# Formal Loop Merging for Signal Transforms

**Franz Franchetti**

Yevgen S. Voronenko

Markus Püschel

Department of Electrical & Computer Engineering

Carnegie Mellon University

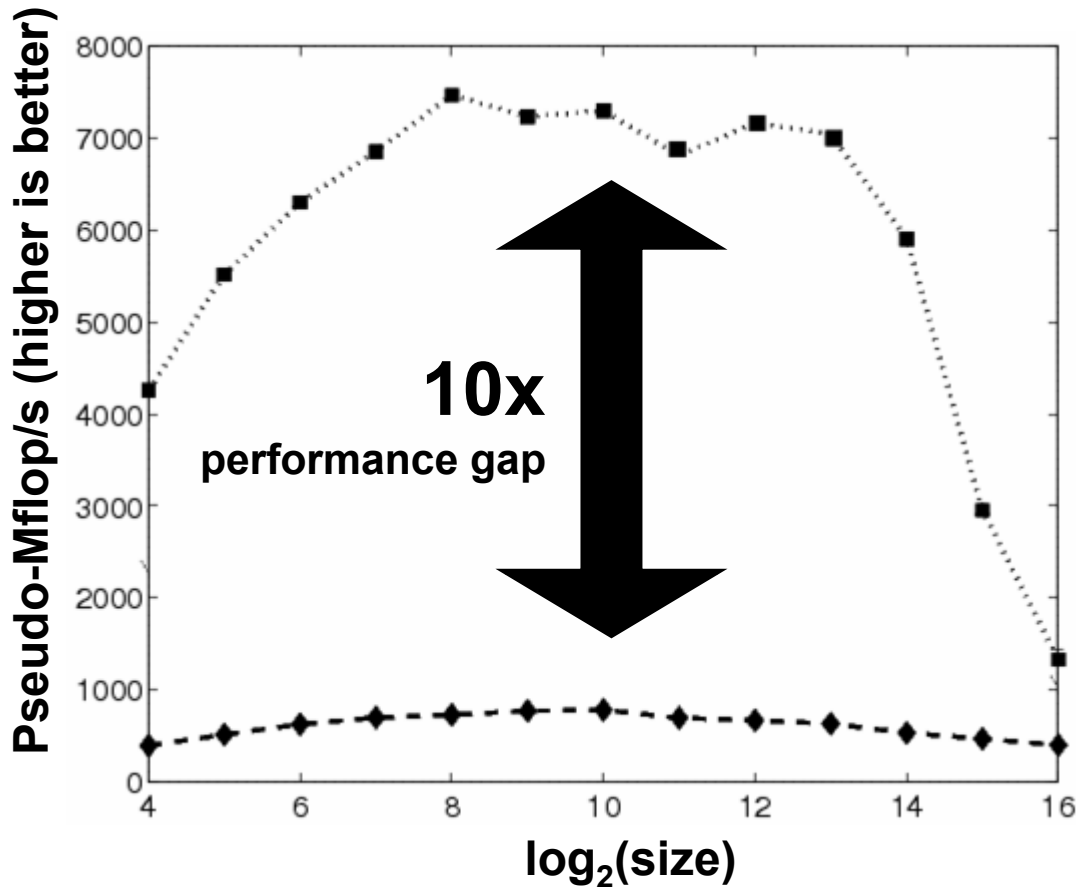SPIRAL

`http://www.spiral.net`

# Problem

- **Runtime of (uniprocessor) numerical applications typically dominated by few compute-intensive kernels**
    - Examples: discrete Fourier transform, matrix-matrix multiplication

- **These kernels are hand-written for every architecture (open-source and commercial libraries)**

- **Writing fast numerical code is becoming increasingly difficult, expensive, and platform dependent, due to:**
    - Complicated memory hierarchies

    - Special purpose instructions
      (short vector extensions, fused multiply-add)

    - Other microarchitectural features
      (deep pipelines, superscalar execution)

# Example: Discrete Fourier Transform (DFT)

**Performance on Pentium 4 @ 3 GHz**



**vendor library**

**roughly the same operations count**

**GNU Scientific Library**

**Writing fast code is hard.  Are there alternatives?**

# Automatic Code Generation and Adaptation

- **ATLAS:** Code generator for basic linear algebra subroutines (BLAS)
  [Whaley, et. al., 1998] [Yotov, et al., 2005]

- **FFTW:** Adaptive library for computing the discrete Fourier transform (DFT) and its variants
  [Frigo and Johnson, 1998]

- **SPIRAL:** Code generator for linear signal transforms (including DFT)
  [Püschel, et al., 2004]

- **See also:** Proceedings of the IEEE special issue on "Program Generation, Optimization, and Adaptation," Feb. 2005.

- **Focus of this talk:**
  A new approach to automatic loop merging in SPIRAL

SPIRAL

# Talk Organization

- **SPIRAL Background**

- **Automatic loop merging in SPIRAL**

- **Experimental Results**

- **Conclusions**

# Talk Organization

- **SPIRAL Background**

- **Automatic loop merging in SPIRAL**

- **Experimental Results**
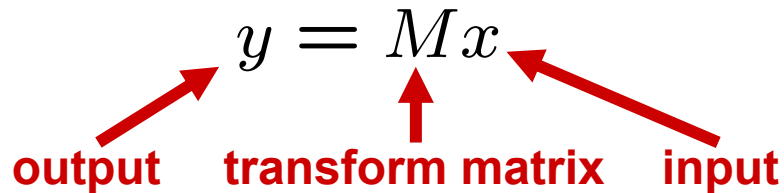
- **Conclusions**

SPIRAL

# SPIRAL: DSP Transforms

- **SPIRAL generates optimized code for linear signal transforms, such as discrete Fourier transform (DFT), discrete cosine transforms, FIR filters, wavelets, and many others.**

- **Linear transform = matrix-vector product:**

$$y = Mx$$

**output**     **transform matrix**     **input**

- **Example: DFT of input vector $x$**

$$y = \mathrm{DFT}_n\, x$$

$$\mathrm{DFT}_n = \left[\omega_n^{k\ell}\right]_{0 \le k,\ell < n}, \quad \omega_n = e^{-2\pi\sqrt{-1}/n}$$

# SPIRAL: Fast Transform Algorithms

- **Reduce computation cost from $O(n^2)$ to $O(n \log n)$**

- **For every transform there are many fast algorithms**

- **Algorithm = sparse matrix factorization**

$$\begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & i & -1 & -i \\ 1 & -1 & 1 & -1 \\ 1 & -i & -1 & i \end{pmatrix} x \rightarrow \begin{pmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & -1 & 0 \\ 0 & 1 & 0 & -1 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & i \end{pmatrix} \begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & -1 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} x$$

**12 adds**　　　　　　**4 adds**　　　**1 mult**　　　**4 adds**
**4 mults**　　　　*(when multiplied with input vector $x$)*

$$\mathsf{DFT}_4 \rightarrow (\mathsf{DFT}_2 \otimes \mathsf{I}_2) D (\mathsf{I}_2 \otimes \mathsf{DFT}_2) P$$

