

Spiral

Computer Synthesis of Computational Programs

Joint work with ...

Franz Franchetti

Yevgen Voronenko

Srinivas Chellappa

Frédéric de Mesmay

Daniel McFarlin

José Moura

James Hoe

...



Markus Püschel

Computer Science

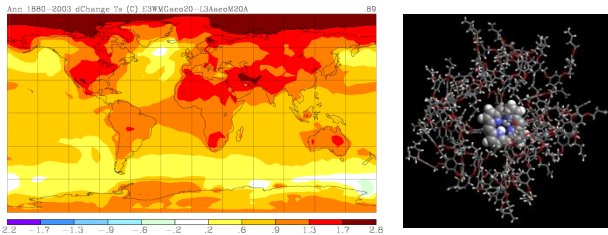
ETH

Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

SPIRAL 
www.spiral.net

Supported by DARPA, ONR, NSF, Intel, Mercury

Scientific Computing



Physics/biology simulations

Consumer Computing



Audio/image/video processing

Embedded Computing



Signal processing, communication, control

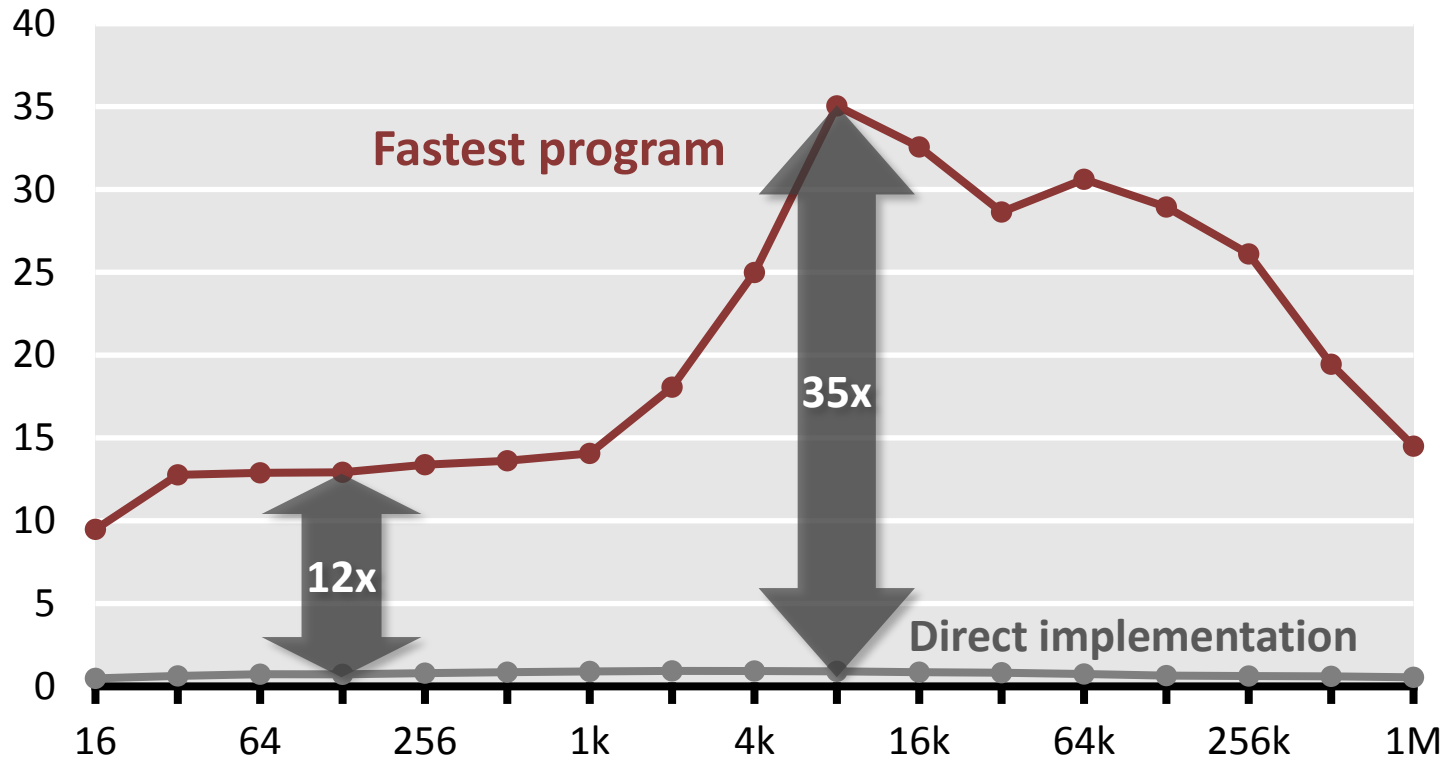
Computing

- **Unlimited need for performance**
- **Large set of applications, but ...**
- **Relatively small set of critical components (100s to 1000s)**
 - Matrix multiplication
 - Discrete Fourier transform
 - Viterbi decoder
 - Filter/stencil
 -

The Problem: Example DFT

DFT on Intel Core i7 (4 Cores, 2.66 GHz)

Performance [Gflop/s]

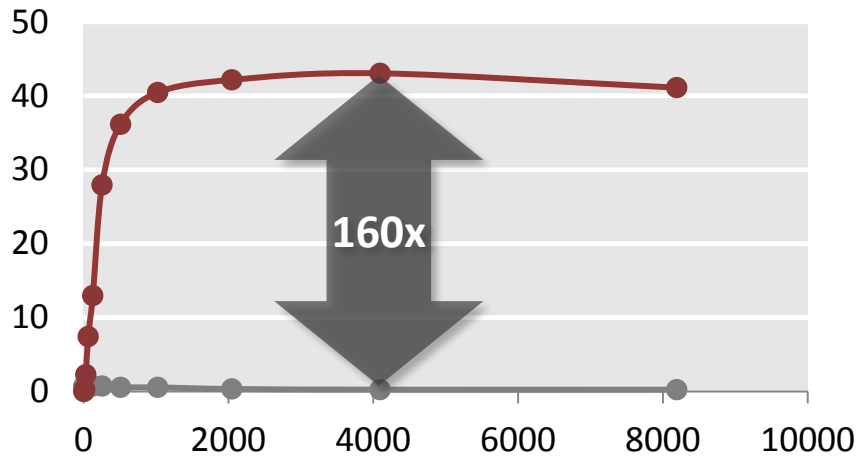


- Same number of operations
- Best compiler

The Problem Is Everywhere

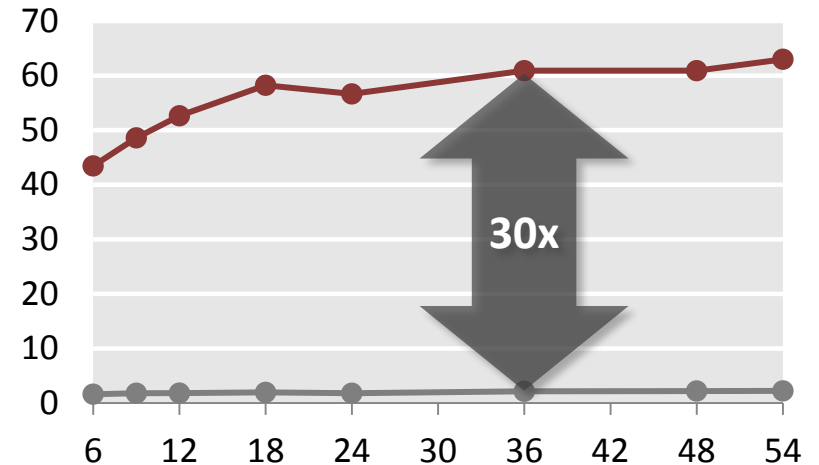
Matrix multiplication

Performance [Gflop/s]



WiFi Receiver

Performance [Mbit/s]

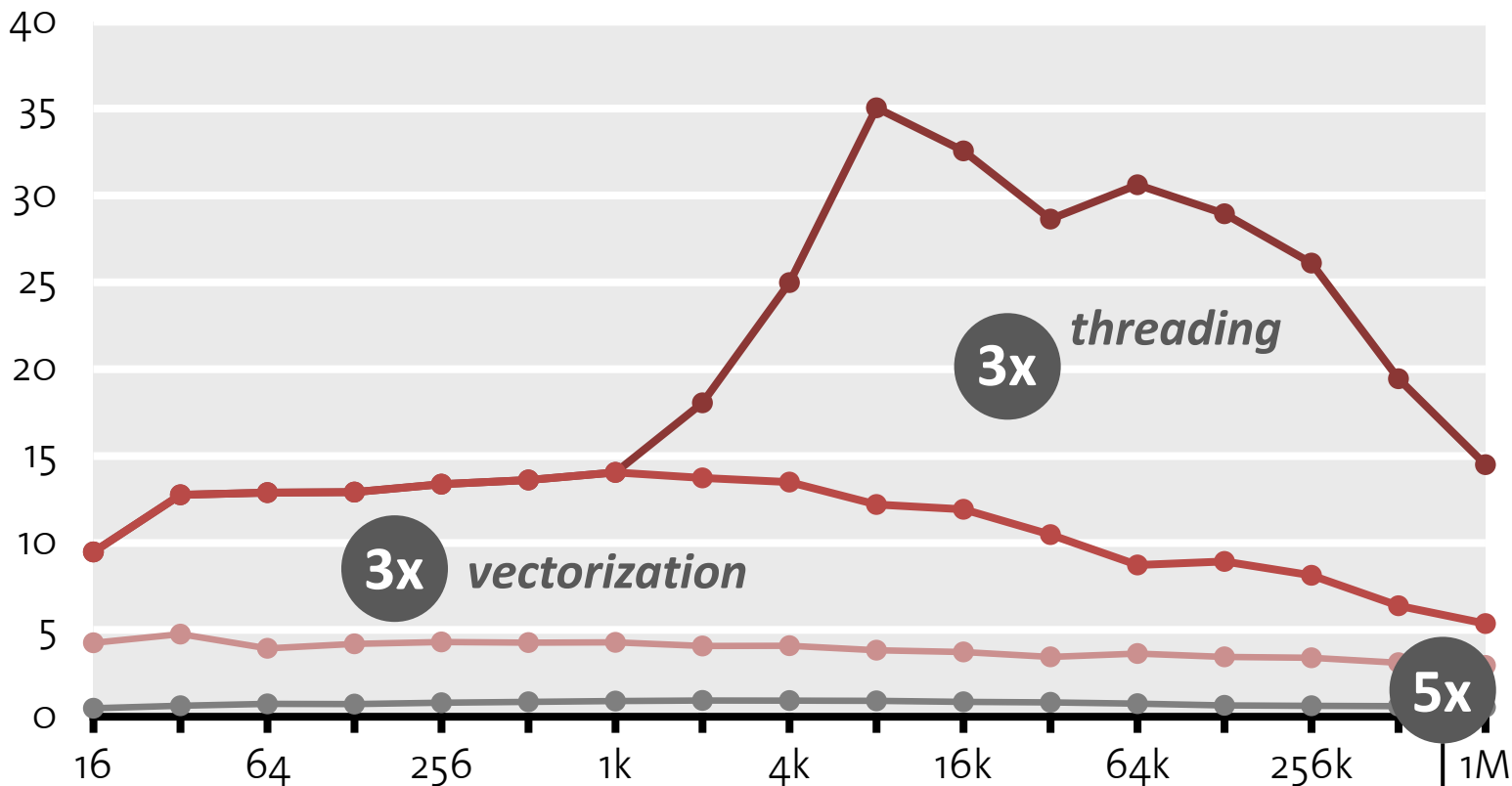


Why is that?

