

# ***SPIRAL-Generated Modular FFTs\****

**Jeremy Johnson Lingchuan Meng**  
**Drexel University**

\* The work is supported by DARPA DESA, NSF, and Intel.  
Material for SPIRAL overview provided by Franz Francheti,  
Yevgen Voronenko, Markus Püschel, Martin Barbella and  
the rest of the SPIRAL team.

# Outline

## ■ SPIRAL Overview

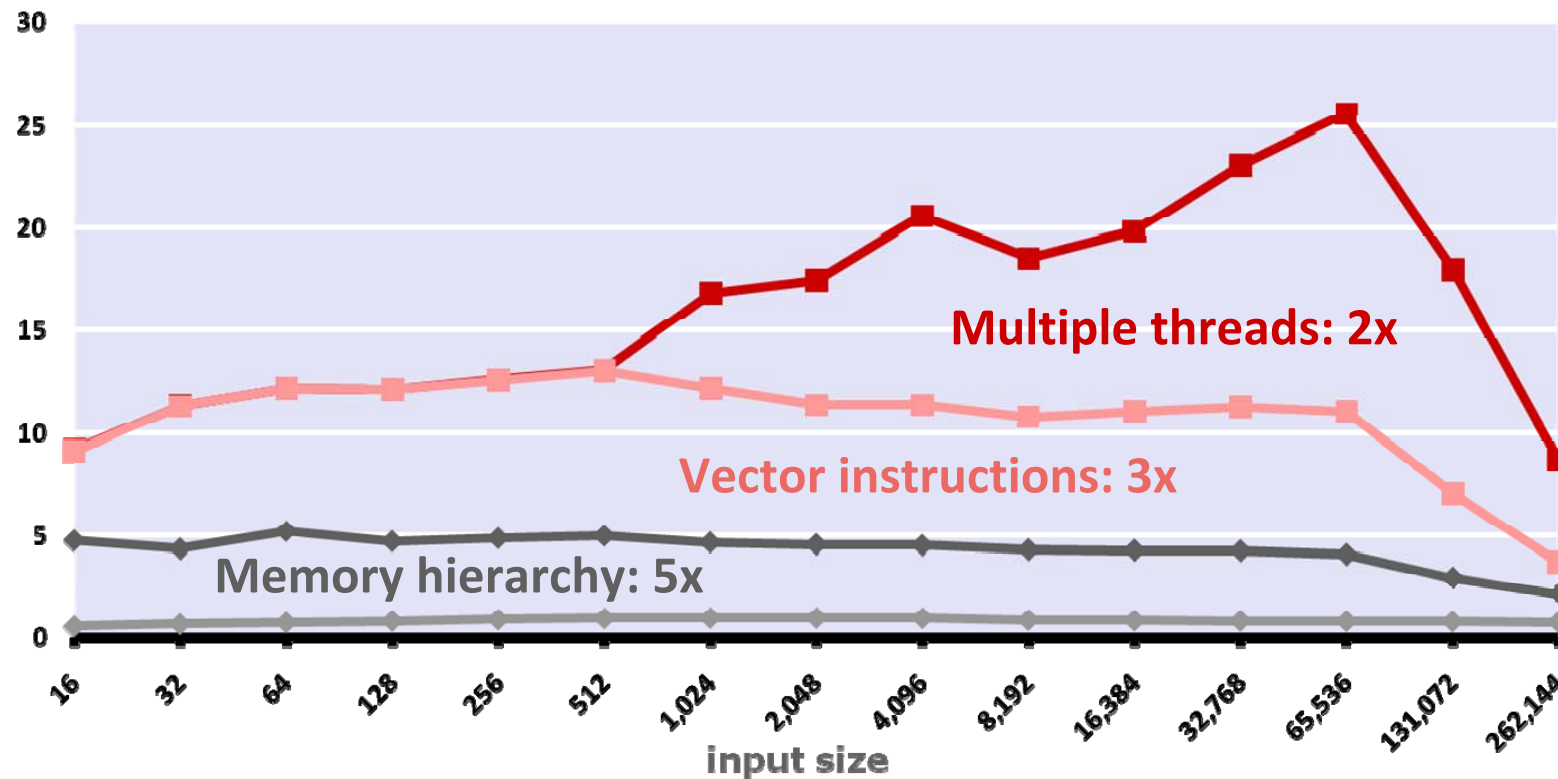
- Challenge of Obtaining Efficient Code
- Transforms and Rules
- SPIRAL Architecture
- SPIRAL Abstraction Levels
- Example
- Performance

## ■ Modular DFT in SPIRAL

- Addition of Modular DFT Transform and Modular FFT Rules
- Code Generation for Modular FFT
- Performance
- Future Work

# Challenge of Obtaining Efficient Code

Discrete Fourier Transform (DFT) on 2xCore2Duo 3 GHz  
Performance [Gflop/s]



*High performance library development  
has become a nightmare*

# Spiral Overview

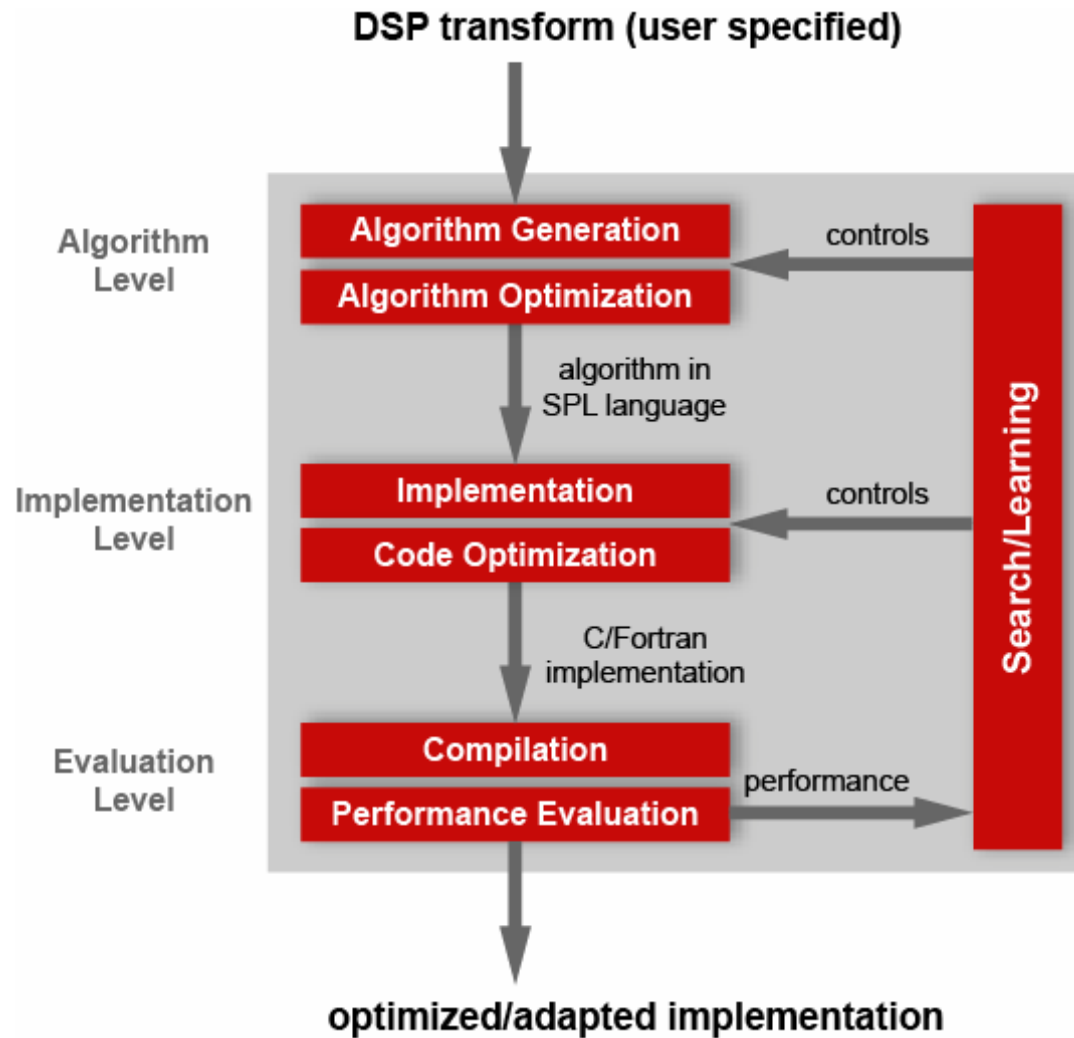
- **Research Goal: “Teach” computers to write fast libraries**
  - Complete automation of implementation and optimization
  - Including vectorization, parallelization
  
