

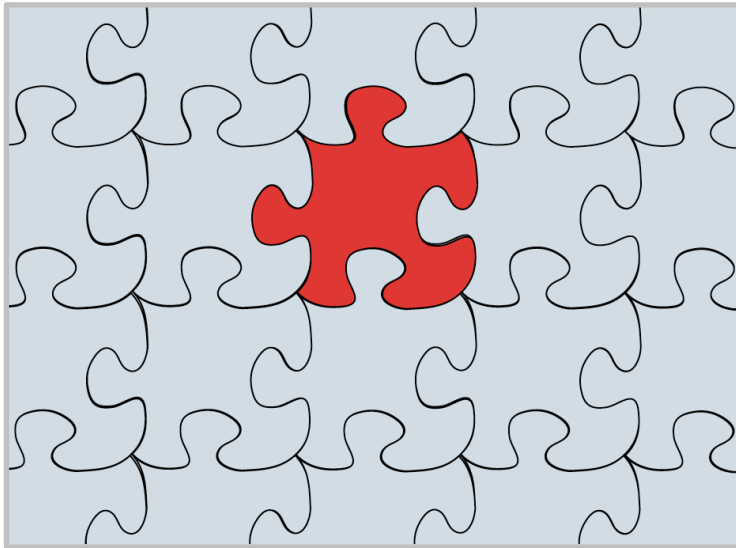
General-Size Library Generation With Spiral

Yevgen Voronenko

Carnegie Mellon University

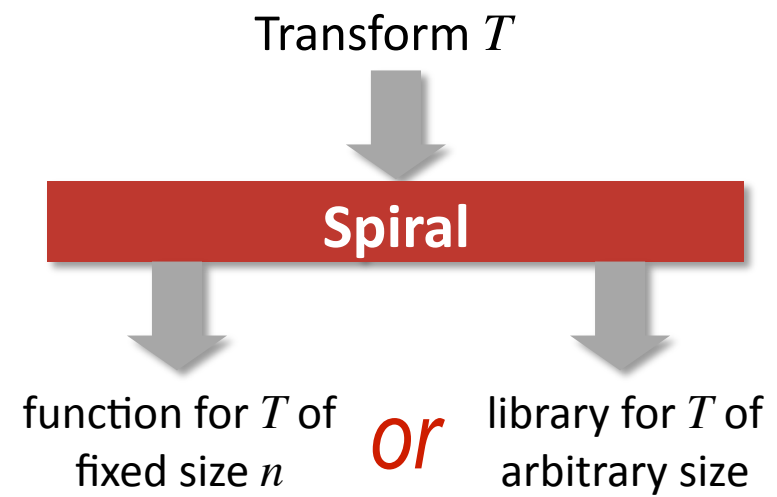
Overview

Spiral: Automating Library Development for Linear Transforms



Automated:

- Loop optimizations
- Vectorization
- Parallelization
- **Derivation of general size libraries**
- Algorithm restructuring and exploration
- Implementation



Fundamentally different problems

Generating Fixed Size Transform Code

Transform
user specified

DFT₈



Fast algorithm
in SPL
many choices

$$(DFT_2 \otimes I_4) T_4^8 (I_2 \otimes ((DFT_2 \otimes I_2) \cdot T_2^4 (I_2 \otimes DFT_2) L_2^4)) L_2^8$$



Σ-SPL:

$$\sum (S_j DFT_2 G_j) \sum \left(\sum (S_{k,l} \text{diag}(t_{k,l}) DFT_2 G_l) \sum (S_m \text{diag}(t_m) DFT_2 G_{k,m}) \right)$$



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

.....

Generating General Size Transform Libraries

Input:

- Transform: DFT_n

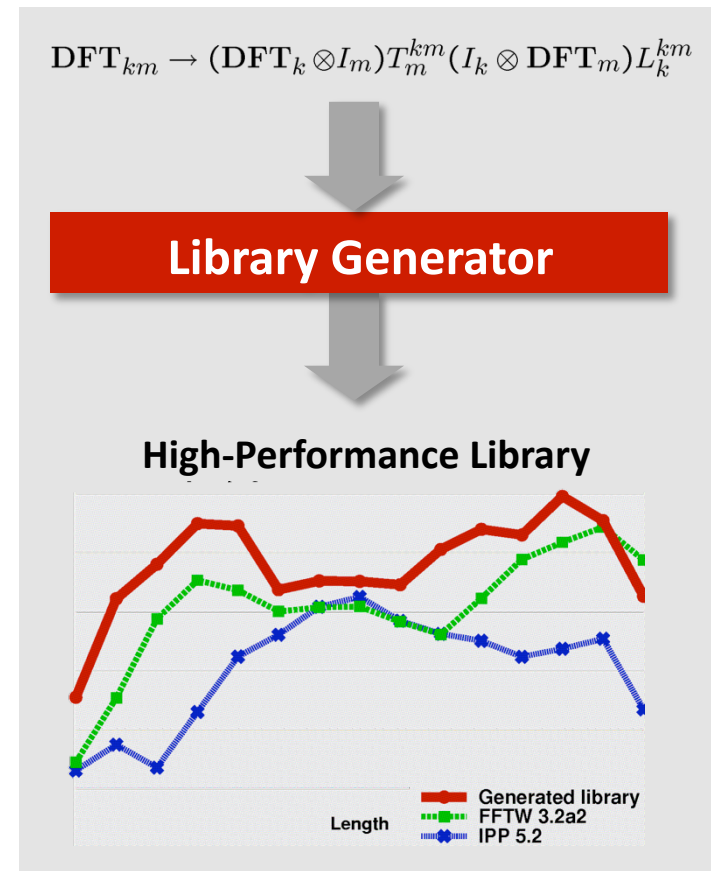
- Algorithms:

$$\begin{aligned} \text{DFT}_{km} &\rightarrow (\text{DFT}_k \otimes I_m) T_m^{km} (I_k \otimes \text{DFT}_m) L_k^{km} \\ \text{DFT}_2 &\rightarrow \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \end{aligned}$$

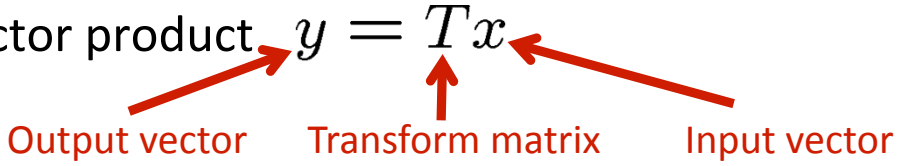
- Vectorization: 2-way SSE
- Threading: Yes

Output:

- Optimized library (10,000 lines of C++)
- For general input size
- Vectorized and multithreaded
- With runtime adaptation mechanism
- Performance competitive with hand-written code



Linear Transforms

- Mathematically:** matrix-vector product $y = Tx$

- Examples:**

$$\begin{array}{ll}
 \text{DFT}_n = \left[\omega_n^{kl} \right]_{0 \leq k, l < n}, & \omega_n = e^{-2\pi j/n} \\
 \text{RDFT}_n = \begin{bmatrix} I''_n \\ \cos \frac{2\pi kl}{n} \\ -\sin \frac{2\pi kl}{n} \end{bmatrix} & \text{WHT}_n = \begin{bmatrix} \text{WHT}_{n/2} & \text{WHT}_{n/2} \\ \text{WHT}_{n/2} & -\text{WHT}_{n/2} \end{bmatrix} \\
 \text{DCT-1}_n = \left[\cos \frac{kl\pi}{n-1} \right] & \text{DHT}_n = \begin{bmatrix} I''_n \\ \cos \frac{2\pi kl}{n} \\ \text{cms} \frac{2\pi kl}{n} \end{bmatrix} \\
 \text{DCT-2}_n = \left[\cos \frac{k(2l+1)\pi}{2n} \right] & \text{DST-1}_n = \left[\sin \frac{(k+1)(l+1)\pi}{n+1} \right] \\
 \text{DCT-3}_n = \left[\cos \frac{(2k+1)l\pi}{2n} \right] & \text{DST-2}_n = \left[\sin \frac{(k+1)(2l+1)\pi}{2n} \right] \\
 \text{DCT-4}_n = \left[\cos \frac{(2k+1)(2l+1)\pi}{4n} \right] & \text{DST-3}_n = \left[\sin \frac{(2k+1)(l+1)\pi}{2n} \right] \\
 \text{MDCT}_n = \left[\cos \frac{(2k+1)(2l+1+n)\pi}{4n} \right] & \text{DST-4}_n = \left[\sin \frac{(2k+1)(2l+1)\pi}{4n} \right] \\
 \text{Filt}_n(t) = \begin{bmatrix} t_0 & \dots & t_{k-1} \\ & \ddots & \dots & \ddots \\ & & t_0 & \dots & t_{k-1} \end{bmatrix} & \text{IMDCT}_n = \left[\cos \frac{(2l+1)(2k+1+n)\pi}{4n} \right] \\
 & \downarrow 2 \text{Filt}_n(t) = \begin{bmatrix} t_0 & t_1 & \dots & t_{k-1} & & & \\ & & t_0 & t_1 & \dots & t_{k-1} & \\ & & & & \ddots & \dots & \dots & \ddots \\ & & & & & t_0 & \dots & t_{k-1} \end{bmatrix}
 \end{array}$$

Transform Algorithms: Example 4-point FFT

Cooley/Tukey fast Fourier transform (FFT):

$$\begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & j & -1 & -j \\ 1 & -1 & 1 & -1 \\ 1 & -j & -1 & j \end{bmatrix} = \begin{bmatrix} 1 & \cdot & 1 & \cdot \\ \cdot & 1 & \cdot & 1 \\ 1 & \cdot & -1 & \cdot \\ \cdot & 1 & \cdot & -1 \end{bmatrix} \begin{bmatrix} 1 & \cdot & \cdot & \cdot \\ \cdot & 1 & \cdot & \cdot \\ \cdot & \cdot & 1 & \cdot \\ \cdot & \cdot & \cdot & j \end{bmatrix} \begin{bmatrix} 1 & 1 & \cdot & \cdot \\ 1 & -1 & \cdot & \cdot \\ \cdot & \cdot & 1 & 1 \\ \cdot & \cdot & 1 & -1 \end{bmatrix} \begin{bmatrix} 1 & \cdot & \cdot & \cdot \\ \cdot & \cdot & 1 & \cdot \\ \cdot & 1 & \cdot & \cdot \\ \cdot & \cdot & \cdot & 1 \end{bmatrix}$$

