

# How to Write Fast Code

18-645, spring 2008





23<sup>rd</sup> Lecture, Apr. 9<sup>th</sup>

**Instructor:** Markus Püschel

**TAs:** Srinivas Chellappa (Vas) and Frédéric de Mesmay (Fred)

# Research Project

## ■ Current status

April						
<i>Sun</i>	<i>Mon</i>	<i>Tue</i>	<i>Wed</i>	<i>Thu</i>	<i>Fri</i>	<i>Sat</i>
		1	2	3	4	5
6	7	8	9 	10	11	12
13	14	15	16	17	18	19
20	21	22	23	24	25 	26
27	28	29	30 			
						
						2008

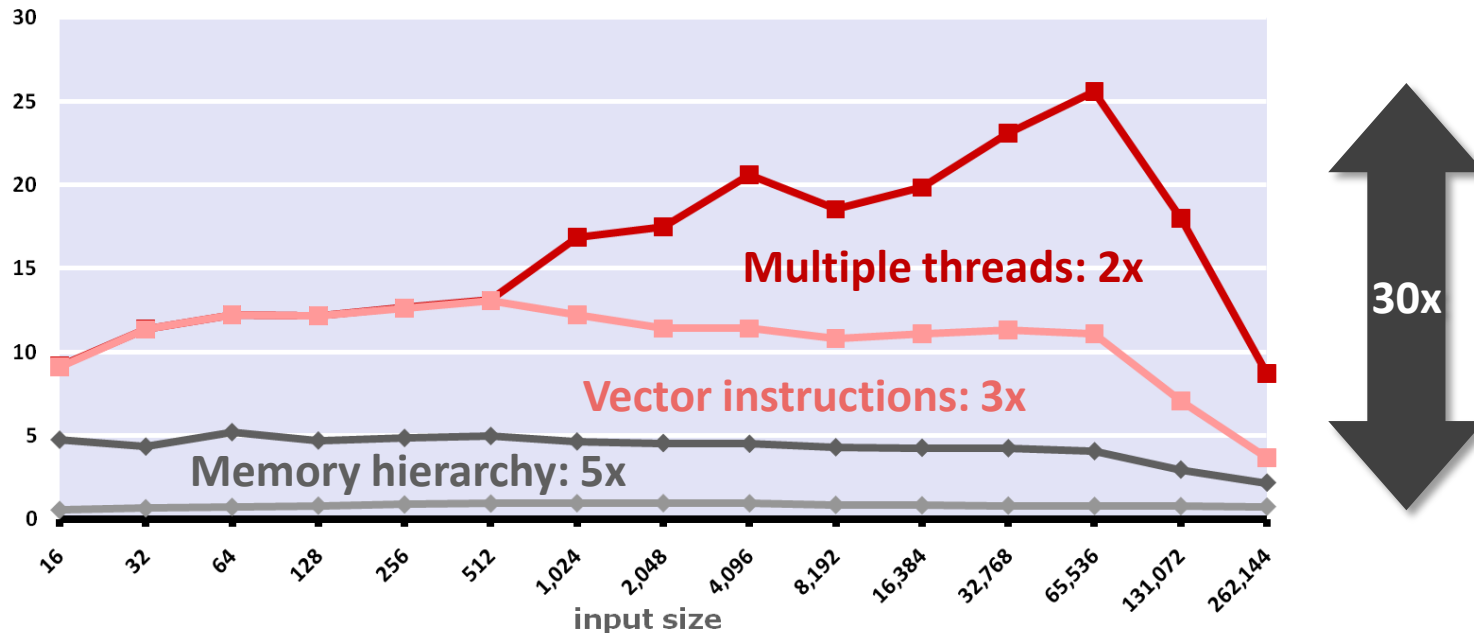
- **Today**
- **Papers due**
- **Last class:  
poster session  
5:30 – 8:30**
- **Final papers due**

# Today

- **Spiral: Can we teach computers to write fast libraries?**

# You Already Know the Problem ...

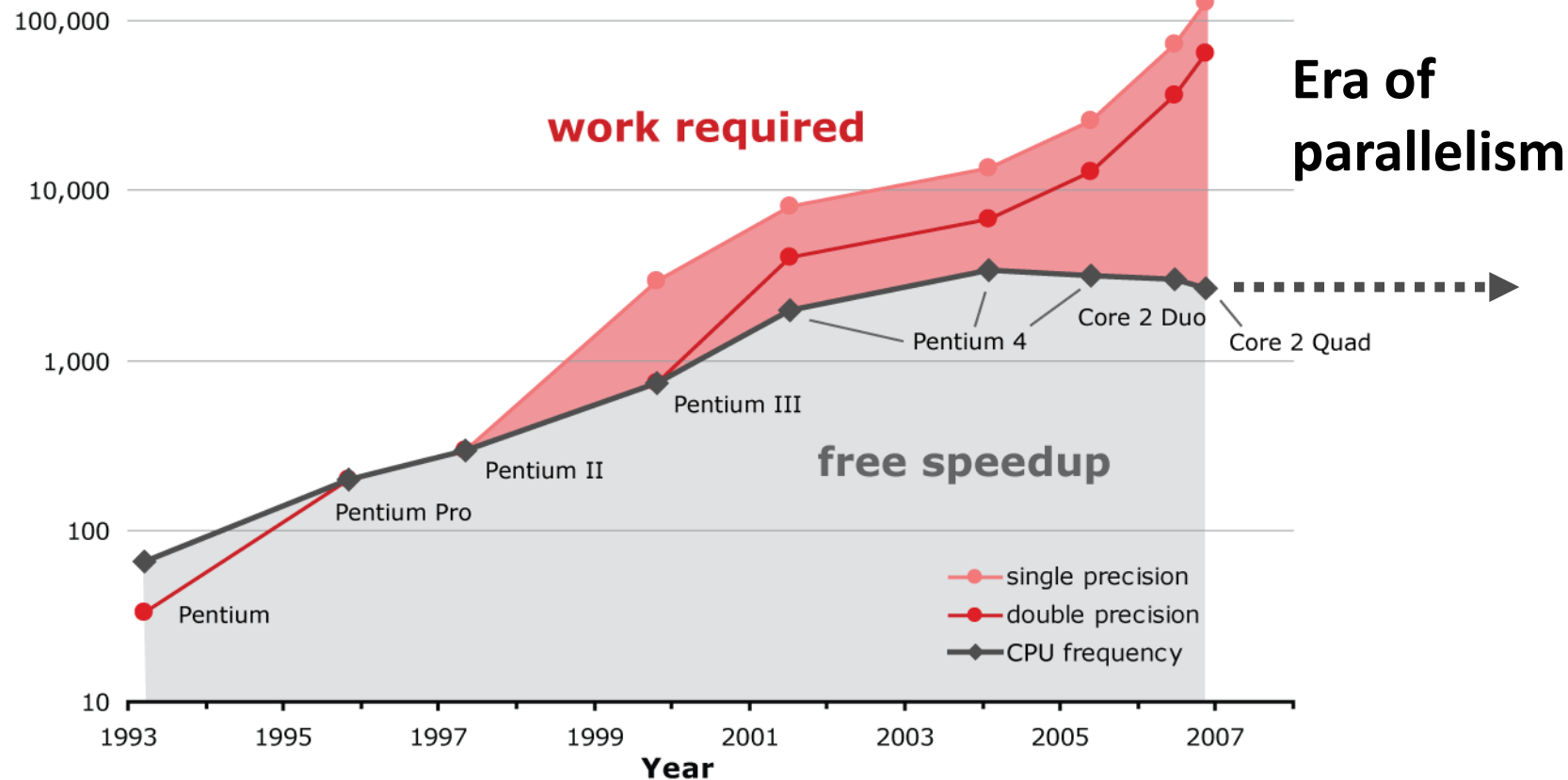
Discrete Fourier Transform (DFT) on 2xCORE2Duo 3 GHz  
Performance [Gflop/s]



- A straightforward implementation of a numerical algorithm is typically **10-100x slower** than the fastest available
- Minimizing operations count  $\neq$  minimizing runtime
- Developing high performance numerical code has become a nightmare
  - Expertise in algorithms, programming, architecture
  - Optimization for the memory hierarchy, vectorization, parallelization

