

Parallelism in Spiral

Franz Franchetti

Electrical and
Computer Engineering
Carnegie Mellon University

Joint work with
Yevgen Voronenko
Markus Püschel

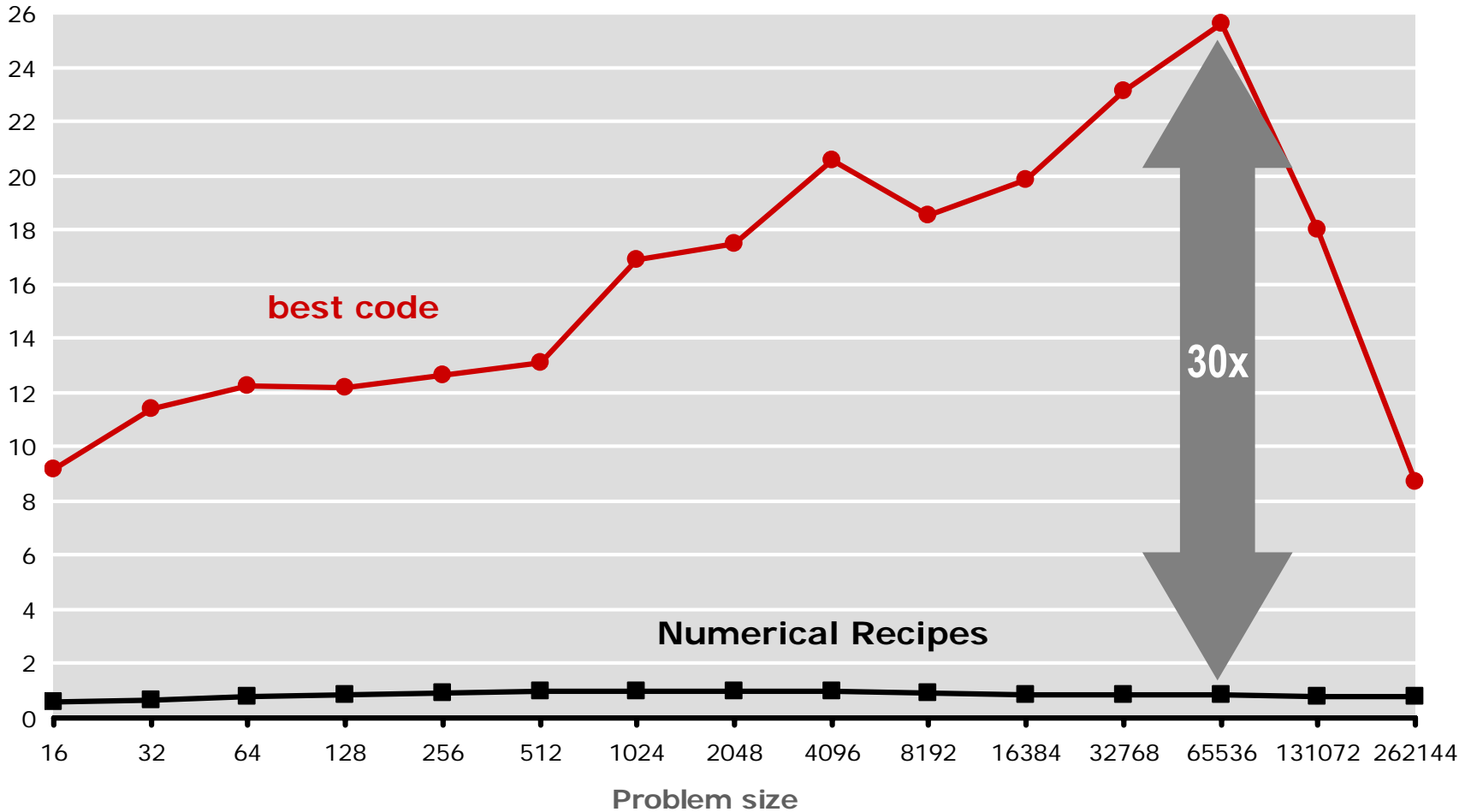
... and the Spiral team (only part shown)



