# How to Write Fast Numerical Code

Spring 2015
*Lecture:* Spiral (Computer generation of FFT code)

**Instructor:** Markus Püschel

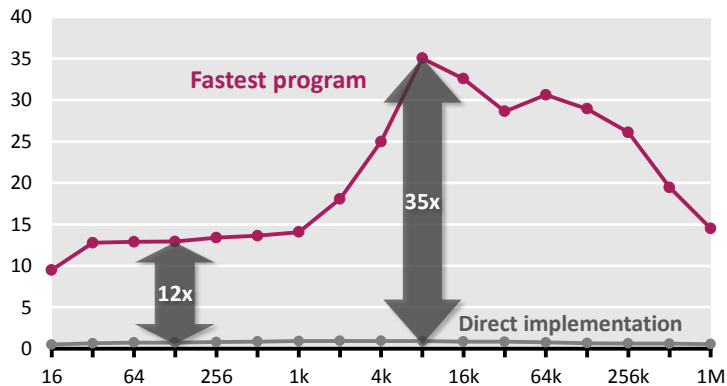**TA:** Gagandeep Singh, Daniele Spampinato, Alen Stojanov

**ETH**
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

---

# The Problem: Example DFT

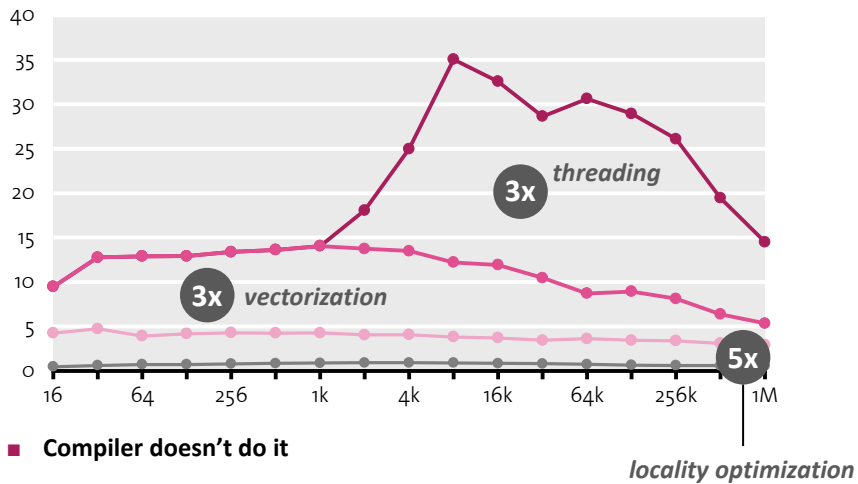**DFT on Intel Core i7 (4 Cores, 2.66 GHz)**
Performance [Gflop/s]



Fastest program

35x

12x

Direct implementation

- **Same number of operations**
- **Best compiler**

## DFT: Analysis

**DFT (single precision) on Intel Core i7 (4 cores, 2.66 GHz)**
Performance [Gflop/s]



**3x** *vectorization*

**3x** *threading*

**5x**

*locality optimization*

- **Compiler doesn't do it**
- **Doing by hand: Very tough**

## Our Goal:

Computer writes high performance library code

Generate Code

*"click"*

## Viterbi Decoder

**Select convolutional code**
Select a preset code or customize parameters

- *custom*
- Voyager
- NASA-DSN
- CCSDS/NASA-GSFC
- WiMax
- CDMA IS-95A
- LTE (3GPP - Long Term Evolution)
- UWB (802.15)
- CDMA 2000
- Cassini
- Mars Pathfinder & Stereo

| | | |
|---|---|---|
| rate | 1 / 2 | code rate (?) |
| K | 7 | constraint length (?) |
| polynomials | 109 | polynomials for the code in decimal notation |
| | 79 | (?) |

**Select implementation options**

| | | |
|---|---|---|
| frame length | 2048 | unpadded frame length |
| Vectorization level | scalar C | type of code (?) |

Generate Code   Reset

**@ www.spiral.net**

### DFT IP Cores

| parameter | value | range | explanation |
|---|---|---|---|
| **Problem specification** | | | |
| transform size | 64 | 4–32768 | Number of samples (?) |
| direction | forward | | forward or inverse DFT (?) |
| data type | fixed point | | fixed or floating point (?) |
| | 16 bits | 4–32 bits | fixed point precision (?) |
| | unscaled | | scaling mode (?) |
| **Parameters controlling implementation** | | | |
| architecture | fully streaming | | iterative or fully streaming (?) |
| radix | 2 | 2, 4, 8, 16, 32, 64 | size of DFT basic block (?) |
| streaming width | 2 | 2–64 | number of complex words per cycle (?) |
| data ordering | natural in / natural out | | natural or digit-reversed data order (?) |
| BRAM budget | 1000 | | maximum # of BRAMs to utilize (-1 for no limit) (?) |