DFT: Analysis

DFT (single precision) on Intel Core i7 (4 cores, 2.66 GHz)

Performance [Gflop/s]



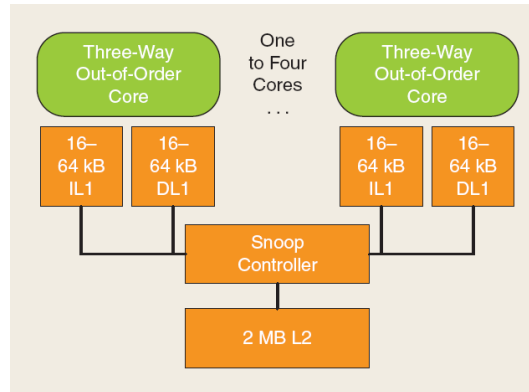
■ Compiler doesn't do it

■ Doing by hand: Very tough

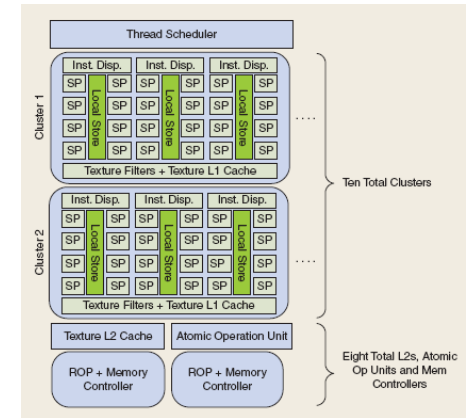
locality optimization

And There Is Not Only Intel ...

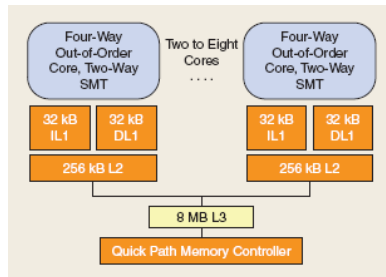
Arm Cortex A9



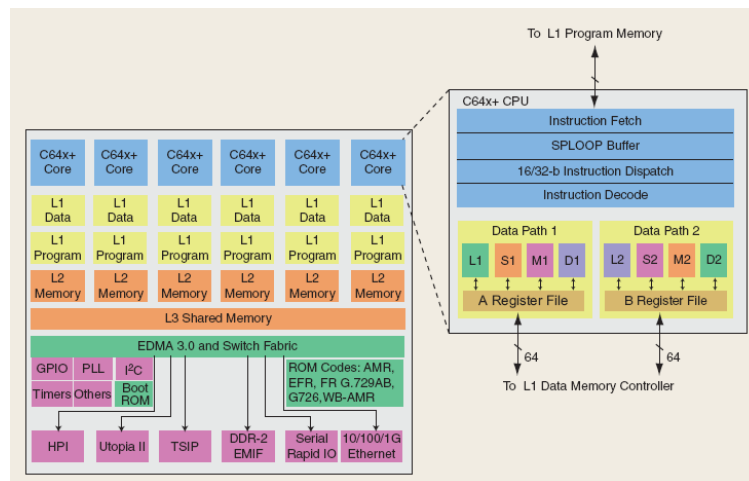
Nvidia G200



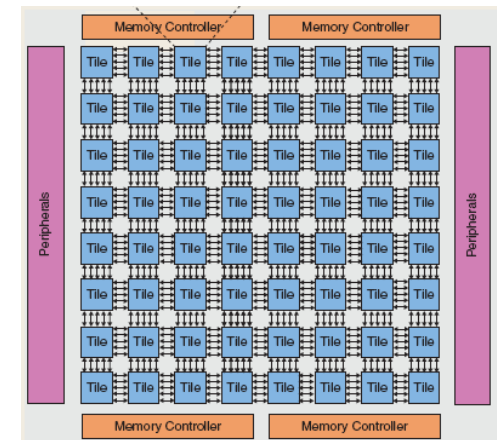
Core i7



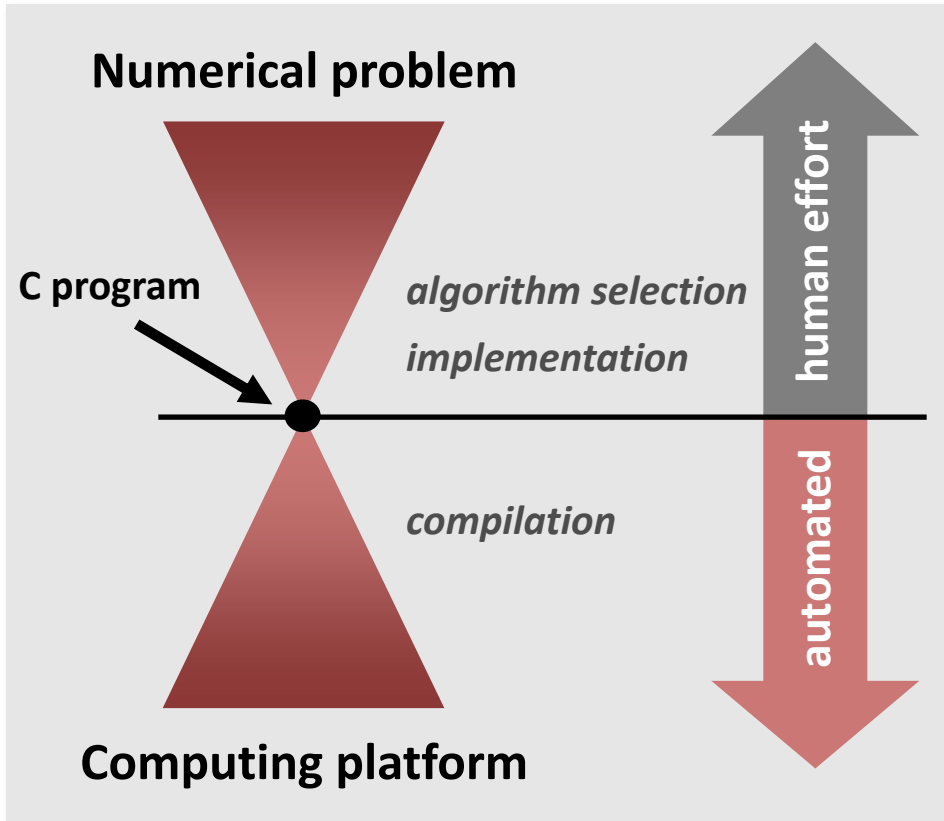
TI TNETV3020



Tilera Tile64



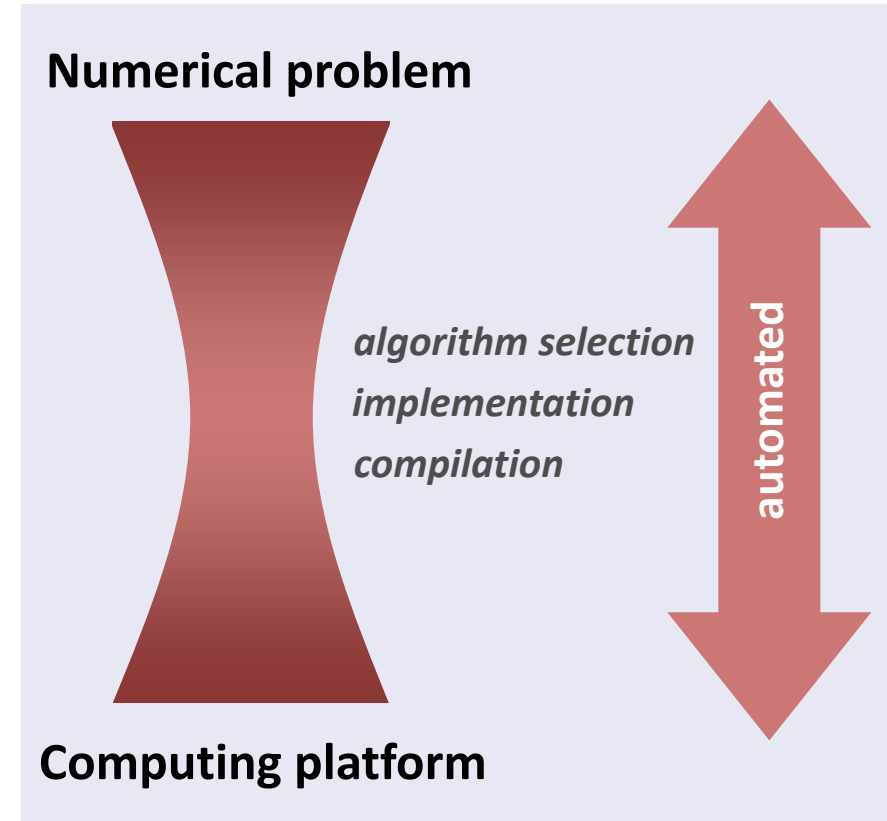
Current



C code a singularity:

- Compiler has no access to high level information
- No structural optimization
- No evaluation of choices

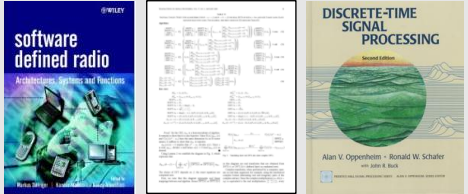
Future



Challenge: conquer the high abstraction level for *complete automation*

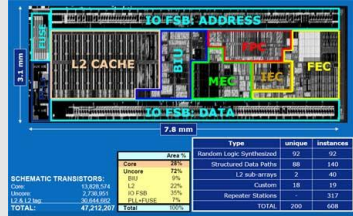
Current Solution

Algorithm knowledge



The 'Algorithm knowledge' box contains three items: the cover of the book 'software defined radio: Architectures, Systems and Functions' by Wiley; a technical document with a block diagram; and the cover of the book 'DISCRETE-TIME SIGNAL PROCESSING: Second Edition' by Alan V. Oppenheim and Ronald W. Schaefer.