This work was supported by
DARPA DESA program, NSF-NGS/ITR, NSF-ACR, and Intel

The Problem

Discrete Fourier Transform (single precision): 2 x Core2 Extreme 3 GHz
Performance [Gflop/s]

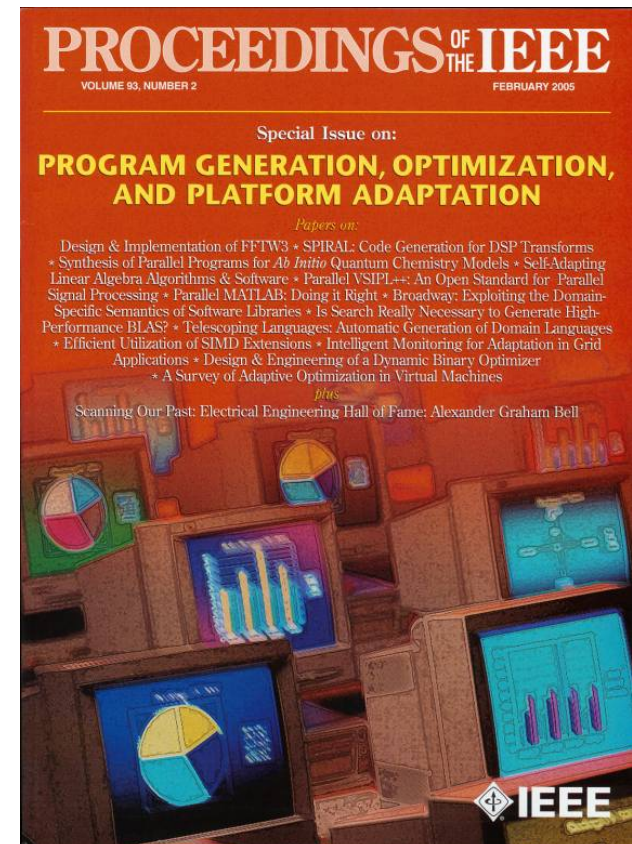


What's going on?

Automatic Performance Tuning

- **Current vicious circle:** Whenever a new platform comes out, the same functionality needs to be rewritten and reoptimized
- **Automatic Performance Tuning**
 - BLAS: ATLAS
 - Linear algebra: Bebop, Spike, Flame
 - Sorting
 - Fourier transform: FFTW
 - **Linear transforms: Spiral**
 - ...others
 - New compiler techniques

