

# Spiral

## Automating Library Development

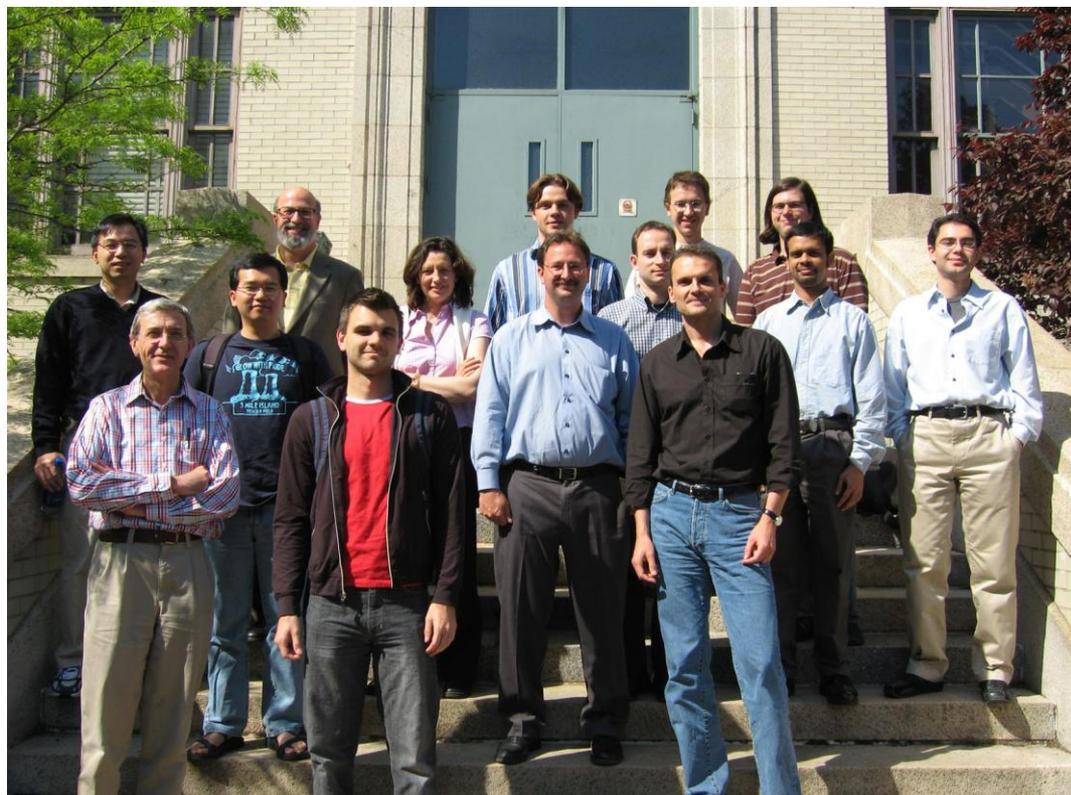
Markus Püschel

*... and the Spiral team (only part shown)*

***With:***

Srinivas Chellappa  
Frédéric de Mesmay  
Franz Franchetti  
Daniel McFarlin  
Yevgen Voronenko

Electrical and  
Computer Engineering  
Carnegie Mellon University



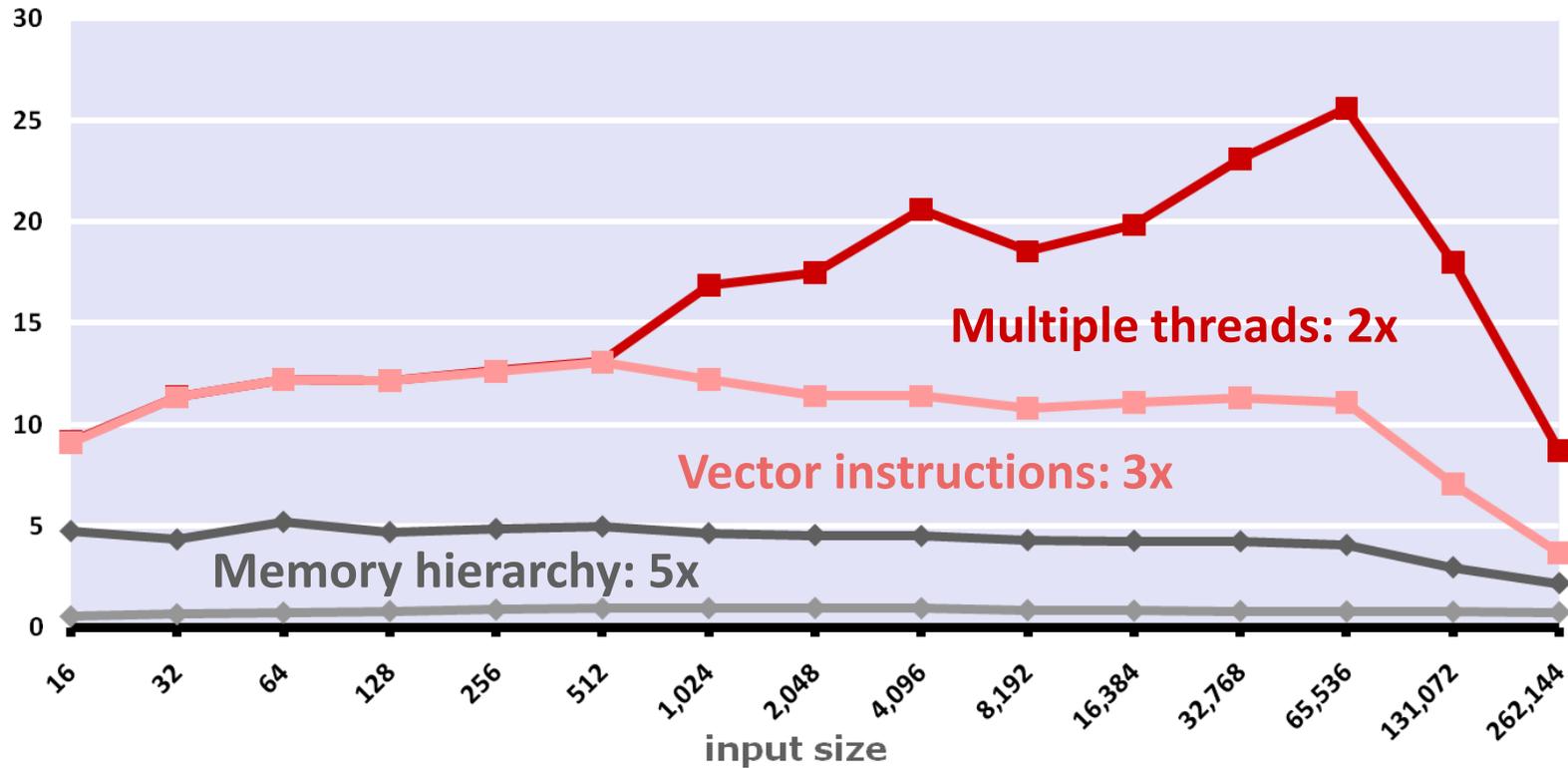
This work was supported by  
DARPA, NSF-NGS/ITR, ACR, CPA, Intel, Mercury, National Instruments

