# Exploiting Domain-Specific Knowledge: Spiral
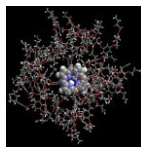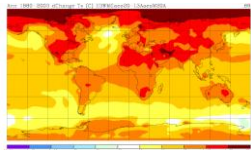
Markus Püschel & Georg Ofenbeck

*Computer Science*
**ETH** *zürich*

SPIRAL
www.spiral.net

DSLs

---

# Mathematical Computing
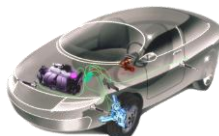
Science simulations

Audio, image, Video processing

Signal processing, communication, control
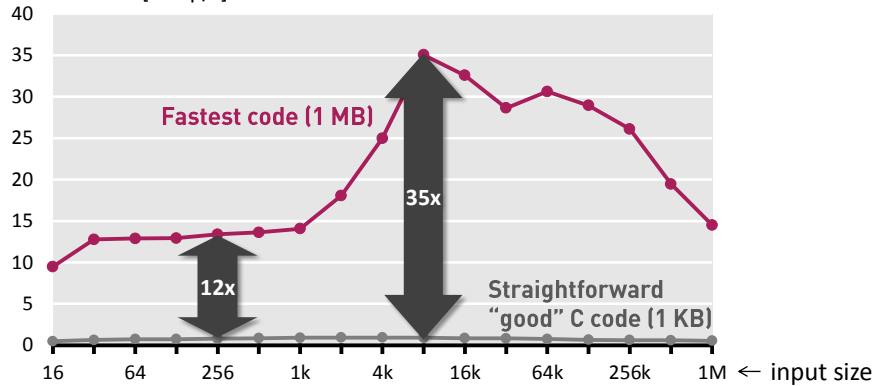
Security

Machine learning, data analytics

Optimization

*Highest performance is often crucial*

---

**Exploiting Domain–Specific Knowledge: Spiral**
EPFL/ETH Summer School on DSL Design and Implementation
Markus Püschel & Georg Ofenbeck, Computer Science, ETH Zürich

1

# Example: Discrete Fourier Transform

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

Performance [Gflop/s]

**Fastest code (1 MB)**

**35x**

**12x**

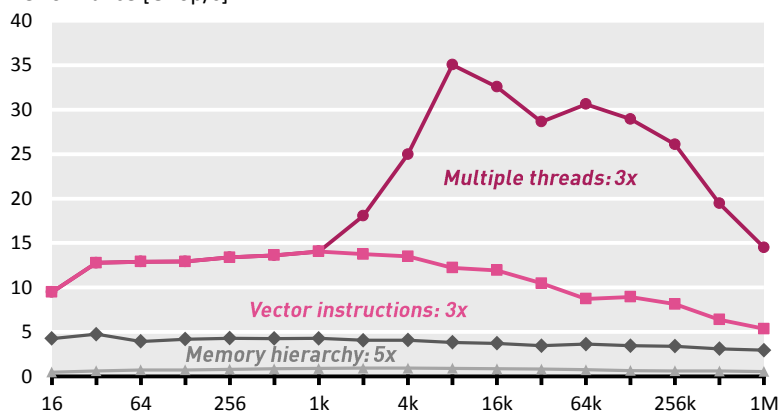Straightforward
"good" C code (1 KB)

input size

Vendor compiler, best flags

Roughly same operations count

---

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

Performance [Gflop/s]

*Multiple threads: 3x*

*Vector instructions: 3x*

*Memory hierarchy: 5x*

Compiler doesn't do the job

**Doing by hand =** restructure algorithm for locality & parallelism,
handle choices, choose proper code style, use vector intrinsics, ….
**= nightmare**

**Exploiting Domain-Specific Knowledge: Spiral**
EPFL/ETH Summer School on  DSL Design and Implementation
Markus Püschel & Georg Ofenbeck, Computer Science, ETH Zürich

2

# Optimization: Register Locality and ILP

```
// straightforward code
for(i = 0; i < N; i += 1)
  for(j = 0; j < N; j += 1)
    for(k = 0; k < N; k += 1)
      c[i][j] += a[i][k]*b[k][j];
```