- **SPIRAL generates the space of algorithms using breakdown rules in the domain-specific Signal Processing Language (SPL)**

$$\mathsf{DFT}_{mn} \rightarrow (\mathsf{DFT}_m \otimes \mathsf{I}_n) D (\mathsf{I}_m \otimes \mathsf{DFT}_n) P$$

# SPL (Signal Processing Language)

■ **SPL expresses transform algorithms as structured sparse matrix factorization**

■ **Examples:**

$$F_2 = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \qquad\qquad A \otimes B = \begin{bmatrix} a_{k,\ell} B \end{bmatrix}_{k,\ell}$$

$$A \oplus B = \begin{bmatrix} A & \\ & B \end{bmatrix} \qquad\qquad I \otimes B = \begin{bmatrix} B & & \\ & \ddots & \\ & & B \end{bmatrix}$$

■ **SPL grammar in Backus-Naur form**

$$
\begin{aligned}
\langle \text{spl} \rangle \quad ::= \quad & \langle \text{generic} \rangle \mid \langle \text{symbol} \rangle \mid \langle \text{transform} \rangle \mid \\
& \langle \text{spl} \rangle \cdots \cdot \langle \text{spl} \rangle \mid \\
& \langle \text{spl} \rangle \oplus \ldots \oplus \langle \text{spl} \rangle \mid \\
& \langle \text{spl} \rangle \otimes \cdots \otimes \langle \text{spl} \rangle \mid \\
& \cdots \\
\langle \text{generic} \rangle \quad ::= \quad & \text{diag}(a_0, \ldots, a_{n-1}) \mid \ldots \\
\langle \text{symbol} \rangle \quad ::= \quad & I_n \mid J_n \mid L_k^n \mid R_\alpha \mid F_2 \mid \ldots \\
\langle \text{transform} \rangle \quad ::= \quad & \mathbf{DFT}_n \mid \mathbf{WHT}_n \mid \mathbf{DCT\text{-}2}_n \mid \mathbf{Filt}_n(h[z]) \mid \ldots
\end{aligned}
$$

# Compiling SPL to Code Using Templates

$$y = \mathsf{L}_n^{mn}\, x$$

```
for i=0..n-1
    for j=0..m-1
        y[i+n*j]=x[m*i+j]
```

$$y = (A_n \oplus B_m)x$$

```
y[0:1:n-1] = call A(x[0:1:n-1])
y[n:1:n+m-1] = call B(x[n:1:n+m-1])
```
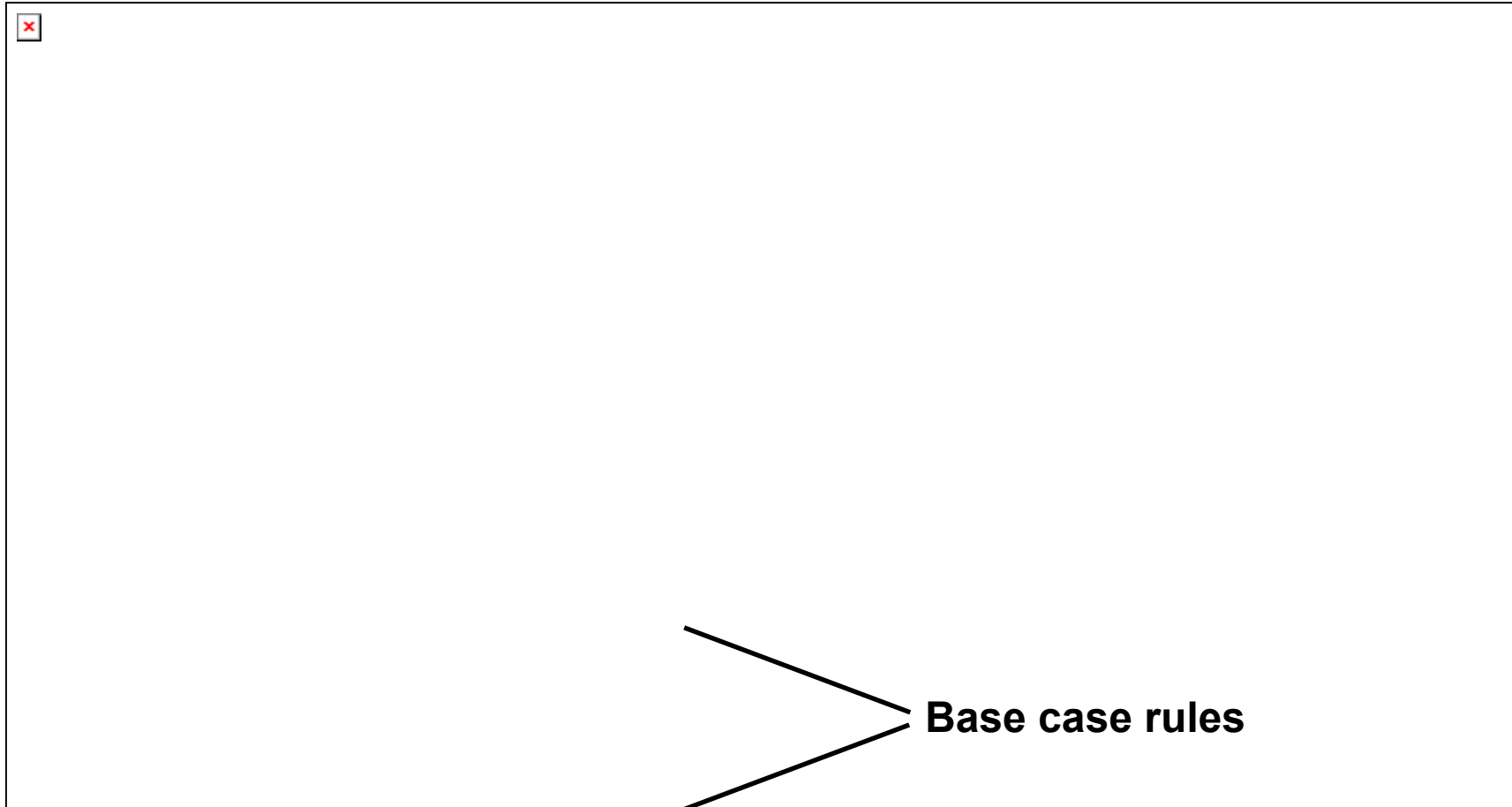
$$y = (\mathrm{I}_n \otimes B_m)x$$

```
for i=0..n-1
    y[im:1:im+m-1] = call B(x[im:1:im+m-1])
```