# Positions and Thoughts

- **Autotuning definition**
  - Search over space of alternatives and
  - Parameter-based tuning are very important
  - but fails to address some key problems; we need to think about
- **Raising the level of abstraction: Enables**
  - Use of domain knowledge
  - Difficult optimizations: parallelization, vectorization, etc.
  - Faster porting to new platforms and platform paradigms
  - Possibly automatic software development
- **We need coarse platform abstractions**
- **We need more interdisciplinary collaborations**
- **Metrics**
  - Time for code development, porting to new platforms
  - Performance

# DFT Plot: Analysis

Discrete Fourier Transform (DFT) on 2xCore2Duo 3 GHz  
Performance [Gflop/s]



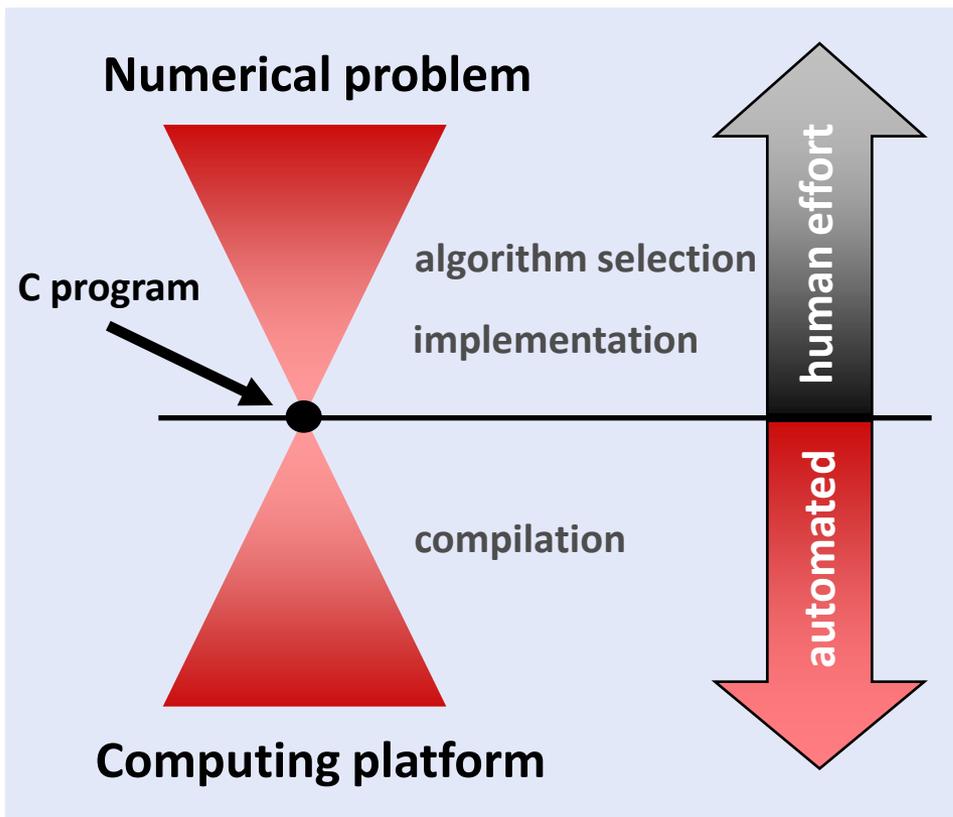
*High performance library development has become a nightmare*

# Spiral

- **Research Goal: “Teach” computers to write fast libraries**
  - Complete automation of implementation and optimization
  - Including vectorization, parallelization
  
- **Functionality:**
  - Linear transforms (discrete Fourier transform, filters, wavelets)
  - BLAS
  - SAR imaging
  - En/decoding (Viterbi, Ebcot in JPEG2000)
  - ... more
  
- **Platforms:**
  - Desktop (vector, SMP), FPGAs, GPUs, distributed, hybrid
  
- **Collaboration with Intel (Kuck, Tang, Sabanin)**
  - Parts of MKL/IPP generated with Spiral
  - IPP 6.0: ippg domain for Spiral generated code

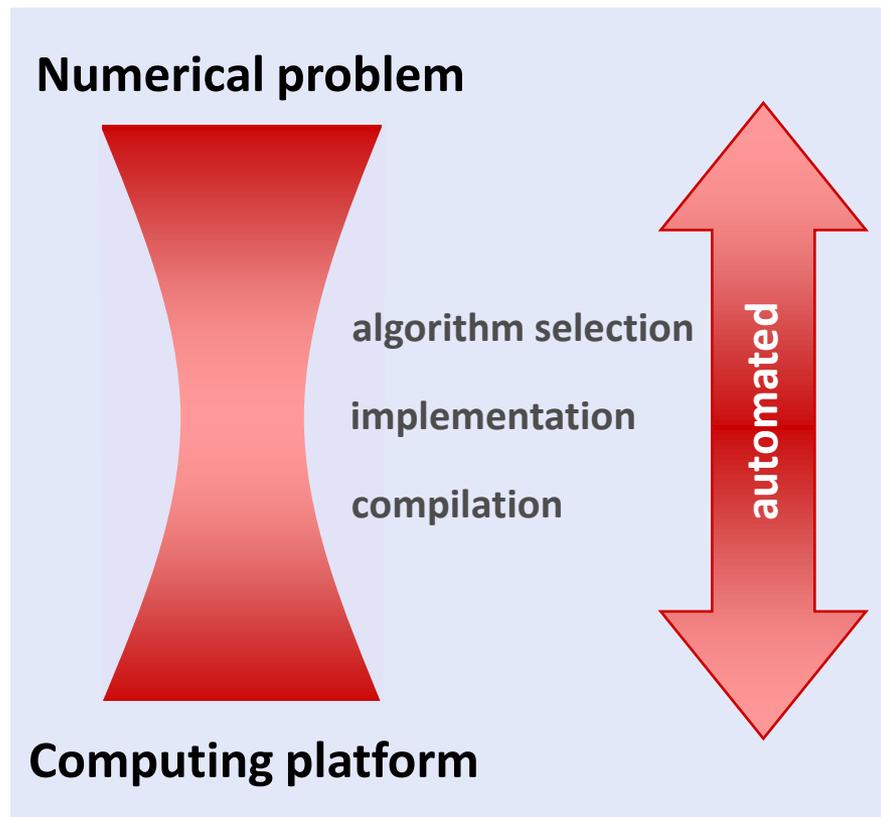
# Vision Behind Spiral

## Current



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

## Future



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

# Organization

- **Spiral's framework: Example transforms**
  - Complete automation achieved
  
- Beyond transforms
  
  
- Conclusions and thoughts

# Linear Transforms

- 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):