Fourier transform

Diagonal matrix (twiddles)

$$\text{DFT}_4 \rightarrow (\text{DFT}_2 \otimes \text{I}_2) \text{T}_2^4 (\text{I}_2 \otimes \text{DFT}_2) \text{L}_2^4$$

Kronecker product
Identity
Permutation

- Algorithms are divide-and-conquer: **Breakdown rules**
- Mathematical, declarative representation: **SPL (signal processing language)**
- SPL describes the structure of the dataflow

Fast Algorithms as “Breakdown Rules”

$$\text{DFT}_n \longrightarrow (\text{DFT}_k \otimes I_m) D_{k,m} (I_k \otimes \text{DFT}_m) L_k^n \quad (2.1)$$

DFT
Cooley-Tukey

$$\text{DFT}_n \longrightarrow V_{m,k}^{-1} (\text{DFT}_k \otimes I_m) (I_k \otimes \text{DFT}_m) V_{m,k} \quad (2.2)$$

$$\text{DFT}_n \longrightarrow W_n^{-1} (I_1 \oplus \text{DFT}_{n-1}) E_n (I_1 \oplus \text{DFT}_{n-1}) W_n \quad (2.3)$$

$$\text{DFT}_n \longrightarrow B_{n,m}^\top D_m \text{DFT}_m D'_m \text{DFT}_m D''_m B_{n,m}, \quad m \geq 2n - 1 \quad (2.4)$$

$$\text{DFT}_n \longrightarrow P_{k/2,2m}^\top (\text{DFT}_{2m} \oplus (I_{k/2-1} \otimes_i C_{2m} \text{rDFT}_{2m}((i+1)/k))) (\text{RDFT}'_k \otimes I_m) \quad (2.5)$$

$$\begin{aligned} \begin{pmatrix} \text{RDFT}_n \\ \text{RDFT}'_n \\ \text{DHT}_n \\ \text{DHT}'_n \end{pmatrix} &\longrightarrow (P_{k/2,m}^\top \otimes I_2) \left(\begin{pmatrix} \text{RDFT}_{2m} \\ \text{RDFT}'_{2m} \\ \text{DHT}_{2m} \\ \text{DHT}'_{2m} \end{pmatrix} \oplus \left(I_{k/2-1} \otimes_i M_{2m} \begin{pmatrix} \text{rDFT}_{2m}((i+1)/k) \\ \text{rDFT}'_{2m}((i+1)/k) \\ \text{rDHT}_{2m}((i+1)/k) \\ \text{rDHT}'_{2m}((i+1)/k) \end{pmatrix} \right) \right) \\ &\cdot \left(\begin{pmatrix} \text{RDFT}'_k \\ \text{RDFT}'_k \\ \text{DHT}'_k \\ \text{DHT}'_k \end{pmatrix} \otimes I_m \right) \end{aligned} \quad (2.6)$$

$$\text{RDFT}_n \longrightarrow D_n \cdot \text{DCT-2}_n \cdot P_n, \quad n \text{ odd} \quad (2.7)$$

$$\begin{aligned} \text{DCT-2}_n &\longrightarrow P_{k/2,2m}^\top \left(\text{DCT-2}_{2m} K_2^{2m} \oplus (I_{k/2-1} \otimes N_{2m} \text{RDFT-3}_{2m}^\top) \right) G_n (L_{k/2}^{n/2} \otimes I_2) \\ &\cdot (I_m \otimes \text{RDFT}'_k) Q_{m/2,k} \end{aligned} \quad (2.8)$$

DCT
“Cooley-Tukey”

Capture complicated algorithms concisely

	Fixed Size	General Size
Adaptation	generation-time	runtime (offline)
Code size	100 – 10,000 LOC	1,000 – 10,000,000 LOC
Generated from	ruletree	breakdown rules
Contains	<ul style="list-style-type: none"> ▪ 1 function ▪ loops ▪ arithmetic 	<ul style="list-style-type: none"> ▪ mutually recursive functions ▪ loops ▪ arithmetic ▪ initialization ▪ runtime adaptation mechanism
User code	<code>dft_1024(Y, X);</code>	<code>Env_1 dft(1024); dft.compute(Y, X);</code>

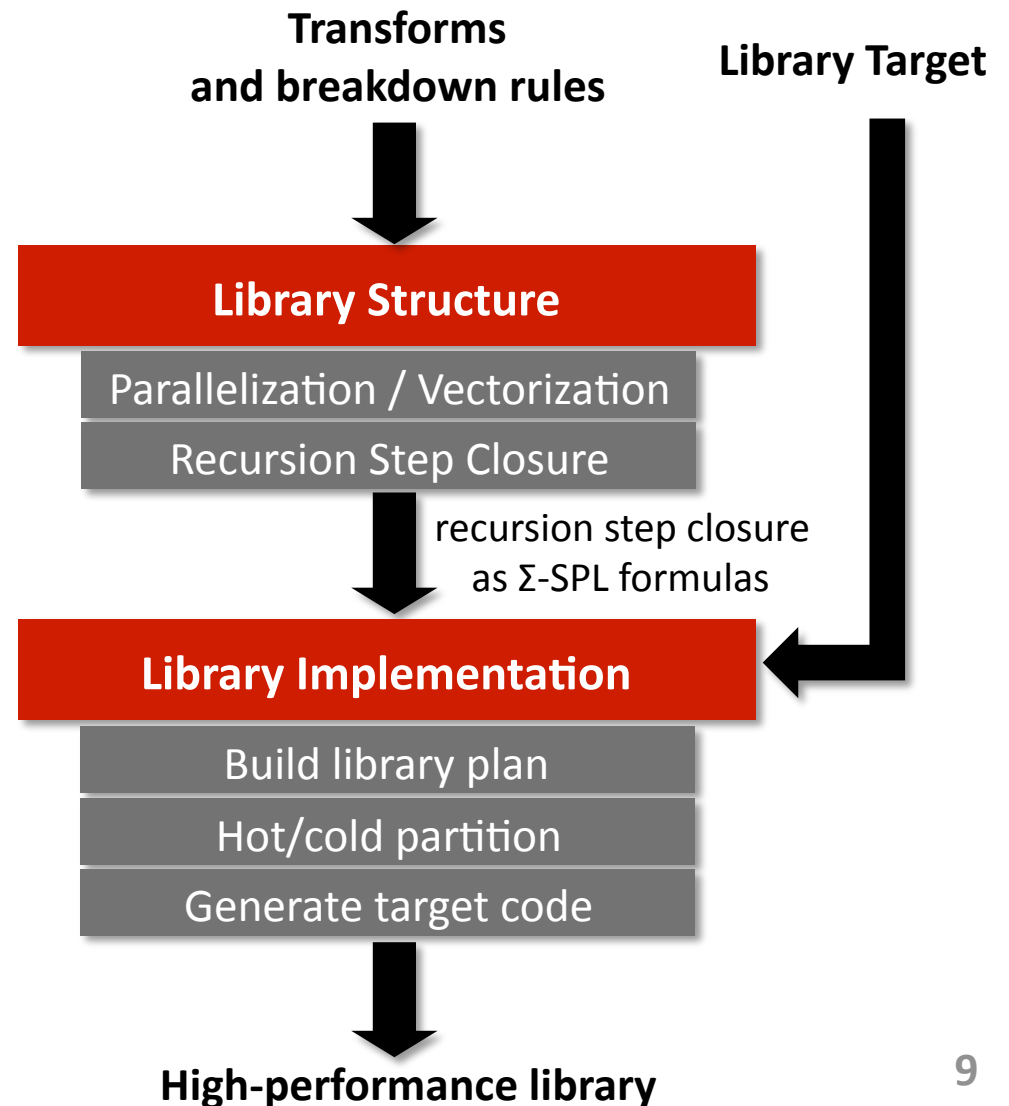
How Library Generation Works

Library Structure : Target Independent

No generation-time search
Fully deterministic

Library Implementation : Target dependent

Generates initialization code
Plugs into the target infrastructure



I. Background

II. Library Generation: Structure

Recursion Step Closure

Base Cases

III. Library Generation: Implementation

IV. Experimental Results

V. Conclusions and Future Work

Library Structure

$$\begin{aligned} & \mathbf{V}_{k,m}^T (\mathbf{DFT}_k \otimes \mathbf{I}_m) (\mathbf{I}_k \otimes \mathbf{DFT}_m) \mathbf{V}_{k,m} \\ & (\mathbf{DFT}_k \otimes \mathbf{I}_m) \mathbf{T}_m^{km} (\mathbf{I}_k \otimes \mathbf{DFT}_m) \mathbf{L}_k^{km} \end{aligned}$$



Library Structure

Parallelization / Vectorization

Recursion Step Closure



DFT

S_z **DFT** G_h

S_h **DFT** G_z

S_h **DFT** $\text{diag}(\text{Dat}) G_h$

S_h **DFT** G_h

S_h **DFT** $G_{h \circ z}$

$S_{h \circ z}$ **DFT** G_h

S_h **DFT** $\text{diag}(\text{Dat}) G_{h \circ z}$

S_z **DFT** $\text{diag}(\text{Dat}) G_h$

$S_{h \circ z}$ **DFT** $\text{diag}(\text{Dat}) G_h$

- **Input:**
 - Breakdown rules
- **Output:**
 - Recursion step closure
 - Σ -SPL Implementation of each recursion step
- **Parallelization/Vectorization**
 - Adds additional breakdown rules
 - Orthogonal to the closure generation
- **Closure**
 - Descend
 - Parametrize
 - Repeat until no new recursion steps found

Recursion Steps

- Cooley-Tukey FFT

$$y = (\text{DFT}_k \otimes I_m) T_m^{km} (I_k \otimes \text{DFT}_m) L_k^{km} x$$