$$y = (\mathrm{I}_n \otimes B_m)\, \mathsf{L}_n^{mn}\, x$$

```
for i=0..n-1
    y[im:1:im+m-1] = call B(x[i:n:i+m-1])
```
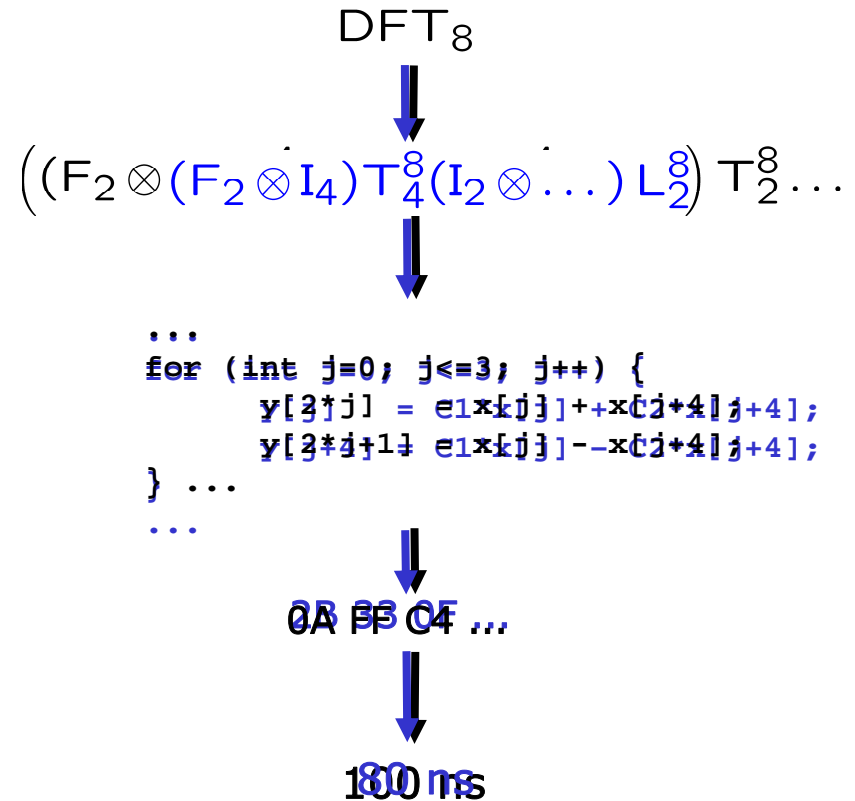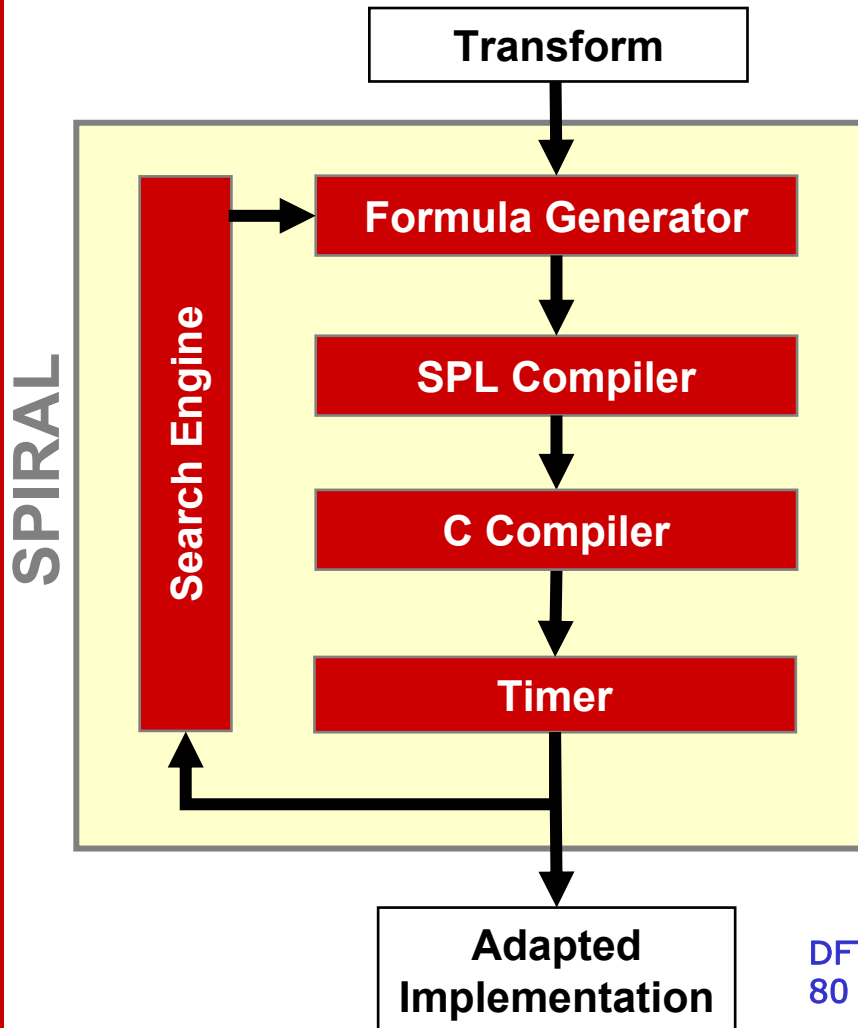
# Some Transforms and Breakdown Rules in SPIRAL



**Base case rules**

**Spiral contains 30+ transforms and 100+ rules**

# SPIRAL Architecture

**Approach:** Empirical search over alternative recursive algorithms



$$DFT_8$$

$$\left(\left(F_2 \otimes (F_2 \otimes I_4) T_4^8 (I_2 \otimes \ldots) L_2^8\right) T_2^8 \ldots\right)$$

```
...
for (int j=0; j<=3; j++) {
    y[2*j] = x[j]+x[j+4];
    y[2*j+1] = x[j]-x[j+4];
} ...
...
```

0A FF C4 ...

80 ns

DFT_8.c
80 ns

# Problem: Fusing Permutations and Loops

$$y = (I_4 \otimes F_2) L_4^8 x$$

**direct mapping**

*hardcoded special case*

State-of-the-art
**SPIRAL:** Hardcoded with templates
**FFTW:** Hardcoded in the infrastructure

**How does hardcoding scale?**

Two passes over the working set
Complex index computation

```
void sub(double *y, double *x) {
   double t[8];
   for (int i=0; i<=7; i++)
      t[(i/4)+2*(i%4)] = 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];
   }
}
```

**C compiler cannot do this**

One pass over the working set
Simple index computation

```
void sub(double *y, double *x) {
   for (int j=0; j<=3; j++){
      y[2*j] = x[j] + x[j+4];
      y[2*j+1] = x[j] - x[j+4];
   }
}
```

13

# General Loop Merging Problem

$\bigcirc$ **= permutations**

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

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

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

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

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

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

- **Combinatorial explosion: Implementing templates for all rules and all recursive combinations is unfeasible**

- **In many cases even theoretically not understood**

SPIRAL

14

# Our Solution in SPIRAL

- **Loop merging at C code level: impractical**

- **Loop merging at SPL level: not possible**

- **Solution:**
  - **New language $\Sigma$-SPL – an abstraction level between SPL and code**