Generate Verilog   Reset

---

# Possible Approach:

Capturing algorithm knowledge:
*Domain-specific languages (DSLs)*

Structural optimization:
*Rewriting systems*

High performance code style:
*Compiler*

Decision making for choices:
*Machine learning*

# Organization

- ***Spiral: Basic system***

- Vectorization

- General input size

- Results

- Final remarks

# 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*

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

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

# Decomposition Rules (>200 for >40 Transforms)

$$\mathbf{DFT}_n \to P^\top_{k/2,2m}\left(\mathbf{DFT}_{2m}\oplus\left(I_{k/2-1}\otimes_i C_{2m}\,\mathbf{rDFT}_{2m}(i/k)\right)\right)\left(\mathbf{RDFT}'_k\otimes I_m\right),\quad k\text{ even,}$$

$$\begin{vmatrix}\mathbf{RDFT}_n\\ \mathbf{RDFT}'_n\\ \mathbf{DHT}_n\\ \mathbf{DHT}'_n\end{vmatrix} \to (P^\top_{k/2,m}\otimes I_2)\left(\begin{vmatrix}\mathbf{RDFT}_{2m}\\ \mathbf{RDFT}'_{2m}\\ \mathbf{DHT}_{2m}\\ \mathbf{DHT}'_{2m}\end{vmatrix}\oplus\left(I_{k/2-1}\otimes_i D_{2m}\begin{vmatrix}\mathbf{rDFT}_{2m}(i/k)\\ \mathbf{rDFT}_{2m}(i/k)\\ \mathbf{rDHT}_{2m}(i/k)\\ \mathbf{rDHT}_{2m}(i/k)\end{vmatrix}\right)\right)\begin{vmatrix}\mathbf{RDFT}'_k\\ \mathbf{RDFT}'_k\\ \mathbf{DHT}'_k\\ \mathbf{DHT}'_k\end{vmatrix}\otimes I_m,\quad k\text{ even,}$$

$$\begin{vmatrix}\mathbf{rDFT}_{2n}(u)\\ \mathbf{rDHT}_{2n}(u)\end{vmatrix}\to L^{2n}_m\left(I_k\otimes_i\begin{vmatrix}\mathbf{rDFT}_{2m}((i+u)/k)\\ \mathbf{rDHT}_{2m}((i+u)/k)\end{vmatrix}\right)\left(\begin{vmatrix}\mathbf{rDFT}_{2k}(u)\\ \mathbf{rDHT}_{2k}(u)\end{vmatrix}\otimes I_m\right),$$

$$\mathbf{RDFT\text{-}3}_n\to(Q^\top_{k/2,2}\otimes I_2)\,(I_k\otimes_i\mathbf{rDFT}_{2m})(i+1/2)/k))\,(\mathbf{RDFT\text{-}3}_k\otimes I_m),\quad k\text{ even,}$$

$$\mathbf{DCT\text{-}2}_n\to P^\top_{k/2,2m}\left(\mathbf{DCT\text{-}2}_{2m}K^{2m}_2\oplus\left(I_{k/2-1}\otimes N_{2m}\mathbf{RDFT\text{-}3}^\top_{2m}\right)\right)B_n(L^{n/2}_{k/2}\otimes I_2)(I_m\otimes\mathbf{RDFT}'_k)Q_{m/2,k},$$

$$\mathbf{DCT\text{-}3}_n\to\mathbf{DCT\text{-}2}_n,$$

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

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

$$\cdots R_p,\quad p\text{ prime}$$

$$\mathbf{DCT\text{-}3}_n\to(I_m\oplus J_m)\,L^n_m(\mathbf{DCT\text{-}3}_m(1/4)\oplus\mathbf{DCT\text{-}3}_m(3/4))$$

$$\cdot(F_2\otimes I_m)\begin{vmatrix}I_m & 0\oplus-J_{m-1}\\ & \frac{1}{\sqrt2}(I_1\oplus2I_m)\end{vmatrix},\quad n=2m$$