# ... And the Cause

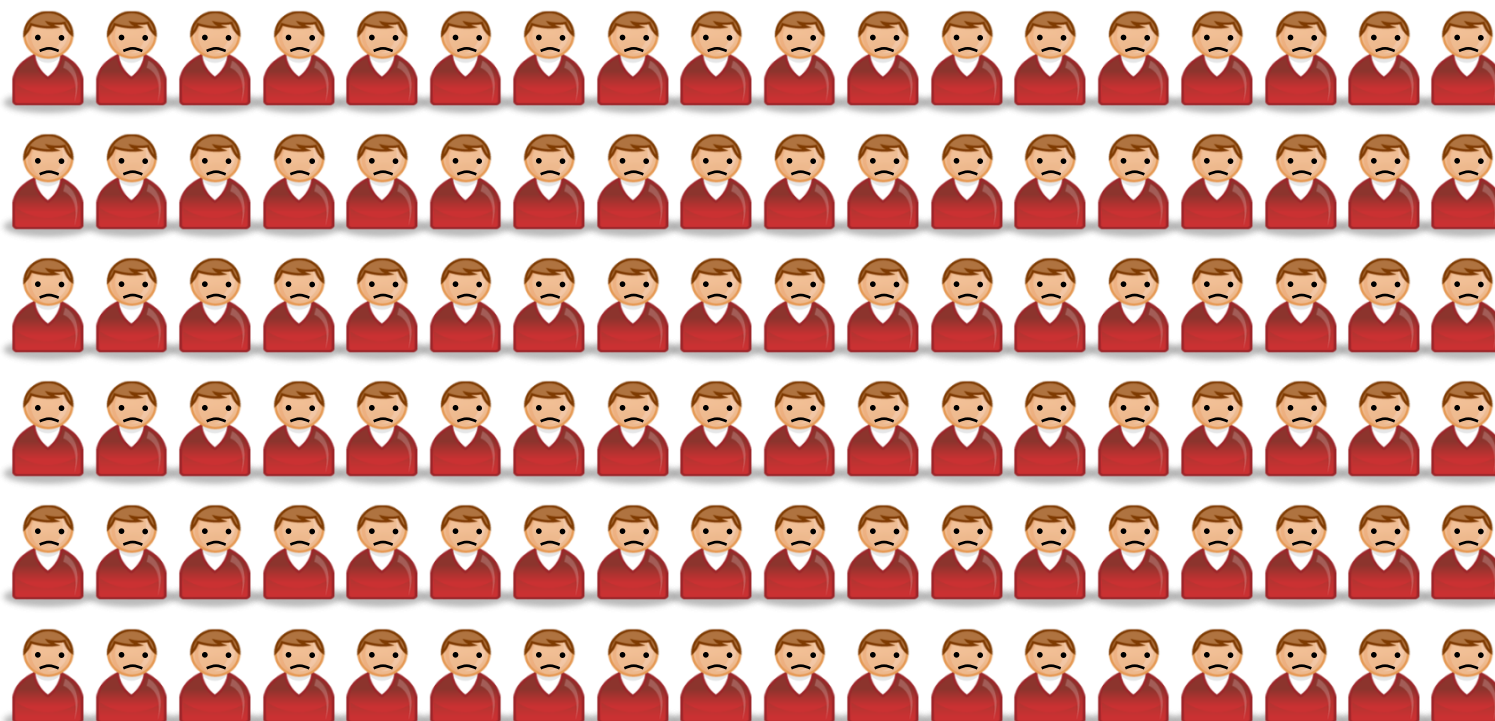
Floating point peak performance [Mflop/s]  
CPU frequency [MHz]



*High performance library development has become a nightmare*

# Current Solution

- **Legions** of programmers implement and optimize the **same** functionality for **every** platform and **whenever** a new platform comes out.



# Better Solution: Automation

- Automate (parts of) the implementation or optimization



- **ATLAS:**

- Automatic generation of BLAS
- Optimization for the memory hierarchy

- **Sparsity/Bebop:**

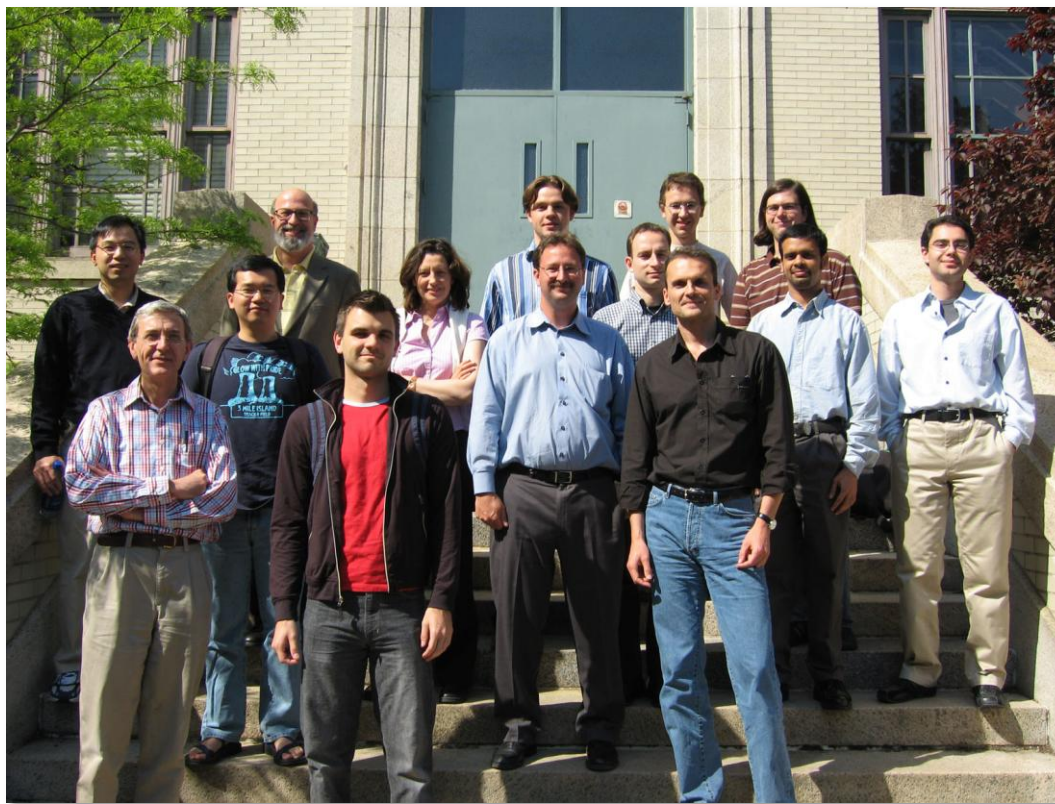
- Automatic platform adaptation of sparse MVM
- Mainly register blocking

- **FFTW:**

- Adaptive DFT library
- Code generator for unrolled basic blocks

- **Spiral (today):**

- Library generator for linear transforms
- Complete automation (computer does everything)



**[www.spiral.net](http://www.spiral.net)**



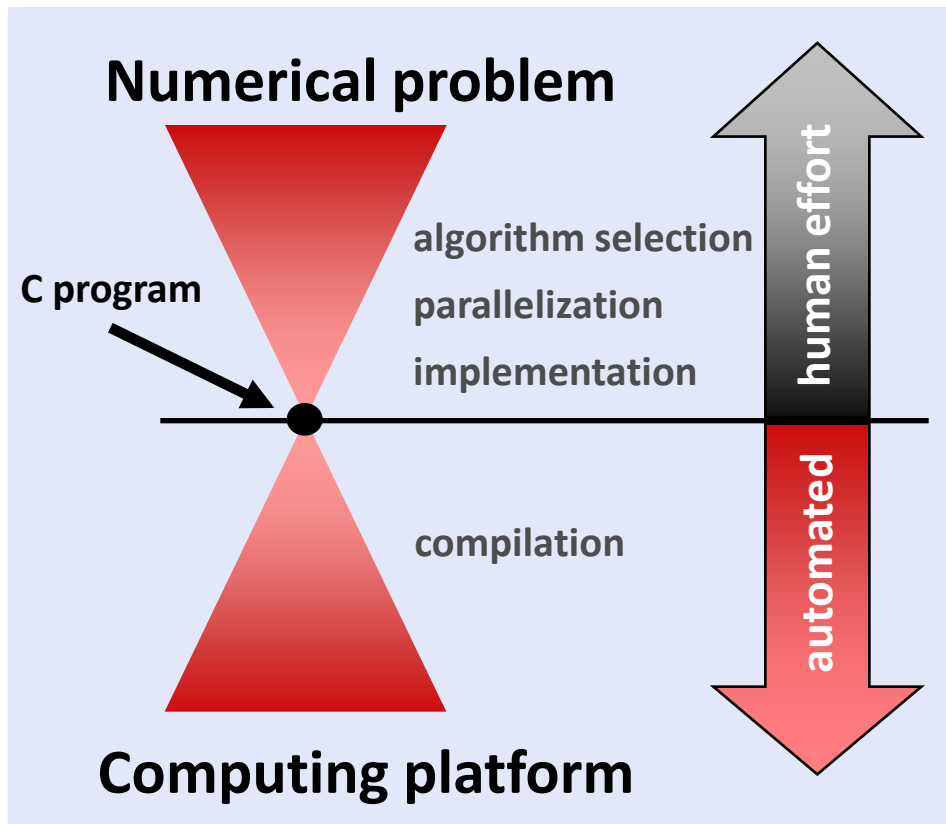
# Organization

- **Spiral: Base system**
- Parallelism
- Libraries (general size)
- Beyond transforms
- Conclusions

Markus Püschel, José M. F. Moura, Jeremy Johnson, David Padua, Manuela Veloso, Bryan Singer, Jianxin Xiong, Franz Franchetti, Aca Gacic, Yevgen Voronenko, Kang Chen, Robert W. Johnson, and Nick Rizzolo, **SPIRAL: Code Generation for DSP Transforms**, Proceedings of the IEEE 93(2), 2005