- **Functionality:**
  - Linear transforms (discrete Fourier transform, filters, wavelets)
  - BLAS
  - SAR imaging
  - En/decoding (Viterbi, Ebcot in JPEG2000)
  - ... more
  
- **Platforms:**
  - Desktop (vector, SMP), FPGAs, GPUs, distributed, hybrid
  
- **Collaboration with Intel (Kuck, Tang, Sabanin)**
  - Parts of MKL/IPP generated with Spiral
  - IPP 6.0: ippg domain for Spiral generated code

# DSP Algorithms: Transforms & Breakdown Rules

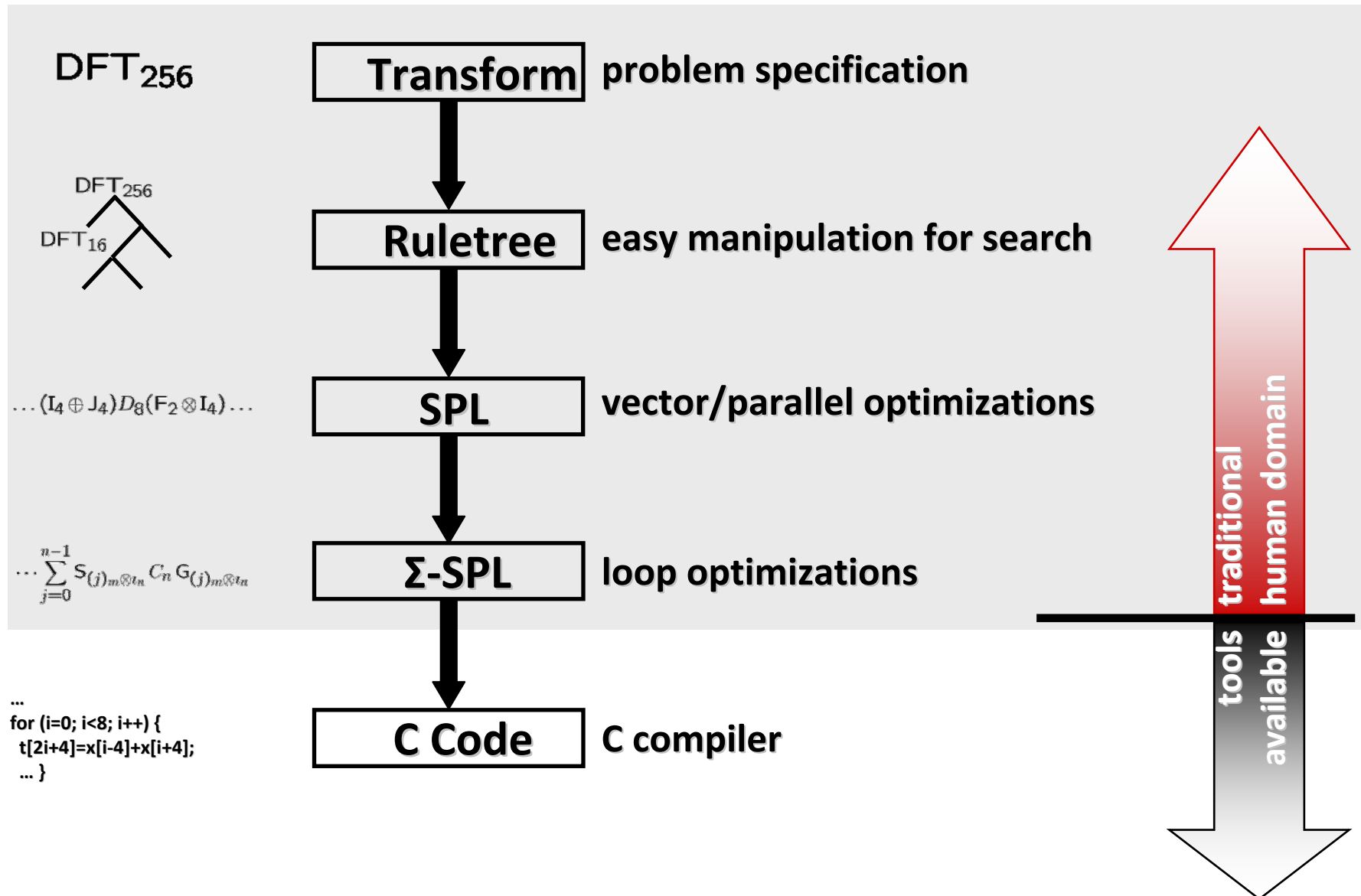
$$\begin{aligned}
 \text{DFT}_n &\rightarrow P_{k/2,2m}^\top (\text{DFT}_{2m} \oplus (I_{k/2-1} \otimes_i C_{2m} \text{rDFT}_{2m}(i/k))) (\text{RDFT}'_k \otimes I_m), \quad k \text{ even,} \\
 \begin{bmatrix} \text{RDFT}'_n \\ \text{RDFT}'_n \\ \text{DHT}'_n \\ \text{DHT}'_n \end{bmatrix} &\rightarrow (P_{k/2,m}^\top \otimes I_2) \left( \begin{bmatrix} \text{RDFT}'_{2m} \\ \text{RDFT}'_{2m} \\ \text{DHT}'_{2m} \\ \text{DHT}'_{2m} \end{bmatrix} \oplus \left( I_{k/2-1} \otimes_i D_{2m} \begin{bmatrix} \text{rDFT}_{2m}(i/k) \\ \text{rDFT}_{2m}(i/k) \\ \text{rDHT}_{2m}(i/k) \\ \text{rDHT}_{2m}(i/k) \end{bmatrix} \right) \right) \left( \begin{bmatrix} \text{RDFT}'_k \\ \text{RDFT}'_k \\ \text{DHT}'_k \\ \text{DHT}'_k \end{bmatrix} \otimes I_m \right), \quad k \text{ even,} \\
 \begin{bmatrix} \text{rDFT}_{2n}(u) \\ \text{rDHT}_{2n}(u) \end{bmatrix} &\rightarrow L_m^{2n} \left( I_k \otimes_i \begin{bmatrix} \text{rDFT}_{2m}((i+u)/k) \\ \text{rDHT}_{2m}((i+u)/k) \end{bmatrix} \right) \left( \begin{bmatrix} \text{rDFT}_{2k}(u) \\ \text{rDHT}_{2k}(u) \end{bmatrix} \otimes I_m \right), \\
 \text{RDFT-3}_n &\rightarrow (Q_{k/2,m}^\top \otimes I_2) (I_k \otimes_i \text{rDFT}_{2m})(i+1/2)/k) (\text{RDFT-3}_k \otimes I_m), \quad k \text{ even,} \\
 \text{DCT-2}_n &\rightarrow P_{k/2,2m}^\top (\text{DCT-2}_{2m} K_2^{2m} \oplus (I_{k/2-1} \otimes N_{2m} \text{RDFT-3}_{2m}^\top)) B_n (L_{k/2}^{n/2} \otimes I_2) (I_m \otimes \text{RDFT}'_k) Q_{m/2,k}, \\
 \text{DCT-3}_n &\rightarrow \text{DCT-2}_n^\top, \\
 \text{DCT-4}_n &\rightarrow Q_{k/2,2m}^\top (I_{k/2} \otimes N_{2m} \text{RDFT-3}_{2m}^\top) B'_n (L_{k/2}^{n/2} \otimes I_2) (I_m \otimes \text{RDFT-3}_k) Q_{m/2,k}. \\
 \text{DFT}_n &\rightarrow (\text{DFT}_k \otimes I_m) \Upsilon_m^n (I_k \otimes \text{DFT}_m) \mathcal{L}_k^n, \quad n = km \\
 \text{DFT}_n &\rightarrow P_n (\text{DFT}_k \otimes \text{DFT}_m) Q_n, \quad n = km, \text{ gcd}(k, m) = 1 \\
 \text{DFT}_p &\rightarrow R_p^\top (I_1 \oplus \text{DFT}_{p-1}) D_p (I_1 \oplus \text{DFT}_{p-1}) R_p, \quad p \text{ prime} \\
 \text{DCT-3}_n &\rightarrow (I_m \oplus J_m) \mathcal{L}_m^n (\text{DCT-3}_m(1/4) \oplus \text{DCT-3}_m(3/4)) \\
 &\quad \cdot (F_2 \otimes I_m) \begin{bmatrix} I_m & 0 \oplus -J_{m-1} \\ \frac{1}{\sqrt{2}}(I_1 \oplus 2I_m) \end{bmatrix}, \quad n = 2m \\
 \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 (J_m \oplus I_m \oplus I_m \oplus J_m) \left( \left( \begin{bmatrix} 1 \\ -1 \end{bmatrix} \otimes I_m \right) \oplus \left( \begin{bmatrix} -1 \\ -1 \end{bmatrix} \otimes I_m \right) \right) J_{2m} \text{DCT-4}_{2m} \\
 \text{WHT}_{2^k} &\rightarrow \prod_{i=1}^t (I_{2^{k_1+\dots+k_{i-1}}} \otimes \text{WHT}_{2^{k_i}} \otimes I_{2^{k_{i+1}+\dots+k_t}}), \quad k = k_1 + \dots + k_t \\
 \text{DFT}_2 &\rightarrow F_2 \\
 \text{DCT-2}_2 &\rightarrow \text{diag}(1, 1/\sqrt{2}) F_2 \\
 \text{DCT-4}_2 &\rightarrow J_2 R_{13\pi/8}
 \end{aligned}$$

