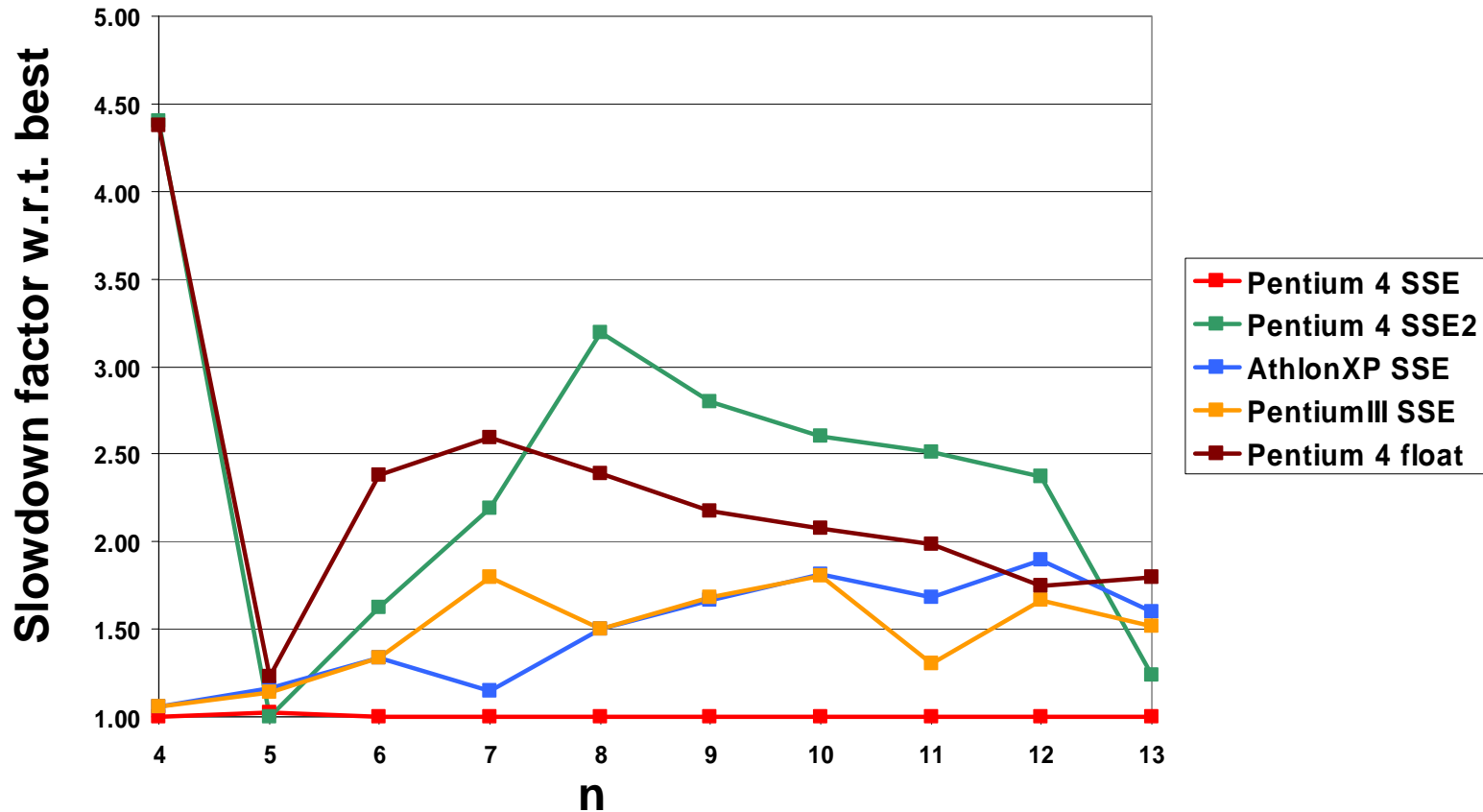


**trees platform/datatype dependent**

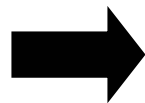
12/5/2002 IBM-Thomas T. J. Watson Res. Center