*Concise and slow*

Removes aliasing

Enables register allocation and instruction scheduling

*Compiler does not do well:*
- often illegal
- many choices

```
// unrolling + scalar replacement
for(i = 0; i < N; i += MU) {
  for(j = 0; j < N; j += NU) {
    for(k = 0; k < N; k += KU) {
      t1 = A[i*N + k];
      t2 = A[i*N + k + 1];
      t3 = A[i*N + k + 2];         load
      t4 = A[i*N + k + 3];
      t5 = A[(i + 1)*N + k];
      <more copies>

      t10 = t1  * t9;
      t17 = t17 + t10;
      t21 = t1  * t8;
      t18 = t18 + t21;
      t12 = t5  * t9;             compute
      t19 = t19 + t12;
      t13 = t5  * t8;
      t20 = t20 + t13;
      <more ops>

      C[i*N + j]         = t17;
      C[i*N + j + 1]     = t18;    store
      C[(i+1)*N + j]     = t19;
      C[(i+1)*N + j + 1] = t20;
    }
  }
}
```
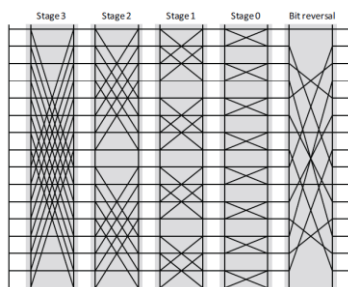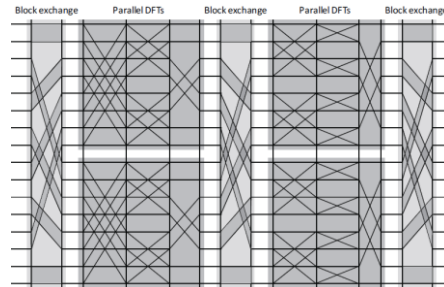
*Ugly and fast*

# Optimization for Parallelism (Threads)



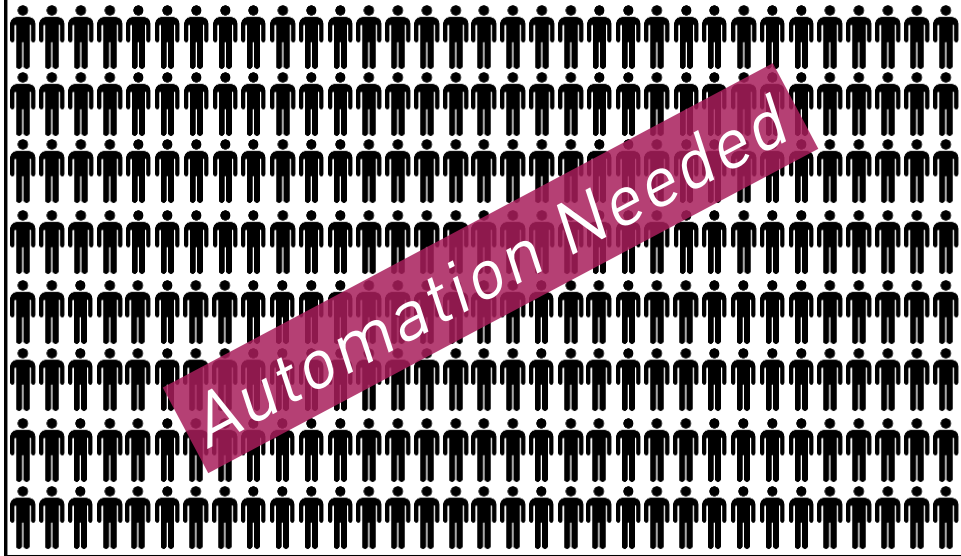**Parallelism is present, but is not in the "right shape"**

**Restructured for locality and parallelism (shared memory, 2 cores, 2 elements per cache line)**

*Compiler usually does not do*
- analysis may be unfeasible
- may require algorithm changes
- may require domain knowledge
- may require processor parameters

**Exploiting Domain-Specific Knowledge: Spiral**
EPFL/ETH Summer School on DSL Design and Implementation
Markus Püschel & Georg Ofenbeck, Computer Science, ETH Zürich

3

**Current practice:** Thousands of programmers re-implement and re-optimize the same functionality for every new processor and for every new processor generation