**Combining these rules yields many algorithms for every given transform**

# SPIRAL: Architecture



# SPIRAL: Abstraction Levels



# DSP Algorithms: E.G. 4-point DFT

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

input vector (signal)  $x$  → output vector (signal)  $y$  = transform = matrix  $M$

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

Fourier transform

Diagonal matrix (twiddles)

Input vector

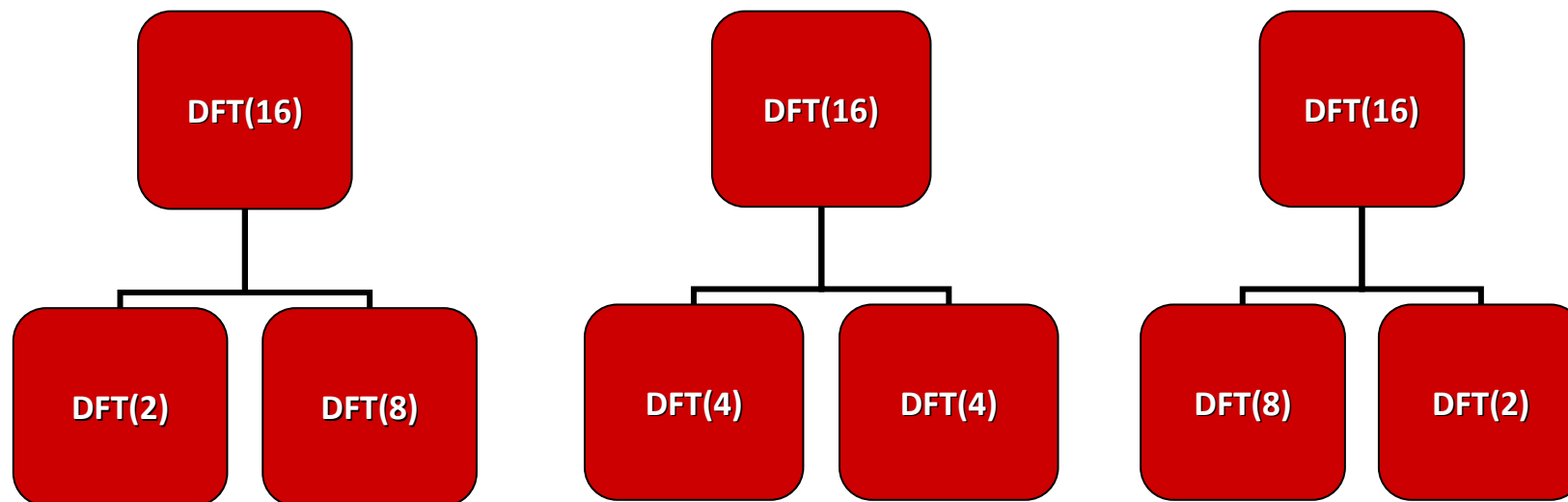
$$y = (\text{DFT}_2 \otimes \text{I}_2) \cdot \text{T}_2^4 \cdot (\text{I}_2 \otimes \text{DFT}_2) \cdot \text{L}_2^4 x$$

Output vector      Kronecker product      Identity      Permutation



# Cooley-Tukey Rule and DP

$$DFT_{mn} = (\underbrace{DFT_m \otimes I_n}_{\text{Kronecker product}}) \cdot \underbrace{T_n^{mn}}_{\text{Diagonal matrix (twiddles)}} \cdot (\underbrace{I_m \otimes DFT_n}_{\text{Identity}}) \cdot \underbrace{L_m^{mn}}_{\text{Permutation}}$$



*Dynamic Programming (DP) searches over different applications of CT Rule using best transform found for recursive calls*

# Program Generation in Spiral (Sketched)

Transform  
*user specified*

DFT<sub>8</sub>



**Fast algorithm**  $(\text{DFT}_2 \otimes I_4) T_4^8 (I_2 \otimes ((\text{DFT}_2 \otimes I_2)$   
**in SPL many choices**  $\cdot T_2^4 (I_2 \otimes \text{DFT}_2) L_2^4)) L_2^8$



$\Sigma$ -SPL:  
*[PLDI 05]*

$$\sum (S_j \text{DFT}_2 G_j) \sum \left( \sum (S_{k,l} \text{diag}(t_{k,l}) \text{DFT}_2 G_l) \right. \\ \left. \sum (S_m \text{diag}(t_m) \text{DFT}_2 G_{k,m}) \right)$$