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

$$\mathbf{IMDCT}_{2m} \to (J_m\oplus I_m\oplus I_m\oplus J_m)\left(\begin{bmatrix}1\\-1\end{bmatrix}\otimes I_m\right)\oplus\left(\begin{bmatrix}-1\\-1\end{bmatrix}\otimes I_m\right)\right)J_{2m}\,\mathbf{DCT\text{-}4}_{2m}$$

$$\mathbf{WHT}_{2^k} \to \prod_{i=1}^{t}(L_{2^{k_1+\cdots+k_{i-1}}}\otimes\mathbf{WHT}_{2^{k_i}}\otimes L_{2^{k_{i+1}+\cdots+k_l}}),\quad k=k_1+\cdots+k_l$$

$$\mathbf{DFT}_2 \to F_2$$

$$\mathbf{DCT\text{-}2}_2 \to \mathrm{diag}(1,1/\sqrt2)\,F_2$$

$$\mathbf{DCT\text{-}4}_2 \to J_2R_{13\pi/8}$$

*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>`<br>`<code for: y = At>` |
| $I_m \otimes A_n$ | `for (i=0; i<m; i++)`<br>`    <code for:`<br>`       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++)`<br>`    <code for:`<br>`       y[i:n:i+m*n-n] = A(x[i:n:i+m*n-n])>` |
| $D_n$ | `for (i=0; i<n; i++)`<br>`    y[i] = D[i]*x[i];` |
| $L^{km}_k$ | `for (i=0; i<k; i++)`<br>`    for (j=0; j<m; j++)`<br>`       y[i*m+j] = x[j*k+i];` |
| $F_2$ | `y[0] = x[0] + x[1];`<br>`y[1] = x[0] - x[1];` |

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

*Correct code:* **easy**      *fast code:* **very difficult**

© *Markus Püschel*
*Computer Science*   ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

*How to write fast numerical code*
*Spring 2015*

## Program Generation in Spiral

| | |
|---|---|
| **Transform** | $\mathrm{DFT}_8$ |

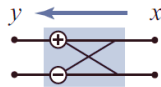*Decomposition rules*

| | |
|---|---|
| **Algorithm** *(SPL)* | $(\mathrm{DFT}_2 \otimes I_4) \, T_4^8 \, (I_2 \otimes ((\mathrm{DFT}_2 \otimes I_2)$ $T_2^4 \, (I_2 \otimes \mathrm{DFT}_2) \, L_2^4)) \, L_2^8$ |

**parallelization vectorization**

| | |
|---|---|
| **Algorithm** *(∑-SPL)* | $\sum \left( S_j \, \mathrm{DFT}_2 \, G_j \right) \sum \left( \sum \left( S_{k,l} \, \mathrm{diag}\left(\mathsf{t}_{k,l}\right) \mathrm{DFT}_2 \, G_l \right) \right.$ $\left. \sum \left( S_m \, \mathrm{diag}\left(\mathsf{t}_m\right) \mathrm{DFT}_2 \, G_{k,m} \right) \right)$ |