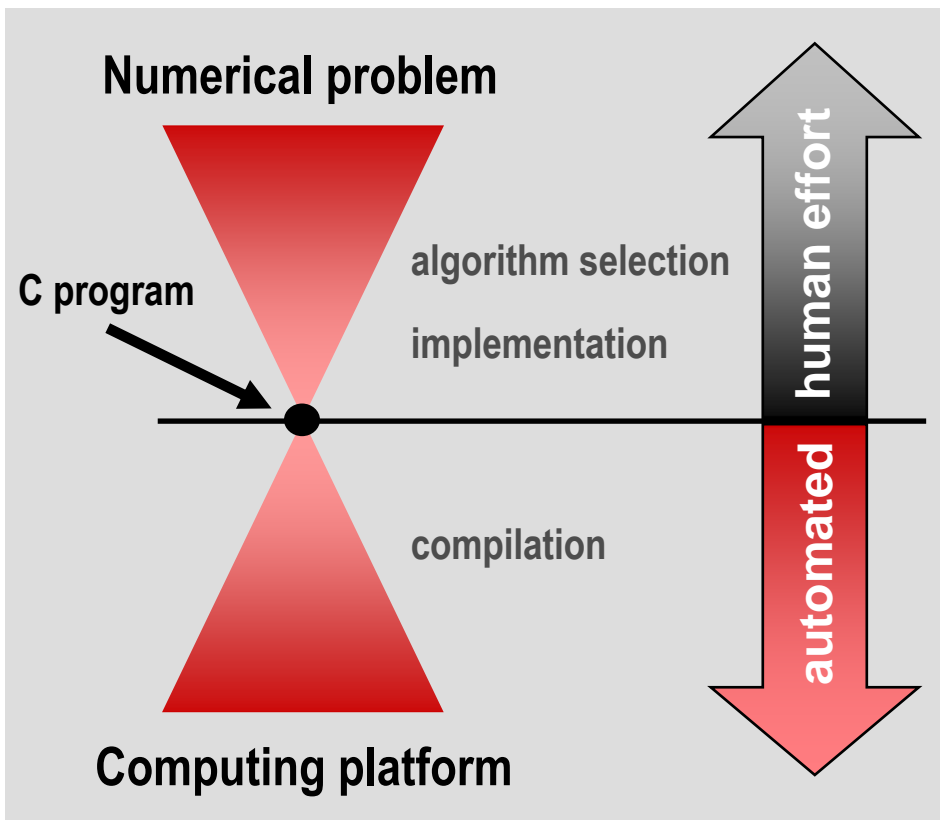
But what about parallelism ... ?



Proceedings of the IEEE special issue, Feb. 2005

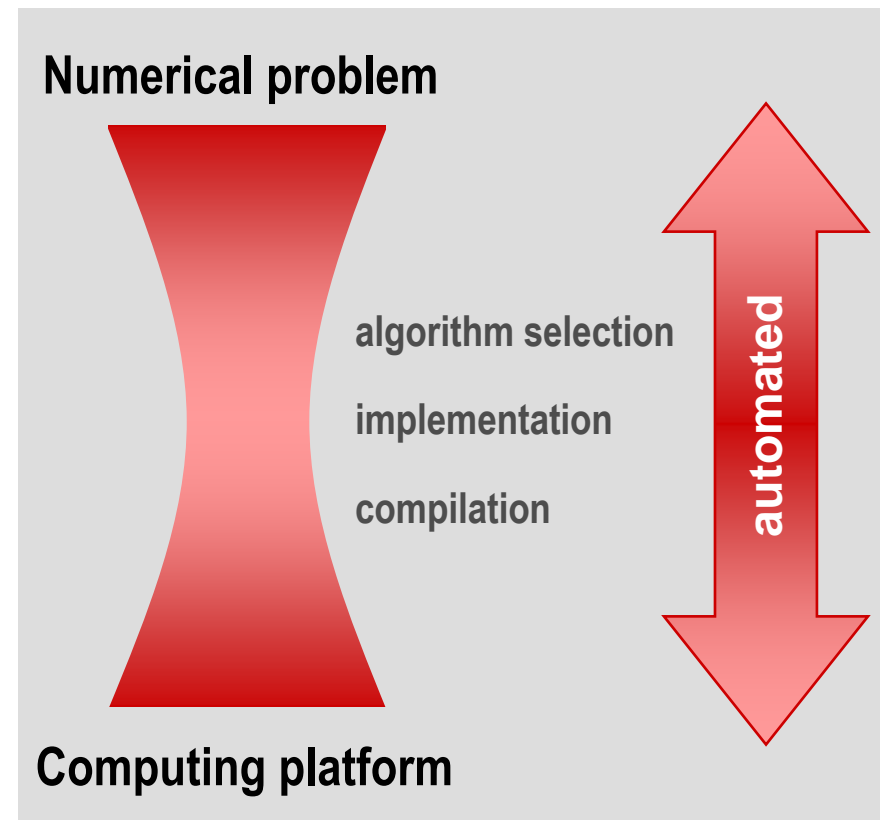
Vision Behind Spiral

Current



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

Future



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

Organization

- **Spiral overview**
- Parallelization in Spiral
- Results
- Concluding remarks

Spiral

- Library generator for linear transforms (DFT, DCT, DWT, filters,) *and recently more ...*
- Wide range of platforms supported: scalar, fixed point, **vector, parallel, Verilog, GPU**
- **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 library