**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;
}
```

*Optimization at all  
abstraction levels*



parallelization  
vectorization



loop  
optimizations



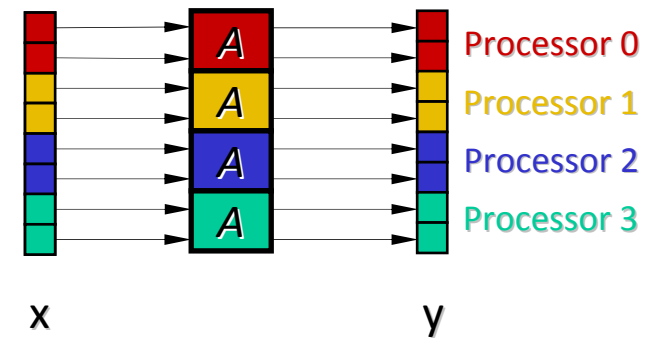
constant folding  
scheduling

.....

# SPL to Shared Memory Code: Basic Idea [SC 06]

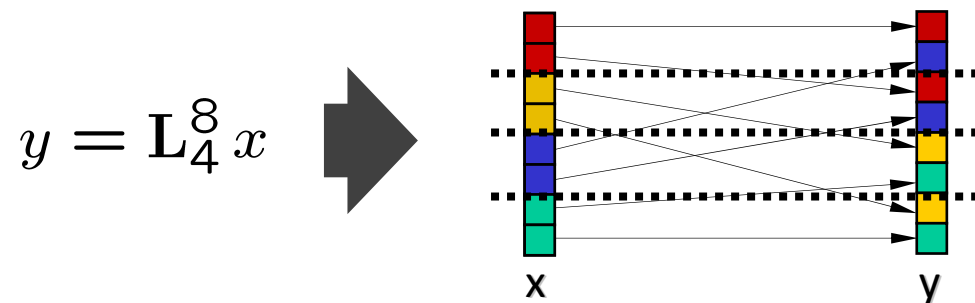
- Governing construct: tensor product

$$y = \left( I_p \otimes A \right) x$$



*p-way embarrassingly parallel, load-balanced*

- Problematic construct: permutations produce false sharing



*Task: Rewrite formulas to extract tensor product + keep contiguous blocks*

# Parallelization by Rewriting

$$\begin{array}{l}
 \underbrace{\text{DFT}_{mn}}_{\text{smp}(p,\mu)} \rightarrow \underbrace{\left( (\text{DFT}_m \otimes \text{I}_n) \text{T}_n^{mn} (\text{I}_m \otimes \text{DFT}_n) \text{L}_m^{mn} \right)}_{\text{smp}(p,\mu)} \\
 \uparrow \\
 \text{coarse platform model} \rightarrow \dots \\
 \rightarrow \underbrace{\left( \text{DFT}_m \otimes \text{I}_n \right)}_{\text{smp}(p,\mu)} \underbrace{\text{T}_n^{mn}}_{\text{smp}(p,\mu)} \underbrace{\left( \text{I}_m \otimes \text{DFT}_n \right)}_{\text{smp}(p,\mu)} \underbrace{\text{L}_m^{nm}}_{\text{smp}(p,\mu)} \\
 \dots \\
 \rightarrow \underbrace{\left( (\text{L}_m^{mp} \otimes \text{I}_{n/p\mu}) \otimes_{\mu} \text{I}_{\mu} \right)}_{\text{smp}(p,\mu)} \underbrace{\left( \text{I}_p \otimes_{\parallel} (\text{DFT}_m \otimes \text{I}_{n/p}) \right)}_{\text{smp}(p,\mu)} \underbrace{\left( (\text{L}_p^{mp} \otimes \text{I}_{n/p\mu}) \otimes_{\mu} \text{I}_{\mu} \right)}_{\text{smp}(p,\mu)} \\
 \underbrace{\left( \bigoplus_{i=0}^{p-1} \text{T}_n^{mn,i} \right)}_{\text{smp}(p,\mu)} \underbrace{\left( \text{I}_p \otimes_{\parallel} (\text{I}_{m/p} \otimes \text{DFT}_n) \right)}_{\text{smp}(p,\mu)} \underbrace{\left( \text{I}_p \otimes_{\parallel} \text{L}_{m/p}^{mn/p} \right)}_{\text{smp}(p,\mu)} \underbrace{\left( (\text{L}_p^{pn} \otimes \text{I}_{m/p\mu}) \otimes_{\mu} \text{I}_{\mu} \right)}_{\text{smp}(p,\mu)}
 \end{array}$$

**Load-balanced**

**No false sharing**

# Same Approach for Other Parallel Paradigms

## Message Passing [ISPA 06]

$$\begin{aligned}
 \underbrace{\text{DFT}_{mn}}_{\text{par}(p)} &\rightarrow \underbrace{(\text{DFT}_m \otimes \text{I}_n)}_{\text{par}(p-q)} \underbrace{\text{T}_n^{mn}}_{\text{par}(q)} \underbrace{(\text{I}_m \otimes \text{DFT}_n)}_{\text{par}(q)} \underbrace{\text{L}_m^{mn}}_{\text{par}(q-p)} \\
 &\dots \\
 &\dots \\
 &\dots \\
 &\rightarrow \underbrace{(\text{I}_p \otimes \text{L}_{m/p}^{mn/p})}_{\text{par}(p)} \underbrace{(\text{L}_p^2 \otimes \text{I}_{mn/p^2})}_{\text{par}(p-q)} \underbrace{(\text{I}_q \otimes (\text{I}_{p/q} \otimes \text{L}_p^n \otimes \text{I}_{m/p}))}_{\text{par}(q)} \underbrace{(\text{I}_q \otimes (\text{I}_{n/q} \otimes \text{DFT}_m))}_{\text{par}(q)} \\
 &\quad \underbrace{(\text{I}_q \otimes \text{L}_{m/q}^{mn/q})}_{\text{par}(q)} \underbrace{(\text{L}_q^2 \otimes \text{I}_{mn/q^2})}_{\text{par}(q)} \underbrace{(\text{I}_q \otimes (\text{L}_q^n \otimes \text{I}_{m/q}))}_{\text{par}(q)} \underbrace{\text{T}_n^{mn}}_{\text{par}(q)} \underbrace{(\text{I}_q \otimes (\text{I}_{m/q} \otimes \text{DFT}_n))}_{\text{par}(q)} \\
 &\quad \underbrace{(\text{I}_q \otimes (\text{I}_{p/q} \otimes \text{L}_{m/p}^{mn/p}))}_{\text{par}(q)} \underbrace{(\text{L}_p^2 \otimes \text{I}_{mn/p^2})}_{\text{par}(q-p)} \underbrace{(\text{I}_p \otimes (\text{L}_p^n \otimes \text{I}_{m/p}))}_{\text{par}(p)}