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

Fourier transform

Diagonal matrix (twiddles)

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

Kronecker product
Identity
Permutation

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

# Breakdown Rules (>200 for >50 Transforms)

$$\begin{aligned}
 \text{DFT}_n &\rightarrow P_{k/2,2m}^\top \left( \text{DFT}_{2m} \oplus \left( I_{k/2-1} \otimes_i C_{2m} \text{rDFT}_{2m}(i/k) \right) \right) \left( \text{RDFT}'_k \otimes I_m \right), \quad k \text{ even,} \\
 \begin{pmatrix} \text{RDFT}_n \\ \text{RDFT}'_n \\ \text{DHT}_n \\ \text{DHT}'_n \end{pmatrix} &\rightarrow \left( P_{k/2,2m}^\top \otimes I_2 \right) \left( \begin{pmatrix} \text{RDFT}_{2m} \\ \text{RDFT}'_{2m} \\ \text{DHT}_{2m} \\ \text{DHT}'_{2m} \end{pmatrix} \oplus \left( I_{k/2-1} \otimes_i D_{2m} \begin{pmatrix} \text{rDFT}_{2m}(i/k) \\ \text{rDFT}'_{2m}(i/k) \\ \text{rDHT}_{2m}(i/k) \\ \text{rDHT}'_{2m}(i/k) \end{pmatrix} \right) \right) \left( \begin{pmatrix} \text{RDFT}'_k \\ \text{RDFT}_k \\ \text{DHT}'_k \\ \text{DHT}_k \end{pmatrix} \otimes I_m \right), \quad k \text{ even,} \\
 \begin{pmatrix} \text{rDFT}_{2n}(u) \\ \text{rDHT}_{2n}(u) \end{pmatrix} &\rightarrow L_m^{2n} \left( I_k \otimes_i \begin{pmatrix} \text{rDFT}_{2m}((i+u)/k) \\ \text{rDHT}_{2m}((i+u)/k) \end{pmatrix} \right) \left( \begin{pmatrix} \text{rDFT}_{2k}(u) \\ \text{rDHT}_{2k}(u) \end{pmatrix} \otimes I_m \right), \\
 \text{RDFT-3}_n &\rightarrow \left( Q_{k/2,2m}^\top \otimes I_2 \right) \left( I_k \otimes_i \text{rDFT}_{2m}(i+1/2/k) \right) \left( \text{RDFT-3}_k \otimes I_m \right), \quad k \text{ even,} \\
 \text{DCT-2}_n &\rightarrow P_{k/2,2m}^\top \left( \text{DCT-2}_{2m} K_2^{2m} \oplus \left( I_{k/2-1} \otimes N_{2m} \text{RDFT-3}_{2m}^\top \right) \right) B_n \left( L_{k/2}^{n/2} \otimes I_2 \right) \left( I_m \otimes \text{RDFT}'_k \right) Q_{m/2,k}, \\
 \text{DCT-3}_n &\rightarrow \text{DCT-2}_n^\top, \\
 \text{DCT-4}_n &\rightarrow Q_{k/2,2m}^\top \left( I_{k/2} \otimes N_{2m} \text{RDFT-3}_{2m}^\top \right) B'_n \left( L_{k/2}^{n/2} \otimes I_2 \right) \left( I_m \otimes \text{RDFT-3}_k \right) Q_{m/2,k}. \\
 \text{DFT}_n &\rightarrow \left( \text{DFT}_k \otimes I_m \right) \Upsilon_m^n \left( I_k \otimes \text{DFT}_m \right) L_k^n, \quad n = km \\
 \text{DFT}_n &\rightarrow P_n \left( \text{DFT}_k \otimes \text{DFT}_m \right) Q_n, \quad n = km, \text{ gcd}(k, m) = 1 \\
 \text{DFT}_p &\rightarrow R_p^\top \left( I_1 \oplus \text{DFT}_{p-1} \right) D_p \left( I_1 \oplus \text{DFT}_{p-1} \right) R_p, \quad p \text{ prime} \\
 \text{DCT-3}_n &\rightarrow \left( I_m \oplus J_m \right) L_m^n \left( \text{DCT-3}_m(1/4) \oplus \text{DCT-3}_m(3/4) \right) \\
 &\quad \cdot \left( F_2 \otimes I_m \right) \begin{bmatrix} I_m & 0 \oplus -J_{m-1} \\ \frac{1}{\sqrt{2}}(I_1 \oplus 2I_m) & \end{bmatrix}, \quad n = 2m \\
 \text{DCT-4}_n &\rightarrow S_n \text{DCT-2}_n \text{diag}_{0 \leq k < n} \left( 1 / (2 \cos((2k+1)\pi/4n)) \right) \\
 \text{IMDCT}_{2m} &\rightarrow \left( J_m \oplus I_m \oplus I_m \oplus J_m \right) \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} \text{DCT-4}_{2m} \\
 \text{WHT}_{2^k} &\rightarrow \prod_{i=1}^t \left( I_{2^{k_1+\dots+k_{i-1}}} \otimes \text{WHT}_{2^{k_i}} \otimes I_{2^{k_{i+1}+\dots+k_t}} \right), \quad k = k_1 + \dots + k_t \\
 \text{DFT}_2 &\rightarrow F_2 \\
 \text{DCT-2}_2 &\rightarrow \text{diag}(1, 1/\sqrt{2}) F_2 \\
 \text{DCT-4}_2 &\rightarrow J_2 R_{13\pi/8}
 \end{aligned}$$

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

# 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:**  
*[PLDI 05]*

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



**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

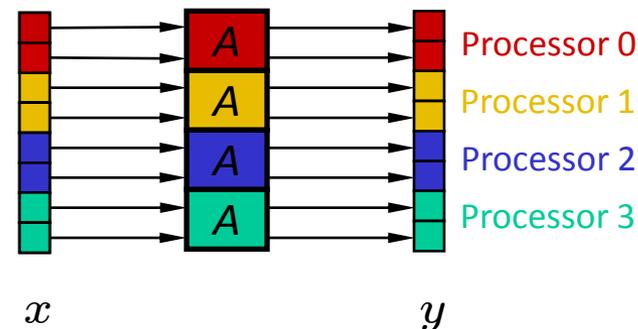


constant folding  
scheduling  
.....