  - **Loop merging through $\Sigma$-SPL formula manipulation**
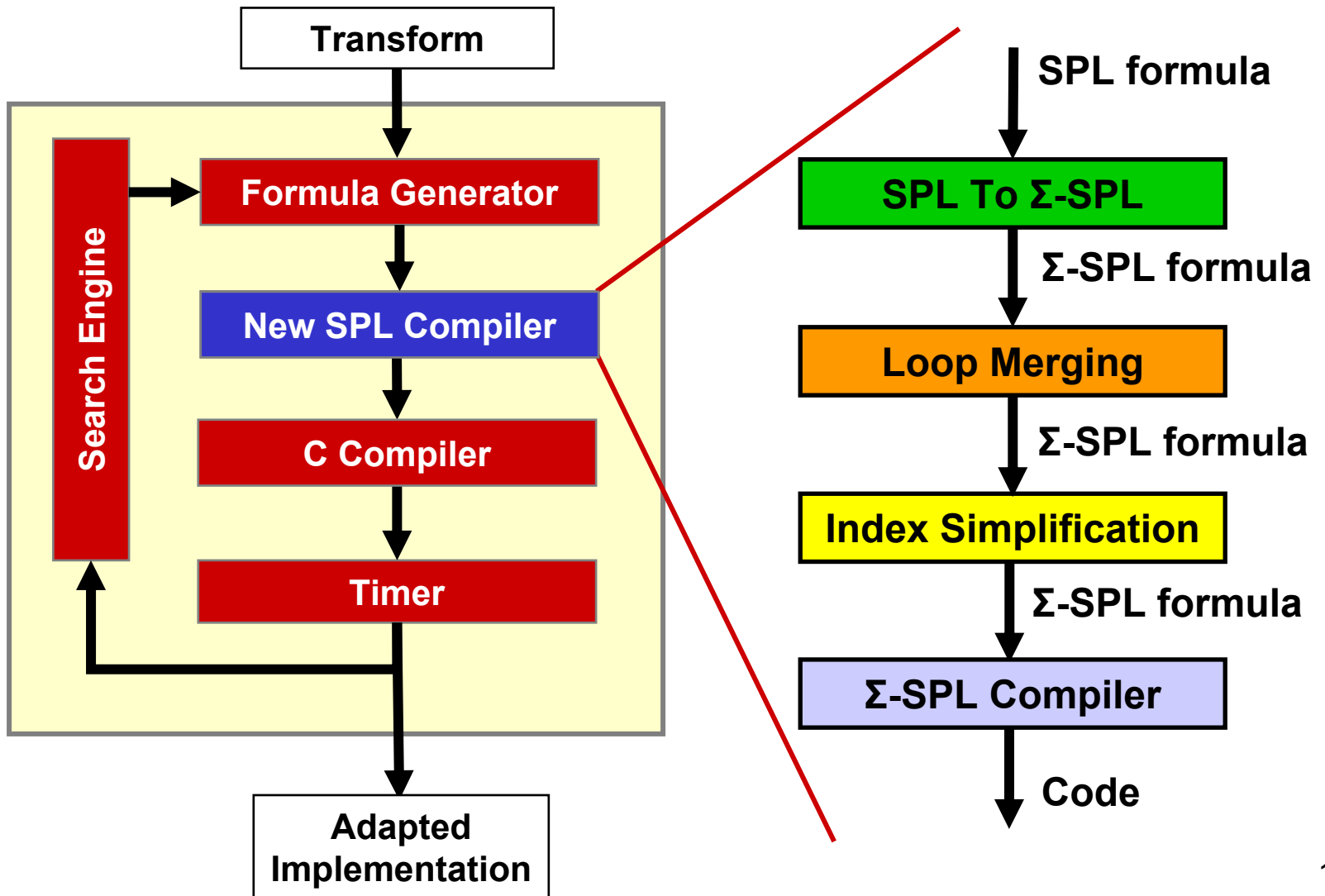
# Talk Organization

■ **SPIRAL Background**

■ **Automatic loop merging in SPIRAL**

■ **Experimental Results**

■ **Conclusions**

# New Approach for Loop Merging

Transform

Search Engine

Formula Generator

New SPL Compiler

C Compiler

Timer

Adapted Implementation

SPL formula

SPL To Σ-SPL

Σ-SPL formula

Loop Merging

Σ-SPL formula

Index Simplification

Σ-SPL formula
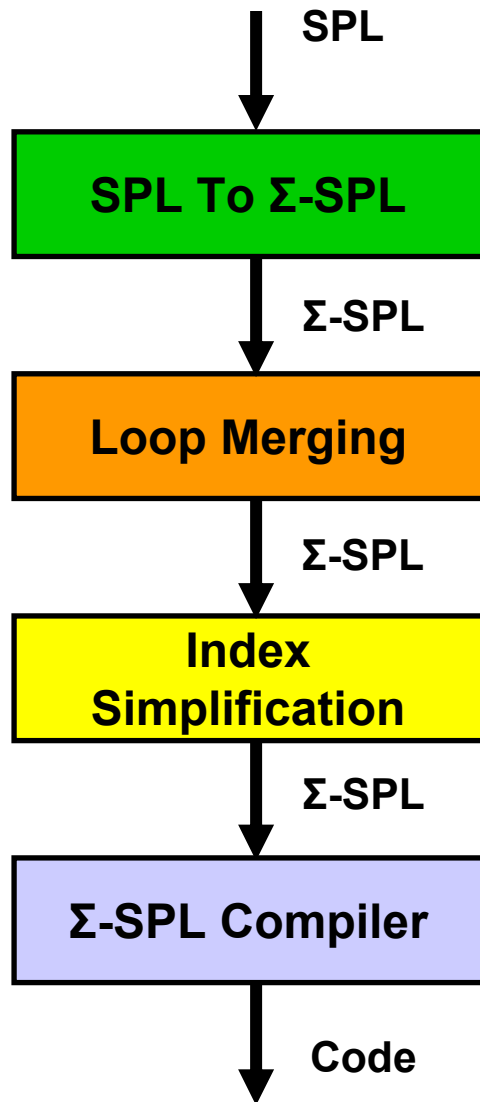
Σ-SPL Compiler

Code

SPIRAL

17

# Σ–SPL

■ **Four central constructs: Σ, G, S, Perm**
  - Σ **(sum)** – **makes loops explicit**
  - $G_f$ **(gather)** – **reads data using the index mapping** $f$
  - $S_f$ **(scatter)** – **writes data using the index mapping** $f$
  - $Perm_f$ – **permutes data using the index mapping** $f$

■ **Every Σ-SPL formula still represents a matrix factorization**

**Example:** $(I_4 \otimes F_2) \rightarrow \sum_{j=0}^{3} S_{f_j} F_2 G_{f_j}$

$$\begin{bmatrix} F_2 & & & \\ & F_2 & & \\ & & F_2 & \\ & & & F_2 \end{bmatrix}$$