 \end{aligned}$$

## Vectorization [VECPAR 06]

$$\begin{aligned}
 \underbrace{(\text{DFT}_{mn})}_{\text{vec}(\nu)} &\rightarrow \underbrace{((\text{DFT}_m \otimes \text{I}_n) \text{T}_n^{mn} (\text{I}_m \otimes \text{DFT}_n) \text{L}_m^{mn})}_{\text{vec}(\nu)} \\
 &\dots \\
 &\dots \\
 &\rightarrow \underbrace{(\text{DFT}_m \otimes \text{I}_n)}_{\text{vec}(\nu)}^\nu \underbrace{(\text{T}_n^{mn})}_{\text{vec}(\nu)}^\nu \underbrace{(\text{I}_m \otimes \text{DFT}_n) \text{L}_m^{mn}}_{\text{vec}(\nu)}^\nu \\
 &\dots \\
 &\rightarrow (\text{I}_{mn/\nu} \otimes \underbrace{\text{L}_\nu^{2\nu}}_{\text{sse}}) (\overline{\text{DFT}_m \otimes \text{I}_{n/\nu}} \otimes \overline{\text{I}_\nu}) \underbrace{(\text{T}_n^{mn})}_{\text{sse}}^\nu \\
 &\quad (\text{I}_{m/\nu} \otimes (\overline{\text{L}_\nu^n} \otimes \overline{\text{I}_\nu})) (\text{I}_{n/\nu} \otimes (\text{L}_\nu^{2\nu} \otimes \overline{\text{I}_\nu})) (\text{I}_2 \otimes \underbrace{\text{L}_\nu^{2\nu}}_{\text{sse}}) (\text{L}_2^{2\nu} \otimes \overline{\text{I}_\nu}) (\overline{\text{DFT}_n} \otimes \overline{\text{I}_\nu}) \\
 &\quad (\text{L}_m^{mn} \otimes \text{I}_2) \otimes \overline{\text{I}_\nu} (\text{I}_{mn/\nu} \otimes \underbrace{\text{L}_\nu^{2\nu}}_{\text{sse}})
 \end{aligned}$$

## Cg/OpenGL for GPUs

$$\begin{aligned}
 \underbrace{(\text{DFT}_{r^k})}_{\text{gpu}(t,c)} &\rightarrow \underbrace{\left( \prod_{i=0}^{k-1} \text{L}_r^{r^k} (\text{I}_{r^{k-1}} \otimes \text{DFT}_r) \left( \text{L}_{r^{k-i-1}}^{r^k} (\text{I}_{r^i} \otimes \text{T}_{r^{k-i-1}}) \underbrace{\text{L}_{r^{i+1}}^{r^k}}_{\text{vec}(c)} \right) \right)}_{\text{gpu}(t,c)} \text{R}_r^{r^k} \\
 &\dots \\
 &\dots \\
 &\rightarrow \left( \prod_{i=0}^{k-1} (\text{L}_r^{r^n/2} \otimes \overline{\text{I}_2}) (\text{I}_{r^{n-1}/2} \otimes \times \underbrace{(\text{DFT}_r \otimes \overline{\text{I}_2}) \text{L}_r^{2r}}_{\text{shd}(t,c)}) \text{T}_i \right) \\
 &\quad (\text{L}_r^{r^n/2} \otimes \overline{\text{I}_2}) (\text{I}_{r^{n-1}/2} \otimes \times \underbrace{\text{L}_r^{2r}}_{\text{shd}(t,c)}) (\text{R}_r^{r^{n-1}} \otimes \overline{\text{I}_r})
 \end{aligned}$$

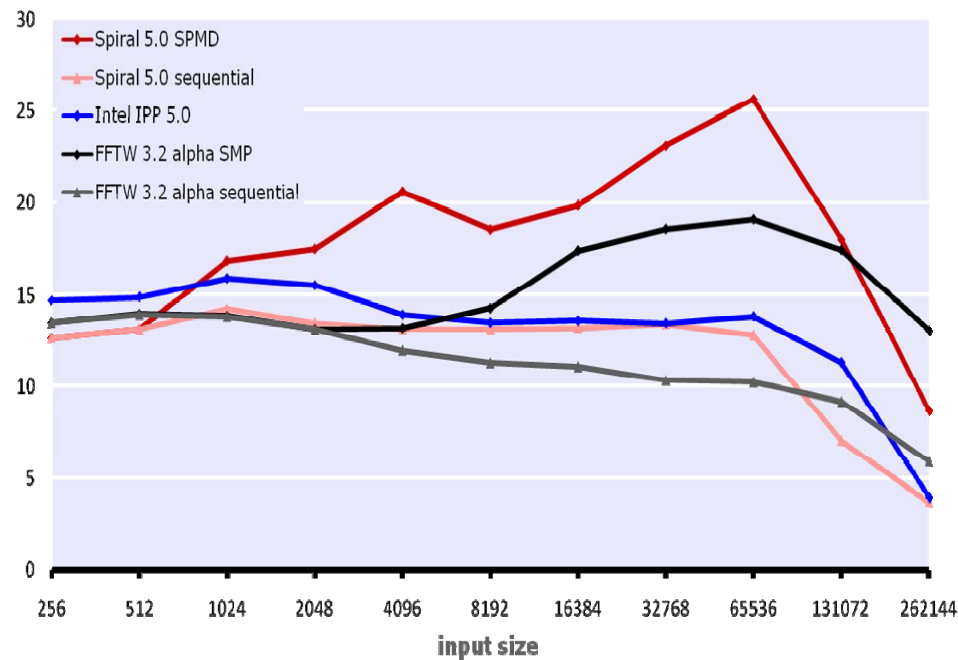
## Verilog for FPGAs [DAC 08]

$$\begin{aligned}
 \underbrace{(\text{DFT}_{r^k})}_{\text{stream}(r^s)} &\rightarrow \underbrace{\left[ \prod_{i=0}^{k-1} \text{L}_r^{r^k} (\text{I}_{r^{k-1}} \otimes \text{DFT}_r) \left( \text{L}_{r^{k-i-1}}^{r^k} (\text{I}_{r^i} \otimes \text{T}_{r^{k-i-1}}) \text{L}_{r^{i+1}}^{r^k} \right) \right]}_{\text{stream}(r^s)} \text{R}_r^{r^k} \\
 &\dots \\
 &\dots \\
 &\rightarrow \left[ \prod_{i=0}^{k-1} \underbrace{\text{L}_r^{r^k}}_{\text{stream}(r^s)} \underbrace{(\text{I}_{r^{k-1}} \otimes \text{DFT}_r)}_{\text{stream}(r^s)} \left( \text{L}_{r^{k-i-1}}^{r^k} (\text{I}_{r^i} \otimes \text{T}_{r^{k-i-1}}) \text{L}_{r^{i+1}}^{r^k} \right) \right]_{\text{stream}(r^s)} \text{R}_r^{r^k} \\
 &\dots \\
 &\rightarrow \left[ \prod_{i=0}^{k-1} \underbrace{\text{L}_r^{r^k}}_{\text{stream}(r^s)} (\text{I}_{r^{k-s-1}} \otimes \text{I}_{r^{s-1}} \otimes \text{DFT}_r) \right]_{\text{stream}(r^s)} \underbrace{\text{T}_i}_{\text{stream}(r^s)} \text{R}_r^{r^k}
 \end{aligned}$$