How Spiral Works

Spiral:

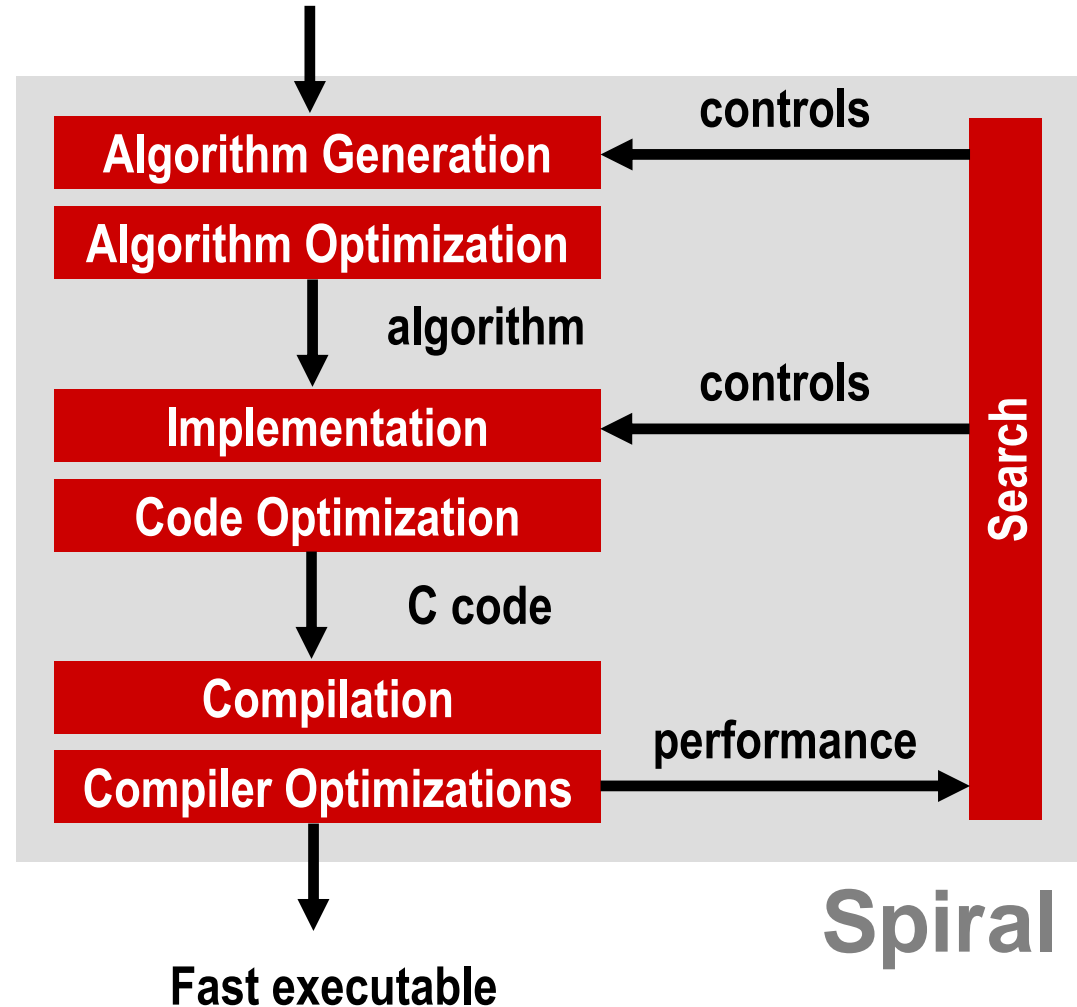
Complete automation of the implementation and optimization task

Basic idea:

Declarative representation of algorithms

Rewriting systems to generate and optimize algorithms

Problem specification (transform)



What is a (Linear) Transform?