Processor knowledge



The 'Processor knowledge' box contains a screenshot of a processor architecture diagram showing components like L2 CACHE, BIU, and FEC. Below the diagram is a table with the following data:

7.6 mm		Type	unique	instances
Random Logic Synthesized	92		92	
Structured Data Paths	88		140	
L2 Data Storage	2		49	
Custom	18		19	
Repeater Structures			317	
TOTAL	200		608	

Below the table is a 'SCHEMATIC TRANSISTORS' table:

Comp.	Area %
Wire	26%
Logic	12%
L2	22%
BIU	19%
BLK2USE	7%
TOTAL	100%

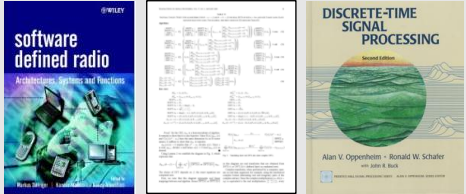


Optimal program
(Repeated for every processor)

Our Research:

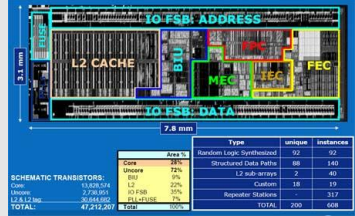
The Computer Writes the Program

Algorithm knowledge



The 'Algorithm knowledge' box contains three items: a book cover for 'software defined radio: Architectures, Systems and Functions', a technical diagram showing a signal flow with various blocks and connections, and a book cover for 'DISCRETE-TIME SIGNAL PROCESSING' by Alan V. Oppenheim and Ronald W. Schaefer.

Processor knowledge



The 'Processor knowledge' box contains a screenshot of a processor architecture diagram. The diagram shows various components like L2 CACHE, BIU, and FEC. Below the diagram is a table with columns for 'Type', 'unique', and 'instances'. The table lists various components and their counts.

Type	unique	instances
Random Logic Synthesizad	92	92
Structured Data Path	88	140
L2 cache	2	40
Custom	18	19
Repeater Structures	—	317
TOTAL	200	608

Automation:
Spiral

Optimal program
(Regenerated for every processor)

“Computer Writes the Program”

Select transform

Transform

DCT2 ▾

transform

Size

2 ▾

number of samples

Pruning

Input

unpruned ▾

number of non-zero input samples

Output

unpruned ▾

number of non-zero output samples

Select implementation options

Data type

double precision ▾

data type

Scaling

-unscaled- ▾

output scaling

Generate Code

Reset



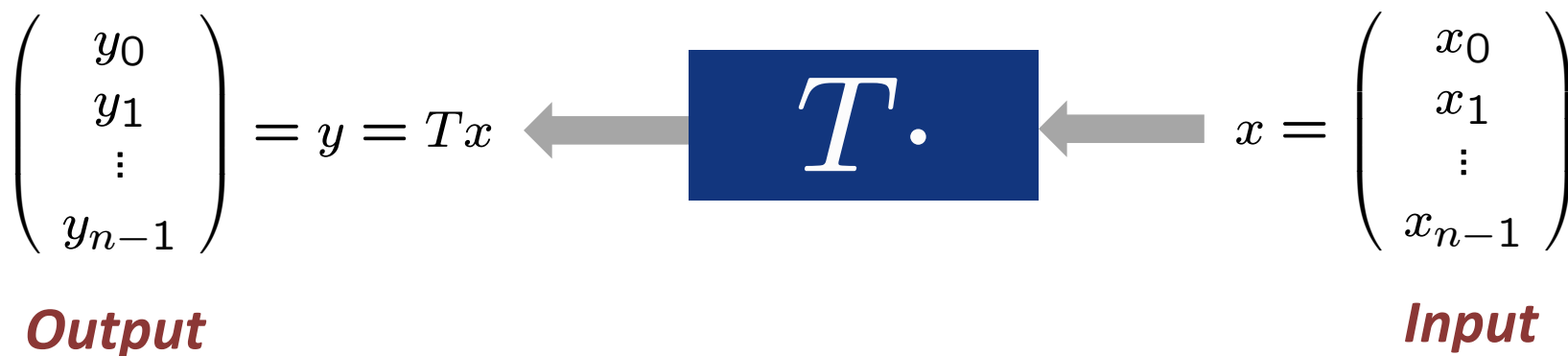
“click”

- Computer writes close to optimal code
- Vectorized, parallelized, etc.

Organization

- Spiral: Basic system
- Parallelism
- General input size
- Results
- Final remarks

Linear Transforms



Example: $T = \text{DFT}_n = [e^{-2k\ell\pi i/n}]_{0 \leq k, \ell < n}$

Algorithms: Example FFT, n = 4

Fast Fourier transform (FFT)

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

Representation using matrix algebra

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

- ***SPL (Signal processing language):*** Mathematical, declarative, point-free
- **Divide-and-conquer algorithms = breakdown rules in SPL**

Decomposition Rules (>200 for >40 Transforms)

$$\begin{aligned}
 \text{DFT}_n &\rightarrow P_{k/2,2m}^\top \left(\text{DFT}_{2m} \oplus \left(I_{k/2-1} \quad i C_{2m} \text{rDFT}_{2m}(i/k) \right) \right) \left(\text{RDFT}'_k \quad I_m \right), \quad k \text{ even,} \\
 \begin{pmatrix} \text{RDFT}_n \\ \text{RDFT}'_n \\ \text{DHT}_n \\ \text{DHT}'_n \end{pmatrix} &\rightarrow \left(P_{k/2,m}^\top \quad I_2 \right) \left(\begin{pmatrix} \text{RDFT}_{2m} \\ \text{RDFT}'_{2m} \\ \text{DHT}_{2m} \\ \text{DHT}'_{2m} \end{pmatrix} \oplus \left(I_{k/2-1} \quad i D_{2m} \begin{pmatrix} \text{rDFT}_{2m}(i/k) \\ \text{rDFT}'_{2m}(i/k) \\ \text{rDHT}_{2m}(i/k) \\ \text{rDHT}'_{2m}(i/k) \end{pmatrix} \right) \right) \left(\begin{pmatrix} \text{RDFT}'_k \\ \text{RDFT}'_k \\ \text{DHT}'_k \\ \text{DHT}'_k \end{pmatrix} \quad I_m \right), \quad k \text{ even,} \\
 \begin{pmatrix} \text{rDFT}_{2n}(u) \\ \text{rDHT}_{2n}(u) \end{pmatrix} &\rightarrow L_m^{2n} \left(I_k \quad i \begin{pmatrix} \text{rDFT}_{2m}((i+u)/k) \\ \text{rDHT}_{2m}((i+u)/k) \end{pmatrix} \right) \left(\begin{pmatrix} \text{rDFT}_{2k}(u) \\ \text{rDHT}_{2k}(u) \end{pmatrix} \quad I_m \right), \\
 \text{RDFT-3}_n &\rightarrow \left(Q_{k/2,m}^\top \quad I_2 \right) \left(I_k \quad i \text{rDFT}_{2m} \right) \left(i + 1/2 \right) / k \left(\text{RDFT-3}_k \quad I_m \right), \quad k \text{ even,} \\
 \text{DCT-2}_n &\rightarrow P_{k/2,2m}^\top \left(\text{DCT-2}_{2m} K_2^{2m} \oplus \left(I_{k/2-1} \quad N_{2m} \text{RDFT-3}_{2m}^\top \right) \right) B_n \left(L_{k/2}^{n/2} \quad I_2 \right) \left(I_m \quad \text{RDFT}'_k \right) Q_{m/2,k}, \\
 \text{DCT-3}_n &\rightarrow \text{DCT-2}_n^\top, \\
 \text{DCT-4}_n &\rightarrow \left(L_{k/2,m}^\top \quad N_{2m} \text{RDFT-3}_{2m}^\top \right) \left(L_{k/2,m}^{n/2} \quad I_2 \right) \left(\text{DCT-3}_k \right) Q_{m/2,k}, \\
 \text{DFT}_n &\rightarrow \left(\text{DFT}_k \quad I_m \right) \left(I_m \quad \text{DFT}_m \right) L_n, \quad n = km \\
 \text{DFT}_n &\rightarrow P_n \left(\text{DFT}_k \quad \text{DFT}_m \right) Q_n, \quad n = km, \text{ gcd}(k, m) = 1 \\
 \text{DFT}_n &\rightarrow \left(I_1 \oplus \text{DFT}_{p-1} \right) R_p, \quad p \text{ prime} \\
 \text{DCT-3}_n &\rightarrow \left(I_m \oplus J_m \right) L_n^n \left(\text{DCT-3}_m(1/4) \oplus \text{DCT-3}_m(3/4) \right) \\
 &\quad \cdot \left(F_2 \quad I_m \right) \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} \left(1 / \left(2 \cos \left((2k+1)\pi/4n \right) \right) \right) \\
 \text{IMDCT}_{2m} &\rightarrow \left(J_m \oplus I_m \oplus I_m \oplus J_m \right) \left(\begin{pmatrix} 1 \\ -1 \end{pmatrix} \quad I_m \right) \oplus \left(\begin{pmatrix} -1 \\ -1 \end{pmatrix} \quad I_m \right) J_{2m} \text{DCT-4}_{2m} \\
 \text{WHT}_{2^k} &\rightarrow \prod_{i=1}^t \left(I_{2^{k_1+\dots+k_{i-1}}} \text{WHT}_{2^{k_i}} \quad I_{2^{k_{i+1}+\dots+k_t}} \right), \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}$$

Decomposition rules = Algorithm knowledge in Spiral
 (from ≈100 publications)

Combining these rules yields many algorithms for every given transform

SPL to Code

SPL S Pseudo code for $y = Sx$

$A_n B_n$ <code for: $t = Bx$ >
<code for: $y = At$ >

$I_m \otimes A_n$ for (i=0; i<m; i++)
<code for:
 $y[i*n:1:i*n+n-1] = A(x[i*n:1:i*n+n-1])$ >

$A_m \otimes I_n$ for (i=0; i<n; i++)
<code for:
 $y[i:n:i+m*n-n] = A(x[i:n:i+m*n-n])$ >