$F_2$

**Output**          **Input**

j=0

j=1

j=2

j=3

18

# Loop Merging With Rewriting Rules

**SPL**

| |
|---|
| **SPL To Σ-SPL** |

**Σ-SPL**

| |
|---|
| **Loop Merging** |

**Σ-SPL**

| |
|---|
| **Index Simplification** |

**Σ-SPL**

| |
|---|
| **Σ-SPL Compiler** |

**Code**

$$y = (\mathsf{I}_4 \otimes \mathsf{F}_2)\,\mathsf{L}_4^8\,x$$

$$\left( \sum_{j=0}^{3} \mathsf{S}_{(j)_4 \otimes \imath_2}\,\mathsf{F}_2\,\mathsf{G}_{(j)_4 \otimes \imath_2} \right) \mathrm{Perm}_{\ell_4^8}$$

$$\left( \sum_{j=0}^{3} \mathsf{S}_{(j)_4 \otimes \imath_2}\,\mathsf{F}_2\,\mathsf{G}_{\ell_4^8 \circ ((j)_4 \otimes \imath_2)} \right)$$

$$\left( \sum_{j=0}^{3} \mathsf{S}_{(j)_4 \otimes \imath_2}\,\mathsf{F}_2\,\mathsf{G}_{\imath_2 \otimes (j)_4} \right)$$

```
for (int j=0; j<=3; j++) {
    y[2*j]   = x[j] + x[j+4];
    y[2*j+1] = x[j] - x[j+4];
}
```



Y        T        X

**F₂**

Y        X

**F₂**

**Rules:**

$$\mathsf{G}_r\,\mathrm{Perm}_p \;\rightarrow\; \mathsf{G}_{p \circ r}$$

$$\ell_m^{mn} \circ \big((j)_m \otimes \imath_n\big) \;\rightarrow\; \imath_n \otimes (j)_m$$

19

# Application: Loop Merging For FFTs

**DFT breakdown rules:**

**Cooley-Tukey FFT** $\quad \mathbf{DFT}_{km} \;\;\rightarrow\;\; (\mathbf{DFT}_k \otimes \mathbf{I}_m)\, \mathsf{T}_m^{km}(\mathbf{I}_k \otimes \mathbf{DFT}_m)\, \mathsf{L}_k^{km}$

**Prime factor FFT** $\quad \mathbf{DFT}_{km} \;\;\rightarrow\;\; \mathsf{V}_{k,m}^T (\mathbf{DFT}_k \otimes \mathbf{I}_m)(\mathbf{I}_k \otimes \mathbf{DFT}_m)\, \mathsf{V}_{k,m}$
$$\gcd(k,m) = 1$$

**Rader FFT** $\quad \mathbf{DFT}_p \;\;\rightarrow\;\; \mathsf{W}_p^T (\mathbf{I}_1 \oplus \mathbf{DFT}_{p-1}) D_p (\mathbf{I}_1 \oplus \mathbf{DFT}_{p-1})\, \mathsf{W}_p$
$$p \text{ - prime}$$

**Index mapping functions are non-trivial:**

$$\mathsf{L}_k^{km} \;\rightarrow\; \mathrm{Perm}_{\ell_k^{km}} \qquad\qquad \ell_k^{km}(i) \;=\; \left\lfloor \frac{i}{m} \right\rfloor + k(i \bmod m)$$

$$\mathsf{V}_{k,m} \;\rightarrow\; \mathrm{Perm}_{v_{k,m}} \qquad\qquad v_{k,m}(i) = \left( m\left\lfloor \frac{i}{m} \right\rfloor + k(i \bmod m) \right) \bmod km$$

$$\mathsf{W}_p \;\rightarrow\; \mathrm{Perm}_{w_{1,g}^p} \qquad\qquad w_{\phi,g}^p(i) \;=\; \begin{cases} 0, & i = 0, \\ \phi g^i \;\bmod\; p, & \text{else.} \end{cases}$$

# Example

## Given DFT$_{pq}$

p – prime
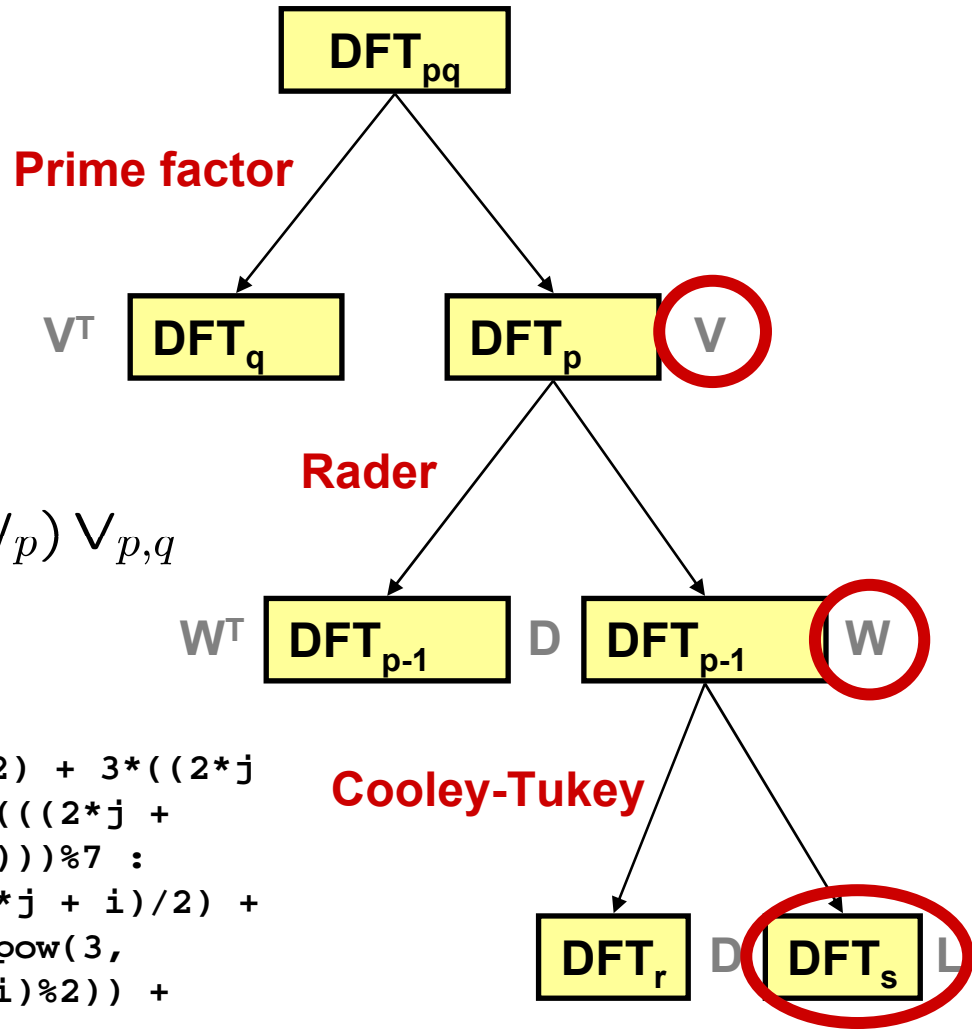p-1 = rs

**Prime factor**

**Rader**

**Cooley-Tukey**

**Formula fragment**

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

**Code for one memory access**

```
p=7; q=4; r=3; s=2;
t=x[((21*((7*k + (((((2*j + i)/2) + 3*((2*j
   + i)%2)) + 1)) ? (5*pow(3, ((((2*j +
   i)/2) + 3*((2*j + i)%2)) + 1)))%7 :
   (0)))/7) + 8*((7*k + (((((2*j + i)/2) +
   3*((2*j + i)%2)) + 1)) ? (5*pow(3,
   ((((2*j + i)/2) + 3*((2*j + i)%2)) +
   1)))%7 : (0)))%7))%28)];
```

