# How to Write Fast Numerical Code

Spring 2013
*Lecture:* Spiral (Computer generation of FFT code)

**Instructor:** Markus Püschel
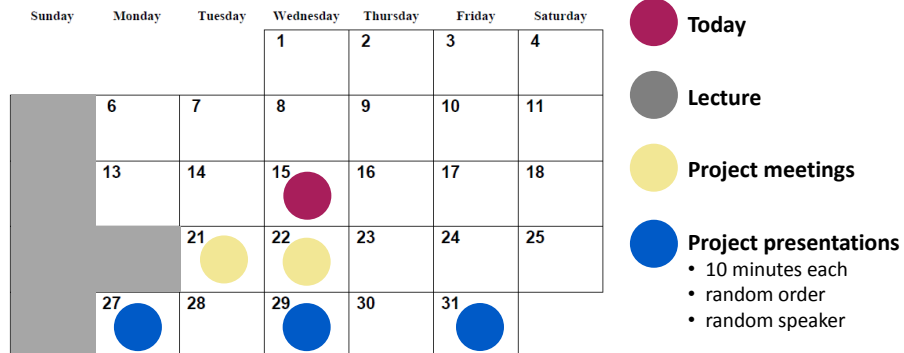
**TA:** Georg Ofenbeck & Daniele Spampinato

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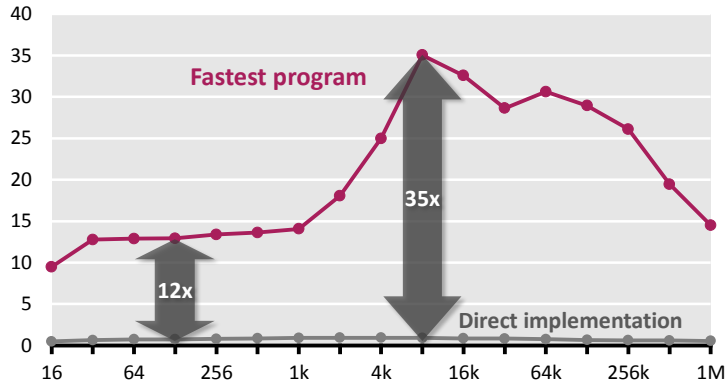
---

# Rest of Semester

**May 2013**

| Sunday | Monday | Tuesday | Wednesday | Thursday | Friday | Saturday |
|--------|--------|---------|-----------|----------|--------|----------|
|        |        |         | 1         | 2        | 3      | 4        |
|        | 6      | 7       | 8         | 9        | 10     | 11       |
|        | 13     | 14      | 15 ●      | 16       | 17     | 18       |
|        |        | 21 ●    | 22 ●      | 23       | 24     | 25       |
|        | 27 ●   | 28      | 29 ●      | 30       | 31 ●   |          |

● **Today**

● **Lecture**

● **Project meetings**

● **Project presentations**
- 10 minutes each
- random order
- random speaker

2

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

*How to write fast numerical code*
*Spring 2013*

## The Problem: Example DFT

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

**Fastest program**

**35x**

**12x**

**Direct implementation**

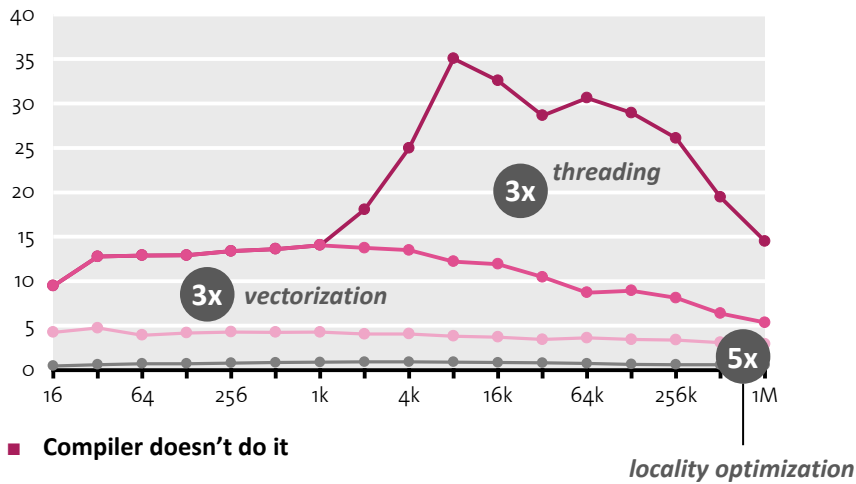16    64    256    1k    4k    16k    64k    256k    1M