# Example Results

## Multicore [SC 06]

**DFT (single precision): on 3 GHz 2 x Core 2 Extreme  
performance [Gflop/s]**



***Code written by the computer is faster (for many sizes)  
than any human-written code***

# Modular DFT in SPIRAL

- **New Transform: ModDFT( $n, p, w, i$ )**
  - $n$ , the size of the transform
  - $p$ , the prime modulus
  - $w$ , primitive  $n^{\text{th}}$  root of unity
  - $i$ , forward/inverse flag
  - *Same as DFT except must store  $p$  and  $w$*
- **Breakdown Rule for ModDFT: Cooley-Tukey in  $Z_p$** 
  - Same as complex DFT except different roots of unity are computed
- **CUnparser for ModDFT: Generate  $Z_p$  C code**
  - C code generated for modular arithmetic operations
  - Alternative approaches explored

# Cooley-Tukey in $Z_p$ : Example

$$\begin{array}{c}
 x \longmapsto y = M \cdot x \\
 \begin{array}{cccc}
 | & | & | & \\
 \text{Input Vector} & \text{Output Vector} & \text{Transform matrix in } Z_p & 
 \end{array}
 \end{array}$$

$$\text{ModDFT}(4, 17, 3) = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 13 & 16 & 4 \\ 1 & 16 & 1 & 16 \\ 1 & 4 & 16 & 13 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & 16 & 0 \\ 0 & 1 & 0 & 16 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 13 \end{bmatrix} \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 16 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 16 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$\begin{array}{c}
 \text{Modular Fourier transform} \quad \text{Diagonal matrix (twiddles)} \quad \text{Input vector} \\
 | \quad \quad \quad | \quad \quad \quad | \\
 y = (\text{ModDFT}_2 \otimes I_2) \cdot T_2^4 \cdot (I_2 \otimes \text{ModDFT}_2) \cdot L_2^4 \cdot x \\
 \begin{array}{cccc}
 | & | & | & | \\
 \text{Output vector} & \text{Kronecker product} & \text{Identity} & \text{Permutation}
 \end{array}
 \end{array}$$



# Modular FFTs: Cooley-Tukey in $Z_p$

## ■ Cooley-Tukey algorithm

A transform of size  $mn$  breaks down into smaller transforms of size  $m$  and  $n$

$$\begin{aligned}
 DFT(m \cdot n) &\rightarrow DFT(m), DFT(n) \\
 &\rightarrow Tensor(DFT(m), I(n)) \cdot T_m^{mn} \\
 &\quad Tensor(I(n), DFT(m)) \cdot L_n^{mn}
 \end{aligned}$$

## ■ Cooley-Tukey algorithm in $Z_p$

Primitive root of  $P$  used in twiddle factors instead of complex root of unity

$$\begin{aligned}
 ModDFT(m \cdot n, p, \omega) &\rightarrow ModDFT(m, p, \omega^m), ModDFT(n, p, \omega^n) \\
 &\rightarrow Tensor(ModDFT(m, p, \omega^m), I(n)) \cdot T_m^{mn}(\omega^m) \cdot \\
 &\quad Tensor(I(m), ModDFT(n, p, \omega^n)) \cdot L_n^{mn}
 \end{aligned}$$

# Modular FFTs: Unparser

- **Unparser**: translate symbolic arithmetic operations to corresponding C code
- **Unparser for Modular FFTs**
  - $\text{add}(a, b) \rightarrow (a + b) \% p$
  - $\text{sub}(a, b) \rightarrow (a - b) \% p$ 
    - All elements are positive and smaller than  $p$
    - $0 < a + b < 2p \rightarrow (a + b) \% p = (a + b < p) ? (a + b) : (a + b - p)$
    - $-p < a - b < p \rightarrow (a - b) \% p = (a - b < 0) ? (a - b + p) : (a - b)$ 
      - Faster than division though may introduce branches
      - Can use Conditional Move instead of branch (avoid pipeline stalls)
  - $\text{mul}(a, b) \rightarrow a * b \% p$

# Experiments

## ■ Platform

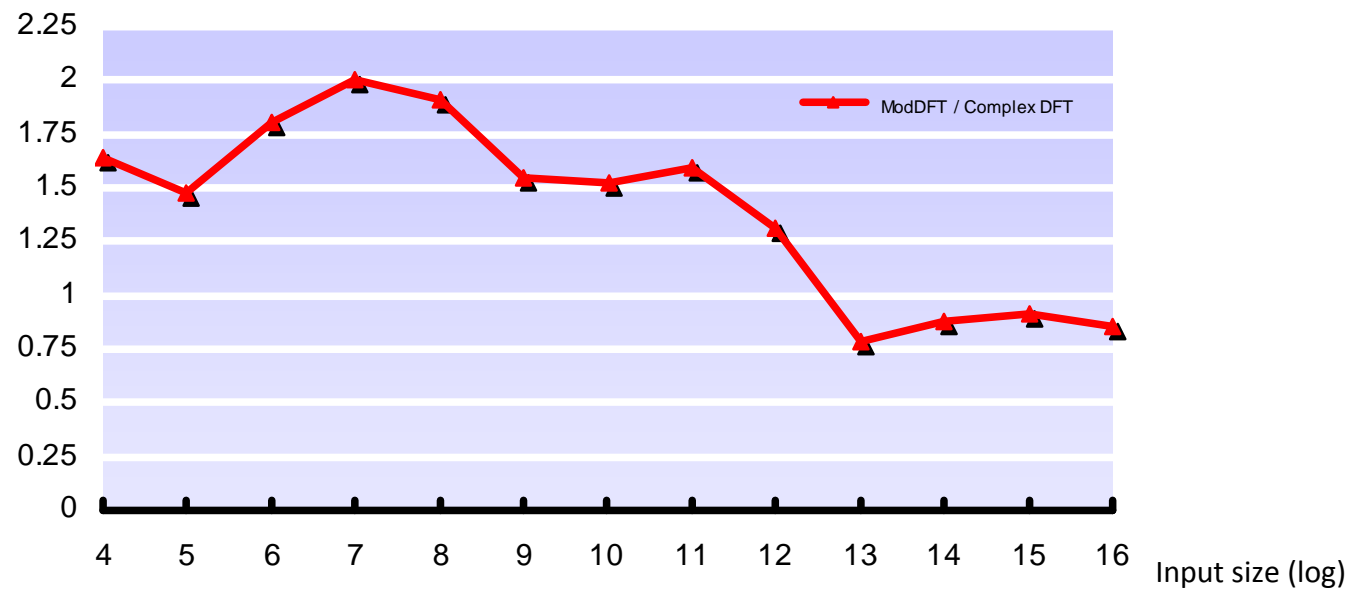
- AMD Phenom II X4 940 quad-core processor
  - **Operating Frequency:** 3.0 GHz
  - **L1 Cache:** 4 X 128 KB (64KB Instruction + 64KB Data)