**locality optimization**

| | |
|---|---|
| **C Program** | ```
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;
< more lines>
``` |

**basic block optimizations**

*+ Search or Learning*

---

## Organization

- Spiral: Basic system

- *Vectorization*

- General input size

- Results

- Final remarks

*© Markus Püschel*
*Computer Science*
ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

*How to write fast numerical code*
*Spring 2015*

# Example: Vectorization in Spiral

- **Relationship SPL expressions ↔ vectorization?**

$$y = \mathrm{DFT}_2\, x \qquad\qquad\qquad y = \left( \mathrm{DFT}_2 \otimes \mathrm{I}_4 \right) x$$



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

---

# Step 1: Identify "Good" Vector Constructs

- **Vector length:** $\nu$

- **Good (= easily vectorizable) SPL constructs:**

    $A \otimes \mathrm{I}_\nu$

    $\mathsf{L}^{\nu^2}_\nu,\ \ \mathsf{L}^{2\nu}_2, \mathsf{L}^{2\nu}_\nu$   *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

$$\mathsf{L}_n^{n\nu} \rightarrow \left(\mathrm{I}_{n/\nu} \otimes \mathsf{L}_\nu^{\nu^2}\right)\left(\mathsf{L}_{n/\nu}^n \otimes \mathrm{I}_\nu\right)$$

$$\mathsf{L}_\nu^{n\nu} \rightarrow \left(\mathsf{L}_\nu^n \otimes \mathrm{I}_\nu\right)\left(\mathrm{I}_{n/\nu} \otimes \mathsf{L}_\nu^{\nu^2}\right)$$

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

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

$$\mathrm{I}_l \otimes \mathsf{L}_n^{kmn} \otimes \mathrm{I}_r \rightarrow \left(\mathrm{I}_l \otimes \mathsf{L}_{kn}^{kmn} \otimes \mathrm{I}_r\right)\left(\mathrm{I}_l \otimes \mathsf{L}_{mn}^{kmn} \otimes \mathrm{I}_r\right)$$

$$\mathrm{I}_l \otimes \mathsf{L}_{kn}^{kmn} \otimes \mathrm{I}_r \rightarrow \left(\mathrm{I}_{kl} \otimes \mathsf{L}_n^{mn} \otimes \mathrm{I}_r\right)\left(\mathrm{I}_l \otimes \mathsf{L}_k^{kn} \otimes \mathrm{I}_{mr}\right)$$

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

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

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

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

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

$$\overline{A^{m\times m} \otimes \mathrm{I}_\nu} \rightarrow \left(\mathrm{I}_m \otimes \mathsf{L}_\nu^{2\nu}\right)\left(\overline{A^{m\times m}} \otimes \mathrm{I}_\nu\right)\left(\mathrm{I}_m \otimes \mathsf{L}_2^{2\nu}\right)$$

$$\overline{\mathrm{I}_m \otimes A^{n\times n}} \rightarrow \mathrm{I}_m \otimes \overline{A^{n\times n}}$$

$$\overline{D} \rightarrow \left(\mathrm{I}_{n/\nu} \otimes \mathsf{L}_\nu^{2\nu}\right)\vec{D}\left(\mathrm{I}_{n/\nu} \otimes \mathsf{L}_2^{2\nu}\right)$$

$$\overline{P} \rightarrow P \otimes \mathrm{I}_2$$

*Manipulation rules = Processor knowledge in Spiral*

# Example

$$\underbrace{\mathbf{DFT}_{mn}}_{\mathrm{vec}(\nu)} \rightarrow \underbrace{(\mathbf{DFT}_m \otimes \mathrm{I}_n)\mathsf{T}_n^{mn}(\mathrm{I}_m \otimes \mathbf{DFT}_n)\mathsf{L}_m^{mn}}_{\mathrm{vec}(\nu)}$$

...

...

...

$$\rightarrow \left(\mathrm{I}_{\frac{mn}{\nu}} \otimes \mathsf{L}_\nu^{2\nu}\right)\left(\overline{\mathbf{DFT}_m \otimes \mathrm{I}_{\frac{n}{\nu}}} \otimes \mathrm{I}_\nu\right)\overline{\mathsf{T}}_n^{\prime mn}$$

$$\left(\mathrm{I}_{\frac{m}{\nu}} \otimes \left(\mathsf{L}_\nu^{2n} \otimes \mathrm{I}_\nu\right)\left(\mathrm{I}_{\frac{2n}{\nu}} \otimes \mathsf{L}_\nu^{\nu^2}\right)\left(\mathrm{I}_{\frac{n}{\nu}} \otimes \mathsf{L}_2^{2\nu} \otimes \mathrm{I}_\nu\right)\left(\overline{\mathbf{DFT}_n} \otimes \mathrm{I}_\nu\right)\right)\left(\mathsf{L}_{\frac{mn}{\frac{m}{\nu}}}^{\frac{mn}{\nu}} \otimes \mathsf{L}_2^{2\nu}\right)$$

**vectorized arithmetic**
**vectorized data accesses**

© *Markus Püschel*
*Computer Science*   ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

*How to write fast numerical code*
*Spring 2015*

# Automatically Generate Base Case Library

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

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

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

```
y = _mm_unpacklo_ps(x0, x1);
```
$$y = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} \vec{x}_0 \\ \vec{x}_1 \end{bmatrix}$$

```
y = _mm_shuffle_ps(x0, x1, _MM_SHUFFLE(1,2,1,2));
```
$$y = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} \vec{x}_0 \\ \vec{x}_1 \end{bmatrix}$$

```
y = _mm_shuffle_ps(x0, x1, _MM_SHUFFLE(3,4,3,4));
```
$$y = \begin{bmatrix} 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \vec{x}_0 \\ \vec{x}_1 \end{bmatrix}$$

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