D_n for (i=0; i<n; i++)
 $y[i] = D[i]*x[i];$

L_k^{km} for (i=0; i<k; i++)
 for (j=0; j<m; j++)
 $y[i*m+j] = x[j*k+i];$

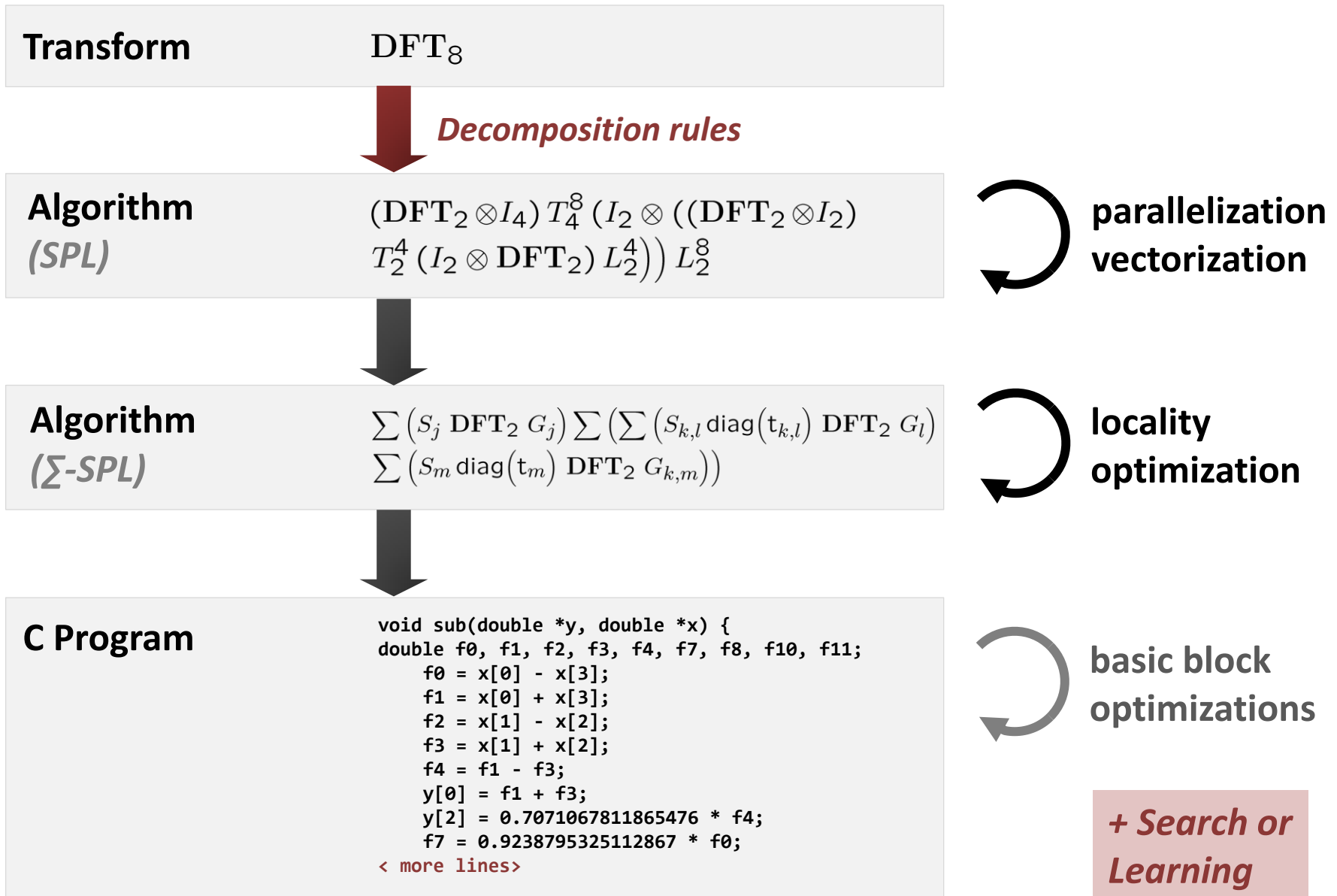
F_2 $y[0] = x[0] + x[1];$
 $y[1] = x[0] - x[1];$

$$I_m \otimes A_n = \begin{bmatrix} A_n & & \\ & \dots & \\ & & A_n \end{bmatrix}$$

Correct code: easy

fast code: very difficult

Program Generation in Spiral



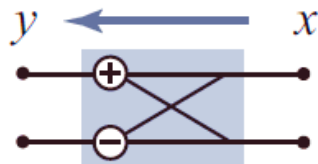
Organization

- Spiral: Basic system
- **Parallelism**
- General input size
- Results
- Final remarks

Example: Vectorization in Spiral

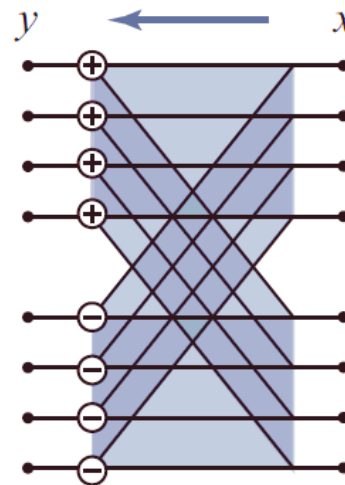
- Relationship SPL expressions \leftrightarrow vectorization?

$$y = \text{DFT}_2 x$$



one addition
one subtraction

$$y = \begin{pmatrix} \text{DFT}_2 & \text{I}_4 \end{pmatrix} x$$



one (4-way) vector addition
one (4-way) vector subtraction

Step 1: Identify “Good” Vector Constructs

- Vector length: ν
- Good (= easily vectorizable) SPL constructs:

$A \quad I_\nu$

$L_\nu^{\nu^2}, L_{\frac{2\nu}{2}}, L_\nu^{2\nu}$ *base cases*

SPL expressions recursively built from those

- **Idea:** Convert a given SPL expression into a “good” SPL expression through rewriting (structural manipulation)

Step 2: Find Manipulation Rules

$$\mathbb{L}_n^{n\nu} \rightarrow \left(\mathbb{I}_{n/\nu} \quad \mathbb{L}_\nu^{\nu^2} \right) \left(\mathbb{L}_{n/\nu}^n \quad \mathbb{I}_\nu \right)$$

$$\mathbb{L}_\nu^{n\nu} \rightarrow \left(\mathbb{L}_\nu^n \quad \mathbb{I}_\nu \right) \left(\mathbb{I}_{n/\nu} \quad \mathbb{L}_\nu^{\nu^2} \right)$$

$$\mathbb{L}_m^{mn} \rightarrow \left(\mathbb{L}_m^{mn/\nu} \quad \mathbb{I}_\nu \right) \left(\mathbb{I}_{mn/\nu^2} \quad \mathbb{L}_\nu^{\nu^2} \right) \left(\mathbb{I}_{n/\nu} \quad \mathbb{L}_{m/\nu}^m \quad \mathbb{I}_\nu \right)$$

$$\mathbb{I}_l \quad \mathbb{L}_n^{kmn} \quad \mathbb{I}_r \rightarrow \left(\mathbb{I}_l \quad \mathbb{L}_n^{kn} \quad \mathbb{I}_{mr} \right) \left(\mathbb{I}_{kl} \quad \mathbb{L}_n^{mn} \quad \mathbb{I}_r \right)$$

$$\mathbb{I}_l \quad \mathbb{L}_n^{kmn} \quad \mathbb{I}_r \rightarrow \left(\mathbb{I}_l \quad \mathbb{L}_{kn}^{kmn} \quad \mathbb{I}_r \right) \left(\mathbb{I}_l \quad \mathbb{L}_{mn}^{kmn} \quad \mathbb{I}_r \right)$$

$$\mathbb{I}_l \quad \mathbb{L}_{km}^{kmn} \quad \mathbb{I}_r \rightarrow \left(\mathbb{I}_{kl} \quad \mathbb{L}_m^{mn} \quad \mathbb{I}_r \right) \left(\mathbb{I}_l \quad \mathbb{L}_k^{kn} \quad \mathbb{I}_{mr} \right)$$

$$\mathbb{I}_l \quad \mathbb{L}_{kn}^{kmn} \quad \mathbb{I}_r \rightarrow \left(\mathbb{I}_l \quad \mathbb{L}_k^{kmn} \quad \mathbb{I}_r \right) \left(\mathbb{I}_l \quad \mathbb{L}_n^{kmn} \quad \mathbb{I}_r \right)$$

Manipulation rules = Processor knowledge in Spiral

$$\left(\mathbb{I}_m \quad A^{n \times n} \right) \mathbb{L}_m^{mn} \rightarrow \left(\mathbb{I}_{m/\nu} \quad \mathbb{L}_\nu^{n\nu} \left(A^{n \times n} \quad \mathbb{I}_\nu \right) \right) \left(\mathbb{L}_{m/\nu}^{mn/\nu} \quad \mathbb{I}_\nu \right)$$

$$\mathbb{L}_n^{mn} \left(\mathbb{I}_m \quad A^{n \times n} \right) \rightarrow \left(\mathbb{L}_n^{mn/\nu} \quad \mathbb{I}_\nu \right) \left(\mathbb{I}_{m/\nu} \quad \left(A^{n \times n} \quad \mathbb{I}_\nu \right) \mathbb{L}_n^{n\nu} \right)$$

$$\left(\mathbb{I}_k \quad \left(\mathbb{I}_m \quad A^{n \times n} \right) \mathbb{L}_m^{mn} \right) \mathbb{L}_k^{kmn} \rightarrow \left(\mathbb{L}_k^{km} \quad \mathbb{I}_n \right) \left(\mathbb{I}_m \quad \left(\mathbb{I}_k \quad A^{n \times n} \right) \mathbb{L}_k^{kn} \right) \left(\mathbb{L}_m^{mn} \quad \mathbb{I}_k \right)$$

$$\mathbb{L}_{mn}^{kmn} \left(\mathbb{I}_k \quad \mathbb{L}_n^{mn} \left(\mathbb{I}_m \quad A^{n \times n} \right) \right) \rightarrow \left(\mathbb{L}_n^{mn} \quad \mathbb{I}_k \right) \left(\mathbb{I}_m \quad \mathbb{L}_n^{kn} \left(\mathbb{I}_k \quad A^{n \times n} \right) \right) \left(\mathbb{L}_m^{km} \quad \mathbb{I}_n \right)$$

$$\overline{AB} \rightarrow \overline{A} \overline{B}$$

$$\overline{A^{m \times m} \quad \mathbb{I}_\nu} \rightarrow \left(\mathbb{I}_m \quad \mathbb{L}_\nu^{2\nu} \right) \left(\overline{A^{m \times m}} \quad \mathbb{I}_\nu \right) \left(\mathbb{I}_m \quad \mathbb{L}_\nu^{2\nu} \right)$$

$$\overline{\mathbb{I}_m \quad A^{n \times n}} \rightarrow \mathbb{I}_m \quad \overline{A^{n \times n}}$$

$$\overline{D} \rightarrow \left(\mathbb{I}_{n/\nu} \quad \mathbb{L}_\nu^{2\nu} \right) \vec{D} \left(\mathbb{I}_{n/\nu} \quad \mathbb{L}_\nu^{2\nu} \right)$$

$$\overline{P} \rightarrow P \quad \mathbb{I}_2$$

Example

$$\underbrace{\text{DFT}_{mn}}_{\text{vec}(\nu)} \rightarrow \underbrace{(\text{DFT}_m \quad \text{I}_n) \text{T}_n^{mn} (\text{I}_m \quad \text{DFT}_n) \text{L}_m^{mn}}_{\text{vec}(\nu)}$$

...