- Mathematically: Matrix-vector multiplication

$$\begin{array}{ccc} & x \mapsto y = T \cdot x & \\ \text{input vector (signal)} & \uparrow & \uparrow \\ \text{output vector (signal)} & \uparrow & \text{transform = matrix} \end{array}$$

- 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)

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

Kronecker product
Identity
Permutation

- Algorithms reduce arithmetic cost $O(n^2) \rightarrow O(n \log(n))$
- Product of structured sparse matrices
- Mathematical notation exhibits structure: **SPL (signal processing language)**

Examples: Transforms

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

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

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

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

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

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

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

Spiral currently contains 55 transforms

Examples: Breakdown Rules (currently ≈220)

$$\text{DFT}_n \rightarrow (\text{DFT}_k \otimes I_m) \Upsilon_m^n (I_k \otimes \text{DFT}_m) L_k^n, \quad n = km$$

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

DFT

DCT-3

- “Teaches” Spiral about existing algorithm knowledge (~200 journal papers)

DCT-4

- Includes many new ones (algebraic theory, Pueschel, Moura, Voronenko)

IMDCT₂

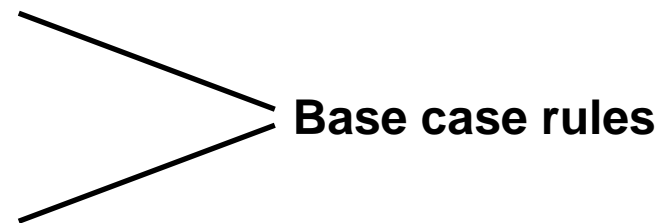
DCT-4_{2m}

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

$$\text{DFT}_2 \rightarrow F_2$$

$$\text{DCT-2}_2 \rightarrow \text{diag}(1, 1/\sqrt{2}) F_2$$

$$\text{DCT-4}_2 \rightarrow J_2 R_{13\pi/8}$$



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<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<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<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<m*n;i++) y[i] = Dmn[i]*x[i];</pre>
$y = L_m^{mn}x$	<pre>for (i=0;i<m;i++) for (j=0;j<n;j++) y[i+m*j]=x[n*i+j];</pre>

Example: tensor product

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

Program Generation in Spiral (Sketched)

Transform
 user specified

DFT₈



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
 scheduling

Organization

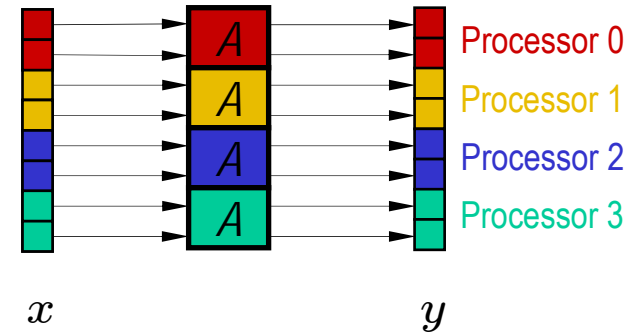
- Spiral overview
- **Parallelization in Spiral**
- Results
- Concluding remarks

SPL to Shared Memory Code: Basic Idea

[SC 06]

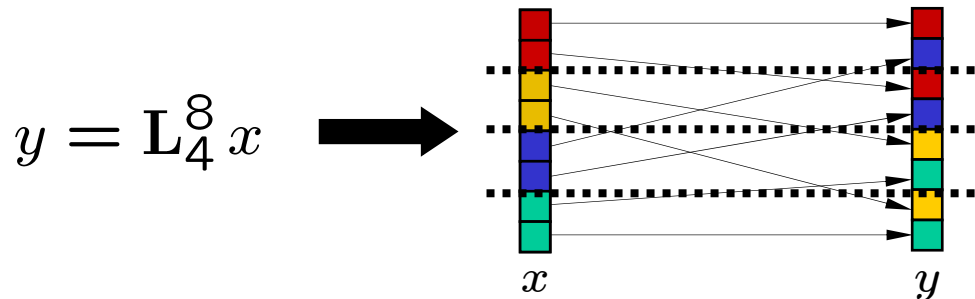
- Governing construct: tensor product

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



Independent operation, 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: μ
- **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 μ