DFT$_{pq}$

V$^T$  DFT$_q$     DFT$_p$  V

W$^T$  DFT$_{p-1}$  D  DFT$_{p-1}$  W

DFT$_r$  D  DFT$_s$  L

**Task: Index simplification**

SPIRAL

# Index Simplification: Basic Idea

- **Example:** Identity necessary for fusing successive Rader and prime-factor step

$$\left(\varphi g^{(b+si) \bmod N'}\right) \bmod N = \left((\varphi g^b)(g^s)^i\right) \bmod N$$

$$s|N', \ N'|N, \ 0 \le i < n$$

- **Performed at the $\Sigma$-SPL level through rewrite rules on function objects:**

$$\overline{w}_{\phi,g}^{N' \to N} \circ \overline{h}_{b,s}^{n \to N'} \quad \longrightarrow \quad \overline{w}_{\phi g^b, g^s}^{n \to N}$$

- **Advantages:**
  - no analysis necessary
  - efficient (or doable at all)

# Index Simplification Rules for FFTs

**Cooley-Tukey**

$$\ell_m^{mn} \circ \left( (j)_m \otimes f^{k \to n} \right) \to f^{k \to n} \otimes (j)_m$$

$$\left( f^{1 \to m} \otimes h \right) \circ g \to f \otimes (h \circ g)$$

$$\left( h \otimes g^{1 \to n} \right) \circ f \to (h \circ f) \otimes g$$

$$(f_0 \otimes f_1) \circ (g_0 \otimes g_1) \to (f_0 \circ g_0) \otimes (f_1 \circ g_1)$$

**Transitional**

$$\iota_n \to h_{0,1}^{n \to n}$$

$$f^{m \to M} \otimes g^{1 \to N} \to h_{g(0),N}^{M \to MN} \circ f$$

$$g^{1 \to N} \otimes f^{m \to M} \to h_{Mg(0),1}^{M \to MN} \circ f$$

**Cooley-Tukey + Prime factor**

$$v_{r,s} \circ h_{b,1}^{s \to rs} \to \overline{h}_{b,r}^{s \to rs}$$

$$h_{b_1,s_1}^{nk \to mnk} \circ h_{b_2,s_2}^{n \to nk} \to h_{b_1+s_1 b_2, s_1 s_2}^{n \to mnk}$$

$$\overline{h}_{b_1,s_1}^{nk \to mnk} \circ h_{b_2,s_2}^{n \to nk} \to \overline{h}_{b_1+s_1 b_2, s_1 s_2}^{n \to mnk}$$

$$\overline{h}_{b_1,s_1}^{nk \to mnk} \circ \overline{h}_{b_2,s_2}^{n \to nk} \to \overline{h}_{b_1+s_1 b_2, s_1 s_2}^{n \to mnk}$$

**Cooley-Tukey + Prime factor + Rader**

$$w_{\phi,g}^N \circ (0)_+^{1 \to N} \to (0)_+^{1 \to N}$$

$$w_{\phi,g}^N \circ (N-1)_+^{N-1 \to N} \to \overline{w}_{\phi,g}^{N-1 \to N}$$

$$\overline{w}_{\phi,g}^{N' \to N} \circ h_{b,s}^{n \to N'} \to \overline{w}_{\phi g^b, g^s}^{n \to N}$$

$$\overline{w}_{\phi,g}^{N' \to N} \circ \overline{h}_{b,s}^{n \to N'} \to \overline{w}_{\phi g^b, g^s}^{n \to N}$$

**These 15 rules cover all combinations.
Some encode novel optimizations.**

SPIRAL

23

# Loop Merging For the FFTs : Example (cont'd)

**SPL**

**SPL To Σ-SPL**

**Σ-SPL**

**Loop Merging**

**Σ-SPL**

**Index Simplification**

**Σ-SPL**

**Σ-SPL Compiler**

**Code**

$$\left(I_p \otimes \left(I_1 \oplus \left(I_r \otimes \mathsf{DFT}_s\right) \mathsf{L}_r^{rs}\right) \mathsf{W}_p\right) \mathsf{V}_{p,q}$$

$$\sum_{j_1=0}^{p-1} \left( \mathsf{S}_{\left((j_1)_p \otimes \imath_q\right) \circ (0)_+^{1 \to q}} \mathsf{G}_{v^{p,q} \circ \left((j_1)_p \otimes \imath_q\right) \circ w_{1,g}^q \circ (0)_+^{1 \to q}} \right.$$

$$+ \sum_{j_0=0}^{r-1} \mathsf{S}_{\left((j_1)_p \otimes \imath_q\right) \circ (1)_+^{q-1 \to q} \circ ((j_0)_r \otimes \imath_s)} \mathsf{DFT}_s$$

$$\left. \mathsf{G}_{v^{p,q} \circ \left((j_1)_p \otimes \imath_q\right) \circ w_{1,g}^q \circ (1)_+^{q-1 \to q} \circ \ell_r^{rs} \circ ((j_0)_r \otimes \imath_s)} \right)$$

$$\sum_{j_1=0}^{p-1} \left( \mathsf{S}_{h_{0,q}^{p \to pq}(j_1)} \mathsf{G}_{\overline{h}_{0,q}^{p \to pq}(j_1)} + \sum_{\substack{j_0=0 \\ \phi_1 = g^{j_0}}}^{r-1} \mathsf{S}_{h_{qj_1+sj_0+1,1}^{s \to pq}} \mathsf{DFT}_s \, \mathsf{G}_{\overline{h}_{b_1,p}^{q \to pq} \circ \overline{w}_{\phi_1,g^s}^{s \to q}} \right)$$

```
// I
in
fo b1=qj1
   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);
   }
}
```