- Implementation in FFTW 2.x

```
DFT(int n, complex *Y, complex *X) {
  k = choose_factor(n); m = n/k;
```

```
  for i=0 to k-1
    DFT_strided(m, k, 1, Y + m*i, T + m*i)
```

```
  for i=0 to m-1
    DFT_scaled(k, m, precomputed_f, Y + i, Y + i)
```

```
}
```

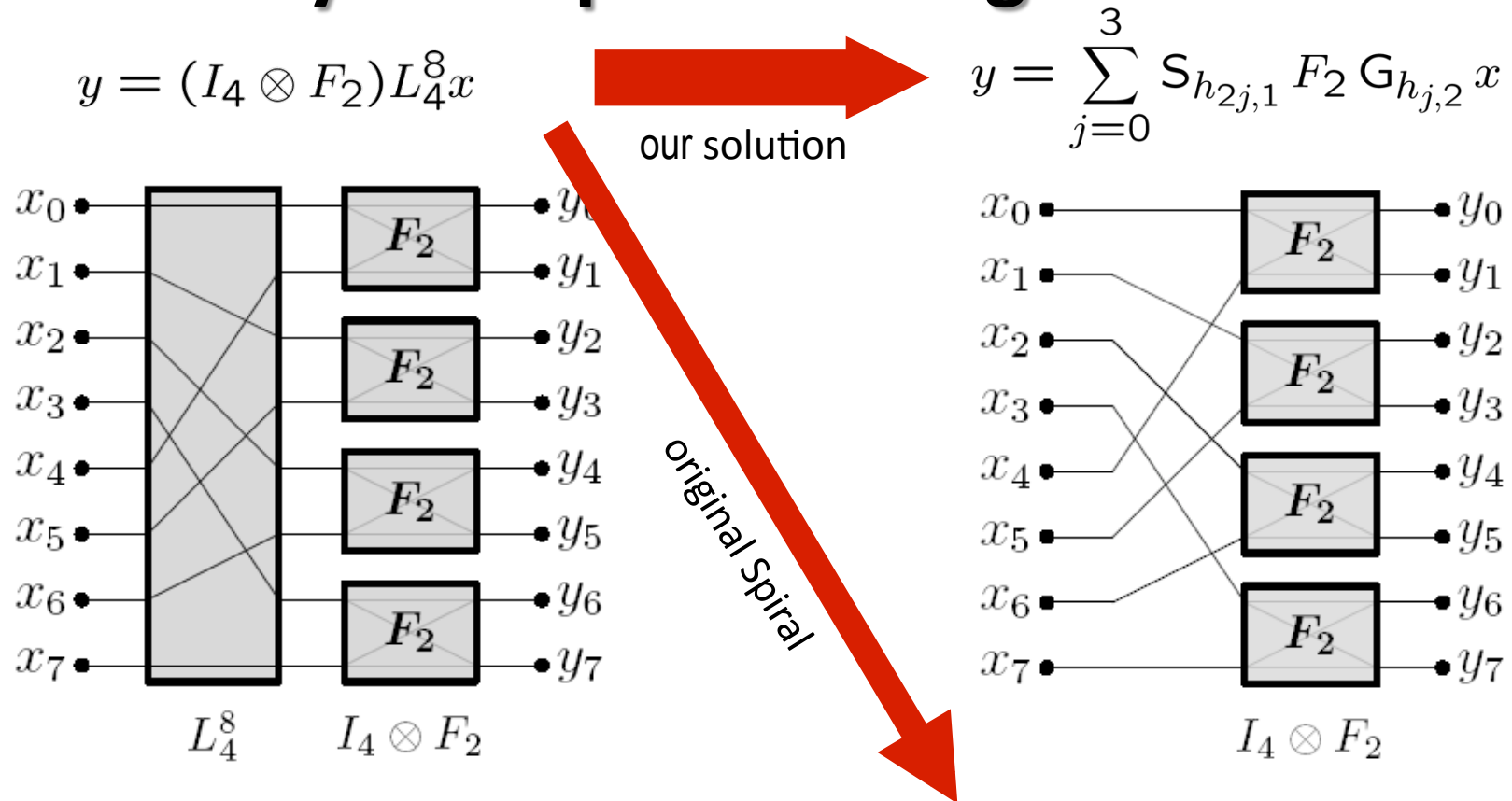
```
DFT_strided(int n, int is, int os, complex *Y, complex *X) {...}
```

```
DFT_scaled(int n, int s, complex *F, complex *Y, complex *X) {...}
```

2 additional functions (recursion steps) needed ...
... and they may spawn more

How to discover automatically?

Σ -SPL: Key to step combining



```

for (int i=0; i<10; i++)
  t[i] = x[(i/5) + 2*(i % 5)];
for (int i=0; i<5; i++) {
  y[2*i] = t[2*i] + t[2*i+1];
  y[2*i+1] = t[2*i] - t[2*i+1];
}

```

impractical

```

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

```

Σ -SPL: Basic Idea

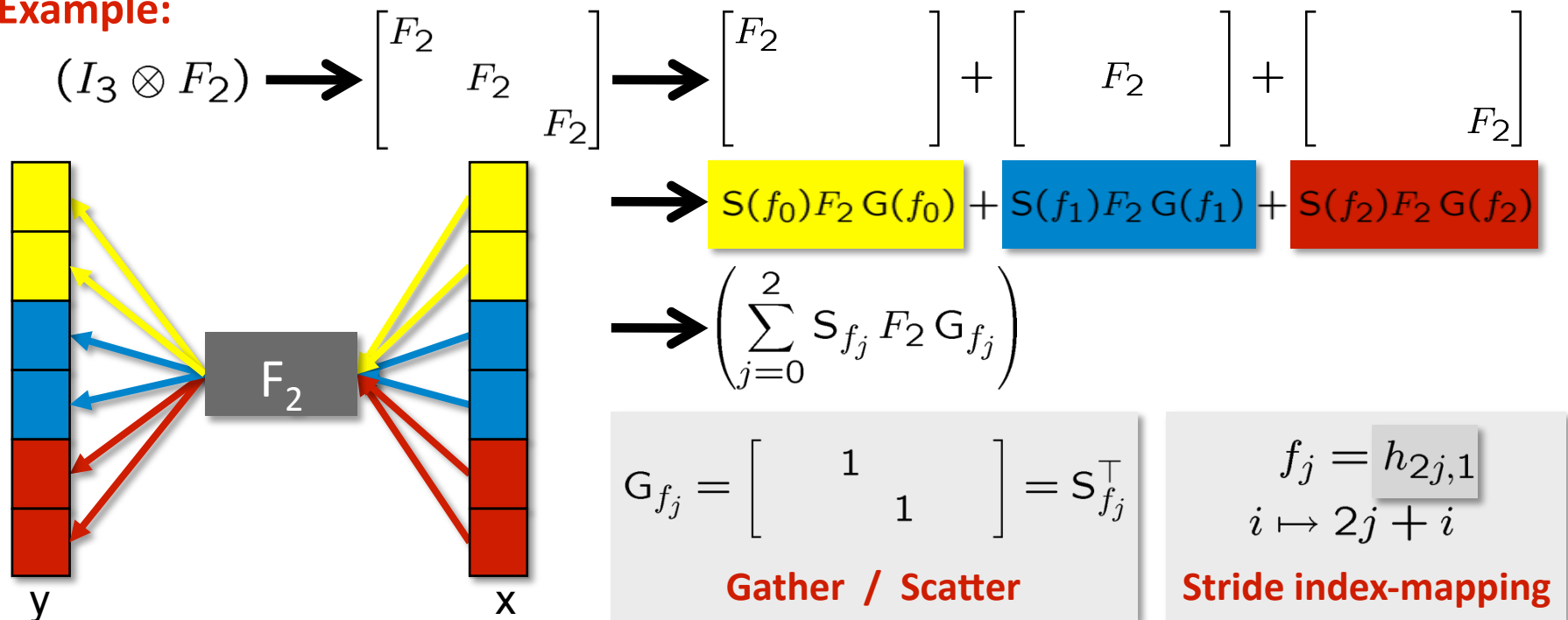
■ Four matrix constructs:

- Σ (sum) – explicit loop
- G_f (gather) – load with index mapping f
- S_f (scatter) – store with index mapping f
- Perm_f – permute with the index mapping f

■ Predefined index-mapping functions

- h – stride
- z – stride mod N (prime-factor algorithms)
- r – reflecting stride (RDFT)
- $1/2$ – exponential

Example:



Recursion Step Closure

- **Input:** transform T and a breakdown rule
- **Output:** spawned recursion steps + Σ -SPL implementation
- **Algorithm:**

1. Apply the breakdown rule

$$\{\text{DFT}_n\}$$

$$\Downarrow$$

$$(\{\text{DFT}_{n/k}\} \otimes I_k) T_k^n (I_{n/k} \otimes \{\text{DFT}_k\}) L_{n/k}^n$$

2. Convert to Σ -SPL

$$\left(\sum_{i=0}^{k-1} S_{h_{i,k}} \{\text{DFT}_{n/k}\} G_{h_{i,k}} \right) \text{diag}(f) \left(\sum_{j=0}^{n/k-1} S_{h_{j,k,1}} \{\text{DFT}_k\} G_{h_{j,k,1}} \right) \text{perm}(\ell_{n/k}^n)$$

3. Apply loop merging + index simplification rules.

$$\sum_{i=0}^{k-1} S_{h_{i,k}} \{\text{DFT}_{n/k}\} \text{diag}(f \circ h_{i,k}) G_{h_{i,k}} \sum_{j=0}^{n/k-1} S_{h_{j,k,1}} \{\text{DFT}_k\} G_{h_{j,n/k}}$$