```
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:**

$$\underset{\text{smp}(p,\mu)}{\underline{\text{DFT}_{ms}}} \rightarrow \underset{\text{smp}(p,\mu)}{\underline{\left( (\text{DFT}_m \otimes I_n) T_n^{mn} (I_m \otimes \text{DFT}_n) L_m^{mn} \right)}}$$

$$\rightarrow \underset{\text{smp}(p,\mu)}{\underline{(\text{DFT}_m \otimes I_n)}} \; \underset{\text{smp}(p,\mu)}{\underline{T_n^{mn}}} \; \underset{\text{smp}(p,\mu)}{\underline{(I_m \otimes \text{DFT}_n)}} \; \underset{\text{smp}(p,\mu)}{\underline{L_m^{mn}}}$$

$$\rightarrow \left( (L_m^{mp} \otimes I_{n/p\mu}) \otimes_\mu I_\mu \right) \left( I_p \otimes_\| (\text{DFT}_m \otimes I_{n/p}) \right) \left( (L_p^{mp} \otimes I_{n/p\mu}) \otimes_\mu I_\mu \right)$$
$$\left( \bigoplus_{i=0}^{p-1} \|^{T_n^{mn,i}} \right) \left( I_p \otimes_\| (I_{m/p} \otimes \text{DFT}_n) \right) \left( I_p \otimes_\| L_{m/p}^{mn/p} \right) \left( (L_p^{pn} \otimes I_{m/p\mu}) \otimes_\mu I_\mu \right)$$

**Vectorization:**

$$\underset{\text{vec}(\nu)}{\underline{(\text{DFT}_{mn})}} \rightarrow \underset{\text{vec}(\nu)}{\underline{\left( (\text{DFT}_m \otimes I_n) T_n^{mn} (I_m \otimes \text{DFT}_n) L_m^{mn} \right)}}$$

$$\rightarrow \underset{\text{vec}(\nu)}{\underline{\overline{(\text{DFT}_m \otimes I_n)}}^\nu} \underset{\text{vec}(\nu)}{\underline{\overline{(T_n^{mn})}^\nu}} \underset{\text{vec}(\nu)}{\underline{\overline{(I_m \otimes \text{DFT}_n)}}^{mn/\nu}}$$

$$\rightarrow (I_{mn/\nu} \otimes \underset{\text{sse}}{\underline{L_\nu^{2\nu}}}) (\overline{\text{DFT}_m \otimes I_{n/\nu}} \bar\otimes I_\nu) \underset{\text{sse}}{\underline{\overline{(T_n^{mn})}^\nu}}$$
$$\left( I_{m/\nu} \otimes (\overline{L_\nu^{2\nu}} \bar\otimes I_\nu)(I_{n/\nu} \otimes (L_\nu^{2\nu} \bar\otimes I_\nu)(I_2 \otimes \underset{\text{sse}}{\underline{L_2^{\nu^2}}})(L_2^{\nu^2} \bar\otimes I_\nu))(\overline{\text{DFT}_n \otimes I_\nu}) \right)$$
$$\left( (L_m^{mn} \otimes I_2) \bar\otimes I_{\nu'} \right)(I_{mn/\nu} \otimes \underset{\text{sse}}{\underline{L_2^{2\nu}}})$$

**GPUs:**

$$\underset{\text{gpu}(t,c)}{\underline{(\text{DFT}_{r^k})}} \rightarrow \underset{\text{gpu}(t,c)}{\underline{\left( \prod_{i=0}^{k-1} L_r^{r^k} \left( I_{r^{k-1}} \otimes \text{DFT}_r \right) \left( L_{r^{k-i-1}}^{r^k} (I_{r^i} \otimes T_{r^{k-i-1}}^{r^{k-i}}) \underset{\text{vec}(c)}{\underline{L_{r^{i+1}}^{r^k}}} \right) R_r^{r^k} \right)}}$$

$$\rightarrow \left( \prod_{i=0}^{k-1} (L_r^{r^n/2} \bar\otimes I_2)(I_{r^{n-1}/2} \otimes_\times \underset{\text{shd}(t,c)}{\underline{(\text{DFT}_r \bar\otimes I_2)}} \underset{}{L_r^{2r}}) \; T_i \right)$$
$$(L_r^{r^n/2} \bar\otimes I_2)(I_{r^{n-1}/2} \otimes_\times \underset{\text{shd}(t,c)}{\underline{L_r^{2r}}})(R_r^{r^{n-1}} \bar\otimes I_r)$$