24

```
// Input: _Complex double x[28], output: y[28]
double t1[28];
for(int i5 = 0; i5 <= 27; i5++)
    t1[i5] = x[(7*3*(i5/7) + 4*2*(i5%7))%28];
for(int i1 = 0; i1 <= 3; i1++) {
    double t3[7], t4[7], t5[7];
    for(int i6 = 0; i6 <= 6; i6++)
        t5[i6] = t1[7*i1 + i6];
    for(int i8 = 0; i8 <= 6; i8++)
        t4[i8] = t5[i8 ? (5*pow(3, i8))%7 : 0];
    {
        double t7[1], t8[1];
        t8[0] = t4[0];
        t7[0] = t8[0];
        t3[0] = t7[0];
    }
    {
        double t10[6], t11[6], t12[6];
        for(int i13 = 0; i13 <= 5; i13++)
            t12[i13] = t4[i13 + 1];
        for(int i14 = 0; i14 <= 5; i14++)
            t11[i14] = t12[(i14/2) + 3*(i14%2)];
        for(int i3 = 0; i3 <= 2; i3++) {
            double t14[2], t15[2];
            for(int i15 = 0; i15 <= 1; i15++)
                t15[i15] = t11[2*i3 + i15];
            t14[0] = (t15[0] + t15[1]);
            t14[1] = (t15[0] - t15[1]);
            for(int i17 = 0; i17 <= 1; i17++)
                t10[2*i3 + i17] = t14[i17];
        }
        for(int i19 = 0; i19 <= 5; i19++)
            t3[i19 + 1] = t10[i19];
    }
    for(int i20 = 0; i20 <= 6; i20++)
        y[7*i1 + i20] = t3[i20];
}
```

```
// 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);
    }
}
```

## After, 2 Loops.

## Before, 11 Loops.

# Talk Organization

- **SPIRAL Background**

- **Automatic loop merging in SPIRAL**

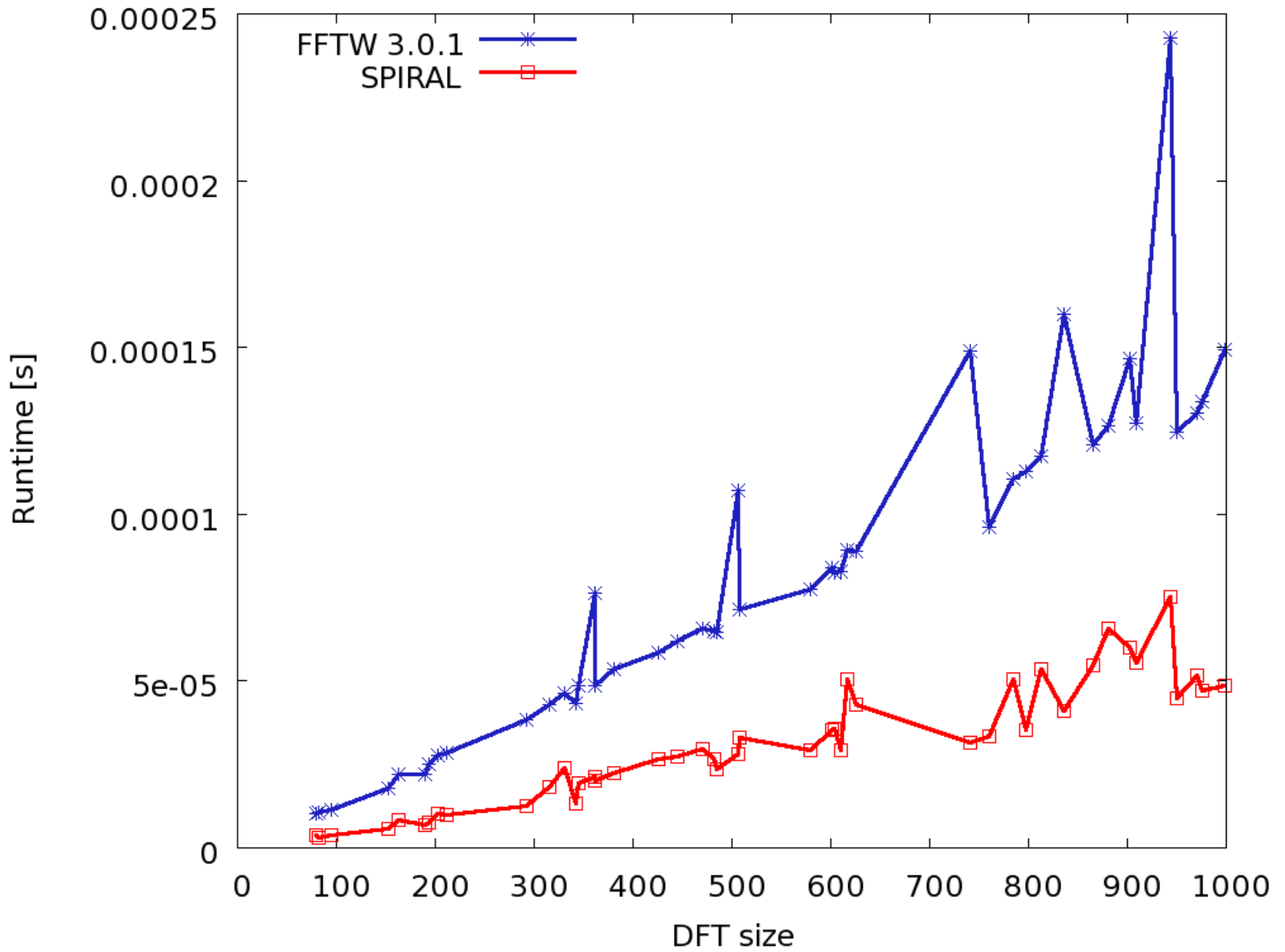- **Experimental Results**

- **Conclusions**

# Benchmarks Setup

- **Comparison against FFTW 3.0.1**

- **Pentium 4 3.6 GHz**

- **We consider sizes requiring at least one Rader step (sizes with large prime factor)**