4. Extract recursion steps

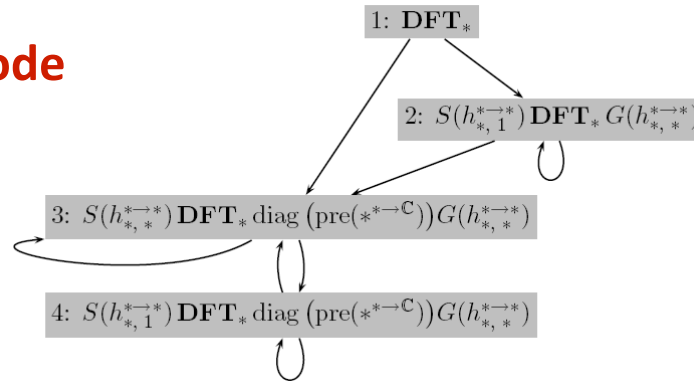
$$\sum_{i=0}^{k-1} \left\{ S_{h_{i,k}} \text{DFT}_{n/k} \text{diag}(f \circ h_{i,k}) G_{h_{i,k}} \right\} \sum_{j=0}^{n/k-1} \left\{ S_{h_{j,k,1}} \text{DFT}_k G_{h_{j,n/k}} \right\}$$

5. Repeat until closure is reached

Parametrization (not shown) derives the independent parameter set for each recursion step

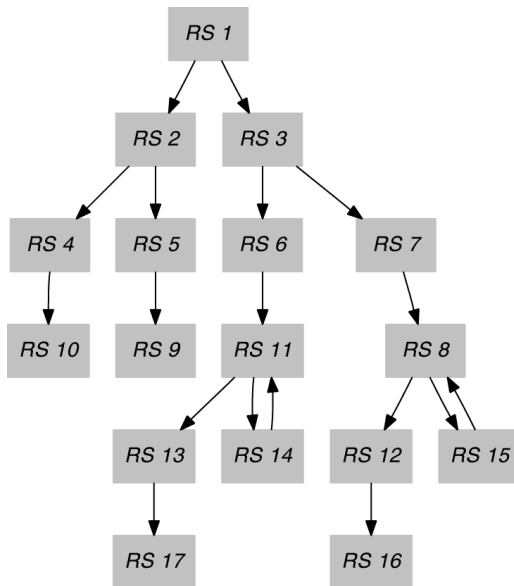
Repeat

DFT: scalar code



**Library =
4 mutually recursive functions**

DCT4: vector code



- 1: $\text{Vec}_2(\text{DCT-4}_{u_1})$
- 2: $\text{Vec}_2(\text{GT}(\text{diag}(N_{2u_8}) \text{RDFT-3}_{2u_8}^\top \text{rcdiag}(\text{pre}(u_4^{\mathbb{Z} \times 2u_8 \rightarrow \mathbb{R}}))), h_{0,1,u_7}^{2u_8 \rightarrow u_6} \circ \ell_{u_8}^{2u_8}, r_{0,u_{11},1,u_{12}}^{2u_8 \rightarrow u_9}, \{u_{13}\})$
- 3: $\text{Vec}_2(\text{GT}(\text{RDFT-3}_{u_1} \text{diag}(N_{u_1}), r_{0,u_5,1,u_6}^{u_1 \rightarrow u_3}, h_{0,u_9,1}^{u_1 \rightarrow u_8}, \{u_{10}\}))$
- 4: $\text{VJam}_2(\text{GT}(\text{diag}(N_{2u_9}) \text{RDFT-3}_{2u_9}^\top \text{rcdiag}(\text{pre}(u_4^{\mathbb{Z} \times 2u_9 \rightarrow \mathbb{R}}))), h_{0,1,u_7,u_8}^{2u_9 \rightarrow u_6} \circ \ell_{u_9}^{2u_9}, r_{0,u_{12},1,2,u_{13}}^{2u_9 \rightarrow u_{10}}, \{2, u_{14}\})$
- 5: $\text{GT}(\text{diag}(N_{2u_9}) \text{RDFT-3}_{2u_9}^\top \text{rcdiag}(\text{pre}(u_4^{\mathbb{Z} \times 2u_9 \rightarrow \mathbb{R}}))), h_{u_7,1,u_8}^{2u_9 \rightarrow u_6} \circ \ell_{u_9}^{2u_9}, r_{u_{12},u_{13},1,u_{14}}^{2u_9 \rightarrow u_{10}}, \{u_{15}\})$
- 6: $\text{VJam}_2(\text{GT}(\text{RDFT-3}_{u_1} \text{diag}(N_{u_1}), r_{0,u_5,1,2,u_6}^{u_1 \rightarrow u_3}, h_{0,u_9,1,2}^{u_1 \rightarrow u_8}, \{2, u_{10}\}))$
- 7: $\text{GT}(\text{RDFT-3}_{u_1} \text{diag}(N_{u_1}), r_{u_5,u_6,1,u_7}^{u_1 \rightarrow u_3}, h_{u_{10},u_{11},1}^{u_1 \rightarrow u_8}, \{u_{12}\})$
- 8: $S(h_{u_3,u_4}^{u_1 \rightarrow u_2}) \text{RDFT-3}_{u_1} \text{diag}(N_{u_1}) G(r_{u_9,u_{10},u_{11}}^{u_1 \rightarrow u_7})$
- 9: $S(r_{u_3,u_4,u_5}^{2u_{13} \rightarrow u_1}) \text{diag}(N_{2u_{13}}) \text{RDFT-3}_{2u_{13}}^\top \text{rcdiag}(\text{pre}(u_9^{2u_{13} \rightarrow \mathbb{R}})) G(h_{u_{12},1}^{2u_{13} \rightarrow u_{11}} \circ \ell_{u_{13}}^{2u_{13}})$
- 10: $\text{VJam}_2(\text{GT}(\text{diag}(N_{2u_9}) \text{RDFT-3}_{2u_9}^\top \text{rcdiag}(\text{pre}(u_4^{\mathbb{Z} \times 2u_9 \rightarrow \mathbb{R}}))), h_{u_7,1,u_8}^{2u_9 \rightarrow u_6} \circ \ell_{u_9}^{2u_9}, r_{u_{12},u_{13},1,u_{14}}^{2u_9 \rightarrow u_{10}}, \{2\})$
- 11: $\text{VJam}_2(\text{GT}(\text{RDFT-3}_{u_1} \text{diag}(N_{u_1}), r_{u_5,u_6,1,u_7}^{u_1 \rightarrow u_3}, h_{u_{10},u_{11},1}^{u_1 \rightarrow u_8}, \{2\}))$
- 12: $\text{GT}(\text{diag}(C_{u_1}) \text{rDFT}_{2u_1}(\lambda\text{-wrap}(\lambda_1^{\mathbb{Z} \times \mathbb{Z} \rightarrow \mathbb{R}}))), h_{0,1,u_5}^{2u_1 \rightarrow u_4}, h_{u_8,u_9}^{2u_{10} \rightarrow u_7} \circ (r_{0,u_{12},1,u_{13}}^{u_1 \rightarrow u_{10}} \otimes v_2), \{u_{14}\})$
- 13: $\text{VJam}_2(\text{GT}(\text{diag}(C_{u_1}) \text{rDFT}_{2u_1}(\lambda\text{-wrap}(\lambda_1^{\mathbb{Z} \times \mathbb{Z} \rightarrow \mathbb{R}}))), h_{u_5,u_6,1,u_7}^{2u_1 \rightarrow u_4}, h_{u_{10},u_{11},1}^{2u_{12} \rightarrow u_9} \circ (r_{0,u_{14},0,1,u_{15}}^{u_1 \rightarrow u_{12}} \otimes v_2), \{2, u_{16}\})$
- 14: $\text{VJam}_2(\text{GT}(\text{RDFT-3}_{u_1} \text{diag}(N_{u_1}), r_{u_5,u_6,1,u_7,u_8}^{u_1 \rightarrow u_3}, h_{u_{11},u_{12},1,u_{13}}^{u_1 \rightarrow u_{10}}, \{2, u_{14}\}))$
- 15: $\text{GT}(\text{RDFT-3}_{u_1} \text{diag}(N_{u_1}), r_{u_5,u_6,u_7,u_8}^{u_1 \rightarrow u_3}, h_{0,u_{11},1}^{u_1 \rightarrow u_{10}}, \{u_{12}\})$
- 16: $S(h_{u_3,u_4}^{2u_5 \rightarrow u_2} \circ (r_{u_7,u_8,u_9}^{u_6 \rightarrow u_5} \otimes v_2)) \text{diag}(C_{u_6}) \text{rDFT}_{2u_6}(\lambda\text{-wrap}(\lambda_1^{\mathbb{Z} \rightarrow \mathbb{R}})) G(h_{u_{14},1}^{2u_6 \rightarrow u_{13}})$
- 17: $\text{VJam}_2(\text{GT}(\text{diag}(C_{u_1}) \text{rDFT}_{2u_1}(\lambda\text{-wrap}(\lambda_1^{\mathbb{Z} \times \mathbb{Z} \rightarrow \mathbb{R}}))), h_{u_5,u_6,1}^{2u_1 \rightarrow u_4}, h_{u_9,u_{10},1}^{2u_{11} \rightarrow u_8} \circ (r_{u_{13},u_{14},u_{15}}^{u_1 \rightarrow u_{11}} \otimes v_2), \{2\})$

Base Cases

- Base cases are called “codelets” in FFTW, UHFFT
- Why needed:
 - Recursion step closure is translated into mutually recursive functions
 - Recursion must be terminated
 - The larger the base case – the less recursion overhead
- How many:
 - In FFTW 3.2: 183 codelets for complex DFT (21 types)
147 codelets for real DFT (18 types)
 - In the generator: # codelet types · # recursion steps
- Obtained by specializing recursion steps to fixed sizes and using standard Spiral to generate code

$$\left\{ S(h_{u_3,1}^{2 \rightarrow u_2}) \text{DFT}_2 G(h_{u_7,u_8}^{2 \rightarrow u_6}) \right\}$$

$$\left\{ S(h_{u_3,1}^{3 \rightarrow u_2}) \text{DFT}_3 G(h_{u_7,u_8}^{3 \rightarrow u_6}) \right\}$$

...