**Verilog for FPGAs:**

$$\underset{\text{stream}(r^s)}{\underline{(\text{DFT}_{r^k})}} \rightarrow \underset{\text{stream}(r^s)}{\underline{\left[ \prod_{i=0}^{k-1} L_r^{r^k} \left( I_{r^{k-1}} \otimes \text{DFT}_r \right) \left( L_{r^{k-i-1}}^{r^k} (I_{r^i} \otimes T_{r^{k-i-1}}^{r^{k-i}}) L_{r^{i+1}}^{r^k} \right) \right] R_r^{r^k}}}$$

$$\rightarrow \underset{\text{stream}(r^s)}{\underline{\left[ \prod_{i=0}^{k-1} \underset{\text{stream}(r^s)}{\underline{L_r^{r^k}}} \underset{\text{stream}(r^s)}{\underline{\left( I_{r^{k-1}} \otimes \text{DFT}_r \right)}} \underset{\text{stream}(r^s)}{\underline{\left( L_{r^{k-i-1}}^{r^k} (I_{r^i} \otimes T_{r^{k-i-1}}^{r^{k-i}}) L_{r^{i+1}}^{r^k} \right)}} \right]}} \underset{\text{stream}(r^s)}{\underline{R_r^{r^k}}}$$

$$\rightarrow \underset{\text{stream}(r^s)}{\underline{\left[ \prod_{i=0}^{k-1} \underset{\text{stream}(r^s)}{\underline{L_r^{r^k}}} \left( I_{r^{k-s-1}} \otimes_s (I_{r^{s+1}} \otimes \text{DFT}_r) \right) \underset{\text{stream}(r^s)}{\underline{T_s^i}} \right]}} \underset{\text{stream}(r^s)}{\underline{R_r^{r^k}}}$$