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

## DFT: Analysis

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

**3x** *threading*

**3x** *vectorization*

**5x**

16    64    256    1k    4k    16k    64k    256k    1M

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

*locality optimization*

## Our Goal:

Computer writes high performance library code

**Generate Code**

*"click"*

---

**Select convolutional code**

Select a preset code or customize parameters

- *custom*
- Voyager — rate 1 / 2 — code rate (?)
- NASA-DSN — K 7 — constraint length (?)
- CCSDS/NASA-GSFC — polynomials 109 — polynomials for the code in decimal notation (?)
- WiMax — 79
- CDMA IS-95A
- LTE (3GPP - Long Term Evolution)
- UWB (802.15)
- CDMA 2000
- Cassini
- Mars Pathfinder & Stereo

**Select implementation options**

frame length 2048 — unpadded frame length

Vectorization level scalar C — type of code (?)

Generate Code   Reset

### Viterbi Decoder

**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

*@ www.spiral.net*

# 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

*© Markus Püschel*
*Computer Science*

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

*How to write fast numerical code*
*Spring 2013*

# 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**

$$\mathrm{DFT}_4 = (\mathrm{DFT}_2 \otimes I_2)\mathsf{T}_2^4(I_2 \otimes \mathrm{DFT}_2)\mathsf{L}_2^4$$

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

---

# Decomposition Rules (>200 for >40 Transforms)

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

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

$$\begin{vmatrix} \mathrm{rDFT}_{2n}(u) \\ \mathrm{rDHT}_{2n}(u) \end{vmatrix} \to L_m^{2n} \left( I_k \quad i \begin{vmatrix} \mathrm{rDFT}_{2m}((i+u)/k) \\ \mathrm{rDHT}_{2m}((i+u)/k) \end{vmatrix} \right) \left( \begin{vmatrix} \mathrm{rDFT}_{2k}(u) \\ \mathrm{rDHT}_{2k}(u) \end{vmatrix} \quad I_m \right),$$

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

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

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



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

$$\mathrm{DCT\text{-}3}_n \to (I_m \oplus J_m)\, L_m^n (\mathrm{DCT\text{-}3}_m(1/4) \oplus \mathrm{DCT\text{-}3}_m(3/4))$$
$$\cdot (\mathsf{F}_2 \quad I_m) \begin{bmatrix} I_m & \cdot & J_{m-1} \\ \cdot & \frac{1}{\sqrt 2}(I_1 \oplus 2I_m) \end{bmatrix}, \quad n = 2m$$

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

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

$$\mathrm{WHT}_{2^k} \to \prod_{i=1}^{t} (I_{2^{k_1+\cdots+k_{i-1}}} \quad \mathrm{WHT}_{2^{k_i}} \quad I_{2^{k_{i+1}+\cdots+k_t}}), \quad k = k_1 + \cdots + k_t$$

$$\mathrm{DFT}_2 \to \mathsf{F}_2$$
$$\mathrm{DCT\text{-}2}_2 \to \mathrm{diag}(1, 1/\sqrt 2)\, \mathsf{F}_2$$
$$\mathrm{DCT\text{-}4}_2 \to \mathsf{J}_2\, \mathsf{R}_{13\pi/8}$$

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

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

*How to write fast numerical code*
*Spring 2013*

## 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];` |
| $L_k^{km}$ | `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***