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

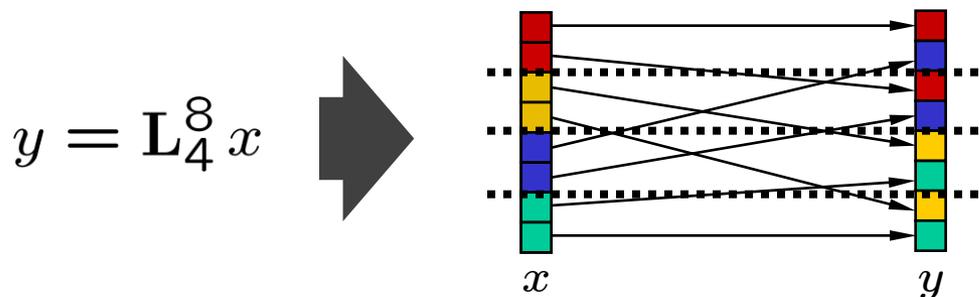
- Governing construct: tensor product

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



*p-way embarrassingly parallel, load-balanced*

- Problematic construct: permutations produce false sharing



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

# Parallelization by Rewriting

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

**Load-balanced**

**No false sharing**

# Same Approach for Other Parallel Paradigms

## Message Passing

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

## Vectorization

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

## Cg/OpenGL for GPUs:

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

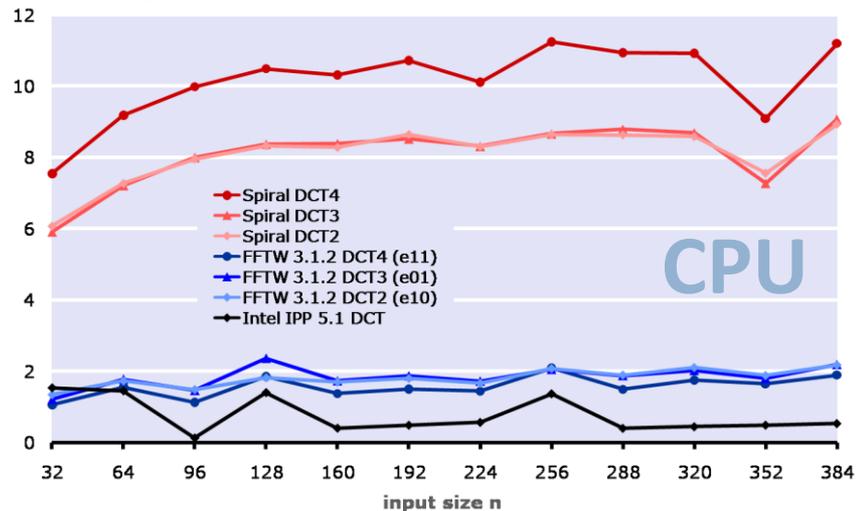
## Verilog for FPGAs:

$$\begin{aligned}
 \underbrace{(\text{DFT}_{r^k})}_{\text{stream}(r^s)} &\rightarrow \underbrace{\left[ \prod_{i=0}^{k-1} \text{L}_r^{r^k} (\text{I}_{r^{k-1}} \otimes \text{DFT}_r) (\text{L}_{r^{k-i-1}}^{r^k} (\text{I}_{r^i} \otimes \text{T}_{r^{k-i-1}}^{r^{k-i}}) \text{L}_{r^{i+1}}^{r^k}) \right]}_{\text{stream}(r^s)} \text{R}_r^{r^k} \\
 &\dots \\
 &\rightarrow \left[ \prod_{i=0}^{k-1} \underbrace{\text{L}_r^{r^k}}_{\text{stream}(r^s)} \underbrace{(\text{I}_{r^{k-1}} \otimes \text{DFT}_r)}_{\text{stream}(r^s)} \underbrace{(\text{L}_{r^{k-i-1}}^{r^k} (\text{I}_{r^i} \otimes \text{T}_{r^{k-i-1}}^{r^{k-i}}) \text{L}_{r^{i+1}}^{r^k})}_{\text{stream}(r^s)} \right] \text{R}_r^{r^k} \\
 &\dots \\
 &\rightarrow \left[ \prod_{i=0}^{k-1} \underbrace{\text{L}_r^{r^k}}_{\text{stream}(r^s)} (\text{I}_{r^{k-s-1}} \otimes \times \underbrace{(\text{I}_{r^{s-1}} \otimes \text{DFT}_r)}_{\text{stream}(r^s)}) \underbrace{\text{T}_i^r}_{\text{stream}(r^s)} \right] \text{R}_r^{r^k}
 \end{aligned}$$

# Example Results

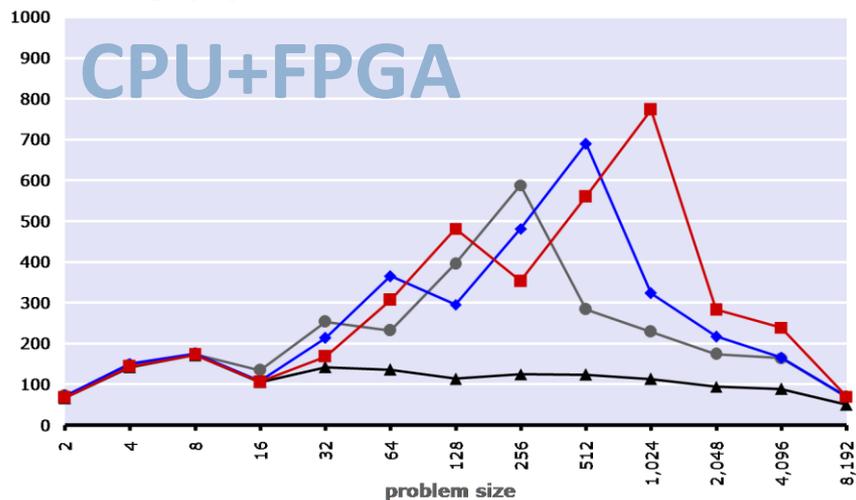
## DCT on 2.66 GHz Core2 (single-precision, 4-way SSE)

performance [Gflop/s]