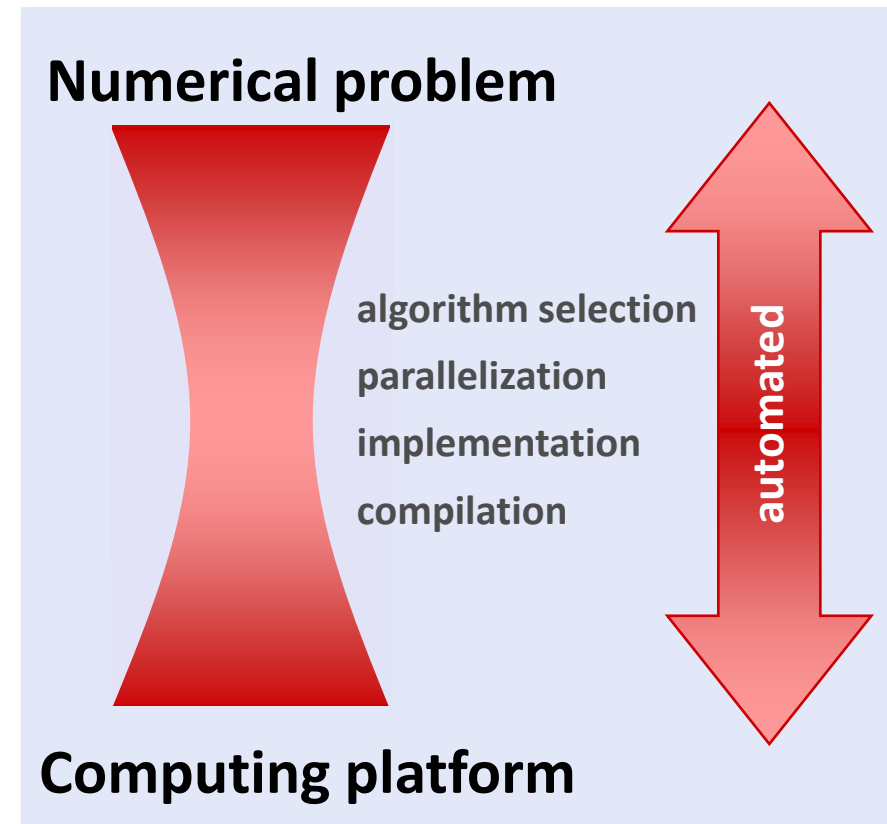
# Vision Behind Spiral

## Current



- C code a singularity: Compiler has no access to high level information

## Future



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

# Spiral

- **Library generator for linear transforms**  
(discrete Fourier/cosine/wavelet transforms, filters, ....) *and recently more ...*
- **Wide range of platforms supported:**  
scalar, fixed point, **vector, parallel, Verilog**
- **Research Goal: “Teach” computers to write fast libraries**
  - Complete automation of implementation and optimization
  - Conquer the “high” algorithm level for automation
- **When a new platform comes out:**  
*Regenerate a retuned library*
- **When a new platform paradigm comes out (e.g., vector or CMPs):**  
*Update the tool rather than rewriting the library*

*Intel has started to use Spiral  
to generate parts of their MKL/IPP library*

# How Spiral Works

## Spiral:

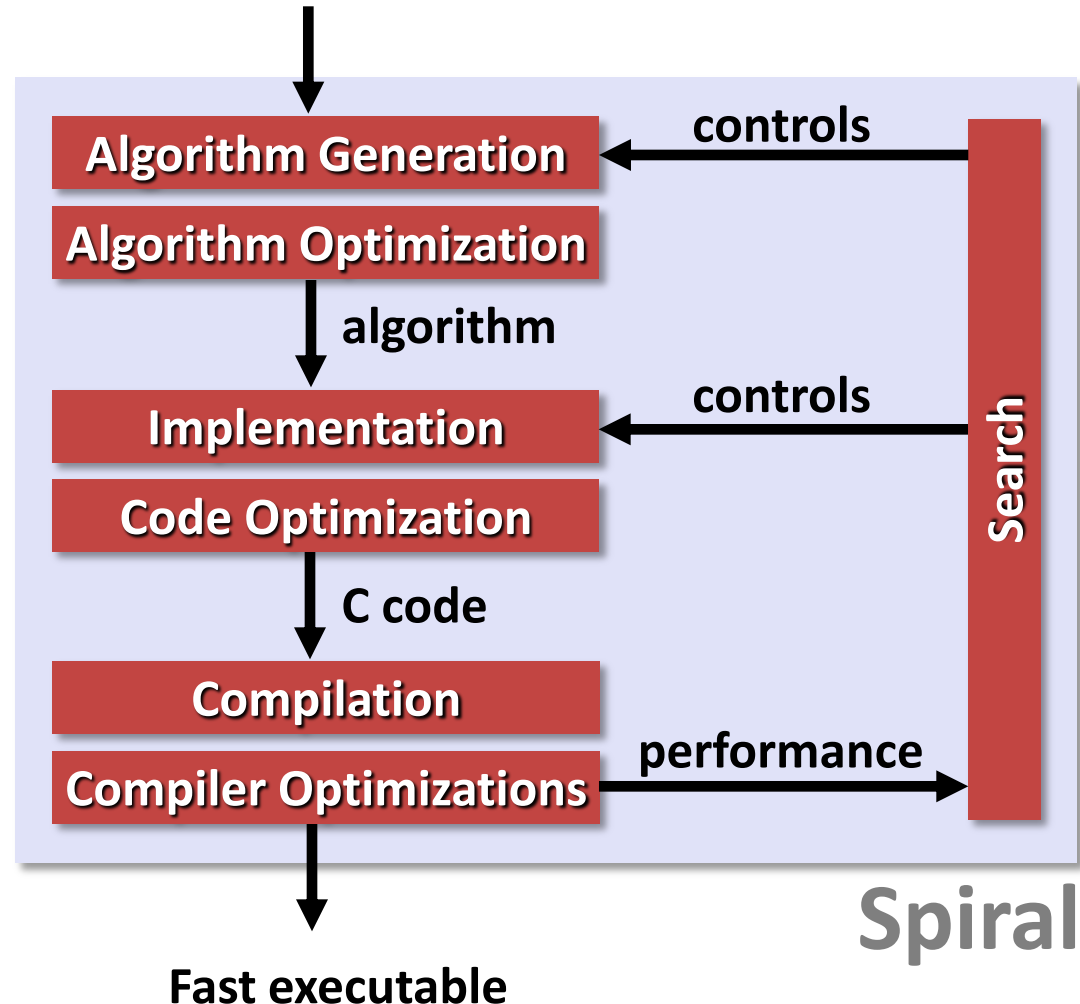
Complete automation of the implementation and optimization task

## Key ideas:

Declarative, mathematical representation of algorithms

Rewriting systems to generate and optimize algorithms at a high level of abstraction

Problem specification (transform)



# Web Interface (Beta Version, @spiral.net)

Program Generation Interface collapse

---

Target platform for optimization:

2x Intel Xeon 3.6 GHz with 2048K L2 cache

parameter	value	explanation
Transform	DCT2 (Discrete Cosine Transform, type 2)	The transform for which you want to request C code
Data type	double	The data type of input and output vector
Transform size	6	The size of the transform = the length of the input vector
Optimize for	runtime	What you want to optimize the code for
Search method	Dynamic Programming	The search method SPIRAL uses (Dynamic Programming is a good choice)
Compiler profile	gcc -O3	Compiler and compiler options used for compilation

Generate code

Browse Archive expand

---

Filter by Platform: All Platforms Selected

Filter by Transform: All Transforms Selected

Filter by Size: All Sizes Selected

Query Database

# What is a (Linear) Transform?

- Mathematically: Matrix-vector multiplication

$$y = \begin{pmatrix} y_0 \\ y_1 \\ \vdots \\ y_{n-1} \end{pmatrix}$$

**Output vector**

$$y = Tx$$



**Transform  
= matrix**

$$x = \begin{pmatrix} x_0 \\ x_1 \\ \vdots \\ x_{n-1} \end{pmatrix}$$

**Input vector**

- Example: Discrete Fourier transform (DFT)

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

# Transform Algorithms: Example 4-point FFT

## Cooley/Tukey fast Fourier transform (FFT):

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

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

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

# Examples: Transforms (currently ≈60)

$$\text{DCT-2}_n = \left[ \cos(k(2l + 1)\pi/2n) \right]_{0 \leq k, l < n},$$

$$\text{DCT-3}_n = \text{DCT-2}_n^T \quad (\text{transpose}),$$

$$\text{DCT-4}_n = \left[ \cos((2k + 1)(2l + 1)\pi/4n) \right]_{0 \leq k, l < n},$$

$$\text{IMDCT}_n = \left[ \cos((2k + 1)(2l + 1 + n)\pi/4n) \right]_{0 \leq k < 2n, 0 \leq l < n},$$

$$\text{RDFT}_n = [r_{kl}]_{0 \leq k, l < n}, \quad r_{kl} = \begin{cases} \cos \frac{2\pi kl}{n}, & k \leq \lfloor \frac{n}{2} \rfloor \\ -\sin \frac{2\pi kl}{n}, & k > \lfloor \frac{n}{2} \rfloor \end{cases},$$

$$\text{WHT}_n = \begin{bmatrix} \text{WHT}_{n/2} & \text{WHT}_{n/2} \\ \text{WHT}_{n/2} & -\text{WHT}_{n/2} \end{bmatrix}, \quad \text{WHT}_2 = \text{DFT}_2,$$

$$\text{DHT} = \left[ \cos(2kl\pi/n) + \sin(2kl\pi/n) \right]_{0 \leq k, l < n}.$$