---

## Program Generation in Spiral

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

*Decomposition rules*

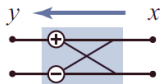| **Algorithm** *(SPL)* | $(\mathrm{DFT}_2 \otimes I_4)\, T_4^8\, \big(I_2 \otimes ((\mathrm{DFT}_2 \otimes I_2)$ <br> $T_2^4\,(I_2 \otimes \mathrm{DFT}_2)\, L_2^4)\big)\, L_2^8$ | **parallelization** **vectorization** |
|---|---|---|
| **Algorithm** *(∑-SPL)* | $\sum \big(S_j\, \mathrm{DFT}_2\, G_j\big) \sum \big(\sum \big(S_{k,l}\, \mathrm{diag}(\mathrm{t}_{k,l})\, \mathrm{DFT}_2\, G_l\big)$ <br> $\sum \big(S_m\, \mathrm{diag}(\mathrm{t}_m)\, \mathrm{DFT}_2\, G_{k,m}\big)\big)$ | **locality** **optimization** |
| **C Program** | ```void sub(double *y, double *x) {```<br>```double f0, f1, f2, f3, f4, f7, f8, f10, f11;```<br>```    f0 = x[0] - x[3];```<br>```    f1 = x[0] + x[3];```<br>```    f2 = x[1] - x[2];```<br>```    f3 = x[1] + x[2];```<br>```    f4 = f1 - f3;```<br>```    y[0] = f1 + f3;```<br>```    y[2] = 0.7071067811865476 * f4;```<br>```    f7 = 0.9238795325112867 * f0;```<br>```< more lines>``` | **basic block** **optimizations** |

*+ Search or Learning*

# Organization

- Spiral: Basic system

- *Vectorization*

- General input size

- Results
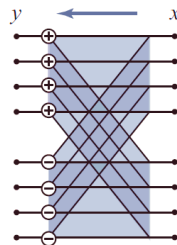
- Final remarks

---

# Example: Vectorization in Spiral

- **Relationship SPL expressions ↔ vectorization?**

$$y = \mathbf{DFT}_2\, x$$



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



**one addition**
**one subtraction**

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

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

*How to write fast numerical code*
*Spring 2013*

## Step 1: Identify "Good" Vector Constructs

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

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

  $$A \quad \mathrm{I}_\nu$$
  $$\mathsf{L}_\nu^{\nu^2}, \ \mathsf{L}_2^{2\nu}, \mathsf{L}_\nu^{2\nu} \quad \textit{base cases}$$

  *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

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

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

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

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

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

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

$$\mathrm{I}_l \quad \mathsf{L}_m^{kmn} \quad \mathrm{I}_r \rightarrow \left(\mathrm{I}_l \quad \mathsf{L}_m^{kmn} \quad \mathrm{I}_r\right)\left(\mathrm{I}_l \quad \mathsf{L}_m^{kmn} \quad \mathrm{I}_r\right)$$

*Manipulation rules = Processor knowledge in Spiral*

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

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

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

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

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

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

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

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

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

## Example

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

$$\dots$$
$$\dots$$
$$\dots$$

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

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

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

---

## Automatically Generate Base Case Library

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

  $\nu = 4:$ $\quad \left\{ \mathsf{L}_2^4, \mathrm{I}_2 \quad \mathsf{L}_2^4, \mathsf{L}_2^4 \quad \mathrm{I}_2, \mathsf{L}_2^8, \mathsf{L}_4^8 \right\}$