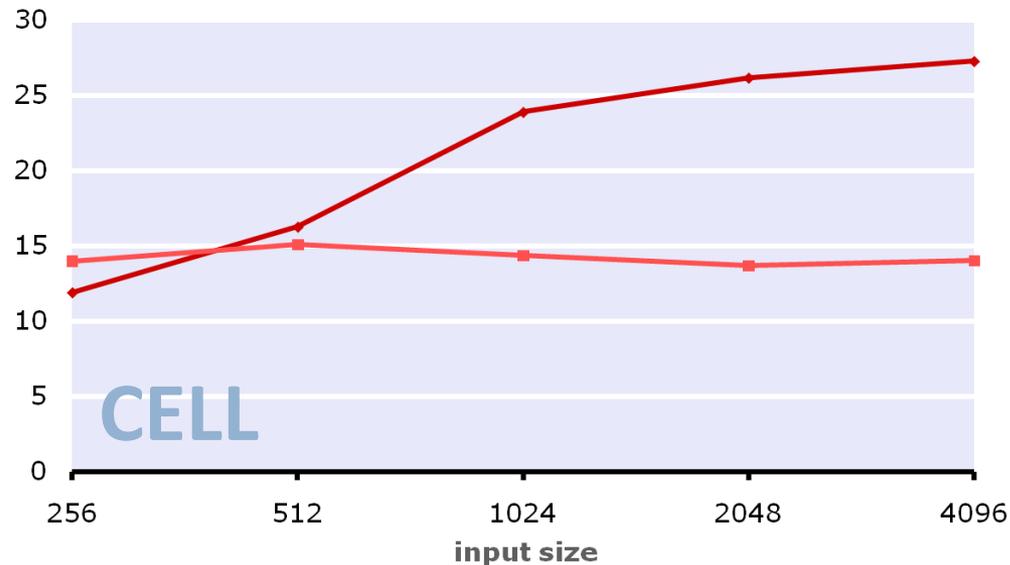
## DFT (16 bit fixed point): Hardware Accelerated Software on Xilinx Virtex-II Pro FPGA and Embedded PowerPC 405 Processor

Performance [Mop/s]



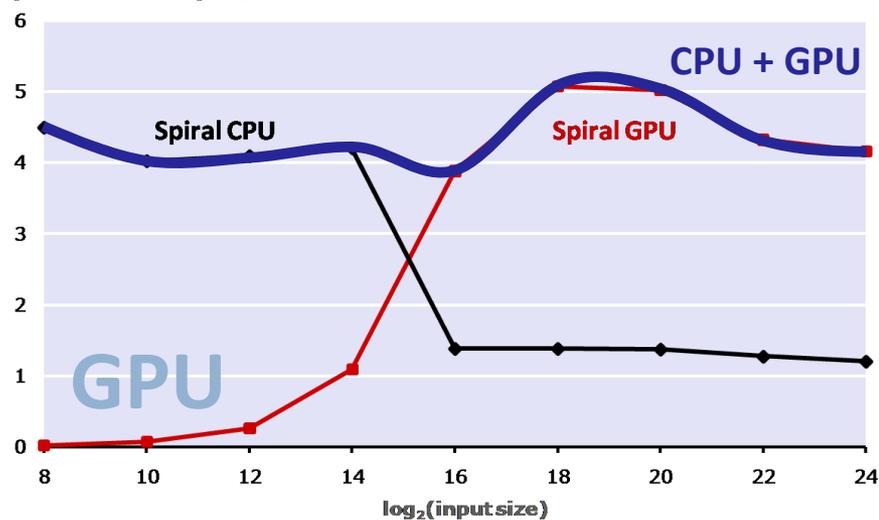
## Spiral-generated FFT on 3.2 GHz Cell BE (PlayStation 3)

performance [Gflop/s], single-precision, block-cyclic format, data resident on SPU



## WHT (single precision) on 3.6 GHz Pentium 4 with Nvidia 7900 GTX

performance [Gflops/s]



# Summary: Complete Automation for Transforms

- Platform: Off-the-shelf desktop
- Often: generated code faster than competition (if exists)

- **Memory hierarchy optimization**

Rewriting and search for algorithm selection  
 Rewriting for loop optimizations

- **Vectorization**

Rewriting

- **Parallelization**

Rewriting

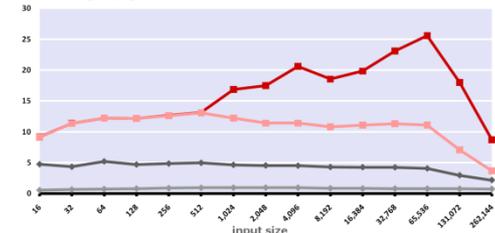
*fixed input size code*

- **Derivation of library structure**

Rewriting

Other methods *general input size library*

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



Numerical problem

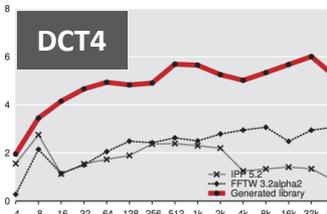
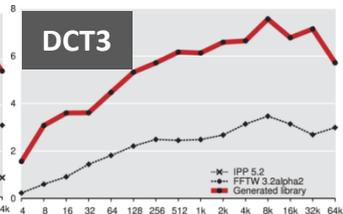
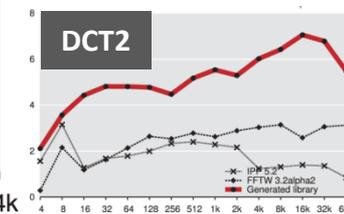
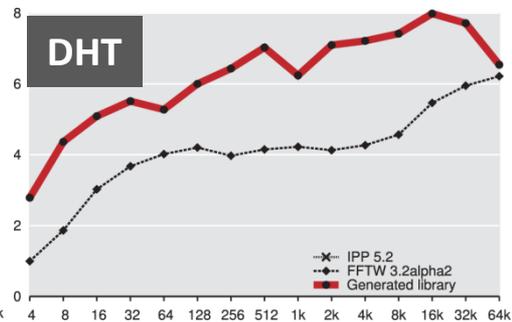
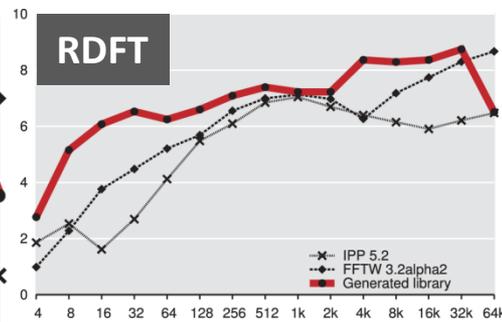
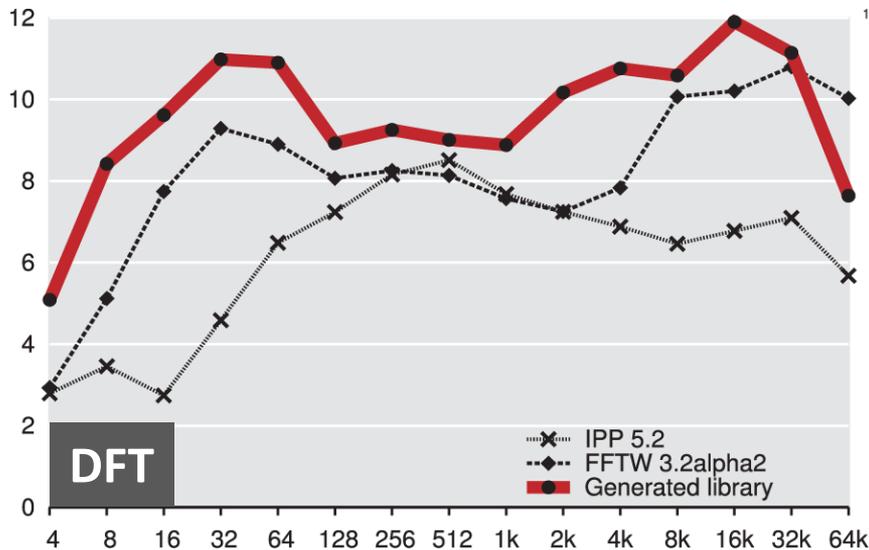