I. Background

II. Library Generation: Structure

Recursion Step Closure

Base Cases

III. Library Generation: Implementation

IV. Experimental Results

V. Conclusions and Future Work

Library Implementation

DFT
 $S_z DFT G_h$
 $S_h DFT G_z$
 $S_h DFT \text{diag} (Dat) G_h$
 $S_h DFT G_h$
 $S_h DFT G_{h_{oz}}$
 $S_{h_{oz}} DFT G_h$
 $S_h DFT \text{diag} (Dat) G_{h_{oz}}$
 $S_z DFT \text{diag} (Dat) G_h$
 $S_{h_{oz}} DFT \text{diag} (Dat) G_h$



Library Implementation

Build library plan

Hot/cold partition

Generate target code



High-performance library

Input:

- Recursion step closure
- Σ -SPL implementation of each recursion step (base cases + recursions)

Output:

- High-performance library
- Target language: C++, Java, etc.



Process:

- Build library plan
- Perform hot/cold partitioning
- Generate target language code

What Kind of Code to Generate: Functions?

- Example: computing DFT_{32} with Intel IPP

```

IppsDFTSpec_C_64fc *f;
ippDFTInitAlloc_C_64fc (f, 32, ...);  initialize
ippDFTFwd_CToC_64fc(X, Y, f, NULL );  compute

```

- In fact, DFT is a higher order function (HOF):

$$N \rightarrow (X \rightarrow Y)$$

- Descriptors = standard way to implement HOFs, also called “closures”
 - “Closures” also equivalent to objects
 - We will implement recursion steps as C++ objects
- Outer parameters = **cold**
- Inner parameters = **hot**

Need to determine hot/cold parameters in recursion steps

Step 1: Build Library Plan

RS 1

- **Formula:** DFT_{u_1}
- **Parameters:** u_1
- **Children:** RS 2, RS 3
- **Implementation 1:** DFT-LibBase
 - *Applicable:* $u_1 = 2$
 - *Formula:* $\begin{bmatrix} 1 & -1 \\ 1 & -1 \end{bmatrix}$
- **Implementation 2:** DFT-CTA
 - *Applicable:* $isPrime(u_1) = 0$
 - *Freedoms:* $f \in divisors(u_1)$
 - *Formula:* <...omitted...>

RS 2

- **Formula:** $S_{h_{u_3,1}}^{u_1 \rightarrow u_2} DFT_{u_1} G_{h_{u_7, u_8}}^{u_1 \rightarrow u_6}$
- **Parameters:** $u_1, u_2, u_3, u_6, u_7, u_8$
- **Children:** RS 2, RS 3
- **Implementation 1:** DFT-LibBase
 - *Applicable:* $u_1 = 2$
 - *Formula:* $S_{h_{u_3,1}}^{2 \rightarrow u_2} \begin{bmatrix} 1 & -1 \\ 1 & -1 \end{bmatrix} G_{h_{u_7, u_8}}^{2 \rightarrow u_6}$
- **Implementation 2:** DFT-CTA
 - *Applicable:* $isPrime(u_1) = 0$
 - *Freedoms:* $f \in divisors(u_1)$
 - *Formula:* <...omitted...>

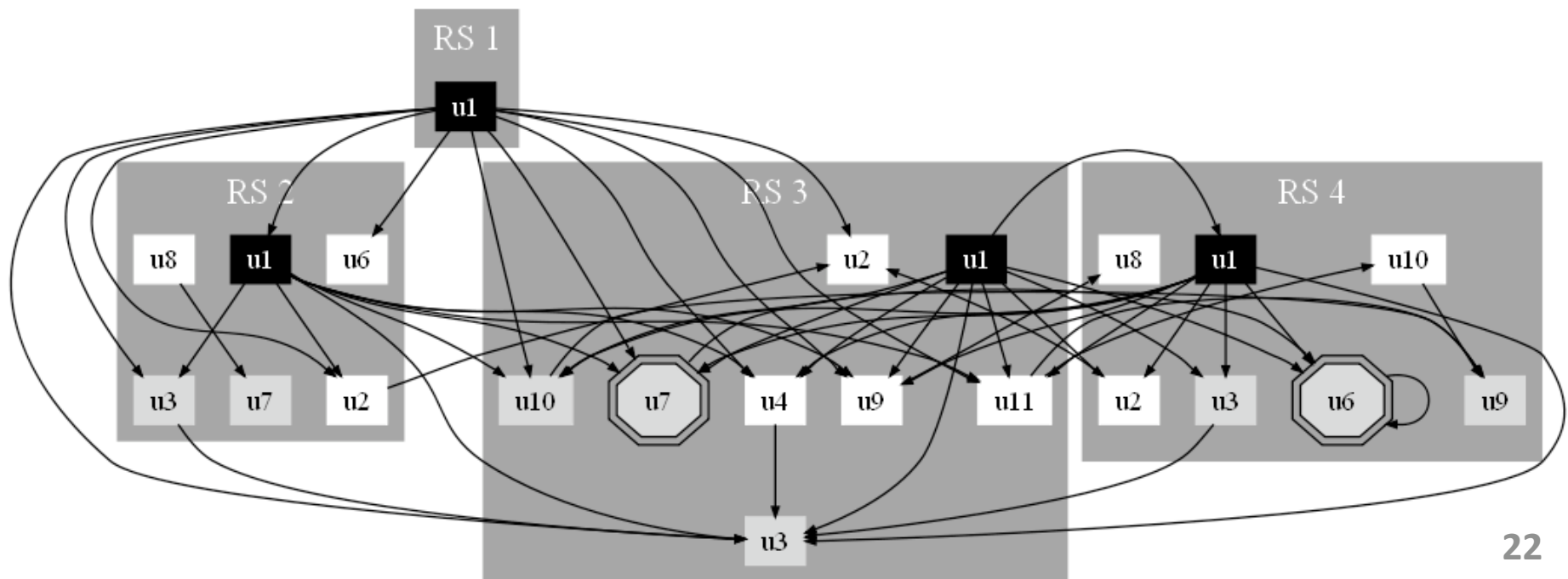
$$\sum_{i=0}^{u_1/f_1-1} RS3 \left((u_2, u_1, u_3, u_4, u_7^{\mathbb{Z} \rightarrow \mathbb{C}}, u_9, u_{10}, u_{11}) \leftarrow (u_1, f_1, i, u_1/f_1, \Omega_{u_1}^1 \circ dt \circ h_{i, u_1/f_1}^{f_1 \rightarrow u_1}, u_1, i, u_1/f_1) \right)$$

Parameters have dependencies

$$\sum_{j=0}^{f_1-1} RS2 \left((u_2, u_1, u_3, u_6, u_7, u_8) \leftarrow (u_1, u_1/f_1, u_1 j/f_1, u_1, j, f_1) \right)$$

Step 2: Hot/Cold Partition

- Build the “parameter flow graph”
- Phase 1 (Backward IDA):
 - Mark mandatory *cold* parameters
 - Propagate coldness
 - none ! cold → cold ! cold
- Phase 2 (Forward IDA):
 - Mark mandatory *hot* parameters
 - Propagate hotness
 - hot ! none → hot ! hot



Step 3: Generate Target Language Code

Generated Code (fragment)

```
class Env_1 : public Env {
    int _rule, f_1, u1;
    char *_dat;
    Env *child1, *child2;
    Env_1(int u1);
    void compute(double *Y, double *X);
};

class Env_2 : public Env {
    int _rule, u8;
    Func_Int_Int_Real *u4;
    char *_dat;
    Env *child1;
    Env_2(int u8, Func_Int_Int_Real *u4);
    void compute(double *Y, double *X, int
        u6, int u7, int u9, int u11, int u12,
        int u13);
};
```

User Code

```
Env_1 dft(1024);
dft.compute(Y, X);
```

- I. Background
- II. Library Generation: Structure
 - Recursion Step Closure
 - Base Cases
- III. Library Generation: Implementation
- IV. Experimental Results**
- V. Conclusions and Future Work

Experimental Results : Setup

- **Generated libraries are in C++, compiled with Intel C++ Compiler 10.1**
- **Use SSE intrinsics for vectorization**
 - 2-way for double precision
 - 4-way for single precision
- **Compared against**
 - FFTW 3.2 alpha 2
 - Intel IPP 5.3
- **Platform**
 - Intel Xeon 5160 3 GHz (2x2 cores, 8 MB L2\$)