  - **L2 Cache:** 4 X 512 KB
  - **L3 Cache:** 6 MB

## ■ Experiments

- Compare SPIRAL performance of modular FFT and complex floating point FFT
  - Optimal algorithms found may be different
- Compare SPIRAL modular FFT to numeric recipes (NR) code modified to perform modular DFT [NR uses iterative FFT]
- Compare optimal SPIRAL modular FFT found by DP to radix 2 recursive FFT

# SPRIAL Code Performance

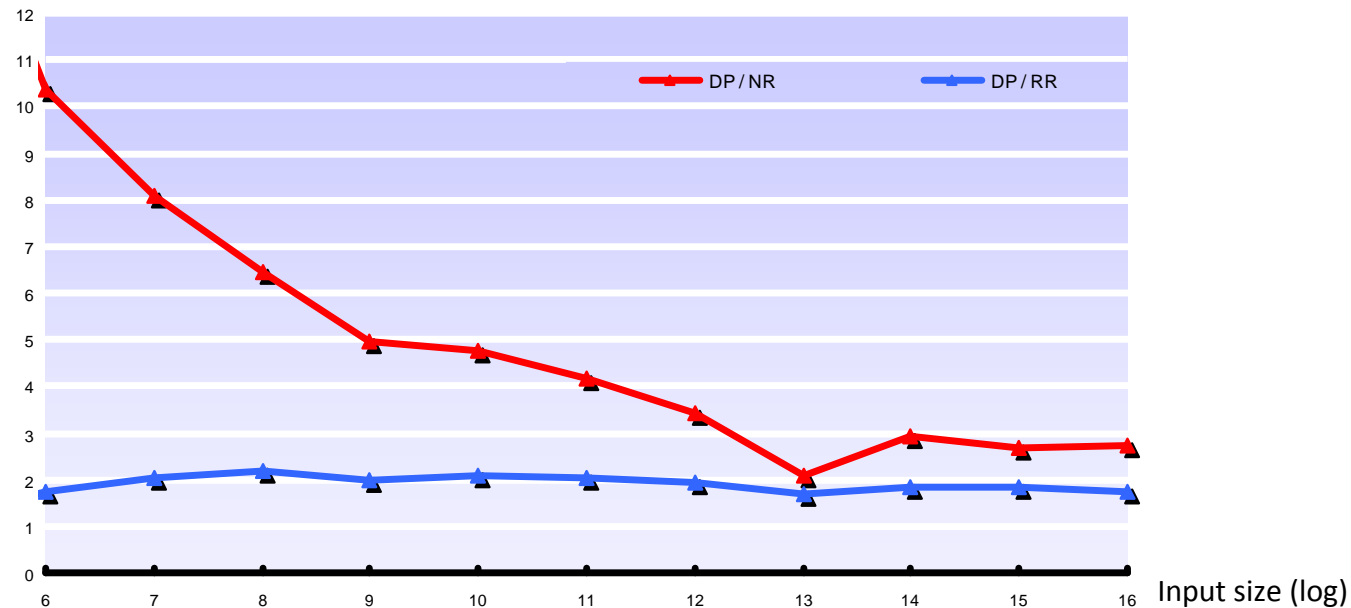
## Best Code found by DP Modular DFT vs. Complex DFT



- Expensive modular operations slow down ModDFT within the cache boundary
- Memory accesses dominate performance when input size crosses cache boundary ( $\geq 2^{13}$ )
  - Cache boundary reached one size earlier for complex DFT

# SPIRAL Code Performance

**Complex DFT Code**  
**Best SPIRAL code (DP) vs. Radix 2 Recursive Code & NR**

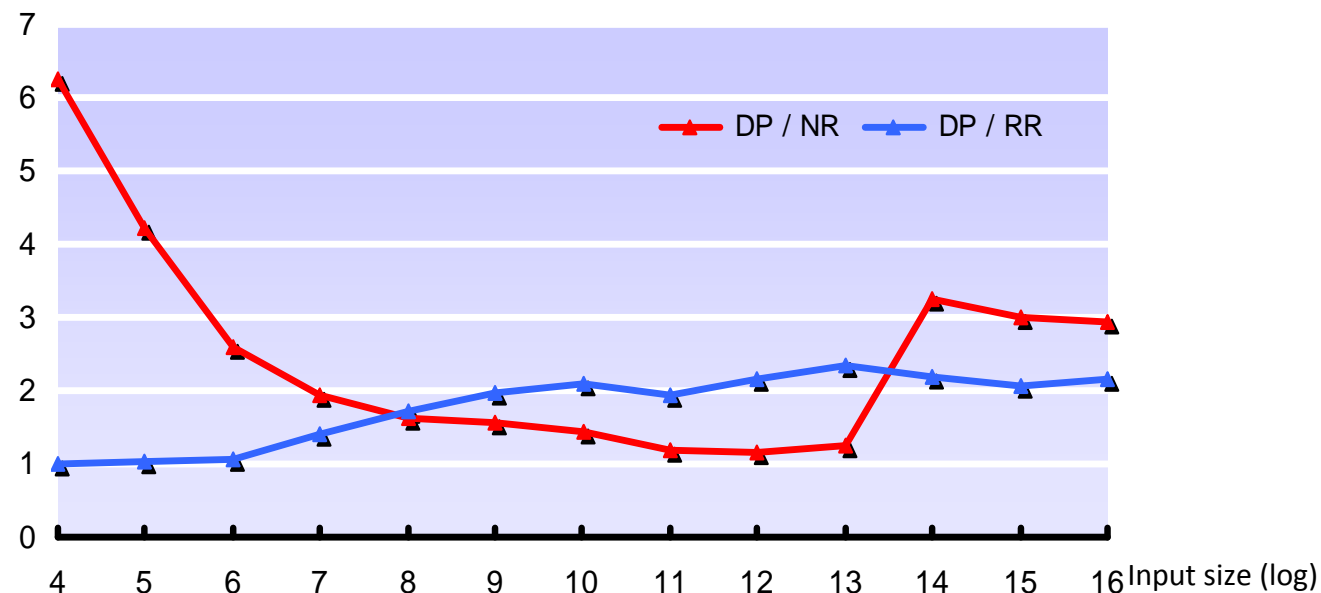


- **SPIRAL best code about twice faster than its right recursive code**
  - Better algorithm found by the search engine
- **SPIRAL best code averagely 4.8 times faster the numerical recipes code.**
  - Better cache/memory utilization by SPIRAL
- **Performance grouped into 3 regions (unrolled, in-cache, outside L1 cache)**

# SPIRAL Code Performance

## Modular FFT

### Best SPIRAL code (DP) vs. Radix 2 Recursive Code & NR



- **Best SPIRAL code twice as fast as right recursive code**
  - Better algorithm found by the search engine
- **Best SPIRAL code on average 3 times faster than NR code.**
  - Better cache/memory utilization by SPIRAL
- **Similar performance trends for modular and complex DFTs**
  - Larger jumps across unrolling and cache boundaries

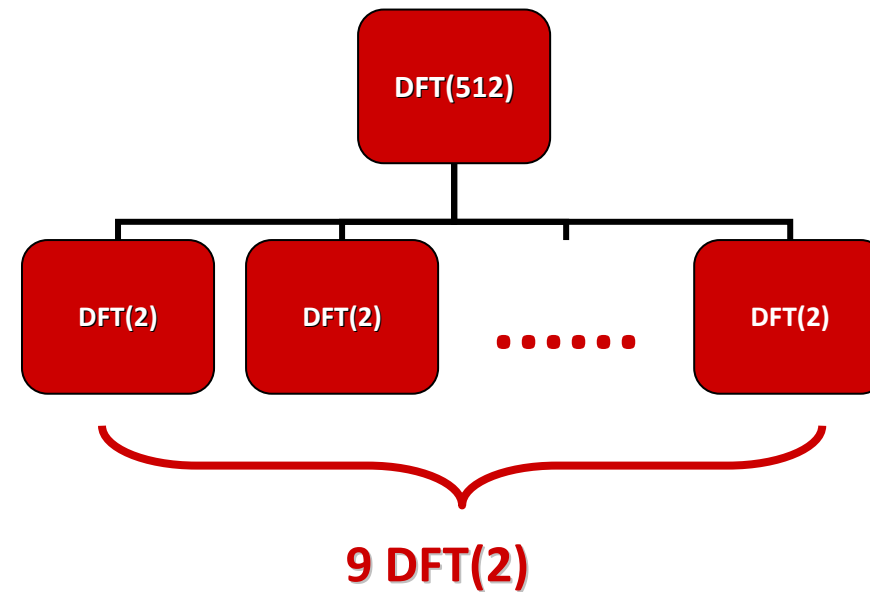
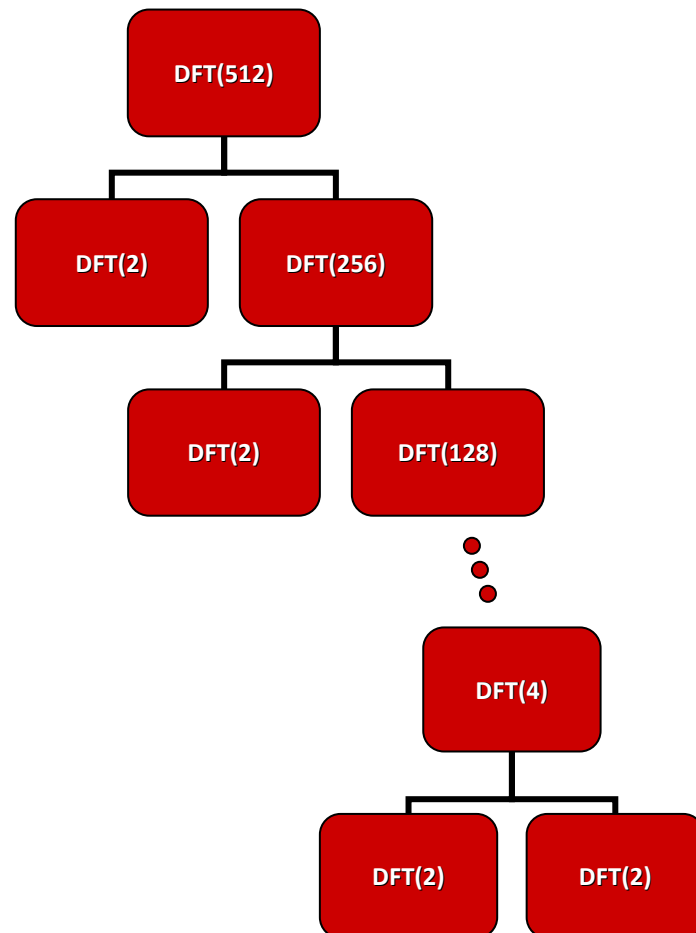
# Experimental Data

Input Size	Radix-2 Right Recursive		Dynamic Programming		Numerical Recipes	
	ModDFT	DFT	ModDFT	DFT	ModDFT	DFT
$2^3$	*107	73	106	68	1852	5202
$2^4$	271	183	271	165	3002	6073
$2^5$	701	559	684	465	5402	7636
$2^6$	2146	1950	2033	1130	10130	11713
$2^7$	7183	5113	5054	2528	19194	20382
$2^8$	19491	12769	11290	5909	39687	38049
$2^9$	46169	30097	23490	15214	81820	75729
$2^{10}$	109924	71651	52297	34368	164575	163198
$2^{11}$	256097	170499	132688	83046	363737	345289
$2^{12}$	611067	424499	284737	217446	800674	748075
$2^{13}$	1440588	1316808	613370	778572	5555759	1628003
$2^{14}$	3814025	3648241	1737673	1979379	13952347	5782618
$2^{15}$	9987370	9698445	4866960	5301481	34076665	14269816
$2^{16}$	24608317	22981890	11386662	13381989	86195220	36172945

\* Number of cycles on AMD Phenom II 940 3.0 GHz Quad-core processor

# Standard FFT Ruletrees

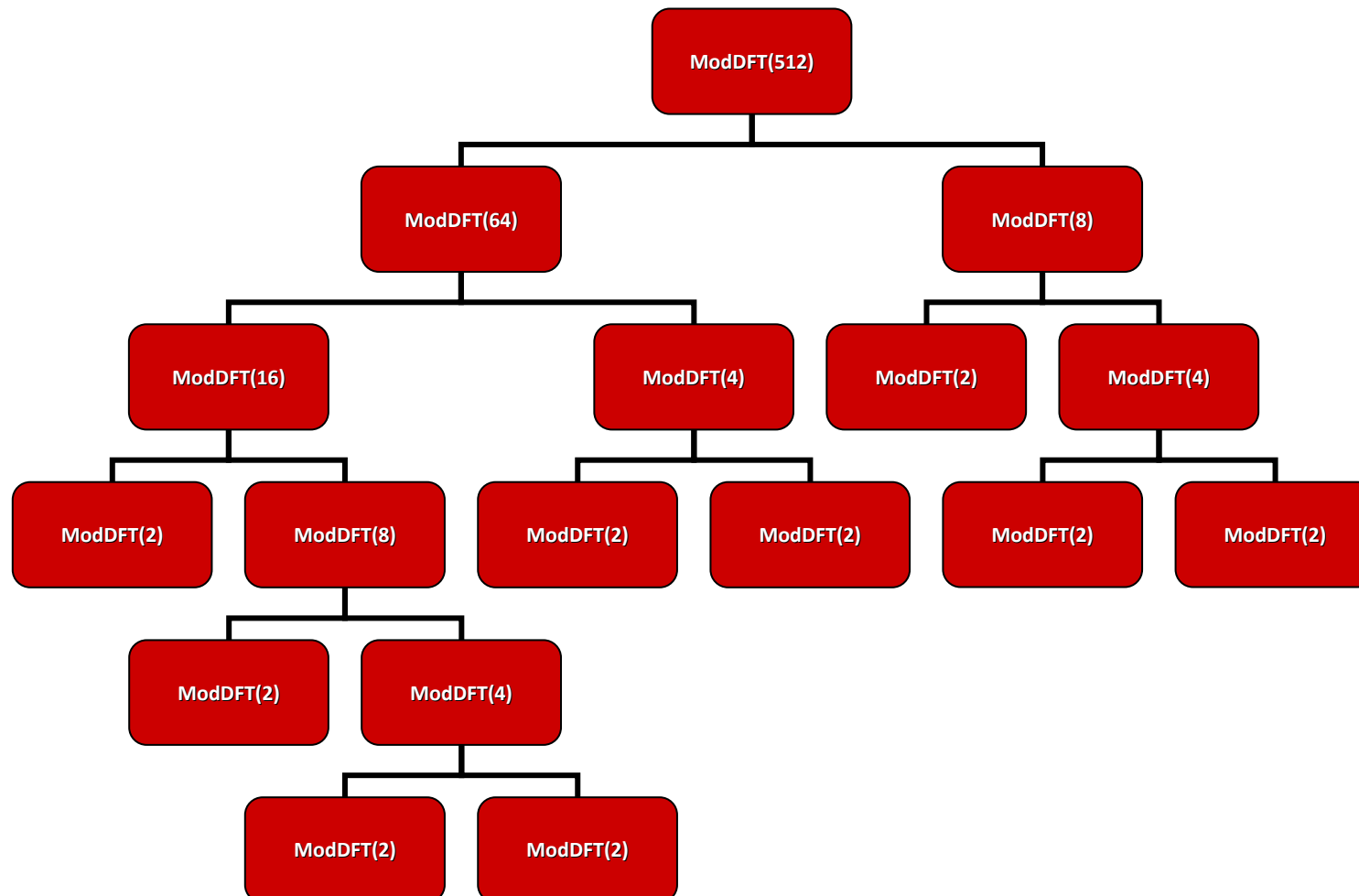
- Radix-2 Right Recursive Rule Tree
- Iterative Tree in Numerical Recipes





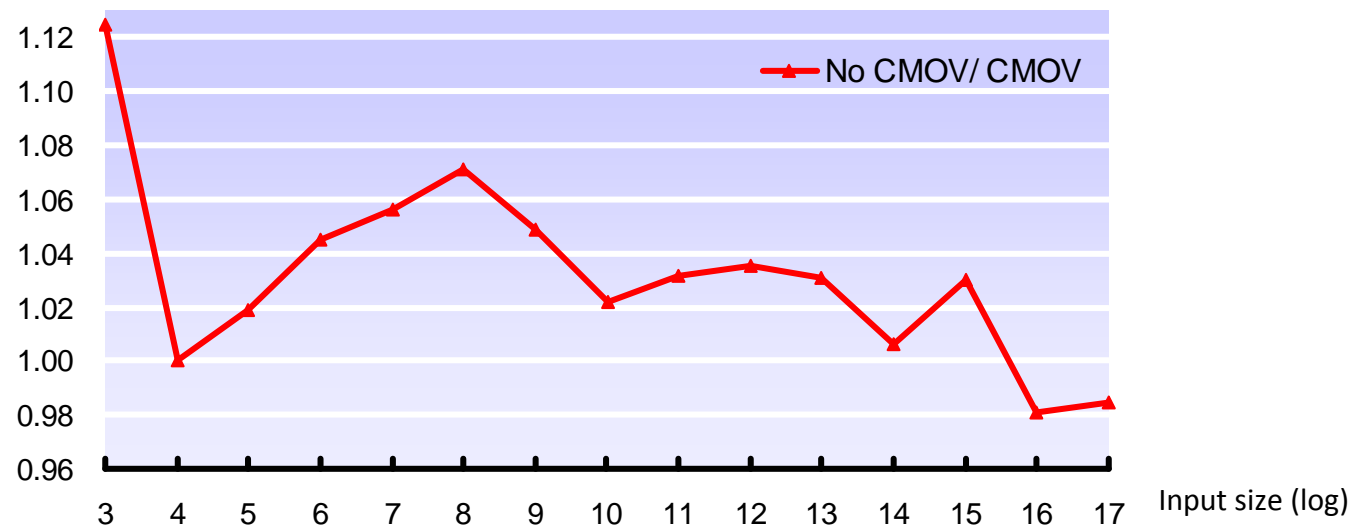
# Optimal Modular FFT Ruletree

- Best Rule Tree found by Dynamic Programming



# Conditional Move

## ■ Conditional Moves



- CMOVcc theoretically reduces pipeline penalty by converting branches to branchless codes in ternary operations of add/sub
- Experiments show improvement only for large sizes
  - Short branches of ternary operations optimized by CPU's multi-instruction fetch and conditional execution
  - Branch prediction table entries overwritten in large-sized programs
  - Instructions fetches insufficient for conditional execution of larger programs