algorithm selection  
 implementation  
 compilation

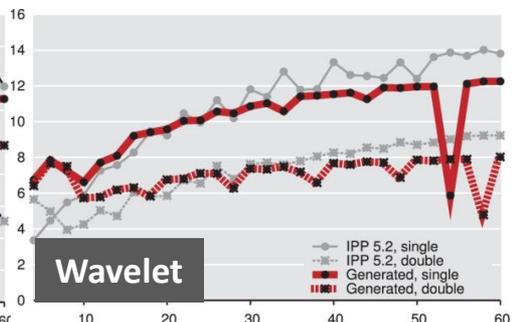
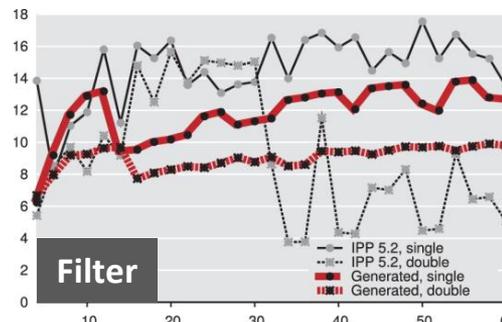
automated

Computing platform

# Generated 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



# Organization

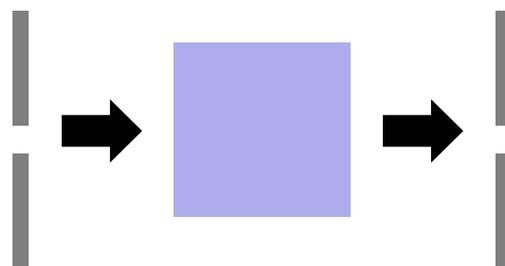
- **Spiral's framework: Example transforms**
  - Complete automation achieved
  
- **Beyond transforms**
  - Operator language
  - BLAS, Viterbi decoding, SAR imaging, Ebcot encoding
  
- **Conclusions and thoughts**

# 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



# Operator Language

| name                   | definition  |
|------------------------|---|
| <i>basic operators</i> |   |
| projection             | $\pi_{\mathbf{x}} : \mathbb{C}^m \times \mathbb{C}^n \rightarrow \mathbb{C}^m; (\mathbf{x}, \mathbf{y}) \mapsto \mathbf{x}$             |
| linear transform       | $M : \mathbb{C}^n \rightarrow \mathbb{C}^m; \mathbf{x} \mapsto M\mathbf{x}$   |
| stride                 | $L_m^{mn} : \mathbb{C}^{mn} \rightarrow \mathbb{C}^{mn}; \mathbf{x} \mapsto L_m^{mn} \mathbf{x}$  |
| vector sum             | $\Sigma_n : \mathbb{C}^n \rightarrow \mathbb{C}; \mathbf{x} \mapsto \sum_{i=0}^{n-1} x_i$   |
| vector minimum         | $\min_n : \mathbb{C}^n \rightarrow \mathbb{C}; \mathbf{x} \mapsto \min(x_0, \dots, x_{n-1})$  |
| constant vector        | $C_c : \emptyset \rightarrow \mathbb{C}^n; () \mapsto \mathbf{c}$   |
| <i>operations</i>      |   |
| addition               | $(M + N)(\mathbf{x}, \mathbf{y}) = M(\mathbf{x}, \mathbf{y}) + N(\mathbf{x}, \mathbf{y})$   |
| multiplication         | $(M \cdot N)(\mathbf{x}, \mathbf{y}) = M(\mathbf{x}, \mathbf{y}) \cdot N(\mathbf{x}, \mathbf{y})$                                       |
| direct sum             | $(M \oplus N)(\mathbf{x} \oplus \mathbf{u}, \mathbf{y} \oplus \mathbf{v}) = M(\mathbf{x}, \mathbf{y}) \oplus N(\mathbf{u}, \mathbf{v})$ |
| cartesian product      | $(M \times N)(\mathbf{x}, \mathbf{y}, \mathbf{u}, \mathbf{v}) = M(\mathbf{x}, \mathbf{y}) \times N(\mathbf{u}, \mathbf{v})$             |
| composition            | $(M \circ N)(\mathbf{x}, \mathbf{y}) = M(N(\mathbf{x}, \mathbf{y}))$  |
| iterative composition  | $(\prod_{i=0}^{n-1} M_i)(\mathbf{x}, \mathbf{y}) = (M_0 \circ \dots \circ M_{n-1})(\mathbf{x}, \mathbf{y})$                             |
| tensor product         | $I \otimes M, M \otimes I$  |

# Example: Matrix Multiplication (MMM)

Breakdown rules = algorithm knowledge:

*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

$$\begin{aligned}
 & \underbrace{\text{MMM}_{m,n,k}}_{\text{smp}(p,\mu)} \\
 \rightarrow & \underbrace{\left( \text{I}_m \otimes \text{L}_p^n \right) \circ \left( \text{MMM}_{m,n/p,k} \otimes (\otimes)_{1 \times p \rightarrow p} \right) \circ \left( \text{I}_{km} \times (\text{I}_k \otimes \text{L}_{n/p}^n) \right)}_{\text{smp}(p,\mu)} \\
 \rightarrow & \underbrace{\left( \text{I}_m \otimes \text{L}_p^n \right)}_{\text{smp}(p,\mu)} \circ \underbrace{\left( \text{MMM}_{m,n/p,k} \otimes (\otimes)_{1 \times p \rightarrow p} \right)}_{\text{smp}(p,\mu)} \circ \underbrace{\left( \text{I}_{km} \times (\text{I}_k \otimes \text{L}_{n/p}^n) \right)}_{\text{smp}(p,\mu)} \\
 \rightarrow & \underbrace{\left( \text{I}_m \otimes \text{L}_p^n \right)}_{\text{smp}(p,\mu)} \circ \underbrace{\text{L}_{m/pn}^{mn}}_{\text{smp}(p,\mu)} \circ \underbrace{\left( (\otimes)_{1 \times p \rightarrow p} \otimes_{\parallel} \text{MMM}_{m/p,n,k} \right)}_{\text{smp}(p,\mu)} \circ \underbrace{\left( \text{I}_{km} \times \text{L}_p^{kn} \right)}_{\text{smp}(p,\mu)} \circ \underbrace{\left( \text{I}_{km} \times (\text{I}_k \otimes \text{L}_{n/p}^n) \right)}_{\text{smp}(p,\mu)} \\
 \rightarrow & \underbrace{\left( (\text{L}_m^{mp} \otimes \text{I}_{n/(p\mu)}) \bar{\otimes} \text{I}_\mu \right)}_{\text{smp}(p,\mu)} \circ \underbrace{\left( (\otimes)_{1 \times p \rightarrow p} \otimes_{\parallel} \text{MMM}_{m,n/p,k} \right)}_{\text{smp}(p,\mu)} \circ \underbrace{\left( (\text{I}_{km/\mu} \bar{\otimes} \text{I}_\mu) \times \left( (\text{L}_p^{kp} \otimes \text{I}_{n/(p\mu)}) \bar{\otimes} \text{I}_\mu \right) \right)}_{\text{smp}(p,\mu)}
 \end{aligned}$$