Automation Needed

## Goal:

Computer generation of high performance code for ubiquitous performance-critical components

Generate Code

🖑 "click"

**Exploiting Domain-Specific Knowledge: Spiral**
EPFL/ETH Summer School on DSL Design and Implementation
Markus Püschel & Georg Ofenbeck, Computer Science, ETH Zürich

4

## Possible Approach:

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

Structural optimization:
*Rewriting systems*

High performance code style:
*Compiler*

Decision making for choices:
*Machine learning*

**Exploiting Domain-Specific Knowledge: Spiral**
EPFL/ETH Summer School on  DSL Design and Implementation
Markus Püschel & Georg Ofenbeck, Computer Science, ETH Zürich

5

Spiral: Program Generation for Performance (www.spiral.net)

Franz Franchetti
Yevgen Voronenko
Jianxin Xiong
Bryan Singer
Srinivas Chellappa
Frédéric de Mesmay
Peter Milder
José Moura
David Padua
Jeremy Johnson
James Hoe
<many more>

funding: DARPA, NSF, ONR, Intel

# Linear Transforms

$$\begin{pmatrix} y_0 \\ y_1 \\ \vdots \\ y_{n-1} \end{pmatrix} = y = Tx \longleftarrow \boxed{T \cdot} \longleftarrow x = \begin{pmatrix} x_0 \\ x_1 \\ \vdots \\ x_{n-1} \end{pmatrix}$$

*Output*  *Input*

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

**Exploiting Domain-Specific Knowledge: Spiral**
EPFL/ETH Summer School on DSL Design and Implementation
Markus Püschel & Georg Ofenbeck, Computer Science, ETH Zürich
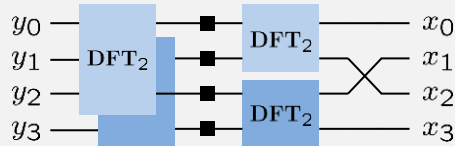
6

# Algorithms: Example FFT, n = 4

## Fast Fourier transform (FFT):

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

**12 adds, 4 mults**          **4 adds**          **1 mult**          **4 adds**          **0 adds/mults**

## Data flow graph



## Description with matrix algebra (SPL)

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

---

# Decomposition Rules (>200 for >40 Transforms)

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

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

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

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



**Rules = algorithm knowledge**
(≈100 journal papers)

$$\mathrm{DCT\text{-}3}_n \to (I_m \oplus J_m) L_m^n (\mathrm{DCT\text{-}3}_m(1/4) \oplus \mathrm{DCT\text{-}3}_m(3/4))$$

$$\cdot (\mathsf{F}_2 \otimes I_m) \begin{bmatrix} I_m & 0 \oplus -J_{m-1} \\ & \frac{1}{\sqrt{2}}(I_1 \oplus 2I_m) \end{bmatrix}, \quad n = 2m$$

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

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

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

$$\mathrm{DFT}_2 \to \mathsf{F}_2$$

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

$$\mathrm{DCT\text{-}4}_2 \to J_2\, \mathsf{R}_{13\pi/8}$$

**Exploiting Domain–Specific Knowledge: Spiral**
EPFL/ETH Summer School on DSL Design and Implementation
Markus Püschel & Georg Ofenbeck, Computer Science, ETH Zürich

7

## 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];` |

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

*Gives reasonable, straightforward code*

---

## Program Generation in Spiral

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

*Decomposition rules (algorithm knowledge)*

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

**parallelization**
**vectorization**

| **Algorithm**<br>*(Σ-SPL)* | $\sum \left( S_j\, \mathrm{DFT}_2\, G_j \right) \sum \left( \sum \left( S_{k,l}\, \mathrm{diag}(t_{k,l})\, \mathrm{DFT}_2\, G_l \right) \right.$<br>$\left. \sum \left( S_m\, \mathrm{diag}(t_m)\, \mathrm{DFT}_2\, G_{k,m} \right) \right)$ |
|---|---|

**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 >``` |
|---|---|

**code style**
**code level**
**optimization**

*+ Search or Learning for Choices*

**Exploiting Domain-Specific Knowledge: Spiral**
EPFL/ETH Summer School on DSL Design and Implementation
Markus Püschel & Georg Ofenbeck, Computer Science, ETH Zürich