# Examples: Breakdown Rules (currently $\approx 220$ )

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

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

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

$$\text{DCT-3}_n \rightarrow (\mathbf{I}_m \oplus \mathbf{J}_m) \mathbf{L}_m^n (\text{DCT-3}_m(1/4) \oplus \text{DCT-3}_m(3/4))$$

$$\cdot (\mathbf{F}_2 \otimes \mathbf{I}_m) \begin{bmatrix} \mathbf{I}_m & 0 \oplus -\mathbf{J}_{m-1} \\ \frac{1}{\sqrt{2}}(\mathbf{I}_1 \oplus 2\mathbf{I}_m) \end{bmatrix}, \quad n = 2m$$

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

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

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

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

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

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



**Base case rules**

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

# Breakdown Rules (220)

BSkewDFT 4  
 BSkewDFT\_Base2  
 BSkewDFT\_Base4  
 BSkewDFT\_Fact  
 BSkewDFT\_Decomp  
 BSkewPRDFT 1  
 BSkewPRDFT\_Decomp  
 BSkewRDFT 1  
 BSkewPRDFT\_Decomp  
 Circulant 10  
 RuleCirculant\_Base  
 RuleCirculant\_toFilt  
 RuleCirculant\_Blocking  
 RuleCirculant\_BlockingDense  
 RuleCirculant\_DiagonalizeStep  
 RuleCirculant\_DFT  
 RuleCirculant\_RDFT  
 RuleCirculant\_PRDFT  
 RuleCirculant\_PRDFT1  
 RuleCirculant\_DHT  
 CosDFT 2  
 CosDFT\_Base  
 CosDFT\_Trigonometric  
 DCT1 3  
 DCT1\_Base2  
 DCT1\_Base4  
 DCT1\_DCT1and3  
 DCT2 5  
 DCT2\_DCT2and4  
 DCT2\_toRDFT  
 DCT2\_toRDFT\_odd  
 DCT2\_PrimePowerInduction  
 DCT2\_PrimeFactor  
 DCT3 5  
 DCT3\_Base2  
 DCT3\_Base3  
 DCT3\_Base5  
 DCT3\_Orth9  
 DCT3\_toSkewDCT3  
 DCT4 8  
 DCT4\_Base2  
 DCT4\_Base3  
 DCT4\_DCT2  
 DCT4\_DCT2t  
 DCT4\_DCT2andDST2  
 DCT4\_DST4andDST2  
 DCT4\_Iterative  
 DCT4\_BSkew\_Decomp

DCT5 1  
 DCT5\_Rader  
 DFT 22  
 DFT\_Base  
 DFT\_Cancellation  
 DFT\_CT  
 DFT\_CT\_MinCost  
 DFT\_CT\_DDL  
 DFT\_CosSinSplit  
 DFT\_Rader  
 DFT\_Bluestein  
 DFT\_GoodThomas  
 DFT\_PFA\_SUMS  
 DFT\_PFA\_RaderComb  
 DFT\_SplitRadix  
 DFT\_DCT1andDST1  
 DFT\_DFT1and3  
 DFT\_CT\_Inplace

DRDFT 1  
 DRDFT\_tSPL\_Pease  
 DSCirculant 2  
 RuleDSCirculant\_Base  
 DST1\_Base2  
 DST1\_DST3and1  
 DST2 3  
 DST2\_Base2  
 DST2\_toDCT2  
 DST2\_DST2and4  
 DST4 2  
 DST4\_Base  
 DST4\_toDCT4  
 DTT 15  
 DTT\_C3\_Base2

## Cooley-Tukey FFT

Filt 10  
 RuleFilt\_Base  
 RuleFilt\_Nest  
 RuleFilt\_OverlapSave  
 RuleFilt\_OverlapSaveFreq  
 RuleFilt\_OverlapAdd  
 RuleFilt\_OverlapAdd2  
 RuleFilt\_KaratsubaSimple  
 RuleFilt\_Circulant  
 RuleFilt\_Blocking  
 RuleFilt\_KaratsubaFast  
 IPRDFT 5  
 IPRDFT1\_Base1  
 IPRDFT1\_Base2  
 IPRDFT1\_CT  
 IPRDFT1\_Complex  
 IPRDFT1\_CT\_Radix2  
 IPRDFT2 4  
 IPRDFT2\_Base1  
 IPRDFT2\_Base2  
 IPRDFT2\_CT  
 IPRDFT2\_CT\_Radix2

PDHT2 2  
 PDHT2\_Base2  
 PDHT2\_CT  
 PDHT3 4  
 PDHT3\_Base2  
 PDHT3\_CT  
 PDHT3\_Trig  
 PDHT3\_CT\_Radix2  
 PDHT4 3  
 PDHT4\_Base2  
 PDHT4\_CT  
 PDHT4\_Trig  
 PDST4 2  
 PDST4\_Base2  
 PDST4\_CT  
 PRDFT 10  
 PRDFT1\_Base1  
 PRDFT1\_Base2  
 PRDFT1\_CT  
 PRDFT1\_Complex  
 PRDFT1\_Complex\_T  
 PRDFT1\_Trig  
 PRDFT1\_PF  
 PRDFT1\_CT\_Radix2

PolyDFT 4  
 PolyDFT\_ToNormal  
 PolyDFT\_Base2  
 PolyDFT\_SkewBase2  
 PolyDFT\_SkewDST3\_CT  
 RDFT 4  
 RDFT\_Base  
 RDFT\_Trigonometric  
 RDFT\_CT\_Radix2  
 RDFT\_toDCT2  
 RHT 2  
 RHT\_Base  
 RHT\_CooleyTukey  
 SinDFT 2  
 SinDFT\_Base  
 SinDFT\_Trigonometric  
 SkewDFT 2  
 SkewDFT\_Base2  
 SkewDFT\_Fact  
 SkewDFT 3  
 SkewDFT\_Base2  
 SkewDCT3\_VarSteidl

**Breakdown rules in Spiral:**

- “Teaches” Spiral about existing algorithm knowledge (≈200 journal papers)
- Includes many new ones (algebraic theory)

DFT4 3  
 DFT4\_Base  
 DFT4\_CT  
 DFT4\_PRDFT4  
 DFTBR 1  
 DFTBR\_HW  
 DHT 3  
 DHT\_Base  
 DHT\_DCT1andDST1  
 DHT\_Radix2

DWT\_Base2  
 DWT\_Base4  
 DWT\_Base8  
 DWT\_Lifting  
 DWTper 4  
 DWTper\_Mallat\_2  
 DWTper\_Mallat  
 DWTper\_Polyphase  
 DWTper\_Lifting

MDDFT\_Tensor  
 MDDFT\_RowCol  
 MDDFT\_Dimless  
 MDDFT\_tSPL\_RowCol  
 PDCT4 2  
 PDCT4\_Base2  
 PDCT4\_CT  
 PDHT 3  
 PDHT1\_Base2  
 PDHT1\_CT  
 PDHT1\_CT\_Radix2

PRDFT4\_Base1  
 PRDFT4\_Base2  
 PRDFT4\_CT  
 PRDFT4\_Trig  
 PolyBDFT 5  
 PolyBDFT\_Base2  
 PolyBDFT\_Base4  
 PolyBDFT\_Fact  
 PolyBDFT\_Trig  
 PolyBDFT\_Decomp

WHT 8  
 WHT\_Base  
 WHT\_GeneralSplit  
 WHT\_BinSplit  
 WHT\_DDL  
 WHT\_Dirsum  
 WHT\_tSPL\_BinSplit  
 WHT\_tSPL\_STerm  
 WHT\_tSPL\_Pease

# SPL to Sequential Code

SPL construct	code
$y = (A_n B_n)x$	<pre>t[0:1:n-1] = B(x[0:1:n-1]); y[0:1:n-1] = A(t[0:1:n-1]);</pre>
$y = (I_m \otimes A_n)x$	<pre>for (i=0;i&lt;m;i++)   y[i*n:1:i*n+n-1] =     A(x[i*n:1:i*n+n-1]);</pre>
$y = (A_m \otimes I_n)x$	<pre>for (i=0;i&lt;m;i++)   y[i:n:i+m-1] =     A(x[i:n:i+m-1]);</pre>
$y = \left(\bigoplus_{i=0}^{m-1} A_n^i\right)x$	<pre>for (i=0;i&lt;m;i++)   y[i*n:1:i*n+n-1] =     A(i, x[i*n:1:i*n+n-1]);</pre>
$y = D_{m,n}x$	<pre>for (i=0;i&lt;m*n;i++)   y[i] = Dmn[i]*x[i];</pre>
$y = L_m^{mn}x$	<pre>for (i=0;i&lt;m;i++)   for (j=0;j&lt;n;j++)     y[i+m*j]=x[n*i+j];</pre>

## Example: tensor product

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

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

# Program Generation in Spiral (Sketched)

**Transform**  
*user specified*

DFT<sub>8</sub>



**Fast algorithm  
in SPL**  
*many choices*

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



**Σ-SPL:**

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



**C Code:**

```
void sub(double *y, double *x) {
    double f0, f1, f2, f3, f4, f7, f8, f10, f11;
    f0 = x[0] - x[3];
    f1 = x[0] + x[3];
    f2 = x[1] - x[2];
    f3 = x[1] + x[2];
    f4 = f1 - f3;
    y[0] = f1 + f3;
    y[2] = 0.7071067811865476 * f4;
    f7 = 0.9238795325112867 * f0;
    f8 = 0.3826834323650898 * f2;
    y[1] = f7 + f8;
    f10 = 0.3826834323650898 * f0;
    f11 = (-0.9238795325112867) * f2;
    y[3] = f10 + f11;
}
```

**Optimization at all  
abstraction levels**



parallelization  
vectorization



loop  
optimizations



constant folding  
scalar replacem.  
scheduling

.....