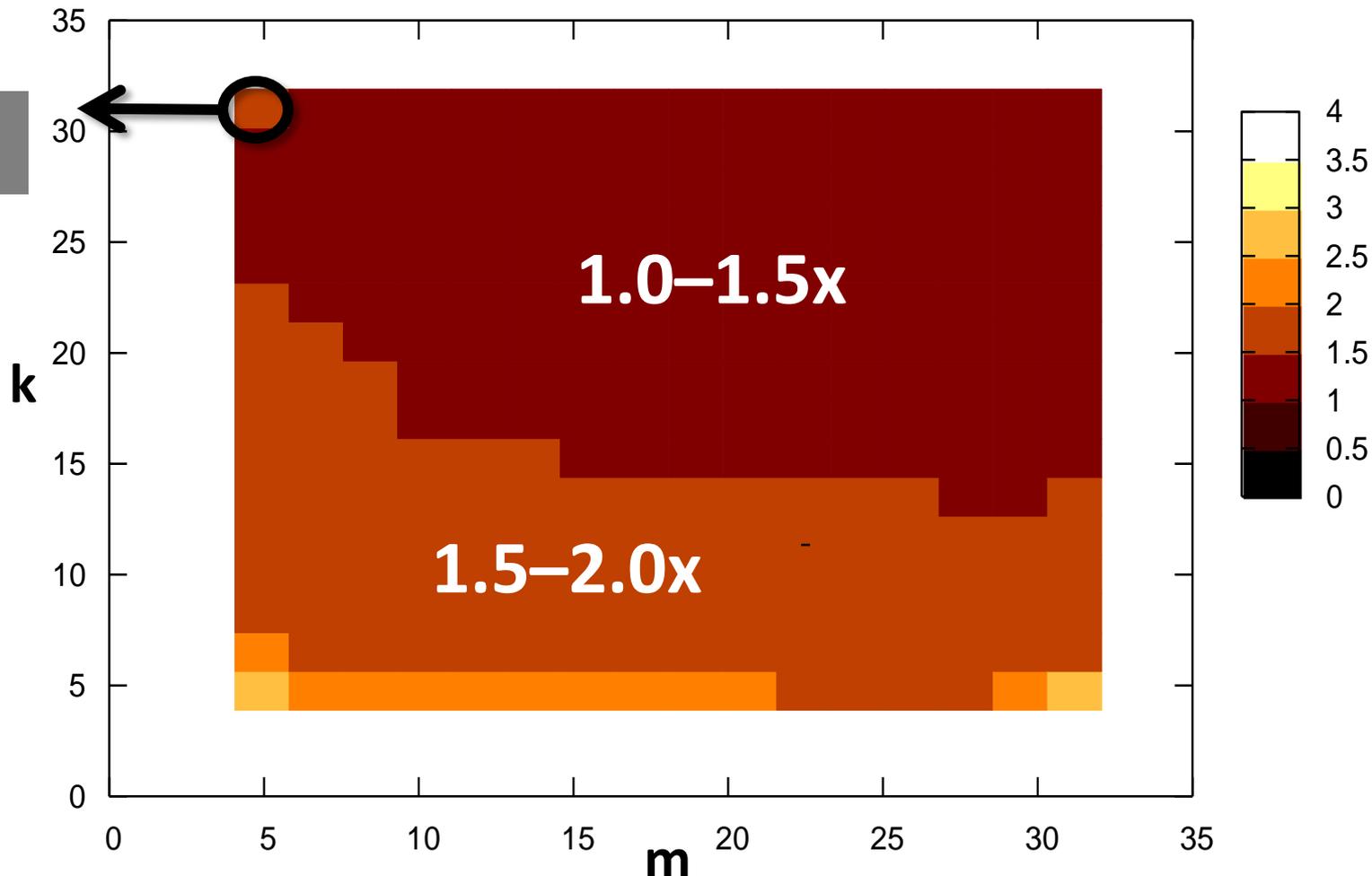
**Load-balanced**

**No false sharing**

# Speed-up ( $m \times k$ ) *times* ( $k \times n$ )

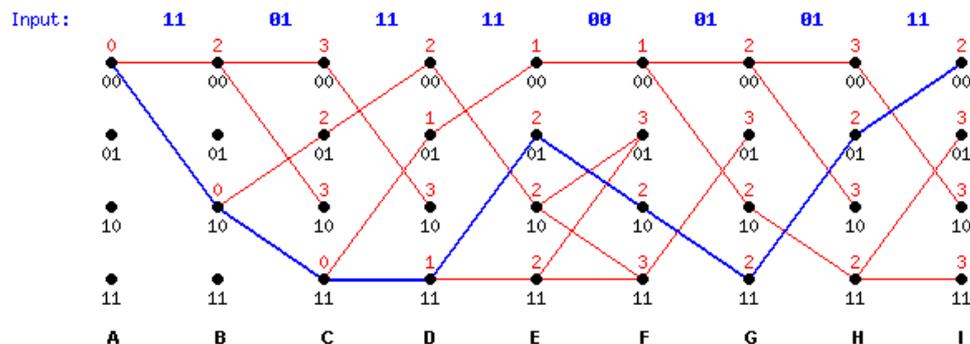
## MMM Speedup over MKL, $n = 32$ , Core 2 Duo

example:



■ Comparison to GotoBLAS similar

# Viterbi Decoding in OL



<http://www.ece.unb.ca/tervo/ee4253/convolution3.htm>

- Operator for Viterbi decoder  $\text{Vit}_{r,K,N,p} : \mathbb{N}^{rN} \rightarrow \mathbb{N}^{2^{K-1}} \times \mathbb{N}^{N2^{K-1}}$
- Breakdown rules