- *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));
```
$$y = \begin{bmatrix} 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}\begin{bmatrix} \vec{x}_0 \\ \vec{x}_1 \end{bmatrix}$$

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

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

$\mathsf{L}_2^4 \otimes \mathrm{I}_2$
~~no base case~~
**Base case**

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

*How to write fast numerical code*
*Spring 2013*

# Same Approach for Different Paradigms

**Threading:**

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

$$\cdots$$

$$\rightarrow \underset{\text{smp}(p,\mu)}{\underline{(\text{DFT}_m \otimes I_n)}}\ \underset{\text{smp}(p,\mu)}{\underline{\mathsf{T}_n^{mn}}}\ \underset{\text{smp}(p,\mu)}{\underline{(I_m \otimes \text{DFT}_n)}}\ \underset{\text{smp}(p,\mu)}{\underline{\mathsf{L}_m^{mn}}}$$

$$\cdots$$

$$\rightarrow \left(\left(\mathsf{L}_m^{mp} \otimes I_{n/p\mu}\right) \otimes_\mu I_\mu\right)\left(I_p \otimes_{\|}(\text{DFT}_m \otimes I_{n/p})\right)\left(\left(\mathsf{L}_p^{mp} \otimes I_{n/p\mu}\right)\otimes_\mu I_\mu\right)$$

$$\left(\bigoplus_{i=0}^{p-1}{}_\| \mathsf{T}_n^{mn,i}\right)\left(I_p \otimes_{\|}(I_{m/p} \otimes \text{DFT}_n)\right)\left(I_p \otimes_{\|} \mathsf{L}_{m/p}^{mn/p}\right)\left((\mathsf{L}_p^{pn} \otimes I_{n/p\mu})\otimes_\mu I_\mu\right)$$

**Vectorization:**

$$\underset{\text{vec}(\nu)}{\underline{(\text{DFT}_{mn})}} \rightarrow \underset{\text{vec}(\nu)}{\underline{\left((\text{DFT}_m \otimes I_n)\,\mathsf{T}_n^{mn}(I_m \otimes \text{DFT}_n)\,\mathsf{L}_m^{mn}\right)}}$$

$$\cdots$$

$$\rightarrow \underset{\text{vec}(\nu)}{\underline{(\text{DFT}_m \otimes I_n)}}^\nu\ \underset{\text{vec}(\nu)}{\underline{\left(\mathsf{T}_n^{mn}\right)}}^\nu\ \underset{\text{vec}(\nu)}{\underline{(I_m \otimes \text{DFT}_n)\,\mathsf{L}_m^{mn}}}^\nu$$

$$\cdots$$

$$\rightarrow (I_{mn/\nu} \otimes \underset{\text{sse}}{\underline{\mathsf{L}_\nu^{2\nu}}})(\overline{\text{DFT}_m} \otimes I_{n/\nu} \tilde{\otimes} I_\nu)(\underset{\text{sse}}{\underline{\mathsf{T}_n^{mn}}})^\nu$$

$$\left(I_{m/\nu} \otimes (\overline{\mathsf{L}_\nu^{\frac{n}{\nu}}}\tilde{\otimes} I_\nu)(I_{n/\nu} \otimes (\mathsf{L}_\nu^{2\nu}\tilde{\otimes}I_\nu)(I_2 \otimes \underset{\text{sse}}{\underline{\mathsf{L}_\nu^{\nu^2}}})(\mathsf{L}_2^{2\nu}\tilde{\otimes}I_\nu))(\overline{\text{DFT}_n}\tilde{\otimes}I_\nu)\right)$$

$$\left((\mathsf{L}_m^{mn} \otimes I_2)\tilde{\otimes} I_\nu\right)(I_{mn/\nu} \otimes \underset{\text{sse}}{\underline{\mathsf{L}_2^{2\nu}}})$$

**GPUs:**

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

$$\cdots$$

$$\rightarrow \left(\prod_{i=0}^{k-1}(\mathsf{L}_r^{r^n/2}\tilde{\otimes}I_2)\left(I_{r^{n-1}/2}\otimes_\times \underset{\text{shd}(t,c)}{\underline{(\text{DFT}_r\tilde{\otimes}I_2)}}\mathsf{L}_r^{2r}\right)\mathsf{T}_i\right)$$

$$(\mathsf{L}_r^{r^n/2}\tilde{\otimes}I_2)(I_{r^{n-1}/2}\otimes_\times \underset{\text{shd}(t,c)}{\underline{\mathsf{L}_r^{2r}}})(\mathsf{R}_r^{r^{n-1}}\tilde{\otimes}I_r)$$

**Verilog for FPGAs:**

$$\underset{\text{stream}(r^s)}{\underline{(\text{DFT}_{r^k})}} \rightarrow \left[\prod_{i=0}^{k-1}\mathsf{L}_r^{r^k}\left(I_{r^{k-1}}\otimes \text{DFT}_r\right)\left(\mathsf{L}_{r^{k-i-1}}^{r^k}(I_{r^i}\otimes \mathsf{T}_{r^{k-i-1}}^{r^{k-i}})\mathsf{L}_{r^{i+1}}^{r^k}\right)\right]\mathsf{R}_r^{r^k}$$

$$\underset{\text{stream}(r^s)}{\underline{\qquad}}$$

$$\rightarrow \left[\prod_{i=0}^{k-1}\underset{\text{stream}(r^s)}{\underline{\mathsf{L}_r^{r^k}}}\left(I_{r^{k-1}}\otimes \text{DFT}_r\right)\underset{\text{stream}(r^s)}{\underline{\left(\mathsf{L}_{r^{k-i-1}}^{r^k}(I_{r^i}\otimes \mathsf{T}_{r^{k-i-1}}^{r^{k-i}})\mathsf{L}_{r^{i+1}}^{r^k}\right)}}\right]\underset{\text{stream}(r^s)}{\underline{\mathsf{R}_r^{r^k}}}$$

$$\cdots$$

$$\rightarrow \left[\prod_{i=0}^{k-1}\underset{\text{stream}(r^s)}{\underline{\mathsf{L}_r^{r^k}}}\left(I_{r^{k-s-1}}\otimes_s(I_{r^{s-1}}\otimes \text{DFT}_r)\right)\underset{\text{stream}(r^s)}{\underline{\mathsf{T}_i^i}}\right]\underset{\text{stream}(r^s)}{\underline{\mathsf{R}_r^{r^k}}}$$

- 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:
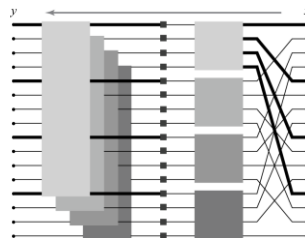*Library for general input size*