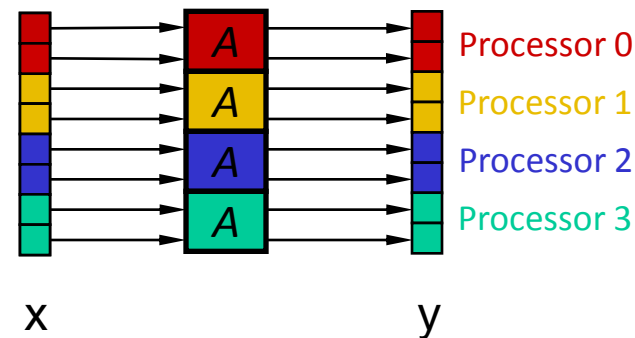
# Organization

- Spiral: Base system
- **Parallelism**
- Libraries (general size)
- Beyond transforms
- Conclusions

# SPL to Shared Memory Code: Basic Idea [SC 06]

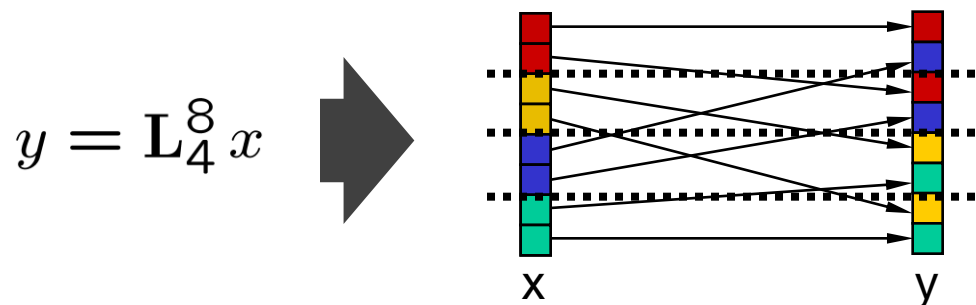
- Governing construct: tensor product

$$y = \left( I_p \otimes A \right) x$$



**p-way embarrassingly parallel, load-balanced**

- Problematic construct: permutations produce false sharing



**Task: Rewrite formulas to extract tensor product + keep contiguous blocks**

# Step 1: Shared Memory Tags

- Identify crucial hardware parameters
  - Number of processors:  $p$
  - Cache line size:  $\mu$
- Introduce them as tags in SPL

$$\overset{A}{\text{smp}(p, \mu)}$$

**This means:** formula  $A$  is to be optimized for  $p$  processors and cache line size  $\mu$

## Step 2: Identify “Good” Formulas

- Load balanced, avoiding false sharing

$$y = (I_p \otimes A)x \quad \text{with } A \in \mathbb{C}^{m\mu \times m\mu}$$

$$y = \left( \bigoplus_{i=0}^{p-1} A_i \right) x \quad \text{with } A_i \in \mathbb{C}^{m\mu \times m\mu}$$

$$y = (P \otimes I_\mu)x \quad \text{with } P \text{ a permutation matrix} \quad \textit{Blackboard}$$

- Tagged operators (no further rewriting necessary)

$$I_p \otimes_{\parallel} A, \quad \bigoplus_{i=0}^{p-1} \parallel A_i, \quad P \bar{\otimes} I_\mu$$

- **Definition:** A formula is **fully optimized** for  $(p, \mu)$  if it is one of the above or of the form

$$I_m \otimes A \quad \text{or} \quad AB$$

where A and B are fully optimized.



# Step 3: Identify Rewriting Rules

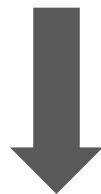
## ■ Goal: Transform formulas into fully optimized formulas

- Formulas rewritten, tags propagated
- There may be choices

$$\begin{aligned}
 \underbrace{AB}_{\text{smp}(p,\mu)} &\rightarrow \underbrace{A}_{\text{smp}(p,\mu)} \underbrace{B}_{\text{smp}(p,\mu)} \\
 \underbrace{A_m \otimes I_n}_{\text{smp}(p,\mu)} &\rightarrow \underbrace{\left( L_m^{mp} \otimes I_{n/p} \right) \left( I_p \otimes (A_m \otimes I_{n/p}) \right) \left( L_p^{mp} \otimes I_{n/p} \right)}_{\text{smp}(p,\mu)} \\
 \underbrace{L_m^{mn}}_{\text{smp}(p,\mu)} &\rightarrow \begin{cases} \underbrace{\left( I_p \otimes L_{m/p}^{mn/p} \right)}_{\text{smp}(p,\mu)} \underbrace{\left( L_p^{pn} \otimes I_{m/p} \right)}_{\text{smp}(p,\mu)} \\ \underbrace{\left( L_m^{pm} \otimes I_{n/p} \right)}_{\text{smp}(p,\mu)} \underbrace{\left( I_p \otimes L_m^{mn/p} \right)}_{\text{smp}(p,\mu)} \end{cases} \\
 \underbrace{I_m \otimes A_n}_{\text{smp}(p,\mu)} &\rightarrow I_p \otimes_{\parallel} \left( I_{m/p} \otimes A_n \right) \\
 \underbrace{(P \otimes I_n)}_{\text{smp}(p,\mu)} &\rightarrow \left( P \otimes I_{n/\mu} \right) \bar{\otimes} I_\mu
 \end{aligned}$$

# Simple Rewriting Example

$$A_m \otimes I_n$$



**Loop splitting + loop exchange**

$$\left( L_m^{mp} \otimes I_{n/p} \right) \left( I_p \otimes_{\parallel} (A_m \otimes I_{n/p}) \right) \left( L_p^{mp} \otimes I_{n/p} \right)$$

**fully optimized**

```
parallel for (i=0; i<p; i++)
  for (j=0; j<n/p; j++)
    y[i*n/p+j:n:i*n/p+j+m-1] =
      A(x[i*n/p+j:n:i*n/p+j+m-1]);
```

# Parallelization by Rewriting

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

**Fully optimized (load-balanced, no false sharing)**  
**in the sense of our definition**

# Same Approach for Other Parallel Paradigms

## Message Passing: [ISPA 06]

$$\begin{aligned}
 \underbrace{\text{DFT}_{mn}}_{\text{par}(p)} &\rightarrow \underbrace{(\text{DFT}_m \otimes \text{I}_n)}_{\text{par}(p-q)} \underbrace{\text{T}_n^{mn}}_{\text{par}(q)} \underbrace{(\text{I}_m \otimes \text{DFT}_n)}_{\text{par}(q)} \underbrace{\text{L}_m^{mn}}_{\text{par}(q-p)} \\
 &\dots \\
 &\dots \\
 &\rightarrow \underbrace{(\text{I}_p \otimes \text{L}_{m/p}^{mn/p})}_{\text{par}(p)} \underbrace{(\text{L}_p^{p^2} \otimes \text{I}_{mn/p^2})}_{\text{par}(p-q)} \underbrace{(\text{I}_q \otimes (\text{I}_{p/q} \otimes \text{L}_p^n \otimes \text{I}_{m/p}))}_{\text{par}(q)} \underbrace{(\text{I}_q \otimes (\text{I}_{n/q} \otimes \text{DFT}_m))}_{\text{par}(q)} \\
 &\dots \\
 &\rightarrow \underbrace{(\text{I}_q \otimes \text{L}_{m/q}^{mn/q})}_{\text{par}(q)} \underbrace{(\text{L}_q^{q^2} \otimes \text{I}_{mn/q^2})}_{\text{par}(q)} \underbrace{(\text{I}_q \otimes (\text{L}_{p/q}^n \otimes \text{I}_{m/q}))}_{\text{par}(q)} \underbrace{\text{T}_n^{mn}}_{\text{par}(q)} \underbrace{(\text{I}_q \otimes (\text{I}_{m/q} \otimes \text{DFT}_n))}_{\text{par}(q)} \\
 &\dots \\
 &\rightarrow \underbrace{(\text{I}_q \otimes (\text{I}_{p/q} \otimes \text{L}_{m/p}^{mn/p}))}_{\text{par}(q)} \underbrace{(\text{L}_p^{p^2} \otimes \text{I}_{mn/p^2})}_{\text{par}(q-p)} \underbrace{(\text{I}_p \otimes (\text{L}_p^n \otimes \text{I}_{m/p}))}_{\text{par}(p)}
 \end{aligned}$$