8

## SPL to Shared Memory Code: Basic Idea

*"Good" SPL structures*

$$y = \left( \mathrm{I}_p \otimes A \right) x$$

#processors $= p = 4$



Processor 0
Processor 1
Processor 2
Processor 3

$x$

$$y = \left( P \otimes \mathrm{I}_\mu \right) x$$



$x$      $y$

} cache block size $= \mu = 2$

*Rewriting:* Bad structures ➡ good structures

---

## Example: SMP Parallelization

*Franchetti, Voronenko & P, SC 2006*

$$\underbrace{\mathbf{DFT}_{mn}}_{\mathsf{smp}(p,\mu)} \;\rightarrow\; \underbrace{\left( (\mathbf{DFT}_m \otimes \mathrm{I}_n) \top_n^{mn} (\mathrm{I}_m \otimes \mathbf{DFT}_n) \mathsf{L}_m^{mn} \right)}_{\mathsf{smp}(p,\mu)}$$

$$\ldots$$

$$\rightarrow\; \underbrace{\left( \mathbf{DFT}_m \otimes \mathrm{I}_n \right)}_{\mathsf{smp}(p,\mu)} \underbrace{\top_n^{mn}}_{\mathsf{smp}(p,\mu)} \underbrace{\left( \mathrm{I}_m \otimes \mathbf{DFT}_n \right)}_{\mathsf{smp}(p,\mu)} \underbrace{\mathsf{L}_m^{nm}}_{\mathsf{smp}(p,\mu)}$$

$$\ldots$$

$$\rightarrow\; \left( (\mathsf{L}_m^{mp} \otimes \mathrm{I}_{n/p\mu}) \otimes \mathrm{I}_\mu \right) \left( \mathrm{I}_p \otimes (\mathbf{DFT}_m \otimes \mathrm{I}_{n/p}) \right) \left( (\mathsf{L}_p^{mp} \otimes \mathrm{I}_{n/p\mu}) \otimes \mathrm{I}_\mu \right)$$

$$\left( \bigoplus_{i=0}^{p-1}{}_{\|} \top_n^{mn,i} \right) \left( \mathrm{I}_p \otimes (\mathrm{I}_{m/p} \otimes \mathbf{DFT}_n) \right) \left( \mathrm{I}_p \otimes \mathsf{L}_{m/p}^{mn/p} \right) \left( (\mathsf{L}_p^{pn} \otimes \mathrm{I}_{m/p\mu}) \otimes \mathrm{I}_\mu \right)$$

**load-balanced, no false sharing**

*One rewriting system for every platform paradigm:*
**SIMD, distributed memory parallelism, FPGA, …**

**Exploiting Domain-Specific Knowledge: Spiral**
EPFL/ETH Summer School on DSL Design and Implementation
Markus Püschel & Georg Ofenbeck, Computer Science, ETH Zürich

9

# Same Approach for Different Paradigms

**Threading:**



**Vectorization:**



**GPUs:**



**Verilog for FPGAs:**



- Rigorous, correct by construction
- Overcomes compiler limitations

---

# 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**

**Exploiting Domain-Specific Knowledge: Spiral**
EPFL/ETH Summer School on  DSL Design and Implementation
Markus Püschel & Georg Ofenbeck, Computer Science, ETH Zürich

10

# Challenge: Recursive Steps Needed

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

    $\Sigma$ (sum)      *matrix 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:** $y = (I_4 \otimes F_2)x \;\rightarrow\; y = \sum_{j=0}^{3} \mathsf{S}_{f_j} F_2 \, \mathsf{G}_{f_j} x$

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

**Exploiting Domain-Specific Knowledge: Spiral**
EPFL/ETH Summer School on DSL Design and Implementation
Markus Püschel & Georg Ofenbeck, Computer Science, ETH Zürich

11

# Find Recursion Step Closure

*Voronenko, thesis 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)*

**Exploiting Domain–Specific Knowledge: Spiral**
EPFL/ETH Summer School on  DSL Design and Implementation
Markus Püschel & Georg Ofenbeck, Computer Science, ETH Zürich