- Rigorous, correct by construction
- Overcomes compiler limitations

# Organization

- Spiral: Basic system

- Vectorization

- *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*

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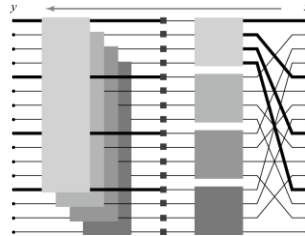DFT(n, x, y) {
  …
  for(i = …) {
    DFT_strided(m, x+mi, y+i, 1, k)
  }
  …
}
```

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

## Challenge: Recursion Steps

- **Cooley-Tukey FFT**

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



- **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 $f$
  - **$\mathrm{Perm}_f$**           permute data with the index mapping $f$

- **Σ-SPL formulas = matrix factorizations** 3
  **Example:** $y = (I_4 \otimes F_2)x \;\rightarrow\; y = \displaystyle\sum_{j=0}^{3} \mathsf{S}_{f_j} F_2 \, \mathsf{G}_{f_j} x$

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

# Find Recursion Step Closure

*Voronenko, 2008*

$$\{\mathbf{DFT}_n\}$$

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

$$\left(\sum_{i=0}^{k-1} \mathsf{S}_{h_{i,k}} \{\mathbf{DFT}_{n/k}\} \, \mathsf{G}_{h_{i,k}}\right) \mathrm{diag}(f) \left(\sum_{j=0}^{n/k-1} \mathsf{S}_{h_{jk,1}} \{\mathbf{DFT}_k\} \, \mathsf{G}_{h_{jk,1}}\right) \mathrm{perm}(\ell_{n/k}^n)$$

$$\sum_{i=0}^{k-1} \mathsf{S}_{h_{i,k}} \{\mathbf{DFT}_{n/k}\} \, \mathrm{diag}(f \circ h_{i,k}) \, \mathsf{G}_{h_{i,k}} \sum_{j=0}^{n/k-1} \mathsf{S}_{h_{jk,1}} \{\mathbf{DFT}_k\} \, \mathsf{G}_{h_{j,n/k}}$$

$$\sum_{i=0}^{k-1} \left\{ \mathsf{S}_{h_{i,k}} \, \mathbf{DFT}_{n/k} \, \mathrm{diag}(f \circ h_{i,k}) \, \mathsf{G}_{h_{i,k}} \right\} \sum_{j=0}^{n/k-1} \left\{ \mathsf{S}_{h_{jk,1}} \, \mathbf{DFT}_k \, \mathsf{G}_{h_{j,n/k}} \right\}$$

*Repeat until closure*

---

# Recursion Step Closure: Examples

*DFT: scalar code*



*DFT: full-fledged (vectorized and parallel code)*

*© Markus Püschel*
**ETH**
Eidgenössische Technische Hochschule Zürich
*Computer Science*  Swiss Federal Institute of Technology Zurich

*How to write fast numerical code*
*Spring 2015*

## 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

- Vectorization

- 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



Spiral generated

FFTW 3.2

Intel IPP 6.0

Spiral generated
(1 thread)

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

**Spiral**

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

---

# Generating 100s of FFTWs

*PhD thesis Voronenko, 2009*

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

$$\begin{vmatrix}\mathbf{RDFT}_n\\\mathbf{RDFT}_n'\\\mathbf{DHT}_n\\\mathbf{DHT}_n'\end{vmatrix} \rightarrow (P_{k/2,m}^\top \otimes I_2) \left(\begin{vmatrix}\mathbf{RDFT}_{2m}\\\mathbf{RDFT}_{2m}'\\\mathbf{DHT}_{2m}\\\mathbf{DHT}_{2m}\end{vmatrix} \oplus \left(I_{k/2-1} \otimes_i D_{2m} \begin{vmatrix}\mathbf{rDFT}_{2m}(i/k)\\\mathbf{rDFT}_{2m}(i/k)\\\mathbf{rDHT}_{2m}(i/k)\\\mathbf{rDHT}_{2m}(i/k)\end{vmatrix}\right)\right) \left(\begin{vmatrix}\mathbf{RDFT}_k'\\\mathbf{RDFT}_k'\\\mathbf{DHT}_k'\\\mathbf{DHT}_k'\end{vmatrix} \otimes I_m\right), \quad k \text{ even,}$$

$$\begin{vmatrix}\mathbf{rDFT}_{2n}(u)\\\mathbf{rDHT}_{2n}(u)\end{vmatrix} \rightarrow L_m^{2n} \left(I_k \otimes_i \begin{vmatrix}\mathbf{rDFT}_{2m}((i+u)/k)\\\mathbf{rDHT}_{2m}((i+u)/k)\end{vmatrix}\right) \left(\begin{vmatrix}\mathbf{rDFT}_{2k}(u)\\\mathbf{rDHT}_{2k}(u)\end{vmatrix} \otimes I_m\right),$$

$$\mathbf{RDFT}\text{-}3_n \rightarrow (Q_{k/2,m}^\top \otimes I_2) (I_k \otimes_i \mathbf{rDFT}_{2m})(i+1/2)/k)) (\mathbf{RDFT}\text{-}3_k \otimes I_m), \quad k \text{ even,}$$

$$\mathbf{DCT}\text{-}2_n \rightarrow P_{k/2,2m}^\top \left(\mathbf{DCT}\text{-}2_{2m} K_2^{2m} \oplus \left(I_{k/2-1} \otimes N_{2m} \mathbf{RDFT}\text{-}3_{2m}^\top\right)\right) B_n(L_{k/2}^{n/2} \otimes I_2)(I_m \otimes \mathbf{RDFT}_k') Q_{m/2,k},$$

$$\mathbf{DCT}\text{-}3_n \rightarrow \mathbf{DCT}\text{-}2_n^\top,$$

$$\mathbf{DCT}\text{-}4_n \rightarrow Q_{k/2,2m}^\top \left(I_{k/2} \otimes N_{2m} \mathbf{RDFT}\text{-}3_{2m}^\top\right) B_n'(L_{k/2}^{n/2} \otimes I_2)(I_m \otimes \mathbf{RDFT}\text{-}3_k) Q_{m/2,k}.$$

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

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

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

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

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

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

$$\mathbf{WHT}_{2^k} \rightarrow \prod_{i=1}^t (\mathsf{I}_{2^{k_1+\cdots+k_{i-1}}} \otimes \mathbf{WHT}_{2^{k_i}} \otimes \mathsf{I}_{2^{k_{i+1}+\cdots+k_l}}), \quad k = k_1 + \cdots + k_l$$

$$\mathbf{DFT}_2 \rightarrow \mathsf{F}_2$$

$$\mathbf{DCT}\text{-}2_2 \rightarrow \operatorname{diag}(1, 1/\sqrt{2}) \mathsf{F}_2$$

$$\mathbf{DCT}\text{-}4_2 \rightarrow \mathsf{J}_2 \mathsf{R}_{13\pi/8}$$

*© Markus Püschel*
*Computer Science*

**ETH**
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

*How to write fast numerical code*
*Spring 2015*

# Generating 100s of FFTWs

*PhD thesis Voronenko, 2009*

| Transform | Code size non-parallelized | parallelized |
|---|---|---|
| *no vectorization* | | |
| 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 | — |
| *2-way vectorization* | | |
| DFT | 14.8 KLOC / 0.73 MB | 15.0 KLOC / 0.74 MB |
| RDFT | 15.6 KLOC / 0.76 MB | 16.0 KLOC / 0.81 MB |
| scaled RDFT | 16.0 KLOC / 0.78 MB | — |
| DHT | 16.9 KLOC / 0.83 MB | 17.2 KLOC / 0.87 MB |
| DCT-2 | 20.7 KLOC / 1.10 MB | 21.0 KLOC / 1.09 MB |
| DCT-3 | 27.9 KLOC / 1.56 MB | 28.2 KLOC / 1.59 MB |
| DCT-4 | 7.8 KLOC / 0.47 MB | 8.1 KLOC / 0.50 MB |
| WHT | 6.9 KLOC / 0.32 MB | 5.8 KLOC / 0.26 MB |
| FIR Filter | 167 KLOC / 7.75 MB | 120 KLOC / 5.12 MB |
| Downsampled FIR Filter | 100 KLOC / 4.2 MB | 68 KLOC / 2.76 MB |
| *4-way vectorization* | | |
| 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 |

# Generating 100s of FFTWs

*PhD thesis Voronenko, 2009*

© Markus Püschel
Computer Science

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

*How to write fast numerical code*
*Spring 2015*

# Computer generated Functions for Intel IPP 6.0

Intel® Integrated Performance Primitives (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.*

---

# Very Large Scale: BG/P



**HPC Challenge Global FFT on BlueGene/P**

[Gflop/s]

theoretical peak

Spiral-generated

UPC coalesced transpose

BlueGene/P node cards and racks

**6.4 Tflop/s**

**32 racks**
**= 32K node cards**
**= 128K cores**

*2010 HPC Challenge Class I Award, Almasi et al.*

# Organization

- Spiral: Basic system

- Vectorization

- General input size

- Results

- *Final remarks*

---

# Spiral: Summary

- **Spiral:**
  Successful approach to automating
  the development of computing software

  Commercial proof-of-concept

  Intel® Integrated Performance Primitives (Intel® IPP) 6.0

$$\mathrm{DFT}_{64}$$