With Bonelli, Lorenz, Ueberhuber, TU Vienna

## Vectorization: [IPDPS 02, VecPar 06]

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

## Cg/OpenGL for GPUs:

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

With Shen, TU Denmark

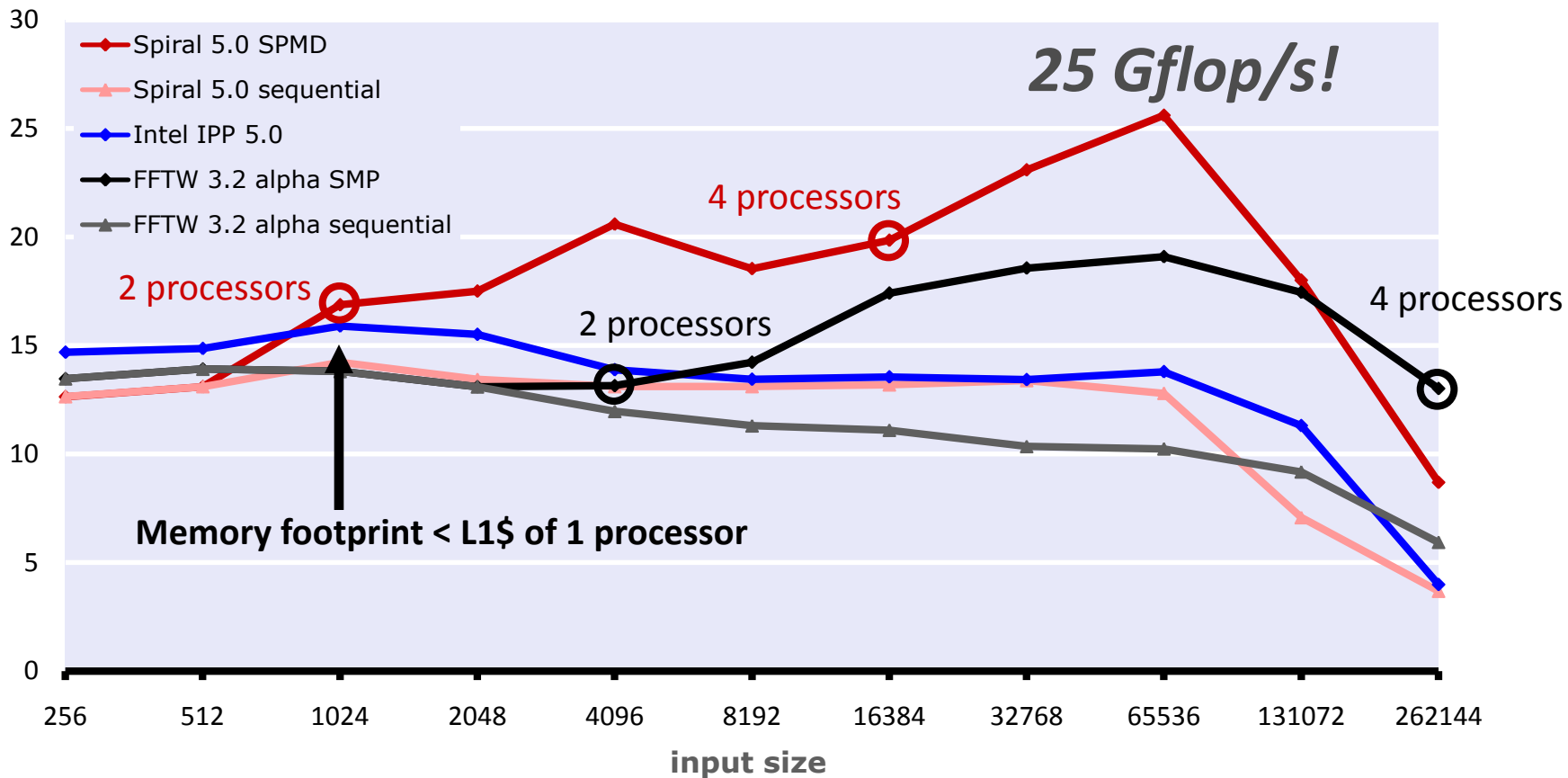
## Verilog for FPGAs: [DAC 05]

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

With Milder, Hoe, CMU

# Benchmark: Vector and SMP

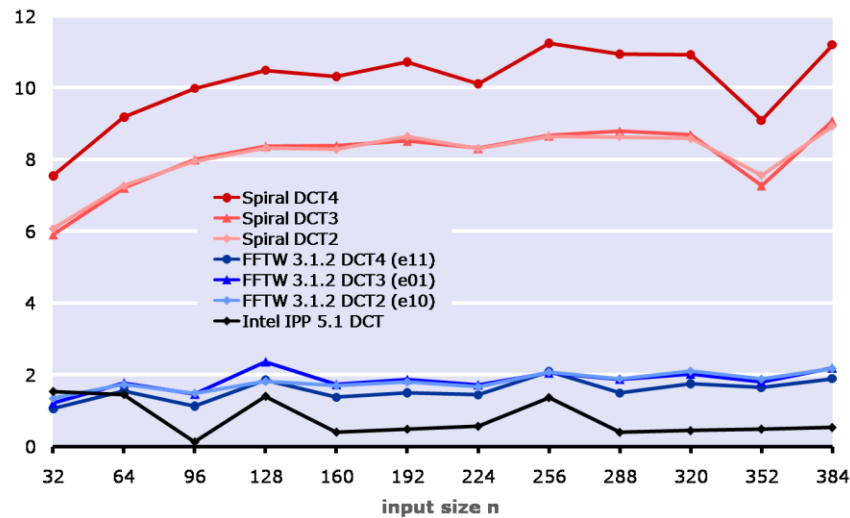
**DFT (single precision): on 3 GHz 2 x Core 2 Extreme  
performance [Gflop/s]**



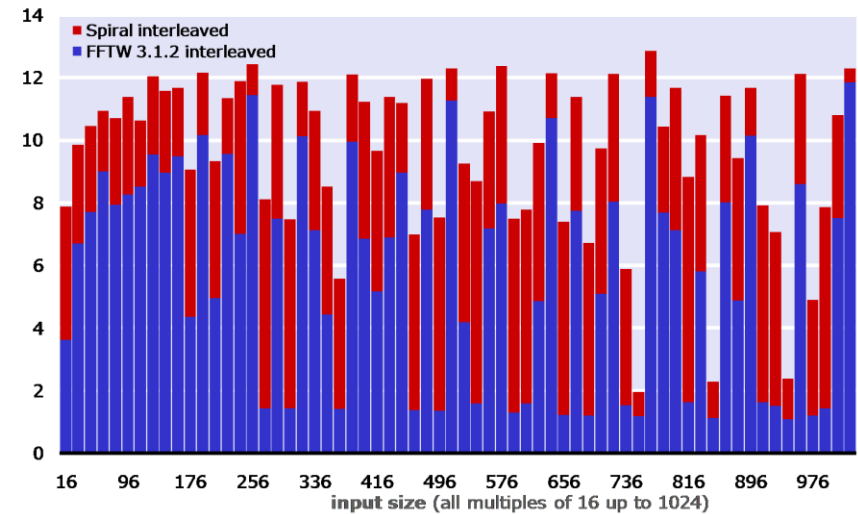
**4-way vectorized + up to 4-threaded + adapted to the memory hierarchy**

# Other Transforms

**DCT on 2.66 GHz Core2 (single-precision, 4-way SSE)**  
performance [Gflop/s]



**FFT on 2.66 GHz Core2 (single-precision, 4-way SSE)**  
performance [Gflop/s]



Less popular transforms or sizes often means suboptimal code

# Organization

- Spiral: Base system
- Parallelism
- **Libraries (general size)**
- Beyond transforms
- Conclusions

# Ultimate Challenge: General Size Libraries

## Spiral so far:

### Fixed size code

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

- Algorithm = recursion fixed at code generation time

## Challenge:

### Recursive general size library (FFTW, MKL)

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

- Algorithm = recursion fixed at time of use
- Requires infrastructure (search, precomputing, base cases = codelets)
- Many other challenges

$$(\text{DFT}_k \otimes \text{I}_m) \text{T}_m^n (\text{I}_k \otimes \text{DFT}_m) \text{L}_k^n$$