$\text{Vit}_{r,K,N,p} \rightarrow$

$$\pi_{2,3} \circ \left( \prod_{0 \leq x < N} (L_{2^{K-2}}^{2^{K-1}} \times I_{r2^{K-1} \times N2^{K-1}}) \circ (I_{2^{K-2} \times 2^{K-2} \times 1} \otimes C_{r,K,p}^x) \right) \circ Id_{(1)_2 \otimes i_{2^{K-2}}}$$

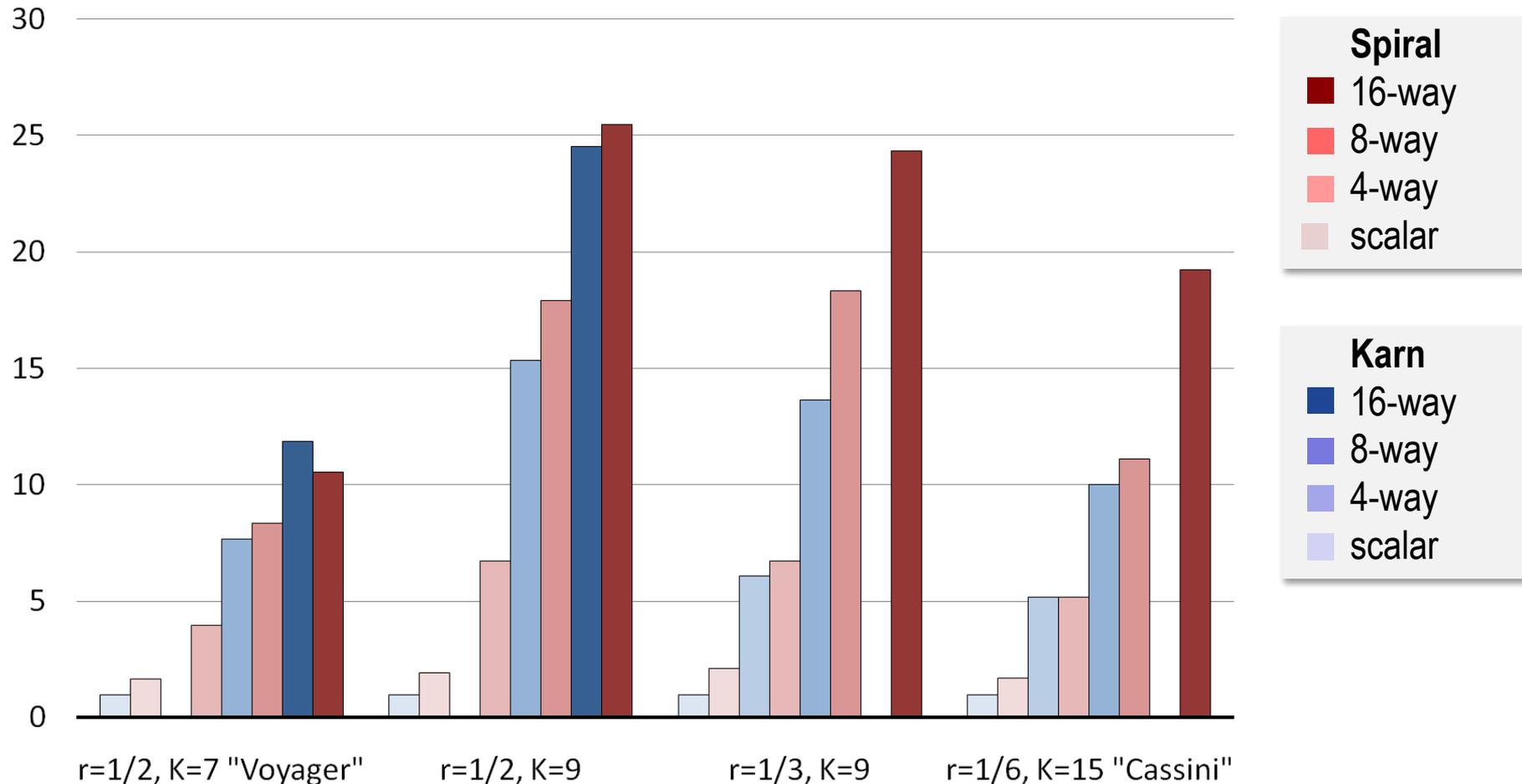
$$C_{r,K,p}^x \circ Id_y \rightarrow B_{r,K,p}^{x,y}$$

# Results

## Karn's implementation: hand-written assembly for 4 Viterbi codes

### Performance Gain of Various Generated Viterbi Decoders

Speedup over Karn's C implementation



# EBCOT Coding in OL

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

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

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

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

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

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

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

$$\text{sppCode}(\sigma_{m,n}) : \mathbb{Z}_2^9 \rightarrow \mathbb{Z}_2$$

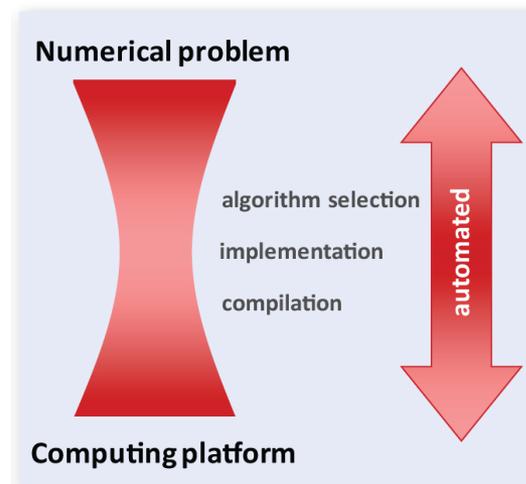
$$\text{and}_2 \circ (\text{eqz} \times \text{nez}) \circ (G_{(4)_9} \times (G_0 + G_1 + G_2 + G_3 + G_5 + G_6 + G_7 + G_8)) \circ \underline{\begin{pmatrix} 1 \\ 1 \end{pmatrix}}$$

# Organization

- **Spiral's framework: Example transforms**
  - Complete automation achieved
- **Beyond transforms**
- **Conclusions and thoughts**

# Raising the Abstraction Level

- **Formally describe and structure algorithms/applications**  
eternally valid
- **In Spiral**
  - Domain-specific, declarative, mathematical language OL
  - Difficult optimizations/transformations by rewriting
  - What it enables
    - Vectorization, parallelization using domain knowledge
    - Efficient retargeting to new platforms and new platform paradigms
    - Complete automation in some cases
- **Other examples**
  - Libraries
  - Identification and definition of BLAS
- **Parameter tuning**
  - Indispensable tool but cannot achieve the above



# Interdisciplinary Research Needed

Programming languages  
Program generation

Symbolic Computation  
Rewriting

*Automating  
High-Performance  
Parallel Library  
Development*

Software  
Scientific Computing

Algorithms  
Mathematics

Compilers

*We Need to Work Together*