12

**Generating Dozens of "FFTWs"**



**It Really Works**

**DFT on Sandybridge (3.3 GHz, 4 Cores, AVX)**
Performance [Gflop/s]



Spiral-generated

FFTW 3.3b1

MKL 10.3.2

IPP 7.0.2

vectorized, threaded, platform-tuned library (5 MB source code)

**Exploiting Domain-Specific Knowledge: Spiral**
EPFL/ETH Summer School on DSL Design and Implementation
Markus Püschel & Georg Ofenbeck, Computer Science, ETH Zürich

13

Computer generated Functions for Intel IPP

**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.*



LGen: Generator for Basic Linear Algebra
*Spampinato & P, CGO 2014*

BLAC $\quad y = x^T(A+B)y + \delta$

Algorithm: **Tiling decision and propagation**
(LL) $\quad \left[ y = x^T(A+B)y + \delta \right]_{2,3}$ → vectorization

Algorithm $\quad \sum_{i,j,i',j'} S_i S_{i'} \left( G_{i'} G_i A G_j G_{j'} \right) \left( G_{j'} G_j x \right) \dots$ → locality optimization
(Σ-LL)

C Program
```
void kernel(float *x, float *A, float *B, …) {
    float t0_64_0, t0_64_1, t0_64_2, t0_64_3 …;
        t0_57_0 = A[0];
        t0_56_0 = A[1];
        …
        t0_59_0 = t0_57_0 + t0_33_0;
        t0_63_0 = t0_59_0 * t0_9_0;
        t0_59_1 = t0_56_0 + t0_32_0;
        t0_60_0 = t0_59_1 * t0_8_0;
        < many more lines>
```
→ code style
code level optimization

**Exploiting Domain-Specific Knowledge: Spiral**
EPFL/ETH Summer School on DSL Design and Implementation
Markus Püschel & Georg Ofenbeck, Computer Science, ETH Zürich

14

# LGen: Sample Results

$$C = \alpha AB + \beta C$$



Performance [f/c]

generated

MKL

BTO

n [Float]

$$A \in \mathbb{R}^{n \times 4}$$

$$B \in \mathbb{R}^{4 \times n}$$

$$C = \alpha(A_0 + A_1)^T B + \beta C$$



Performance [f/c]

generated

MKL

n [Float]

$$A_0 \in \mathbb{R}^{4 \times 4}$$

$$B \in \mathbb{R}^{4 \times n}$$

— LGen
—▽— Handwritten fixed size
—△— Handwritten gen size
—○— MKL 11.0
—□— Eigen 3.1.3
—☆— BTO 1.3
—◇— IPP 7.1

---

# PL Support: Example Code Style
*Ofenbeck*, *Rompf*, *Stojanov*, *Odersky & P*, *GPCE 2012*

**SPL**

$$y = (I_2 \otimes DFT_2)x$$

**Data flow graph**

$$y_0 \quad \text{DFT}_2 \quad x_0$$
$$y_1 \quad \quad x_1$$
$$y_2 \quad \text{DFT}_2 \quad x_2$$
$$y_3 \quad \quad x_3$$

**Scala function**

```scala
def f(x: Array[Double], y: Array[Double]) = {
    for (i <- 0 until 2) {
      y(2*i)   = x(i*2) + x(i*2+1)
      y(2*i+1) = x(i*2) - x(i*2+1)
    }
  }
```

**Exploiting Domain-Specific Knowledge: Spiral**
EPFL/ETH Summer School on  DSL Design and Implementation
Markus Püschel & Georg Ofenbeck, Computer Science, ETH Zürich

15

## Top slide

```
def f(x: Array[Rep[Double]],
      y: Array[Rep[Double]]) = {
    for (i <- 0 until 2) {
      y(2*i)   = x(i*2) + x(i*2+1)
      y(2*i+1) = x(i*2) - x(i*2+1)
    }
}
```



*scalarized*

```
t0 = s0 + s1;
t1 = s0 - s1;
t2 = s2 + s3;
t2 = s2 - s3;
```

```
def f(x: Rep[Array[Double]],
      y: Rep[Array[Double]]) = {
    for (i <- 0 until 2) {
      y(2*i)   = x(i*2) + x(i*2+1)
      y(2*i+1) = x(i*2) - x(i*2+1)
    }
}
```