Step 2: Identify “Good” Formulas

- Load balanced, avoiding false sharing

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

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

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

- 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** 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

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

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

Same Approach for Other Parallel Paradigms

Message Passing: [ISPA 06]

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

With Bonelli, Lorenz, Ueberhuber, TU Vienna

Cg/OpenGL for GPUs:

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

With Shen, TU Denmark

Vectorization: [IPDPS 02, VecPar 06]

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

Verilog for FPGAs: [DAC 05]

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

With Milder, Hoe, CMU

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)

Cooley-Tukey FFT in OL: $DFT \rightarrow (DFT \otimes I) \circ D \circ (I \otimes DFT) \circ L.$

Viterbi in OL: $Vit \rightarrow \pi \circ \left(\prod (I \otimes V) \circ (L \times I) \right) \circ (C \times C \times I)$

Mat-Mat-Mult: $MMM \rightarrow I \otimes MMM$

$MMM \rightarrow (I \otimes L) \circ (MMM \otimes I) \circ (I \times (I \otimes L))$

OL Rewriting Rules

- SPL rules reused
- Only few OL-specific rules required

$$\underbrace{\left(\mathbf{I}_k \otimes \mathbf{L}_n^{mn} \right)}_{\text{smp}(p,\mu)} \circ \underbrace{\mathbf{L}_{km}^{kmn}}_{\text{smp}(p,\mu)} \rightarrow \left(\mathbf{L}_k^{kn} \otimes \mathbf{I}_{m/\mu} \right) \bar{\otimes} \mathbf{I}_\mu$$

$$\underbrace{\mathbf{L}_n^{kmn}}_{\text{smp}(p,\mu)} \circ \underbrace{\left(\mathbf{I}_k \otimes \mathbf{L}_m^{mn} \right)}_{\text{smp}(p,\mu)} \rightarrow \left(\mathbf{L}_n^{kn} \otimes \mathbf{I}_{m/\mu} \right) \bar{\otimes} \mathbf{I}_\mu$$

$$\underbrace{\mathbf{A} \circ \mathbf{B}}_{\text{smp}(p,\mu)} \rightarrow \underbrace{\mathbf{A}}_{\text{smp}(p,\mu)} \circ \underbrace{\mathbf{B}}_{\text{smp}(p,\mu)}$$

$$\underbrace{\mathbf{A}^{k \times m \rightarrow n} \otimes \mathbf{I}^{1 \times p \rightarrow p}}_{\text{smp}(p,\mu)} \rightarrow \underbrace{\mathbf{L}_n^{pn}}_{\text{smp}(p,\mu)} \circ \left(\mathbf{I}_{1 \times p \rightarrow p} \otimes_{\parallel} \mathbf{A}^{k \times m \rightarrow n} \right) \circ \underbrace{\left(\mathbf{I}_k \times \mathbf{L}_p^{pm} \right)}_{\text{smp}(p,\mu)}$$

$$\underbrace{\left(\mathbf{A} \times \mathbf{B} \right)}_{\text{smp}(p,\mu)} \circ \underbrace{\left(\mathbf{C} \times \mathbf{D} \right)}_{\text{smp}(p,\mu)} \rightarrow \underbrace{\left(\mathbf{A} \circ \mathbf{C} \right)}_{\text{smp}(p,\mu)} \times \underbrace{\left(\mathbf{B} \circ \mathbf{D} \right)}_{\text{smp}(p,\mu)}$$