...

...

$$\rightarrow \underbrace{\left(\text{I}_{\frac{mn}{\nu}} \quad \text{L}_{\nu}^{2\nu} \right)}_{\text{red}} \underbrace{\left(\overline{\text{DFT}_m} \quad \text{I}_{\frac{n}{\nu}} \quad \text{I}_{\nu} \right)}_{\text{blue}} \text{T}_n^{mn}$$

$$\underbrace{\left(\text{I}_{\frac{m}{\nu}} \quad \left(\text{L}_{\nu}^{2n} \quad \text{I}_{\nu} \right) \left(\text{I}_{\frac{2n}{\nu}} \quad \text{L}_{\nu}^{\nu^2} \right) \left(\text{I}_{\frac{n}{\nu}} \quad \text{L}_2^{2\nu} \quad \text{I}_{\nu} \right) \right)}_{\text{red}} \underbrace{\left(\overline{\text{DFT}_n} \quad \text{I}_{\nu} \right)}_{\text{blue}} \underbrace{\left(\text{L}_{\frac{m}{\nu}}^{\frac{mn}{\nu}} \quad \text{L}_2^{2\nu} \right)}_{\text{red}}$$

vectorized arithmetic

vectorized data accesses

Automatically Generate Base Case Library

- **Goal:** Given instruction set, generate base cases

$$\nu = 4 : \quad \{ L_2^4, I_2 \quad L_2^4, L_2^4 \quad I_2, L_2^8, L_4^8 \}$$

- **Idea:** Instructions as matrices + search

```
y = _mm_unpacklo_ps(x0, x1);
```

```
y = _mm_shuffle_ps(x0, x1, _MM_SHUFFLE(1,2,1,2));
```

```
y = _mm_shuffle_ps(x0, x1, _MM_SHUFFLE(3,4,3,4));
```

$$y = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} \vec{x}_0 \\ \vec{x}_1 \end{bmatrix}$$

$$y = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} \vec{x}_0 \\ \vec{x}_1 \end{bmatrix}$$

$$y = \begin{bmatrix} 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \vec{x}_0 \\ \vec{x}_1 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ \hline 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

```
y0 = _mm_shuffle_ps(x0, x1,
                    _MM_SHUFFLE(1,2,1,2));
y1 = _mm_shuffle_ps(x0, x1,
                    _MM_SHUFFLE(3,4,3,4));
```



$L_2^4 \otimes I_2$
no base case
Base case

Same Approach for Different Paradigms

Threading:

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

Vectorization:

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

GPUs:

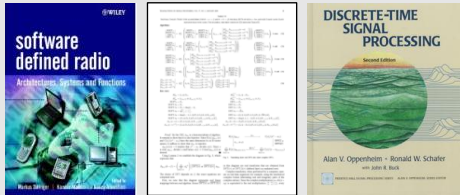
$$\begin{aligned}
 \underbrace{(\text{DFT}_{r^k})}_{\text{gpu}(t,c)} &\rightarrow \underbrace{\left(\prod_{i=0}^{k-1} \text{L}_r^{r^k} (\text{I}_{r^{k-1}} \otimes \text{DFT}_r) (\text{L}_{r^{k-i-1}} (\text{I}_{r^i} \otimes \text{T}_{r^{k-i-1}}) \text{L}_{r^{i+1}}^{r^k}) \right)}_{\text{gpu}(t,c)} \text{R}_r^{r^k}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

Verilog for FPGAs:

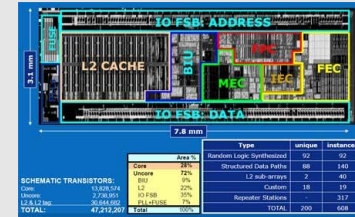
$$\begin{aligned}
 \underbrace{(\text{DFT}_{r^k})}_{\text{stream}(r^s)} &\rightarrow \underbrace{\left[\prod_{i=0}^{k-1} \text{L}_r^{r^k} (\text{I}_{r^{k-1}} \otimes \text{DFT}_r) (\text{L}_{r^{k-i-1}} (\text{I}_{r^i} \otimes \text{T}_{r^{k-i-1}}) \text{L}_{r^{i+1}}^{r^k}) \right]}_{\text{stream}(r^s)} \text{R}_r^{r^k} \\
 &\dots \\
 &\rightarrow \left[\prod_{i=0}^{k-1} \underbrace{\text{L}_r^{r^k}}_{\text{stream}(r^s)} \underbrace{(\text{I}_{r^{k-1}} \otimes \text{DFT}_r)}_{\text{stream}(r^s)} \underbrace{(\text{L}_{r^{k-i-1}} (\text{I}_{r^i} \otimes \text{T}_{r^{k-i-1}}) \text{L}_{r^{i+1}}^{r^k})}_{\text{stream}(r^s)} \right] \text{R}_r^{r^k} \\
 &\dots \\
 &\rightarrow \left[\prod_{i=0}^{k-1} \underbrace{\text{L}_r^{r^k}}_{\text{stream}(r^s)} (\text{I}_{r^{k-s-1}} \otimes_s (\text{I}_{r^{s-1}} \otimes \text{DFT}_r)) \underbrace{\text{T}'_i}_{\text{stream}(r^s)} \right] \text{R}_r^{r^k} \\
 &\quad \text{stream}(r^s)
 \end{aligned}$$

- Rigorous, correct by construction
- Overcomes compiler limitations

Algorithm knowledge



Processor knowledge



Automation:
Spiral

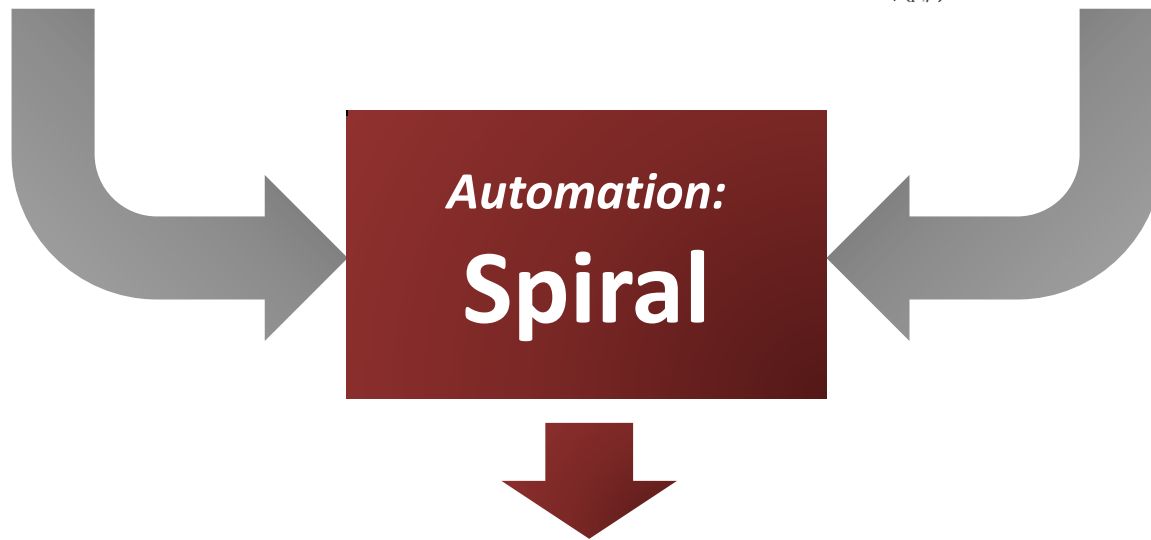
Optimal program
(Regenerated for every processor)

Decomposition rules

$$\begin{aligned} \text{DFT}_n &\rightarrow P_{k/2,2m}^\top (\text{DFT}_{2m} \oplus (I_{k/2-1} \otimes_i C_{2m} \text{rDFT}_{2m}(i/k))) (\text{RDFT}'_k \otimes I_m) \\ \begin{pmatrix} \text{rDFT}_{2n}(u) \\ \text{rDHT}_{2n}(u) \end{pmatrix} &\rightarrow L_m^{2n} \left(I_k \otimes_i \begin{pmatrix} \text{rDFT}_{2m}((i+u)/k) \\ \text{rDHT}_{2m}((i+u)/k) \end{pmatrix} \right) \left(\begin{pmatrix} \text{rDFT}_{2k}(u) \\ \text{rDHT}_{2k}(u) \end{pmatrix} \otimes I_m \right) \\ \text{RDFT-3}_n &\rightarrow (Q_{k/2,m}^\top \otimes I_2) (I_k \otimes_i \text{rDFT}_{2m}(i+1/2/k)) (\text{RDFT-3}_k \otimes I_m) \end{aligned}$$

Manipulation rules

$$\begin{aligned} \underbrace{A_m \otimes I_n}_{\text{smp}(p,\mu)} &\rightarrow \underbrace{(L_m^{mp} \otimes I_{n/p}) (I_p \otimes (A_m \otimes I_{n/p})) (L_p^{mp} \otimes I_{n/p})}_{\text{smp}(p,\mu)} \\ \underbrace{I_m \otimes A_n}_{\text{smp}(p,\mu)} &\rightarrow I_p \otimes_{\parallel} (I_{m/p} \otimes A_n) \\ \underbrace{(P \otimes I_n)}_{\text{smp}(p,\mu)} &\rightarrow (P \otimes I_{n/\mu}) \bar{\otimes} I_\mu \end{aligned}$$



Automation:
Spiral

Optimal program
(Regenerated for every processor)

Organization

- Spiral: Basic system
- Parallelism
- **General input size**
- Results
- Final remarks

Challenge: General Size Libraries

So far:

Code specialized to fixed input size

```
DFT_384(x, y) {  
  ...  
  for(i = ...) {  
    t[2i]    = x[2i] + x[2i+1]  
    t[2i+1] = x[2i] - x[2i+1]  
  }  
  ...  
}
```

- Algorithm fixed
- Nonrecursive code

Challenge:

Library for general input size

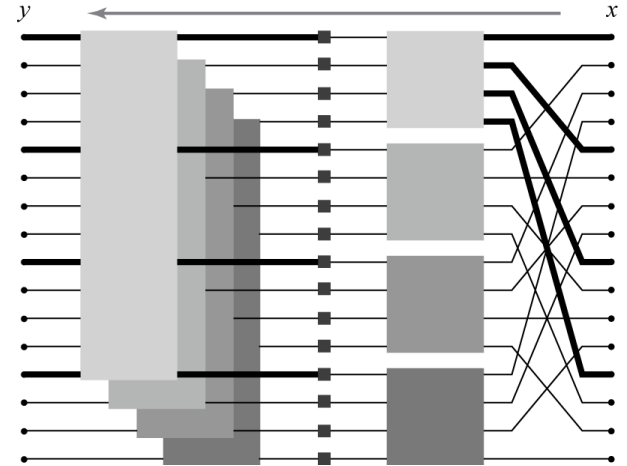
```
DFT(n, x, y) {  
  ...  
  for(i = ...) {  
    DFT_strided(m, x+mi, y+i, 1, k)  
  }  
  ...  
}
```

- Algorithm cannot be fixed
- Recursive code
- Creates many challenges

Challenge: 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 that increases locality (e.g., FFTW 2.x)

```
void DFT(int n, cpx *y, cpx *x) {  
    int k = choose_dft_radix(n);
```

```
    for (int i=0; i < k; ++i)  
        DFTrec(m, y + m*i, x + i, k, 1);  
    for (int j=0; j < m; ++j)  
        DFTscaled(k, y + j, t[j], m);
```

```
}
```

Σ -SPL : Basic Idea

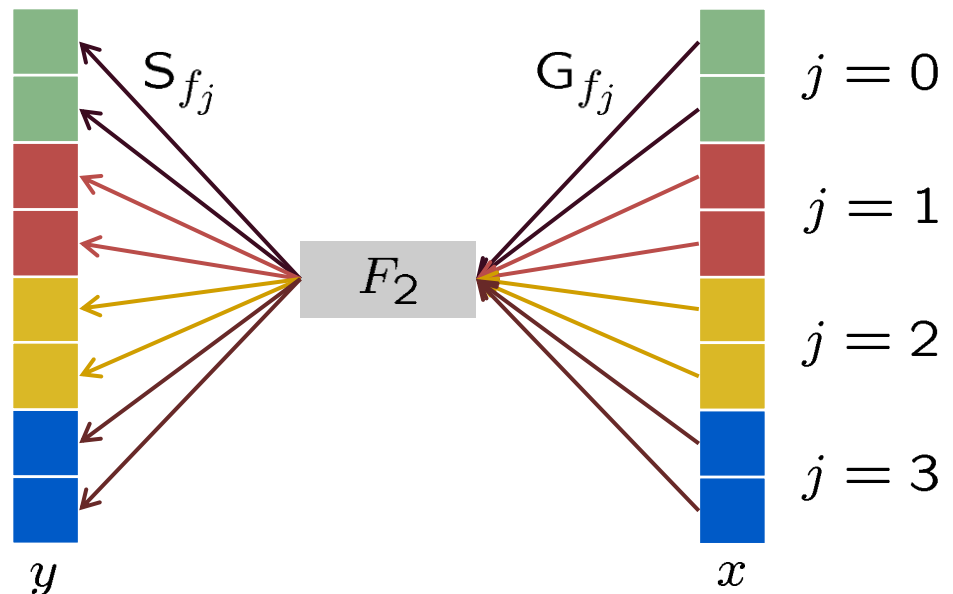
- Four additional matrix constructs: Σ , G , S , Perm

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

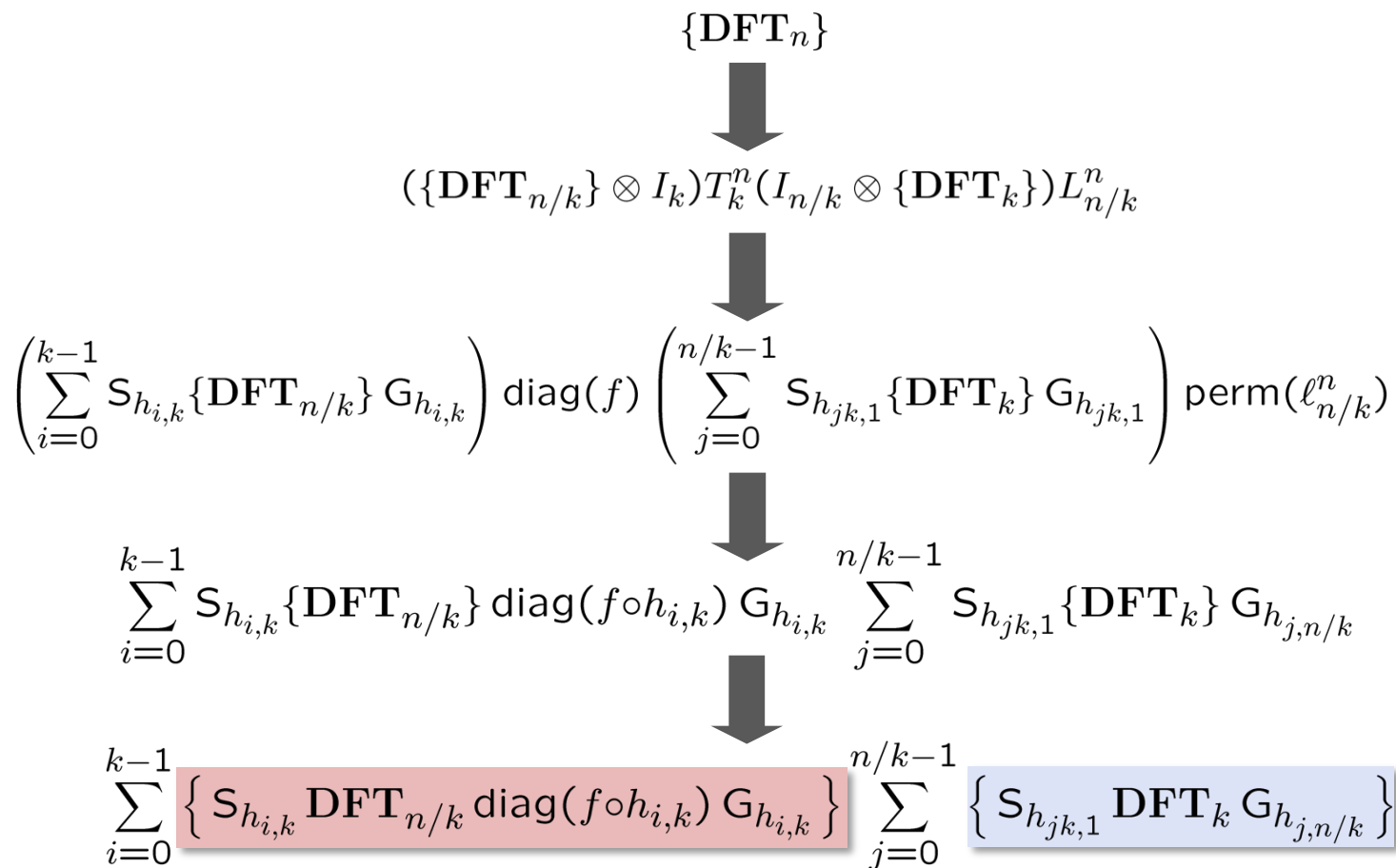
- Σ -SPL formulas = matrix factorizations

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

$$y = \begin{bmatrix} F_2 & & & \\ & F_2 & & \\ & & F_2 & \\ & & & F_2 \end{bmatrix} x$$



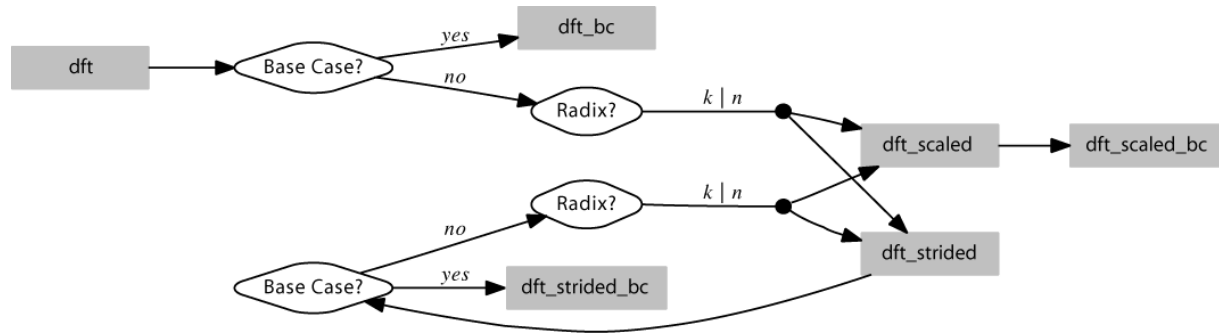
Find Recursion Step Closure



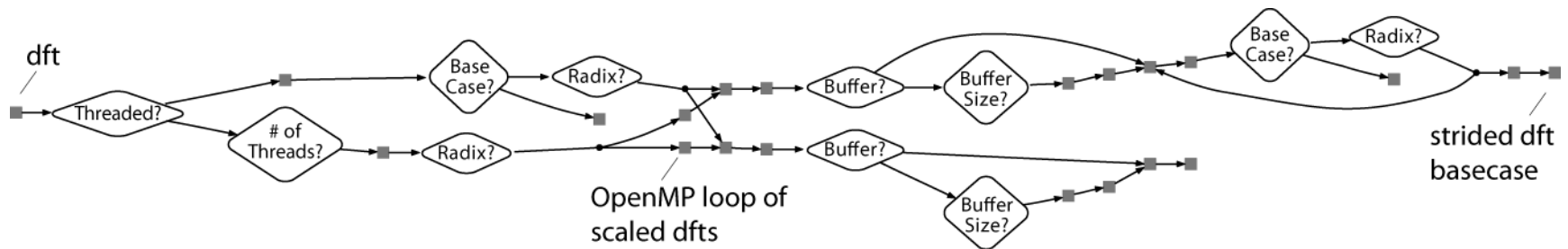
Repeat until closure

Recursion Step Closure: Examples

DFT: scalar code



DFT: full-fledged (vectorized and parallel code)