*unrolled, scalar repl.*

```
t0 = x[0];
t1 = x[1];
t2 = t0 + t1;
y[0] = t2;
t3 = t0 - t1;
y[1] = t3;
t4 = x[0];
t5 = x[1];
t6 = t4 + x5;
y[0] = t6;
t7 = t4 – x5;
y[3] = t7;
```

```
def f(x: Rep[Array[Double]],
      y: Rep[Array[Double]]) = {
    for (i <- 0 until 2: Rep[Range]) {
      y(2*i)   = x(i*2) + x(i*2+1)
      y(2*i+1) = x(i*2) - x(i*2+1)
    }
}
```



*looped, scalar repl.*

```
for (int i=0; i < 2; i++)
{
    t0 = x[i];
    t1 = x[i+1];
    t2 = t0 + t1;
    y[i] = t2;
    t3 = t0 - t1;
    y[i+1] = t3;
}
```

## Bottom slide

*Staging enables program generation*

*Abstracting over code style =
abstracting over staging decisions*

```
def f[L[_],A[_],T](looptype: L, x: A[Array[T]], y: A[Array[T]]) = {
    for (i <- 0 until 2: L[Range]) {
      y(2*i)  = x(i*2) + x(i*2+1)
      y(2*i+1)= x(i*2) - x(i*2+1)
    }
}
```

**Exploiting Domain-Specific Knowledge: Spiral**
EPFL/ETH Summer School on  DSL Design and Implementation
Markus Püschel & Georg Ofenbeck, Computer Science, ETH Zürich

16

# Related Work

**Program generators for performance**
- *FFTW codelet generator (Frigo)*
- *Flame (van de Geijn, Quintana-Orti, Bientinesi, …)*
- *cvxgen (Mattingley, Boyd)*
- *PetaBricks (Ansel, Amarasinghe, …)*
- *Spiral*

**Autotuning**
- *ATLAS/PhiPAC (Whaley, Bilmes, Demmel, Dongarra, …)*
- *FFTW adaptive library (Frigo, Johnson)*
- *OSKI (Vuduc et al.)*
- *Adaptive sorting (Li et al.)*

**Environments for DSLs and program generation**
- *see this workshop*

---

# Automatically from Math to Fast Code

## Principles

Capturing algorithm knowledge:
*Mathematical DSLs*

Structural optimization:
*Rewriting*

Decision making:
*Search and learning*

Generate Code

*"click"*

$$\text{DFT}_n \rightarrow (\text{DFT}_k \otimes \text{I}_m)\, \text{T}_m^n (\text{I}_k \otimes \text{DFT}_m)\, \text{L}_k^n, \quad n = km$$
$$\text{DFT}_n \rightarrow P_n(\text{DFT}_k \otimes \text{DFT}_m)Q_n, \quad n = km,\ \gcd(k,m) = 1$$
$$\text{DCT-4}_n \rightarrow S_n \text{DCT-2}_n \operatorname{diag}_{0 \le k < n}(1/(2\cos((2k+1)\pi/4n)))$$
$$\text{IMDCT}_{2m} \rightarrow (\text{J}_m \oplus \text{I}_m \oplus \text{I}_m \oplus \text{J}_m)\left(\left[\begin{smallmatrix}1\\-1\end{smallmatrix}\right] \otimes \text{I}_m\right) \oplus \left(\left[\begin{smallmatrix}-1\\-1\end{smallmatrix}\right] \otimes \text{I}_m\right)\right) \text{J}_{2m}\, \text{DCT-4}_{2m}$$

$$\underbrace{A_m \otimes \text{I}_n}_{\text{smp}(p,\mu)} \rightarrow \underbrace{\text{L}_m^{mn}}_{\text{smp}(p,\mu)}\ \left(\text{I}_p \otimes_{\parallel}(\text{I}_{n/p} \otimes A_m)\right)\ \underbrace{\text{L}_n^{mn}}_{\text{smp}(p,\mu)}$$

## Key Challenges

New domains (linear algebra, filters, …)

Programming language support (DSLs, staging)

More information: www.spiral.net

**Exploiting Domain-Specific Knowledge: Spiral**
EPFL/ETH Summer School on DSL Design and Implementation
Markus Püschel & Georg Ofenbeck, Computer Science, ETH Zürich

17