```
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);
  }
  ```
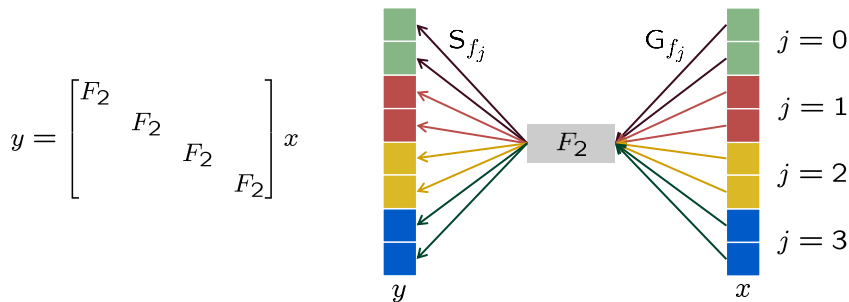
# Σ–SPL : Basic Idea

- **Four additional matrix constructs: $\Sigma$, G, S, Perm**
  - **$\Sigma$ (sum)**      explicit loop
  - **$G_f$ (gather)**      load data with index mapping $f$
  - **$S_f$ (scatter)**      store data with index mapping $f$
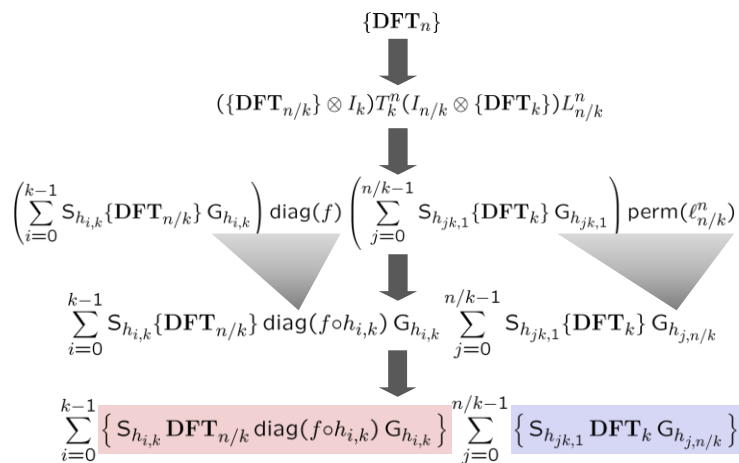  - **$Perm_f$**      permute data with the index mapping $f$

- **$\Sigma$-SPL formulas = matrix factorizations** 3

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

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



---

# Find Recursion Step Closure

*Voronenko, 2008*

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

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

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

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

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

*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*

*© Markus Püschel*
*Computer Science*

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

*How to write fast numerical code*
*Spring 2013*

# 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



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

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

© Markus Püschel

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

Computer Science

*How to write fast numerical code*
*Spring 2013*

# 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.*

---

# 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(0.7071067811865475f),0xAAAA,a2154,U926),t1341),
        _mm512_set_1to16_ps(0.7071067811865475f),_);
        _mm512_swizupconv_r32(t1341, _MM_SWIZ_REG_CDAB));
    U927 = _mm512_swizupconv_r32
        <many more lines>
}
```