New OL rules

Example: Viterbi Decoder in OL

- Viterbi decoder (forward part) as operator

$$\text{Vit}_{m,n,N}^{e,f,x} : \mathbb{R}^{nN} \rightarrow \mathbb{R}^{2^m} \times \mathbb{N}^{2^m n}$$

- Viterbi kernel (butterfly)

$$\mathbf{V}_{i,j}^{e,f} : \mathbb{R}^2 \times \mathbb{R}^{2n} \times \mathbb{R}^{nN} \rightarrow \mathbb{R}^2 \times \mathbb{R}^{2n} \times \mathbb{R}^{nN}; (\mathbf{x}, \mathbf{d}, \mathbf{c}) \mapsto (\mathbf{y}, \mathbf{d}', \mathbf{c}) \quad , \quad 0 \leq i < n, 0 \leq j < 2^{m-1}$$

- Viterbi algorithm as breakdown rule

$$\text{Vit}_{m,n,N}^{e,f,x} \rightarrow \pi_{(\mathbf{x}, \mathbf{d})} \circ \left(\prod_{i=0}^{n-1} \left(\mathbf{I}_{2^{m-1} \times 2^{m-1} \times 1} \otimes_j \mathbf{V}_{i,j}^{e,f} \right) \circ \left(\mathbf{L}_{2^{m-1}}^{2^m} \times \mathbf{I}_{2^m n \times nN} \right) \right) \circ \left(\mathbf{C}_x \times \mathbf{C}_0 \times \mathbf{I}_{nN} \right)$$

First non-transform supported by Spiral

Viterbi: Vectorization Through Rewriting

$$\begin{aligned}
 \underbrace{\text{Vit}_{m,n,N}^{e,f,x}}_{\text{vec}(\nu)} &\rightarrow \underbrace{\pi_{(x,d)} \circ \left(\prod_{i=0}^{n-1} \left(\text{I}_{2^{m-1} \times 2^{m-1} \times 1} \otimes_j \text{V}_{i,j}^{e,f} \right) \circ \left(\text{L}_{2^{m-1}}^{2^m} \times \text{I}_{2^m n \times n N} \right) \right)}_{\text{vec}(\nu)} \circ \left(\text{C}_x \times \text{C}_0 \times \text{I}_{nN} \right) \\
 &\rightarrow \pi_{(x,d)} \circ \left(\prod_{i=0}^{n-1} \underbrace{\left(\text{I}_{2^{m-1} \times 2^{m-1} \times 1} \otimes_j \text{V}_{i,j}^{e,f} \right) \circ \left(\text{L}_{2^{m-1}}^{2^m} \times \text{I}_{2^m n \times n N} \right)}_{\text{vec}(\nu)} \right) \circ \left(\text{C}_x \times \text{C}_0 \times \text{I}_{nN} \right) \\
 &\dots \\
 &\rightarrow \pi_{(x,d)} \circ \left(\prod_{i=0}^{n-1} \left(\text{I}_{2^{m-1}/\nu \times 2^{m-1} \times 1} \otimes_j \underbrace{\left(\left(\underbrace{\text{L}_{\nu}^{2\nu}}_{\text{reg}(\nu)} \right) \times \text{I}_{2^m \times n N} \right)}_{\text{vec}(\nu)} \circ \underbrace{\left(\text{V}_{i,4j+k}^{e,f} \otimes_k \text{I}_{\nu \times 1 \times 1} \right)}_{\text{vec}(\nu)} \right) \right) \circ \left(\underbrace{\left(\text{L}_{2^{m-1}/\nu}^{2^m} \otimes \text{I}_{\nu} \right)}_{\text{vec}(\nu)} \times \text{I}_{2^m n \times n N} \right) \circ \left(\text{C}_x \times \text{C}_0 \times \text{I}_{nN} \right)
 \end{aligned}$$