```
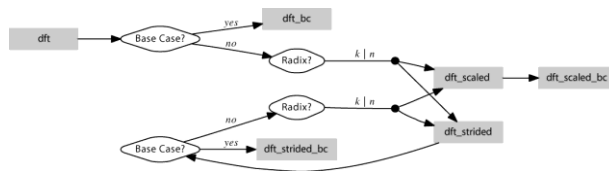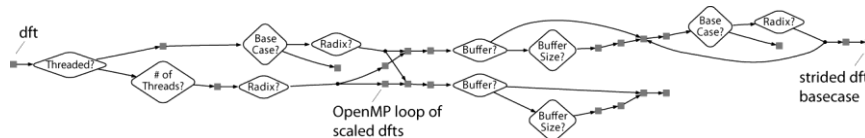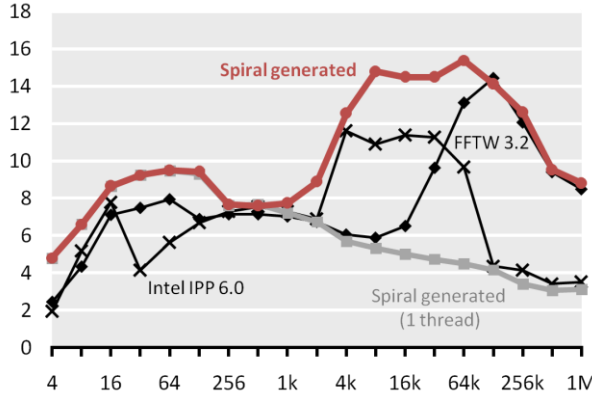void dft64(float  *Y, float *X) {
  __m512 U912, U913, U914, U915,...
  __m512 *a2153, *a2155;
  a2153 = ((__m512 *) X);   s1107 = *(a2153);
  s1108 = *((a2153 + 4));   t1323 = _mm512_add_ps(s1107,s1108);
  t1324 = _mm512_sub_ps(s1107,s1108);
  <many more lines>
  U926 = _mm512_swizupconv_r32(...);
  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>
}
```

- **Key ideas:**

  *Algorithm knowledge:*
  Domain specific symbolic representation

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

  *Platform knowledge:*
  Tagged rewrite rules, SIMD specification

  $$\underbrace{\mathrm{I}_m \otimes A_n}_{\mathsf{smp}(p,\mu)} \rightarrow \mathrm{I}_p \otimes_{\parallel} \left( \mathrm{I}_{m/p} \otimes A_n \right)$$

*© Markus Püschel*
**ETH**
*Computer Science* Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

*How to write fast numerical code*
*Spring 2015*