■ **Key ideas:**

*Algorithm knowledge:*
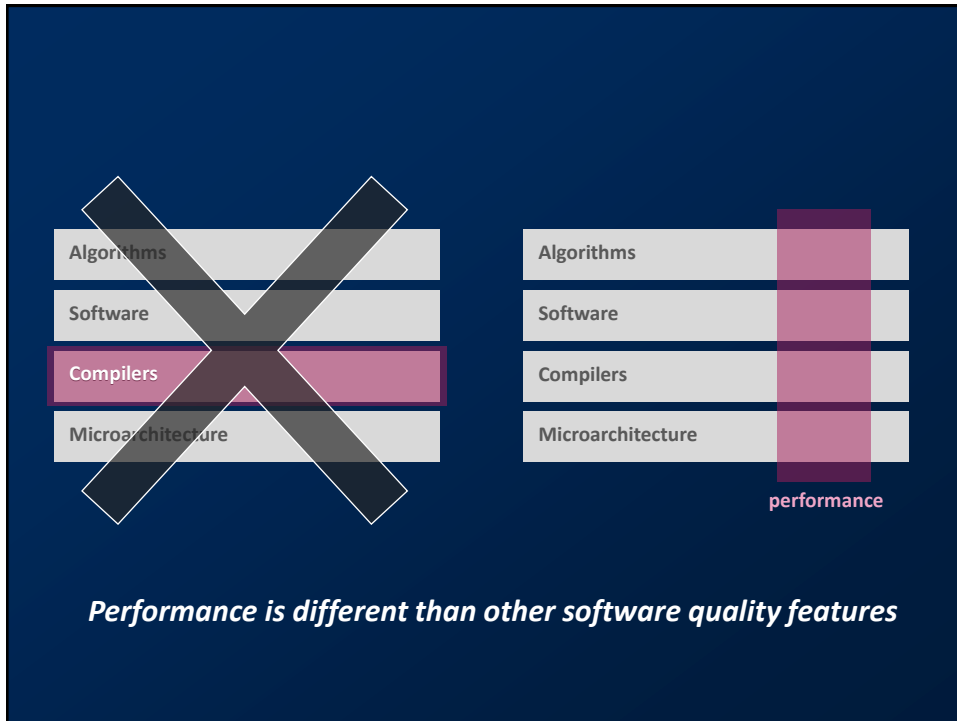Domain specific symbolic representation

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

*Platform knowledge:*
Tagged rewrite rules, SIMD specification

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

*Performance is different than other software quality features*

---

# Research Questions

- **How to automate the production of fastest numerical code?**
    - *Domain-specific languages*
    - *Rewriting*
    - *Compilers*
    - *Machine Learning*
- **What program language features help with program generation?**
- **What environment should be used to build generators?**
- **How to represent mathematical functionality?**
- **How to formalize the mapping to fast code?**
- **How to handle various forms of parallelism?**
- **How to integrate into standard work flows?**