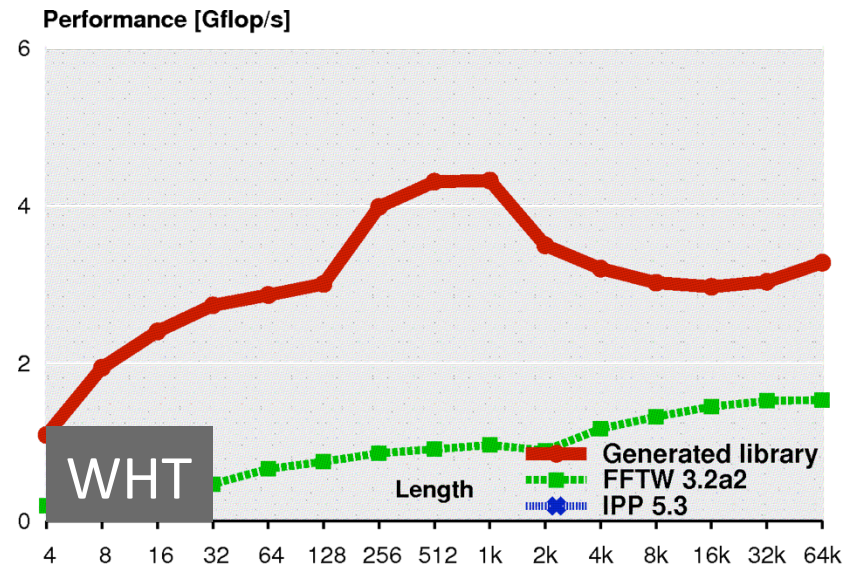
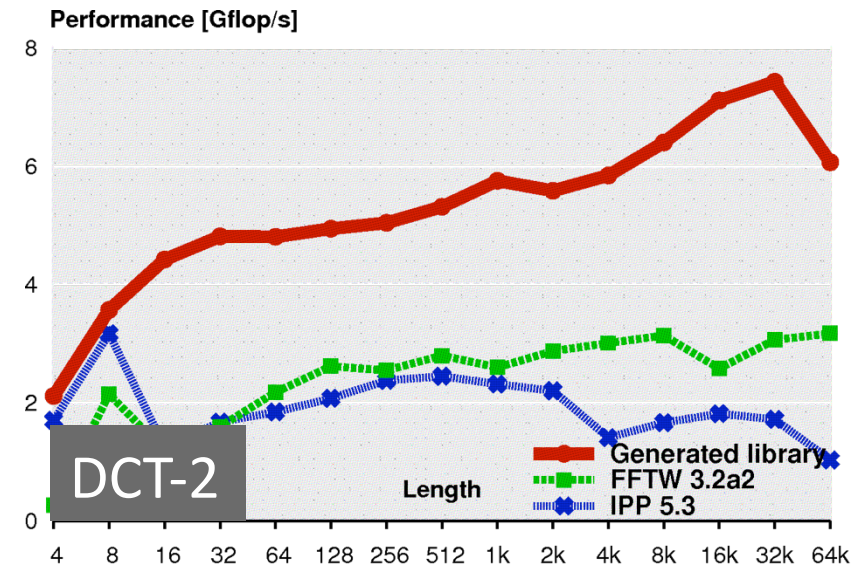
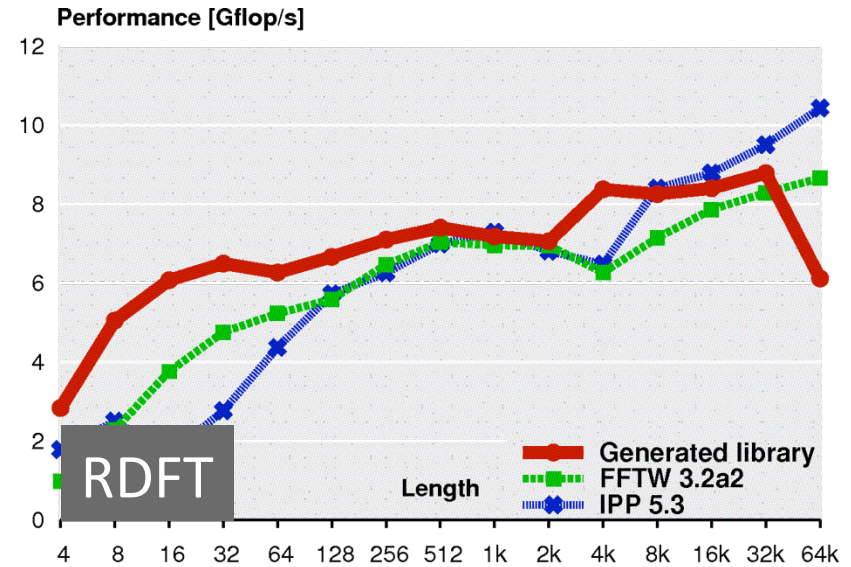
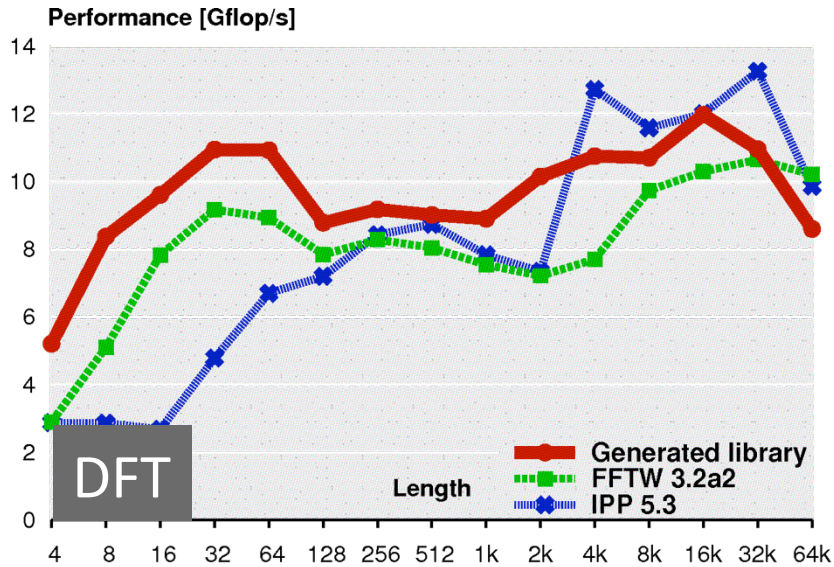
Code Size Snapshot

Transform	Code size	
	non-parallelized	parallelized
<i>no vectorization</i>		
DFT	13.1 KLOC / 0.59 MB	10.3 KLOC / 0.45 MB
RDFT	8.5 KLOC / 0.36 MB	8.8 KLOC / 0.39 MB
DHT	9.1 KLOC / 0.40 MB	9.4 KLOC / 0.39 MB
DCT-2	12.0 KLOC / 0.55 MB	12.4 KLOC / 0.57 MB
DCT-3	12.0 KLOC / 0.56 MB	12.3 KLOC / 0.59 MB
DCT-4	6.8 KLOC / 0.33 MB	7.1 KLOC / 0.35 MB
WHT	5.6 KLOC / 0.21 MB	—
<i>4-way vectorization</i>		
DFT	17.9 KLOC / 1.09 MB	18.2 KLOC / 1.11 MB
RDFT	16.2 KLOC / 0.86 MB	16.5 KLOC / 0.91 MB
scaled RDFT	16.5 KLOC / 0.88 MB	—
DHT	17.9 KLOC / 1.02 MB	18.3 KLOC / 1.04 MB
DCT-2	23.3 KLOC / 1.50 MB	23.6 KLOC / 1.53 MB
DCT-3	32.0 KLOC / 2.17 MB	32.3 KLOC / 2.20 MB
DCT-4	8.3 KLOC / 0.63 MB	8.6 KLOC / 0.66 MB
WHT	8.5 KLOC / 0.53 MB	6.9 KLOC / 0.4 MB
2D DFT	20.6 KLOC / 1.32 MB	20.8 KLOC / 1.33 MB
2D DCT-2	27.0 KLOC / 2.1 MB	27.2 KLOC / 2.11 MB
FIR Filter	109 KLOC / 5.69 MB	74 KLOC / 3.44 MB
Downsampled FIR Filter	151 KLOC / 7.7 MB	92 KLOC / 4.61 MB

Quick turnaround time. **Example:** native algorithm for DCT-2

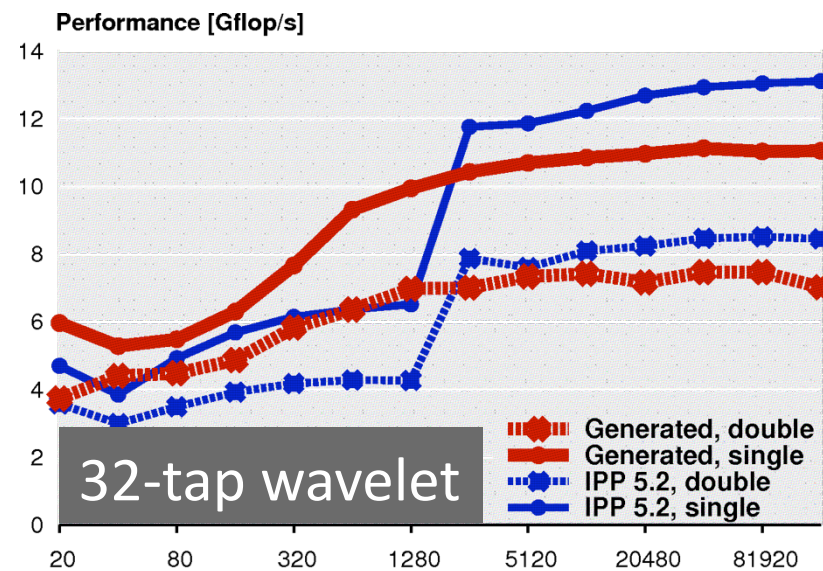
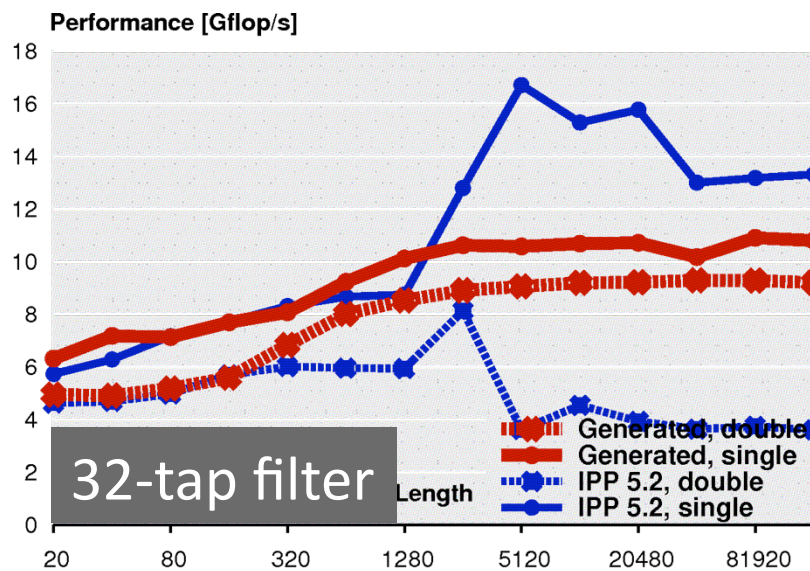
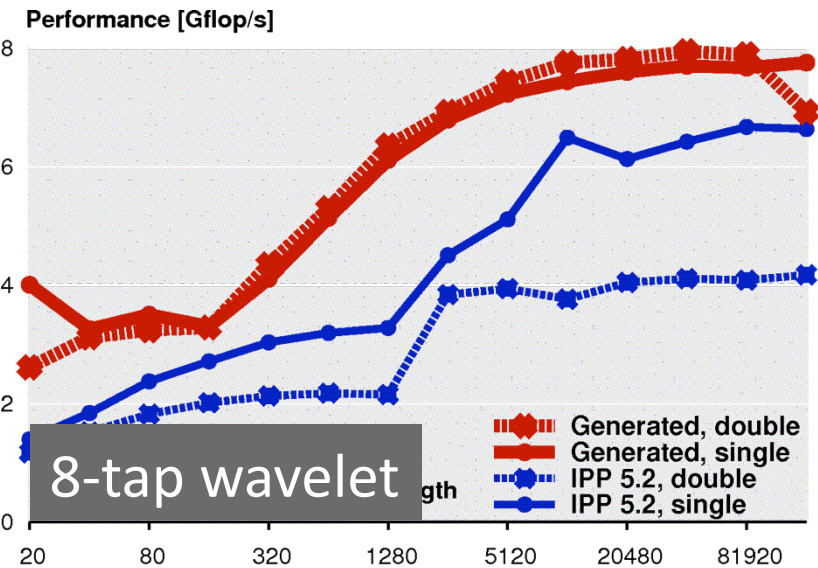
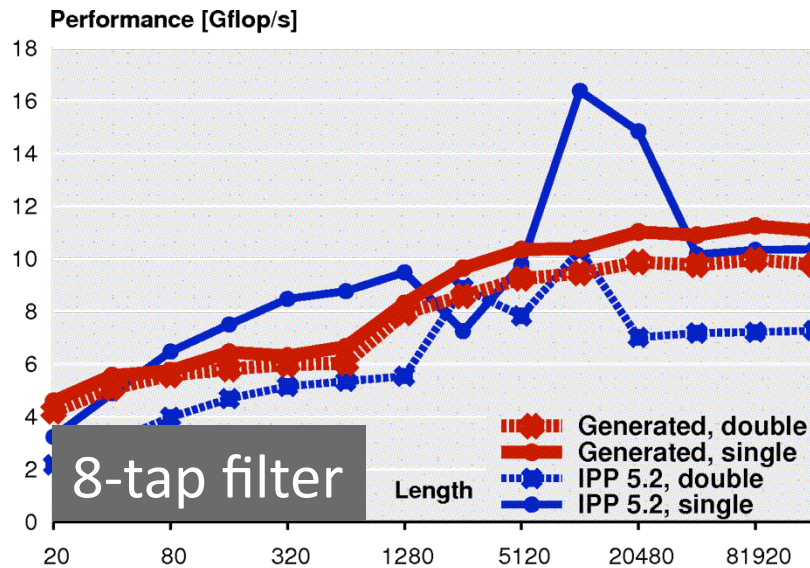
Double Precision Performance: Intel Xeon 5160