Summary: Complete Automation for Transforms

- **Memory hierarchy optimization**

Rewriting and search for algorithm selection

Rewriting for loop optimizations

- **Vectorization**

Rewriting

- **Parallelization**

Rewriting

fixed input size code

- **Derivation of library structure**

Rewriting

Other methods *general input size library*

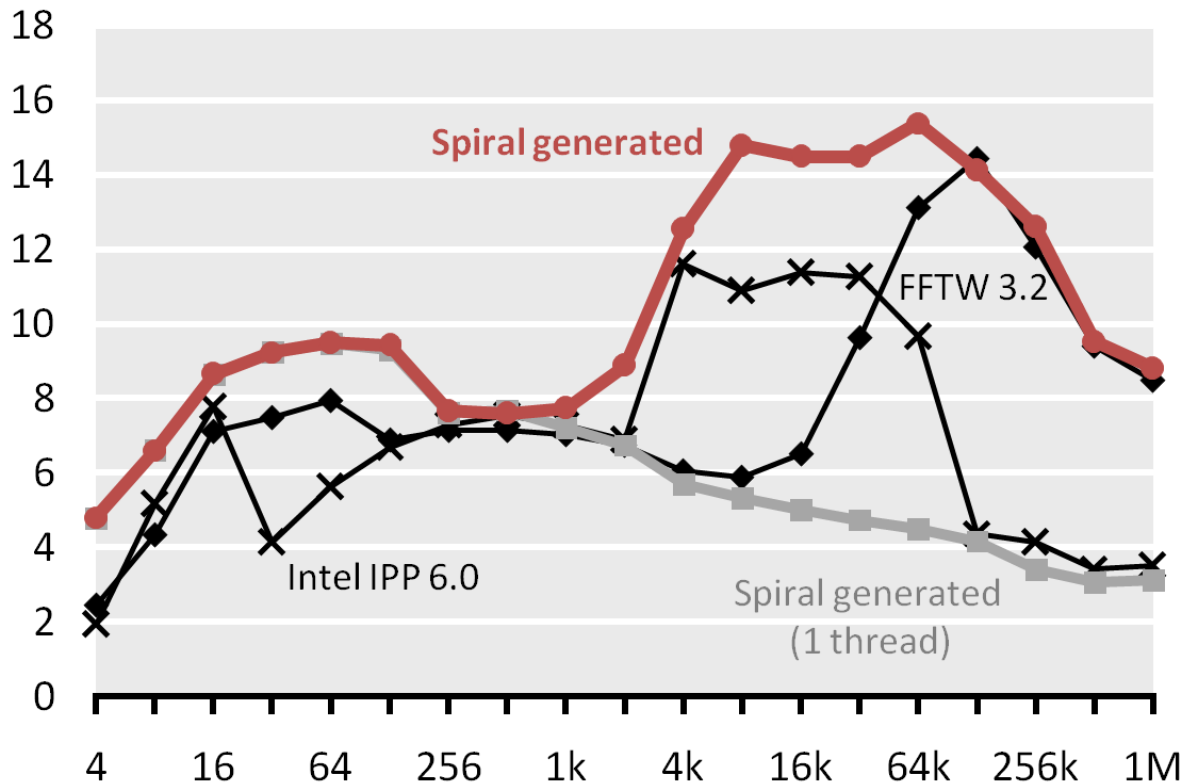
Organization

- Spiral: Basic system
- Parallelism
- General input size
- **Results**
- Final remarks

DFT on Intel Multicore

Complex DFT (Intel Core i7, 2.66 GHz, 4 cores)

Performance [Gflop/s] vs. input size



$$\begin{aligned}
 \text{DFT}_n &\rightarrow (\text{DFT}_k \otimes I_m) T_m^n (I_k \otimes \text{DFT}_m) L_k^n \\
 \text{DFT}_n &\rightarrow P_{k/2, 2m}^\top (\text{DFT}_{2m} \oplus (I_{k/2-1} \otimes_i C_{2m} \text{rDFT}_{2m}(i/k))) (\text{RDFT}_k \otimes I_m) \\
 \text{RDFT}_n &\rightarrow (P_{k/2, m}^\top \otimes I_2) (\text{RDFT}_{2m} \oplus (I_{k/2-1} \otimes_i D_{2m} \text{rDFT}_{2m}(i/k))) (\text{RDFT}_k \otimes I_m) \\
 \text{rDFT}_{2n}(u) &\rightarrow L_m^{2n} (I_k \otimes_i \text{rDFT}_{2m}((i+u)/k)) (\text{rDFT}_{2k}(u) \otimes I_m)
 \end{aligned}$$

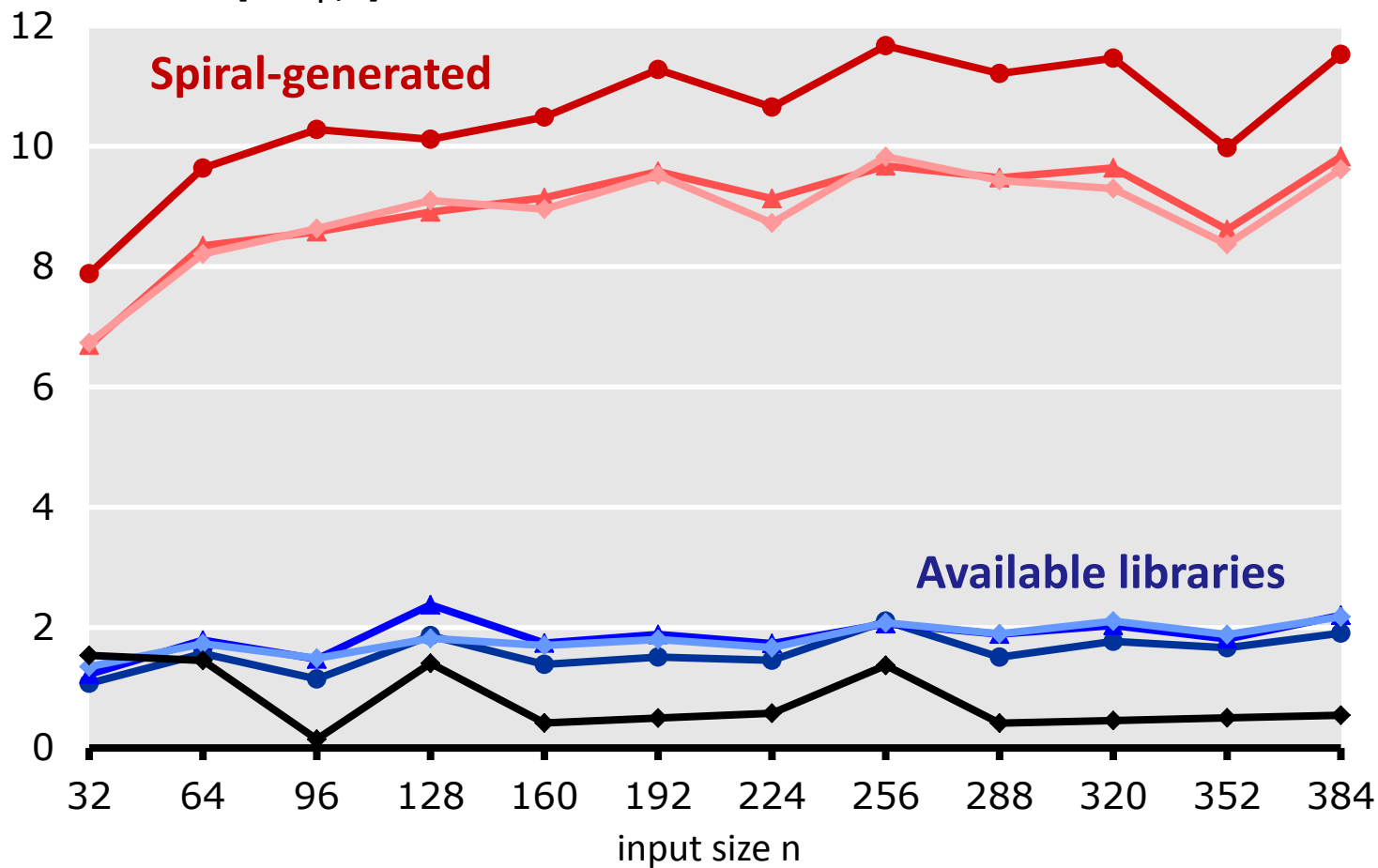


**5MB vectorized, threaded,
general-size, adaptive library**

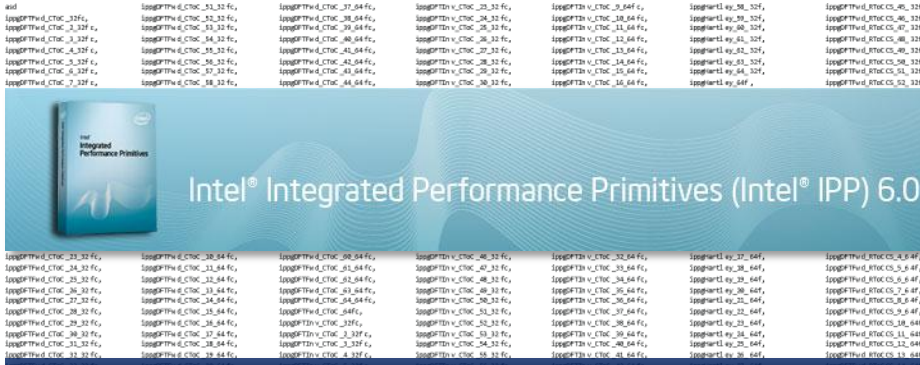
Often it Looks Like That

DCTs on 2.66 GHz Core2 (4-way SSSE3)

Performance [Gflop/s]