**Possible?**



**Vectorized, parallel,  
general-size, adaptive library**



# 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 in FFTW 2.x

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

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

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

```
}
```

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

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

**2 additional functions (recursion steps) needed**

***How to discover automatically?***

# $\Sigma$ -SPL : Basic Idea

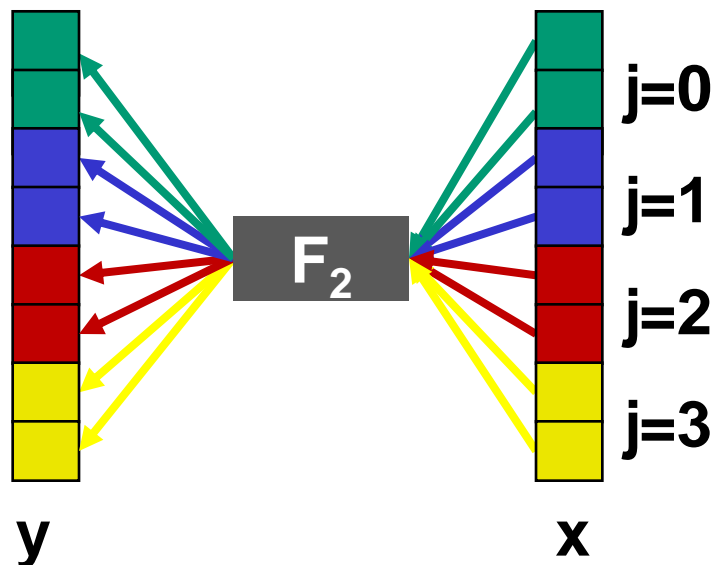
- Four central constructs:  $\Sigma$ ,  $G$ ,  $S$ ,  $\text{Perm}$

- $\Sigma$  (sum) – explicit loop
- $G_f$  (gather) – load data with index mapping  $f$
- $S_f$  (scatter) – store data with index mapping  $f$
- $\text{Perm}_f$  – permute data with the index mapping  $f$

- $\Sigma$ -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

- **Input:** transform T and a breakdown rule
- **Output:** recursion step closure + implementation of each recursion step
- **Algorithm:**

1. Apply the breakdown rule to T

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



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

2. Convert to  $\Sigma$ -SPL



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

3. Apply loop merging + index simplification rules.



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

4. Extract required recursion steps

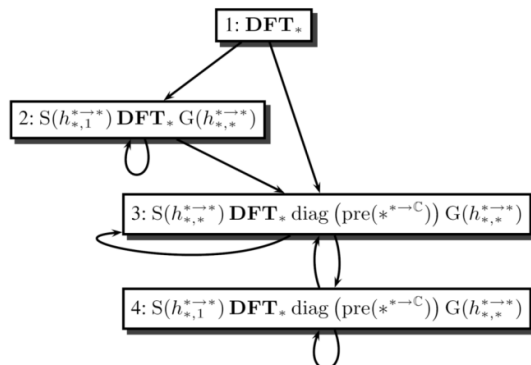


$$\sum_{i=0}^{k-1} \left\{ S_{h_{i,k}} \text{DFT}_{n/k} \text{diag}(f \circ h_{i,k}) G_{h_{i,k}} \right\} \quad \sum_{j=0}^{n/k-1} \left\{ S_{h_{j,k,1}} \text{DFT}_k G_{h_{j,n/k}} \right\}$$

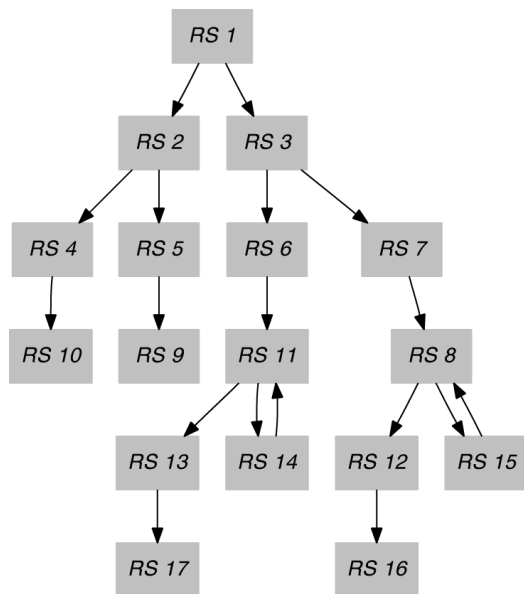
5. Repeat until closure is reached

# Recursion Step Closure: Example

## DFT: scalar code



## DCT4: vector code



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

# Generated Full Libraries

**Input:** breakdown rules (a few are not shown)

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

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

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

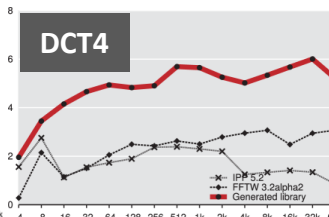
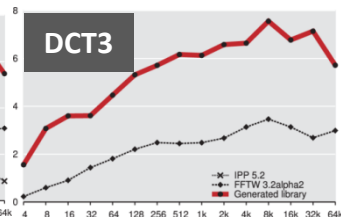
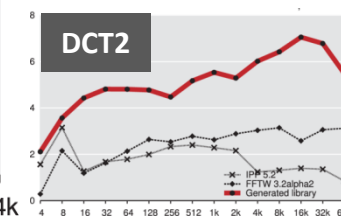
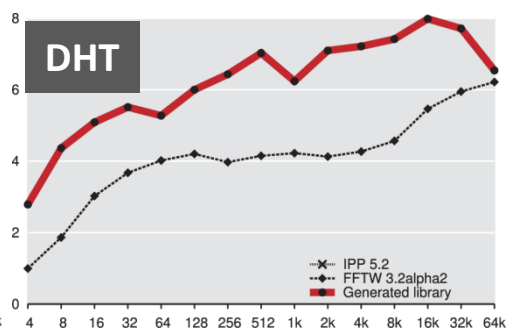
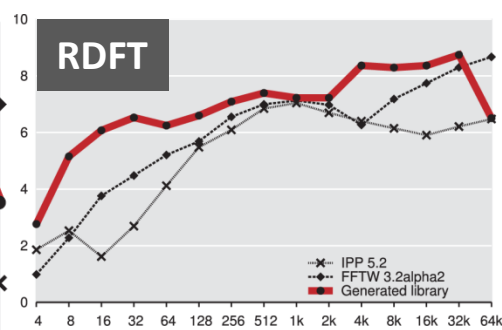
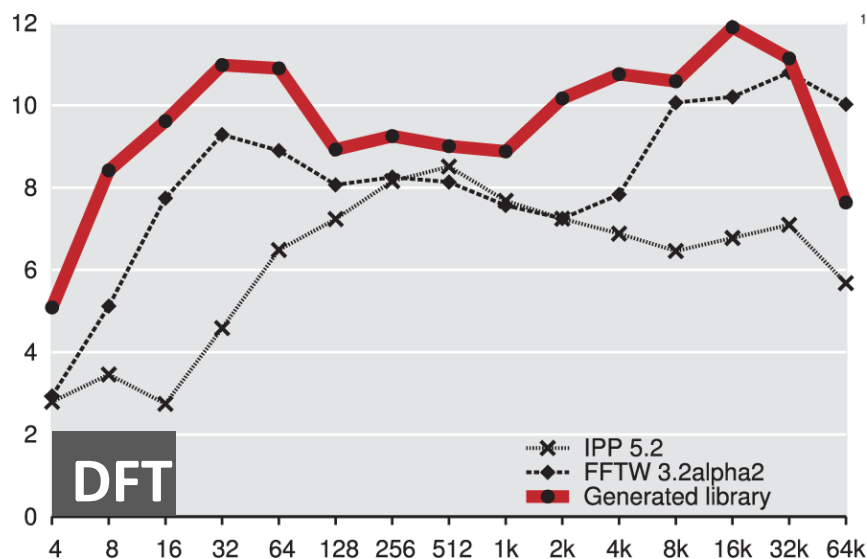
$$\text{RDFT-3}_n = (Q_{k/2,m}^\top \otimes I_2) (I_k \otimes_i \text{rDFT}_{2m}(i+1/2)/k) (\text{RDFT-3}_k \otimes I_m), \quad k \text{ even},$$

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

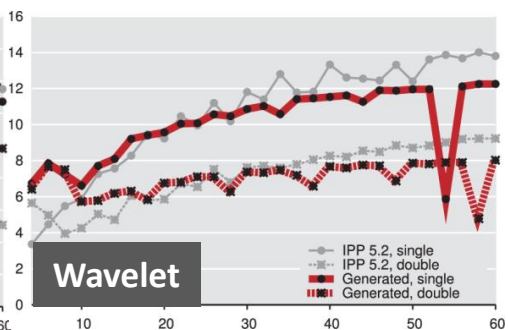
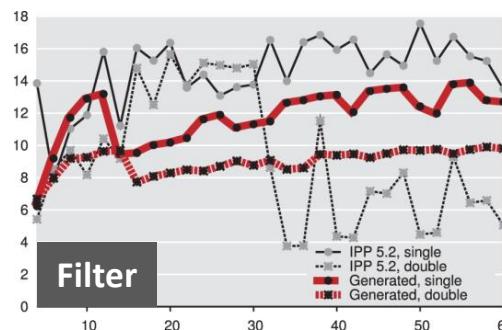
$$\text{DCT-3}_n = \text{DCT-2}_n^\top,$$

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

# Generated Full Libraries



- 2-way vectorized, 2-threaded
- Most are faster than hand-written libs
- Recursion steps: 4–17
- Code size: 8–120 kloc or 0.5–5 MB
- Generation time: 1–3 hours



**Total: 300 kloc or 13.3 MB of code generated in < 20 hours**  
*from a few breakdown rules*

# Organization

- Spiral: Base system
- Parallelism
- Libraries (general size)
- **Beyond transforms**
- Conclusions