2-way vectorization, up to 2 threads



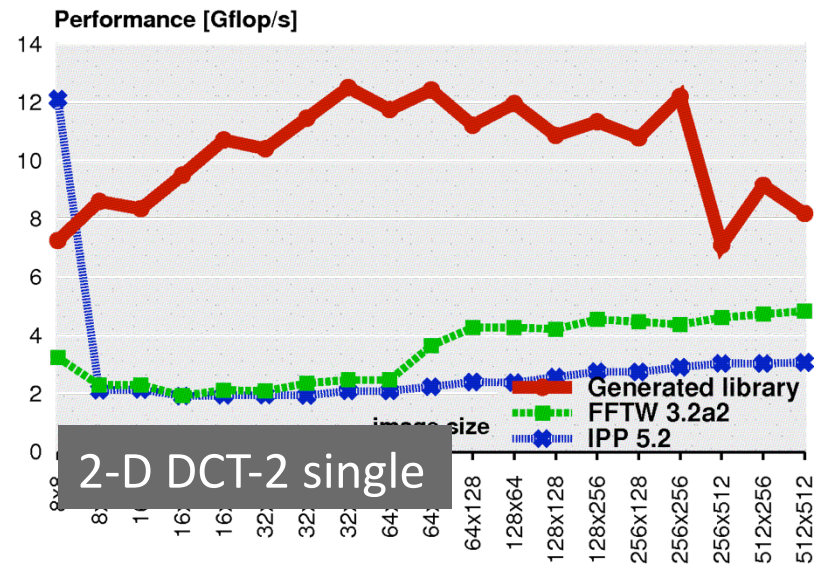
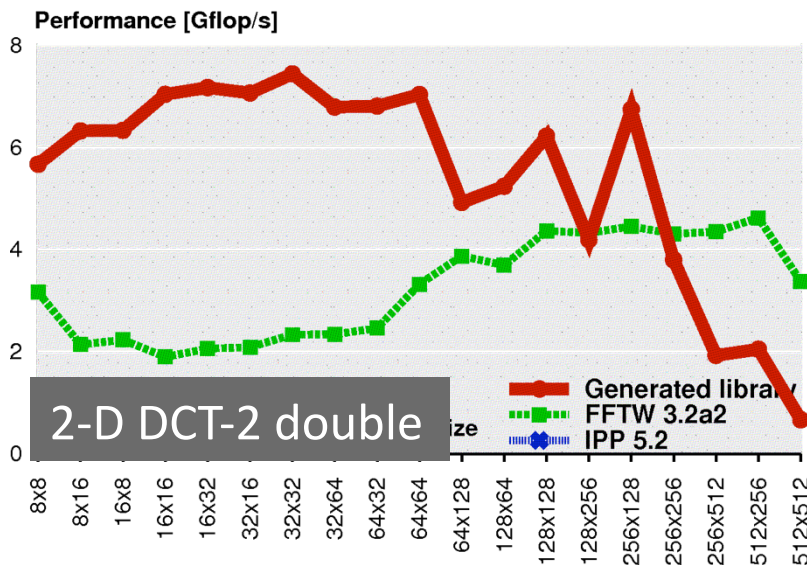
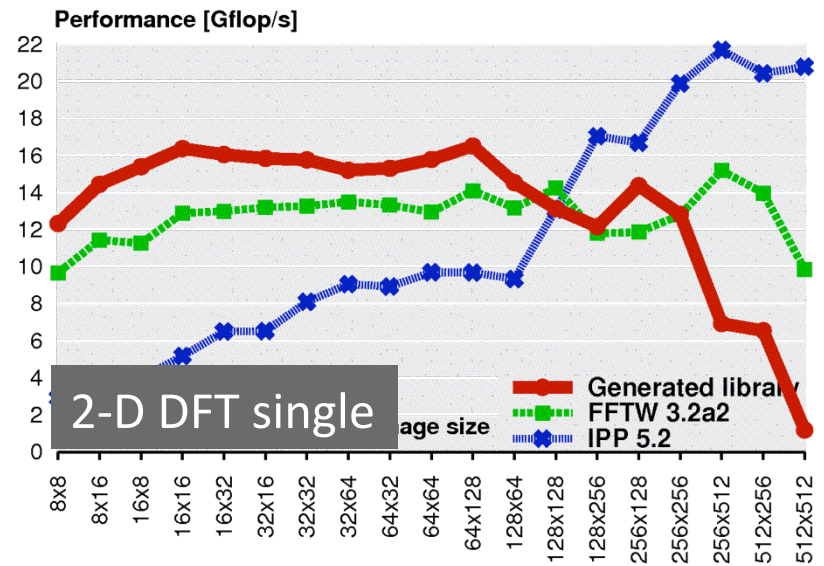
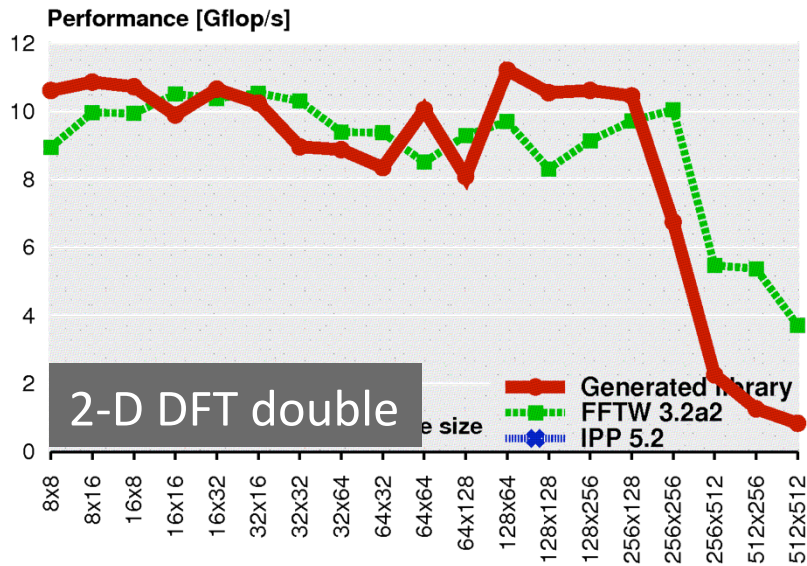
FIR Filter Performance