# Crosstiming of best trees on Pentium 4

Carnegie Mellon



DFT  $2^n$  single precision, runtime of best found of other platforms



**software adaptation is necessary**



SPIRAL

# Organization

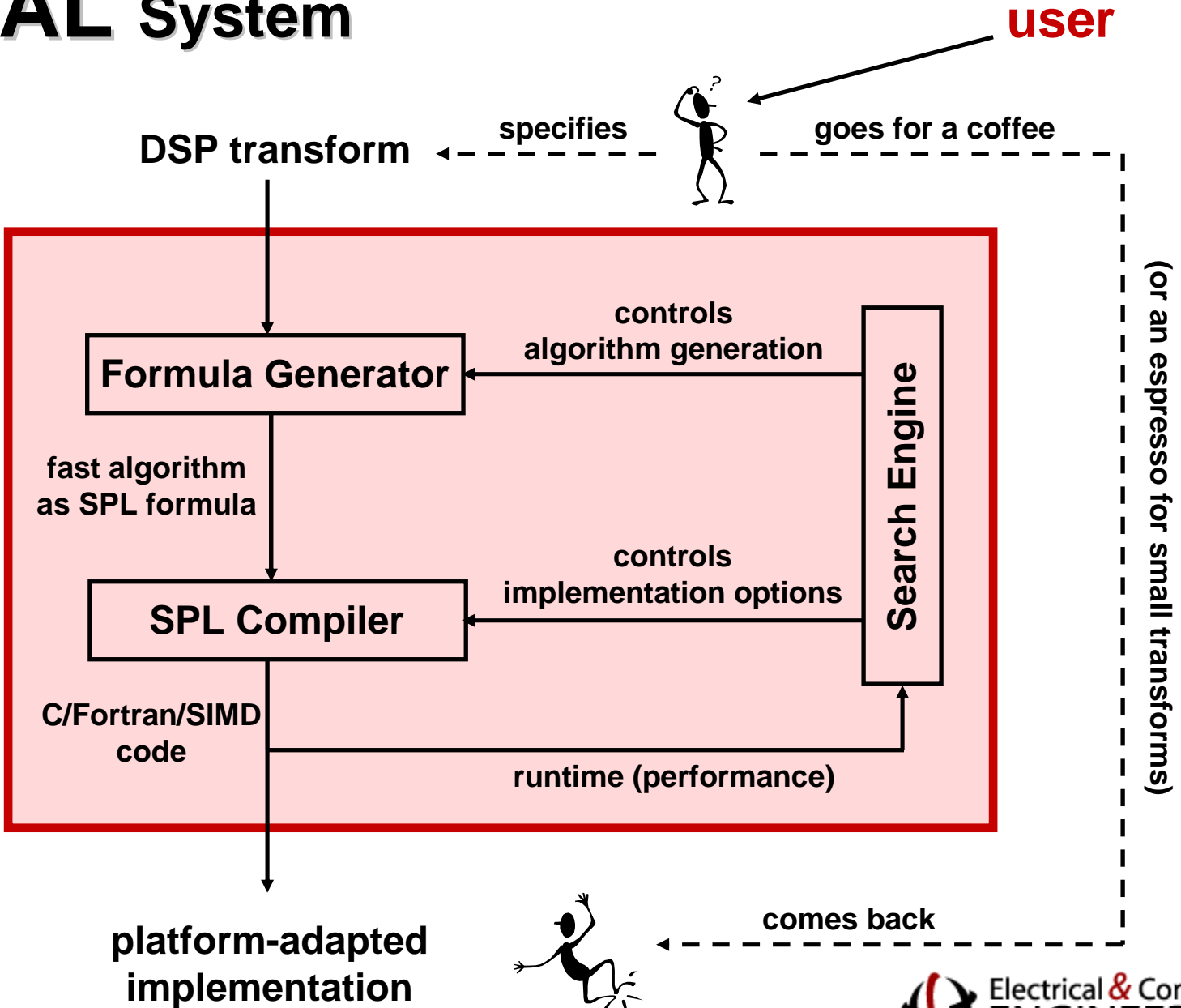


- **Mathematical Framework**  
Transforms, Rules, and Formulas
- **Formula Generator**  
Transform  $\rightarrow$  Algorithm
- **SPL and SPL Compiler**  
Algorithm  $\rightarrow$  Implementation
- **Search Engine**  
How to find the best implementation
- **SPIRAL system**  
Everything taken together
- **Conclusions**

# SPIRAL System

Carnegie Mellon

SPIRAL



platform-adapted implementation

12/5/2002 IBM-Thomas T. J. Watson Res. Center





# SPIRAL System

Carnegie Mellon



- Available for download (v3.1): [www.ece.cmu.edu/~spiral](http://www.ece.cmu.edu/~spiral)
  - Easy installation (Unix: configure/make; Windows: install shield)
  - Unix/Linux and Windows 98/ME/NT/2000/XP
  - Current transforms: DFT, DHT, WHT, RHT, DCT/DST type I – IV, MDCT, Filters, Wavelets, Toeplitz, Circulants
  - Extensible
  - **New version (4.0) in preparation**
- ~ 30 Publications in the areas of Signal Processing, High Performance Computing, Compilers, Machine Learning, Mathematics

# SPIRAL System: Summary

- Available for download: [www.ece.cmu.edu/~spiral](http://www.ece.cmu.edu/~spiral)
- Easy installation (Unix: configure/make; Windows: install shield)
- Unix/Linux and Windows 98/ME/NT/2000/XP
- Current transforms: DFT, DHT, WHT, DCT/DST type I – IV, MDCT, Filters, Wavelets, Toeplitz, Circulants
- Extensible



# Organization



- **Mathematical Framework**  
Transforms, Rules, and Formulas
- **Formula Generator**  
Transform  $\rightarrow$  Algorithm
- **SPL and SPL Compiler**  
Algorithm  $\rightarrow$  Implementation
- **Search Engine**  
How to find the best implementation
- **SPIRAL system**  
Everything taken together
- **Conclusions**

# Conclusions

**SPIRAL closes the gap between math domain (algorithms) and implementation domain (programs)**

- Mathematical computer representation of algorithms
- Automatic translation of algorithms into code

**SPIRAL does automatic optimization by intelligent search/learning in the space of alternatives**

- High level: Mathematical manipulation of algorithms
- Low level: Coding degrees of freedom



**a new paradigm of software development**