Computer generated Functions for Intel IPP 6.0



3984 C functions
1M lines of code

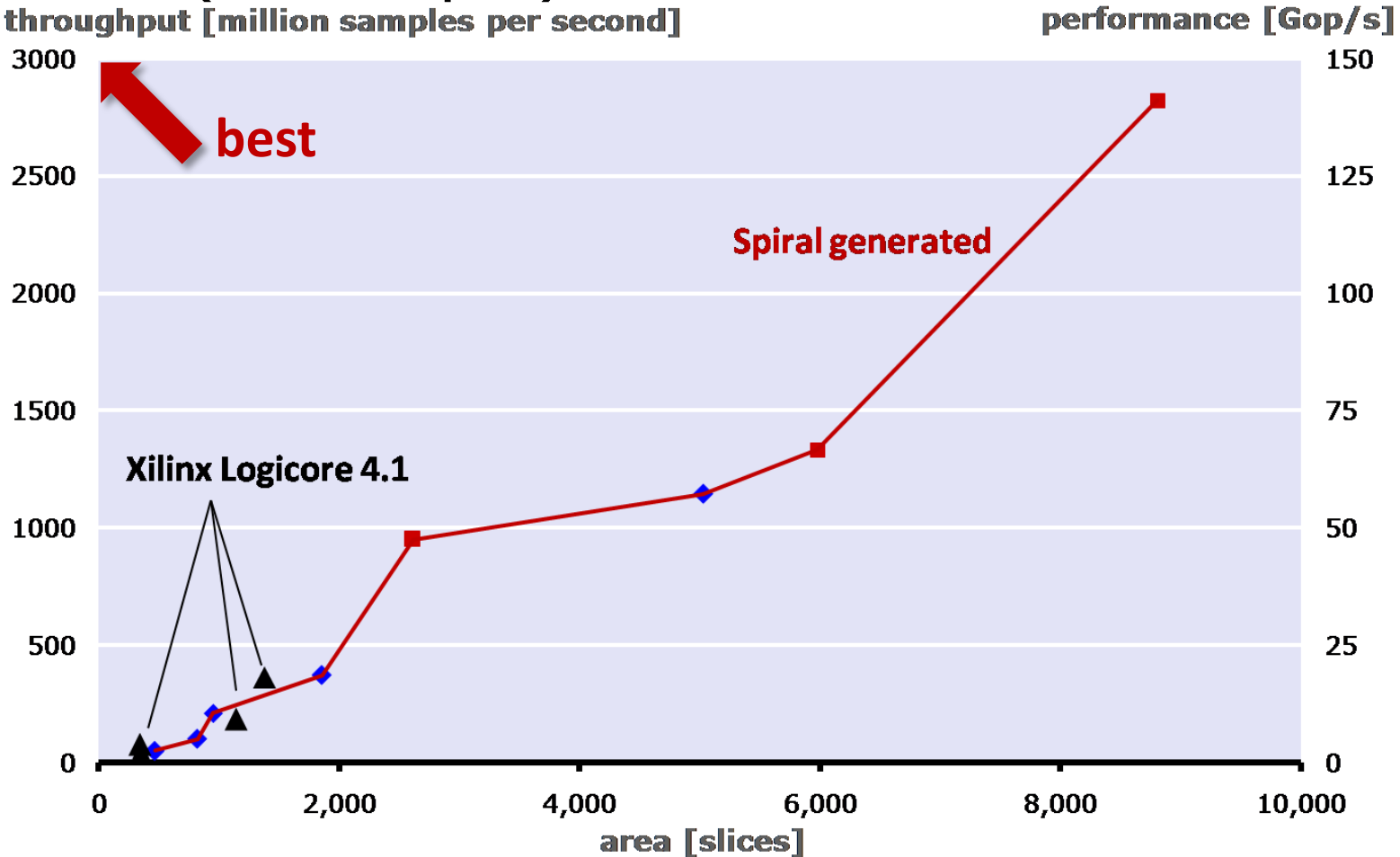
Transforms: DFT (fwd+inv), RDFT (fwd+inv), DCT2, DCT3, DCT4, DHT, WHT
Sizes: 2–64 (DFT, RDFT, DHT); 2-powers (DCTs, WHT)
Precision: single, double
Data type: scalar, SSE, AVX (DFT, DCT), LRB (DFT)

Computer generated

Results: SpiralGen Inc.

Mapping to FPGAs

DFT 1024 (16 bit fixed point) on Xilinx Virtex-5 FPGA
throughput [million samples per second]



Organization

- Spiral: Basic system
- Parallelism
- General input size
- Results
- **Final remarks**

Beyond Transforms

Matrix-Matrix Multiplication



$$\text{MMM}_{1,1,1} \rightarrow (\cdot)_1$$

$$\text{MMM}_{m,n,k} \rightarrow (\otimes)_{m/m_b \times 1} \otimes \text{MMM}_{m_b,n,k}$$

$$\text{MMM}_{m,n,k} \rightarrow \text{MMM}_{m,n_b,k} \otimes (\otimes)_{1 \times n/n_b}$$

$$\text{MMM}_{m,n,k} \rightarrow ((\Sigma_{k/k_b} \circ (\cdot)_{k/k_b}) \otimes \text{MMM}_{m,n,k_b}) \circ ((L_{k/k_b}^{m k/k_b} \otimes I_{k_b}) \times I_{kn})$$

$$\text{MMM}_{m,n,k} \rightarrow (L_m^{mn/n_b} \otimes I_{n_b}) \circ ((\otimes)_{1 \times n/n_b} \otimes \text{MMM}_{m,n_b,k}) \circ (I_{km} \times (L_{n/n_b}^{kn/n_b} \otimes I_{n_b}))$$

JPEG 2000 (Wavelet, EBCOT)



$$\text{SC}(\chi_{m,n}, \sigma_{m,n}) : (\mathbb{Z}_2^9 \times \mathbb{Z}_2^9) \rightarrow (\mathbb{N}, \mathbb{Z}_2)$$

$$(I \times \text{xor}_2) \circ (T_{SC} \times I) \circ (H \times V \times I) \circ (L_4^2 \times G_4) \circ \left(\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix} \right) \times \left(\begin{pmatrix} 1 \\ 1 \end{pmatrix} \right) \right)$$

$$H, V : (\mathbb{Z}_2^9 \times \mathbb{Z}_2^9) \rightarrow \mathbb{N}$$

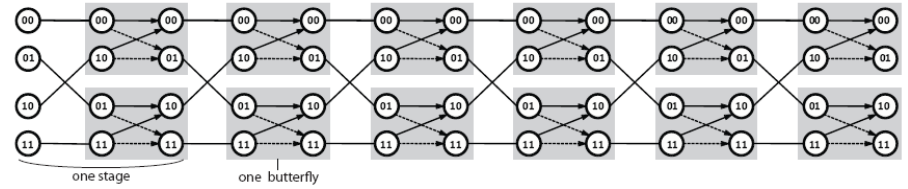
$$H : h \circ (f \times f) \circ (G_1 \times C_{-2} \times G_1 \times G_7 \times C_{-2} \times G_7) \circ L_4^2 \circ \left(\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix} \right) \times \left(\begin{pmatrix} 1 \\ 1 \end{pmatrix} \right) \right)$$

$$V : h \circ (f \times f) \circ (G_3 \times C_{-2} \times G_3 \times G_5 \times C_{-2} \times G_5) \circ L_4^2 \circ \left(\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix} \right) \times \left(\begin{pmatrix} 1 \\ 1 \end{pmatrix} \right) \right)$$

$$f : \text{mul}_2 \circ (I \times \text{sub}_2) \circ (I \times C_1 \times \text{mul}_2)$$

$$h : \text{min}_2 \circ (C_1 \times \text{max}_2) \circ (C_{-1} \times \text{sum}_2)$$

Viterbi Decoder



$$\mathbf{F}_{K,F} \rightarrow \prod_{i=1}^F \left((I_{2^{K-2}} \otimes_j B_{F-i,j}) L_{2^{K-2}}^{2^{K-1}} \right)$$

Synthetic Aperture Radar (SAR)



$$\text{SAR} \rightarrow 2\text{D-iDFT} \circ \text{Interpl} \circ \text{MatchFilt} \circ \text{prep}$$

$$2\text{D-iDFT} \rightarrow \text{iDFT} \otimes \text{iDFT}$$

$$\text{MatchFilt} \rightarrow \text{Filt} \circ (I \times C_f)$$

$$\text{Filt} \rightarrow (I \otimes (\cdot))$$

$$\text{Interpl} \rightarrow (\Sigma \otimes I) \circ (I \otimes_j S_{x_j \otimes y_j}) \circ \text{Filt} \circ ((I \otimes \mathbf{1} \otimes I) \times I) \circ (I \times C_{i \otimes_j g_j})$$

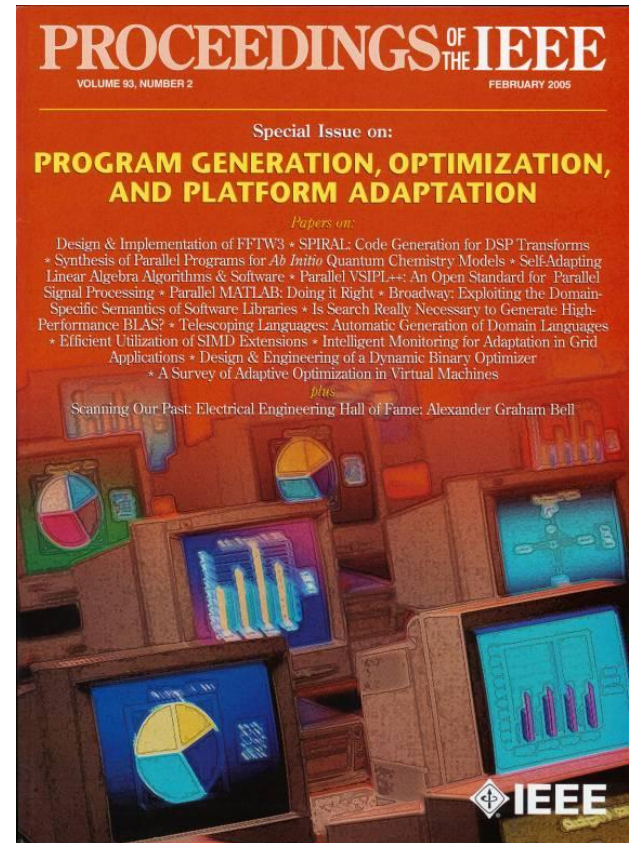
Related Work: Autotuning

- **Goal:** (Partial) automation of runtime optimization

- **Projects**

- ATLAS (U. Tennessee)
- FFTW (MIT)
- BeBop/OSKI (U. Berkeley)
- Spiral (Carnegie Mellon)

- Tensor contractions (Ohio State)
- Fenics (some universities)
- FLAME (UT Austin)
- Adaptive sorting (UIUC)
- **<many others>**



Proceedings of the IEEE special issue, Feb. 2005

Spiral: Summary

■ Spiral:

Successful approach to automating the development of computing software

Commercial proof-of-concept



■ Key ideas:

Algorithm knowledge:

Domain specific symbolic representation

Platform knowledge:

Tagged rewrite rules, SIMD specification

DFT₆₄



```
void dft64(float *Y, float *X) {
  __m512 U912, U913, U914, U915, ...
  __m512 *a2153, *a2155;
  a2153 = ((__m512 *) X);  s1107 = *(a2153);
  s1108 = *((a2153 + 4));  t1323 = _mm512_add_ps(s1107, s1108);
  t1324 = _mm512_sub_ps(s1107, s1108);
  <many more lines>
  U926 = _mm512_swizupconv_r32(...);
  s1121 = _mm512_madd231_ps(_mm512_mul_ps(_mm512_mask_or_pi(
    _mm512_set_1to16_ps(0.70710678118654757), 0xAAAA, a2154, U926), t1341),
    _mm512_mask_sub_ps(_mm512_set_1to16_ps(0.70710678118654757), ...),
    _mm512_swizupconv_r32(t1341, MM_SWIZ_REG_CDAB));
  U927 = _mm512_swizupconv_r32
  <many more lines>
}
```

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

$$\underbrace{\text{I}_m \otimes \text{A}_n}_{\text{smp}(p, \mu)} \rightarrow \text{I}_p \otimes_{\parallel} (\text{I}_{m/p} \otimes \text{A}_n)$$

Compilers

Programming languages
Program generation

Symbolic Computation
Rewriting

*Automating
High-Performance
Numerical Software
Development*

**Search
Learning**

**Algorithms
Mathematics**

**Software
Scientific Computing**