2- and 4-way vectorization, up to 2 threads



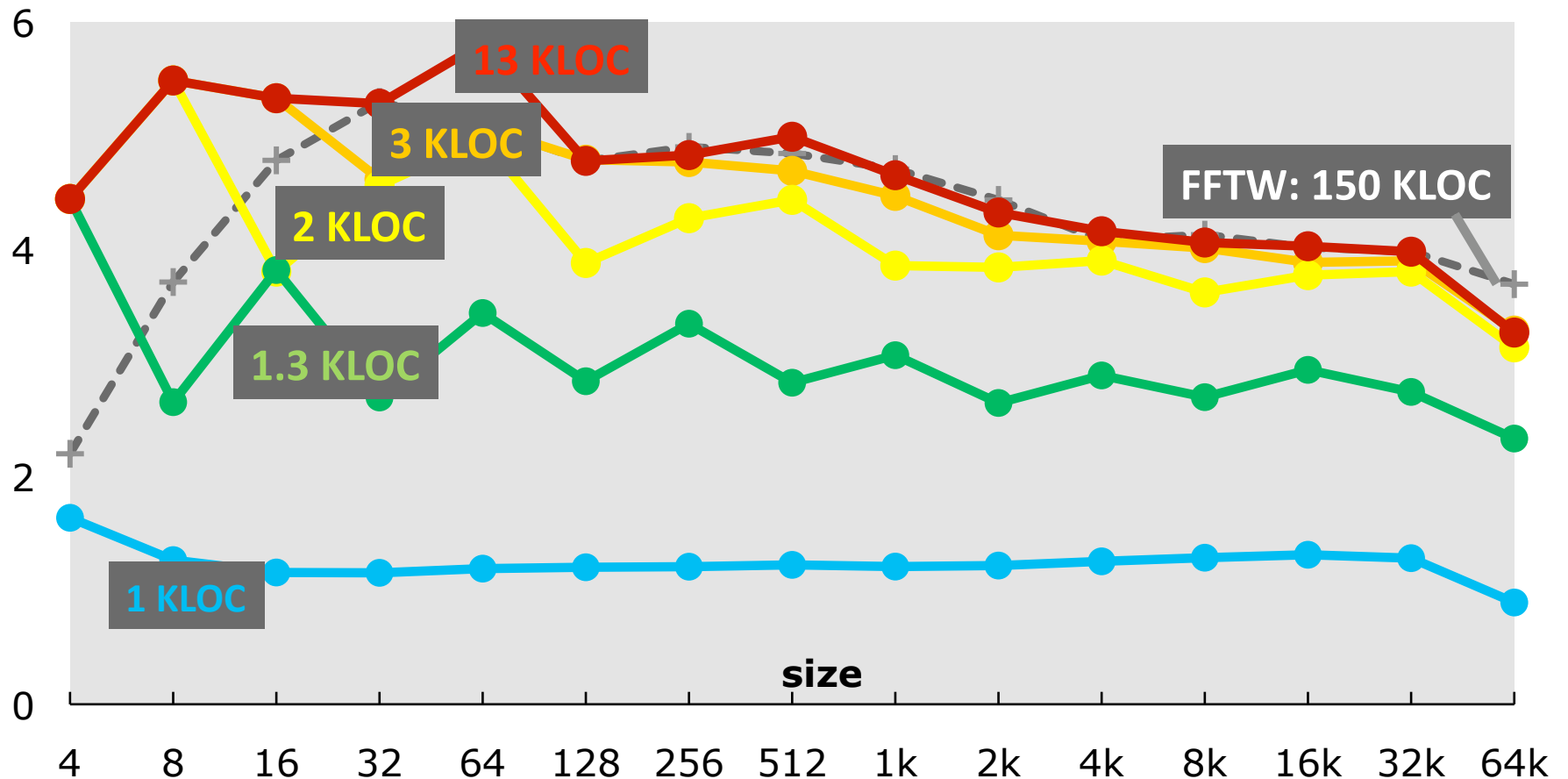
2-D Transforms Performance

2- or 4-way vectorization, up to 2 threads



Qualitative Customization: Code Size

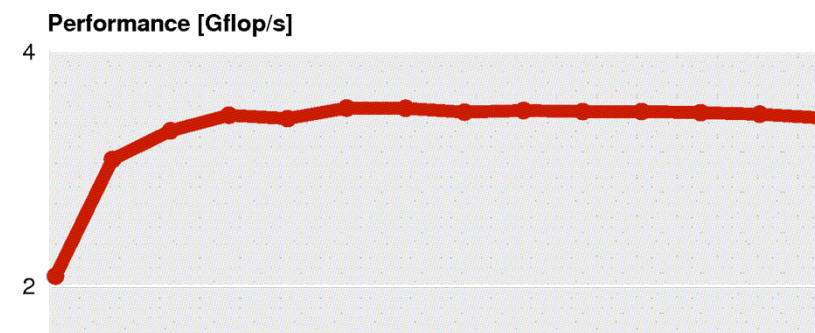
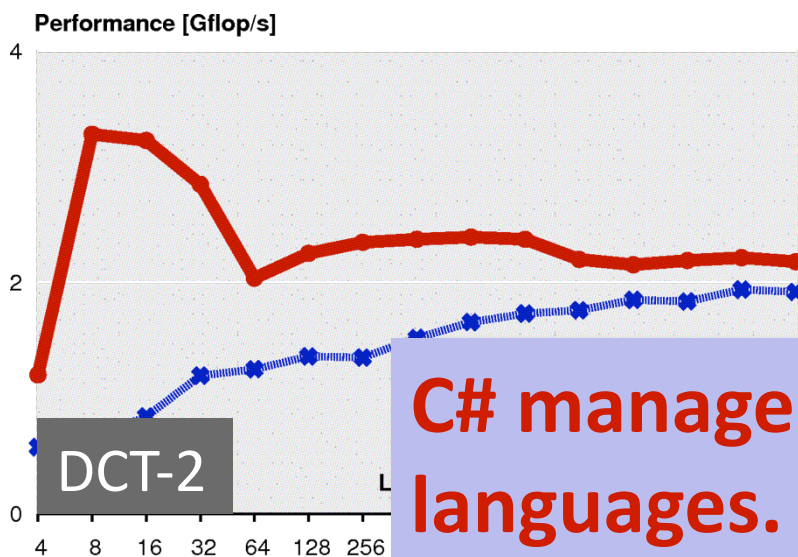
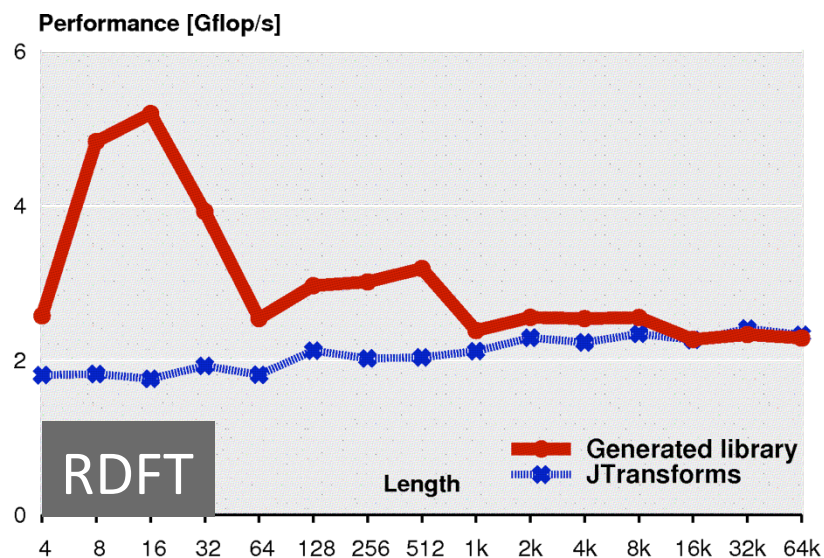
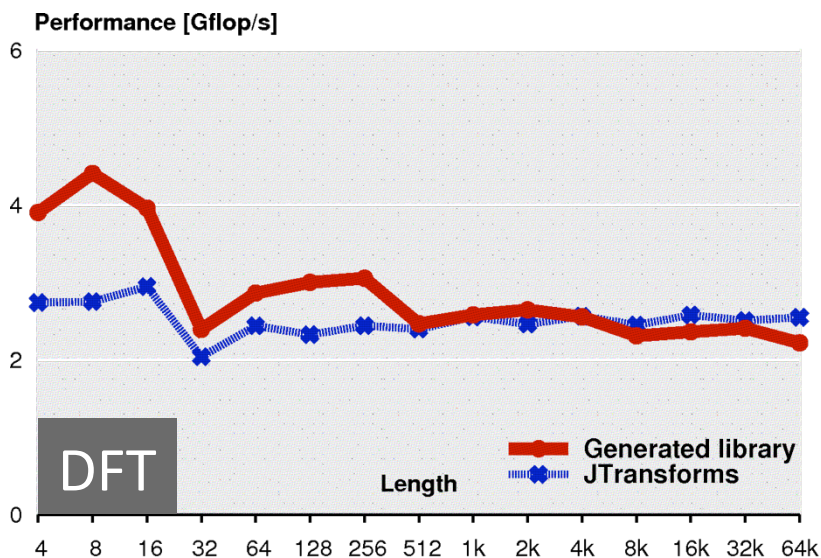
Performance [Gflop/s]



Code size controlled by changing # of base cases...
 ... but other possibilities are also available

Backend Customization: Java

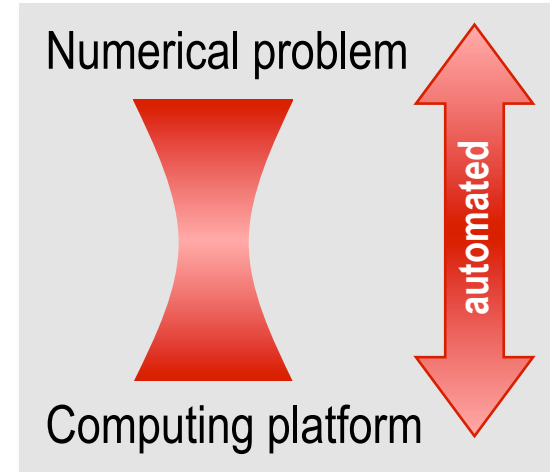
Joint work with Frédéric de Mesmay



C# managed code possible, and other languages.

Future Work and Conclusions

- **Milestone:** Computer generation of high-performance libraries for the entire domain of linear transforms
- **Vertical integration (algorithm ! “FTW”)**
- **Important advantage: customization**
- **Ongoing work: Industry collaboration**
 - Robustness
 - Extra features
 - Completeness (full cross-product)
- **Future: new optimizations for last bits of performance**
 - Software pipelining (Cell, PowerPC)
 - Scheduling and register allocation (target = assembly code)
 - Advanced index mapping optimizations (e.g., sliding pointers)



Question: What can we do to build common tool bases for compiler-based autotuning and for construction of self-tuning or autotuning libraries?

- **Layers of DSLs very powerful idea**
- **Use rewriting to move between levels**
- **Existing generic tools:** Stratego, Maude, (Spiral?), others.
- **In Spiral:**
 - Transform
 - SPL
 - Sigma-SPL
 - Typed Sigma-SPL
 - Code IR 1
 - Code IR 2 -- binary ops, registers