# Future Work

- **Montgomery Algorithm [MATHCOMP 85]**
  - Modular multiplication without division
  - Enable vectorization when vector division not available
  
- **Vectorization**
  - Extend SPIRAL vectorization to support modular DFT
  
- **Multicore Parallelization**
  - Extend SPIRAL parallelization to support modular DFT

# Conclusion

- **Minor modifications to SPIRAL to support generation of efficient modular FFT code**
- **SPIRAL generated code 2-5 times faster than standard implementations**
  - Better memory utilization
  - Reduced loop/recursion overhead
- **Modular code 1.38 times slower than floating point code**
  - Beyond cache boundary modular code becomes faster
- **SPIRAL machinery can be used to obtain vectorized/parallel code**

# References

- **[CSSP 90]** J. R. Johnson, R. W Johnson, D. Rodriguez, and R. Tolimieri, **A methodology for designing, modifying, and implementing Fourier Transform algorithms on various architectures**, Circuits Systems Signal Process., (9)4: 449-500, 1990.
- **[IEEE 05]** 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 Nicholas Rizzolo, **SPIRAL: Code Generation for DSP Transforms**, Proc. of the IEEE, special issue on "**Program Generation, Optimization, and Adaptation**", Vol. 93, No. 2, pp. 232- 275, 2005.
- **[IPDPS 02]** Kang Chen and Jeremy Johnson, **A Prototypical Self-Optimizing Package for Parallel Implementation of Fast Signal Transforms**, Proc. International Parallel and Distributed Processing Symposium (IPDPS), pp. 58-63, 2002.
- **[SC 06]** Franz Franchetti, Yevgen Voronenko and Markus Püschel, **FFT Program Generation for Shared Memory: SMP and Multicore**, Proc. Supercomputing (SC), 2006.
- **[ISPA 06]** Andreas Bonelli, Franz Franchetti, Juergen Lorenz, Markus Püschel and Christoph W. Ueberhuber, **Automatic Performance Optimization of the Discrete Fourier Transform on Distributed Memory Computers**, Proc. International Symposium on Parallel and Distributed Processing and Application (ISPA), Lecture Notes In Computer Science, Springer, Vol. 4330, pp. 818-832, 2006.

# References

- **[DAC 08]** Peter A. Milder, Franz Franchetti, James C. Hoe and Markus Püschel, **Formal Datapath Representation and Manipulation for Implementing DSP Transforms**, Proc. Design Automation Conference (DAC), 2008.
- **[IPDPS 04]** Jeremy Johnson and Kang Chen, **A Self-Adapting Distributed Memory Package for Fast Signal Transforms**, Proc. International Parallel and Distributed Processing Symposium (IPDPS), pp. 44-, 2004.
- **[PLDI 01]** Jianxin Xiong, Jeremy Johnson, Robert W. Johnson and David Padua, **SPL: A Language and Compiler for DSP Algorithms**, Proc. Programming Languages Design and Implementation (PLDI), pp. 298-308, 2001.
- **[PLDI 05]** Franz Franchetti, Yevgen Voronenko and Markus Püschel, **Formal Loop Merging for Signal Transforms**, Proc. Programming Languages Design and Implementation (PLDI), pp. 315-326 , 2005.
- **[MATHCOMP 85]** Peter Montgomery, **Modular Multiplication Without Trial Division**, Mathematics of Computation (MATHCOMP), Vol. 44, pp. 519-521, 1985.
- **[VECPAR 06]** Franz Franchetti, Yevgen Voronenko and Markus Püschel, **A Rewriting System for the Vectorization of Signal Transforms**, Proc. High Performance Computing for Computational Science (VECPAR), Lecture Notes in Computer Science, Springer, Vol. 4395, pp. 363-377, 2006