- **We divide sizes into levels depending on number of Rader steps needed (Rader FFT has most expensive index mapping)**
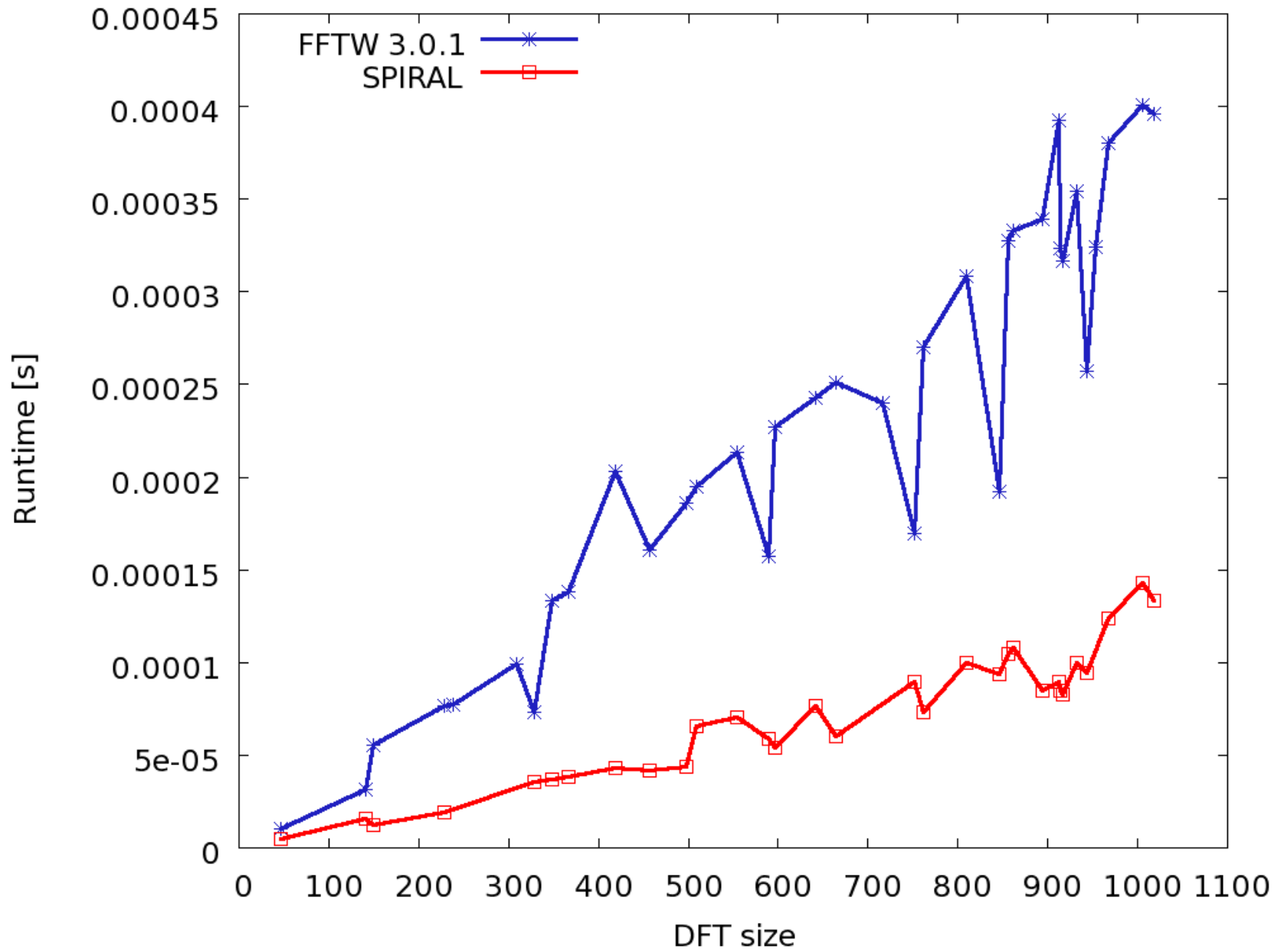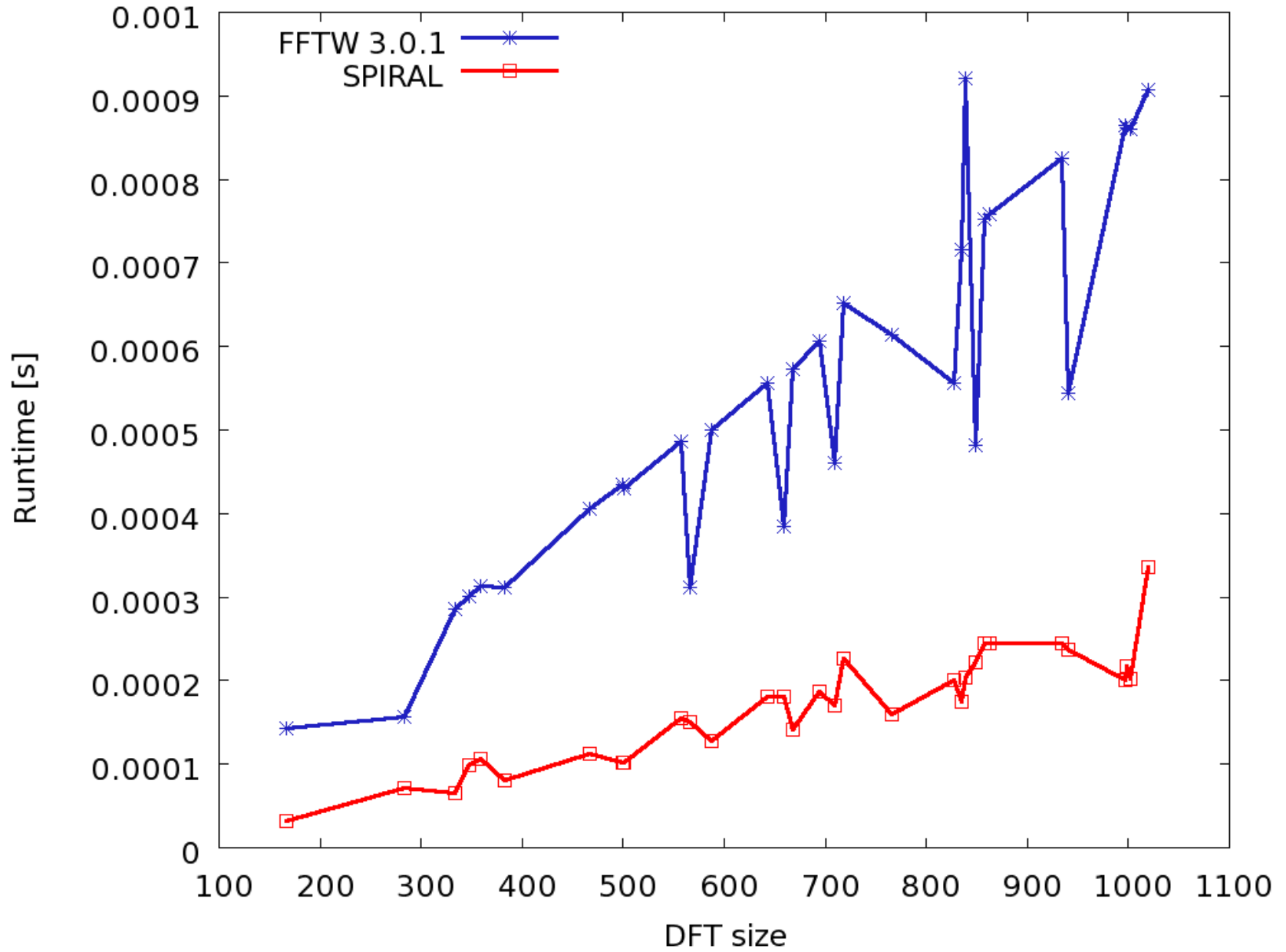
# One Rader Step

Average SPIRAL speedup: factor of 2.7

# Two Rader Steps      Average SPIRAL speedup: factor of 3.3

# **Three** Rader Steps   **Average SPIRAL speedup: factor of 3.4**

# Talk Organization

- **SPIRAL Background**

- **Automatic loop merging in SPIRAL**

- **Experimental Results**

- **Conclusions**

# Conclusion

- **General loop optimization framework for linear DSP transforms in SPIRAL**

- **Loop optimization at the "right" abstraction level: Σ-SPL**

- **Application to FFT: Speedups of a factor of 2-5 over FFTW**

- **Future work: Other Σ-SPL optimizations**
  - **Loop merging for other transforms**
  - **Loop elimination, interchange, peeling**

**http://www.spiral.net**

SPIRAL