# Going Beyond Transforms

- Transform = **linear** operator with **one** vector input and **one** vector output
- Key ideas:
  - Generalize to (**possibly nonlinear**) operators with **several** inputs and **several** outputs
  - Generalize SPL (including tensor product) to OL (operator language)
  - Generalize rewriting systems for parallelizations

## Basic operators

name	definition
stride	$L_n^{mn} : \mathbb{R}^{mn} \rightarrow \mathbb{R}^{mn}; M \mapsto M^T$
vector sum	$\Sigma_m : \mathbb{R}^m \rightarrow \mathbb{R}; \mathbf{x} \mapsto \sum_{i=0}^{m-1} x_i$
point-wise multiplication	$(\cdot)_m : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}^m; (\mathbf{x}, \mathbf{y}) \mapsto (x_0y_0, \dots, x_{m-1}y_{m-1})$
Kronecker product	$(\otimes)_{m \times n} : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}^{mn}; (\mathbf{x}, \mathbf{y}) \mapsto x_0\mathbf{y} \oplus \dots \oplus x_{m-1}\mathbf{y}$

## Operations (higher-order operators)

name	definition
composition	$(\mathcal{B} \circ \mathcal{A})(\mathbf{x}, \mathbf{y}) = \mathcal{B}(\mathcal{A}(\mathbf{x}, \mathbf{y}))$
Cartesian product	$(\mathcal{A} \times \mathcal{B})(\mathbf{x}, \mathbf{y}, \mathbf{z}, \mathbf{t}) = \mathcal{A}(\mathbf{x}, \mathbf{y}) \times \mathcal{B}(\mathbf{z}, \mathbf{t})$
tensor product	$(\mathcal{A} \otimes \mathcal{B})(\mathbf{x}, \mathbf{y}) = \sum_{0 \leq i < p, 0 \leq j < r} \mathcal{A}(\mathbf{e}_{p,i}, \mathbf{e}_{r,j}) \otimes \mathcal{B}(\mathbf{x}^i, \mathbf{y}^j)$



# Example: MMM

Breakdown rules capture various forms of blocking:

breakdown rule	description
$MMM_{1,1,1} \rightarrow (\cdot)_1$	base case
$MMM_{m,n,k} \rightarrow (\otimes)_{m/m_b \times 1} \otimes MMM_{m_b,n,k}$	horizontal blocking
$MMM_{m,n,k} \rightarrow MMM_{m,n_b,k} \otimes (\otimes)_{1 \times n/n_b}$	interleaved blocking
$MMM_{m,n,k} \rightarrow ((\sum_{k/k_b} \circ (\cdot)_{k/k_b}) \otimes MMM_{m,n,k_b}) \circ ((L_{k/k_b}^{mk/k_b} \otimes I_{k_b}) \times I_{kn})$	accumulative blocking
$MMM_{m,n,k} \rightarrow (L_m^{mn/n_b} \otimes I_{n_b}) \circ ((\otimes)_{1 \times n/n_b} \otimes MMM_{m,n_b,k}) \circ (I_{km} \times (L_{n/n_b}^{kn/n_b} \otimes I_{n_b}))$	vertical blocking

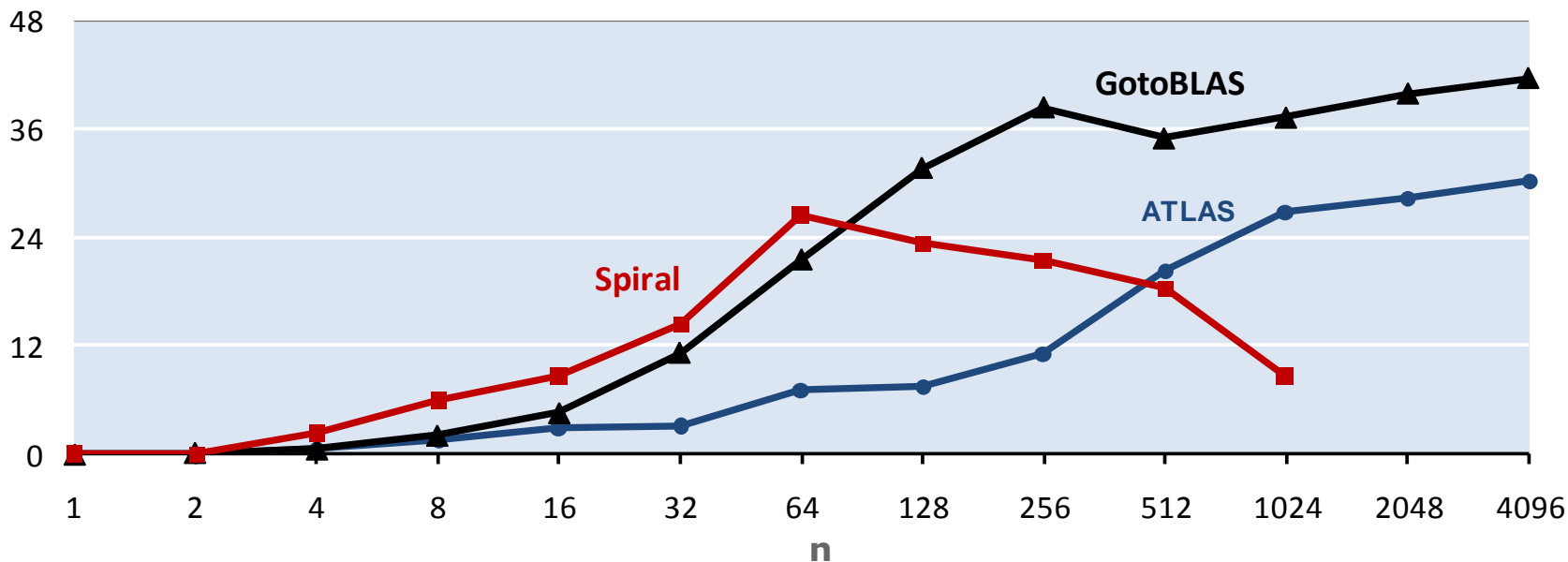
## Parallelization through rewriting

$$\begin{aligned}
 & \underbrace{MMM_{m,n,k}}_{\text{smp}(p,\mu)} \\
 & \rightarrow \underbrace{(I_m \otimes L_p^n)}_{\text{smp}(p,\mu)} \circ \underbrace{(MMM_{m,n/p,k} \otimes (\otimes)_{1 \times p \rightarrow p})}_{\text{smp}(p,\mu)} \circ (I_{km} \times (I_k \otimes L_{n/p}^n)) \\
 & \rightarrow \underbrace{(I_m \otimes L_p^n)}_{\text{smp}(p,\mu)} \circ \underbrace{(MMM_{m,n/p,k} \otimes (\otimes)_{1 \times p \rightarrow p})}_{\text{smp}(p,\mu)} \circ \underbrace{(I_{km} \times (I_k \otimes L_{n/p}^n))}_{\text{smp}(p,\mu)} \\
 & \rightarrow \underbrace{(I_m \otimes L_p^n)}_{\text{smp}(p,\mu)} \circ \underbrace{L_{m/pm}^{mn}}_{\text{smp}(p,\mu)} \circ ((\otimes)_{1 \times p \rightarrow p} \otimes \parallel MMM_{m/p,n,k}) \circ \underbrace{(I_{km} \times L_p^{kn})}_{\text{smp}(p,\mu)} \circ \underbrace{(I_{km} \times (I_k \otimes L_{n/p}^n))}_{\text{smp}(p,\mu)} \\
 & \rightarrow \underbrace{(L_m^{mp} \otimes I_{n/(p\mu)}) \bar{\otimes} I_\mu}_{\text{smp}(p,\mu)} \circ \underbrace{((\otimes)_{1 \times p \rightarrow p} \otimes \parallel MMM_{m,n/p,k})}_{\text{smp}(p,\mu)} \circ \underbrace{((I_{km/\mu} \bar{\otimes} I_\mu) \times ((L_p^{kp} \otimes I_{n/(p\mu)}) \bar{\otimes} I_\mu))}_{\text{smp}(p,\mu)}
 \end{aligned}$$

**Load-balanced**  
**No false sharing**

# First Results

## MMM performance on 2x2 core2 3Ghz (nxn Matrices, float) performance [Gflop/s]



### ■ Other numerical problems we currently study

- SAR imaging
- Viterbi decoding
- JPEG 2000 decoding

# Organization

- Spiral: Base system
- Parallelism
- Libraries (general size)
- Beyond transforms
- **Conclusions**

# Summary

- **Spiral: The computer writes fast transform libraries**
  - Complete automation
  - Including memory hierarchy optimizations, vectorization, parallelization
  
- **What we have learned**
  - Declarative representation of algorithms (mathematical domain-specific language)
  - Optimization at a high, “right” level of abstraction using rewriting
  - It makes sense to use math to represent and optimize math functionality
  - It makes sense to “teach” the computer algorithms and math (does not become obsolete)
  - Domain-specific is necessary to get fastest code
  - One needs techniques from different disciplines .....

**Programming languages**  
Program generation

**Symbolic Computation**  
Rewriting

***Automating  
High-Performance  
Numerical Library  
Development***

**Software**  
Scientific Computing

**Algorithms**  
Mathematics

**Compilers**