Sufficient to vectorize one input

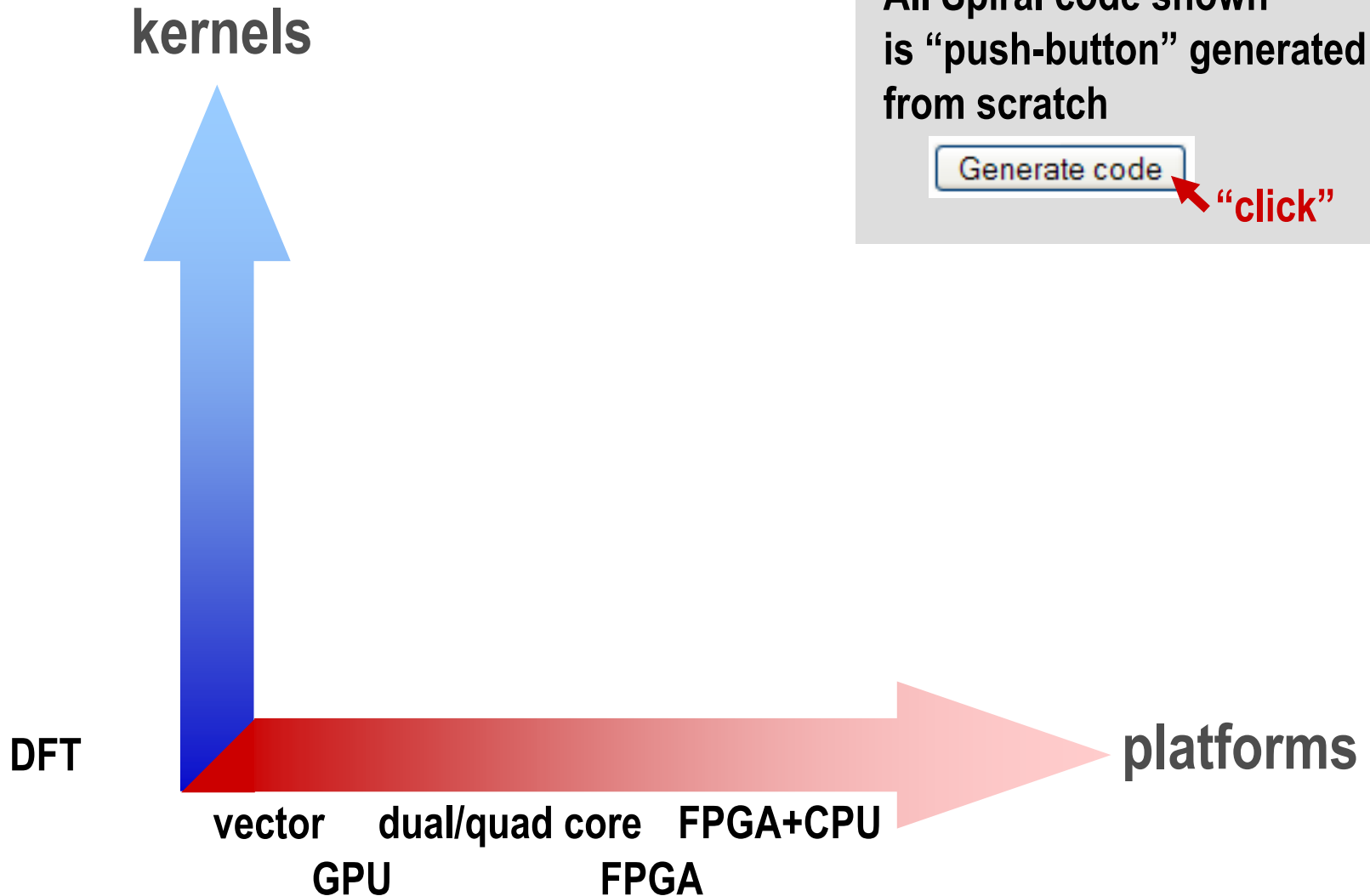
Vectorized kernel

In-register shuffle operation

Organization

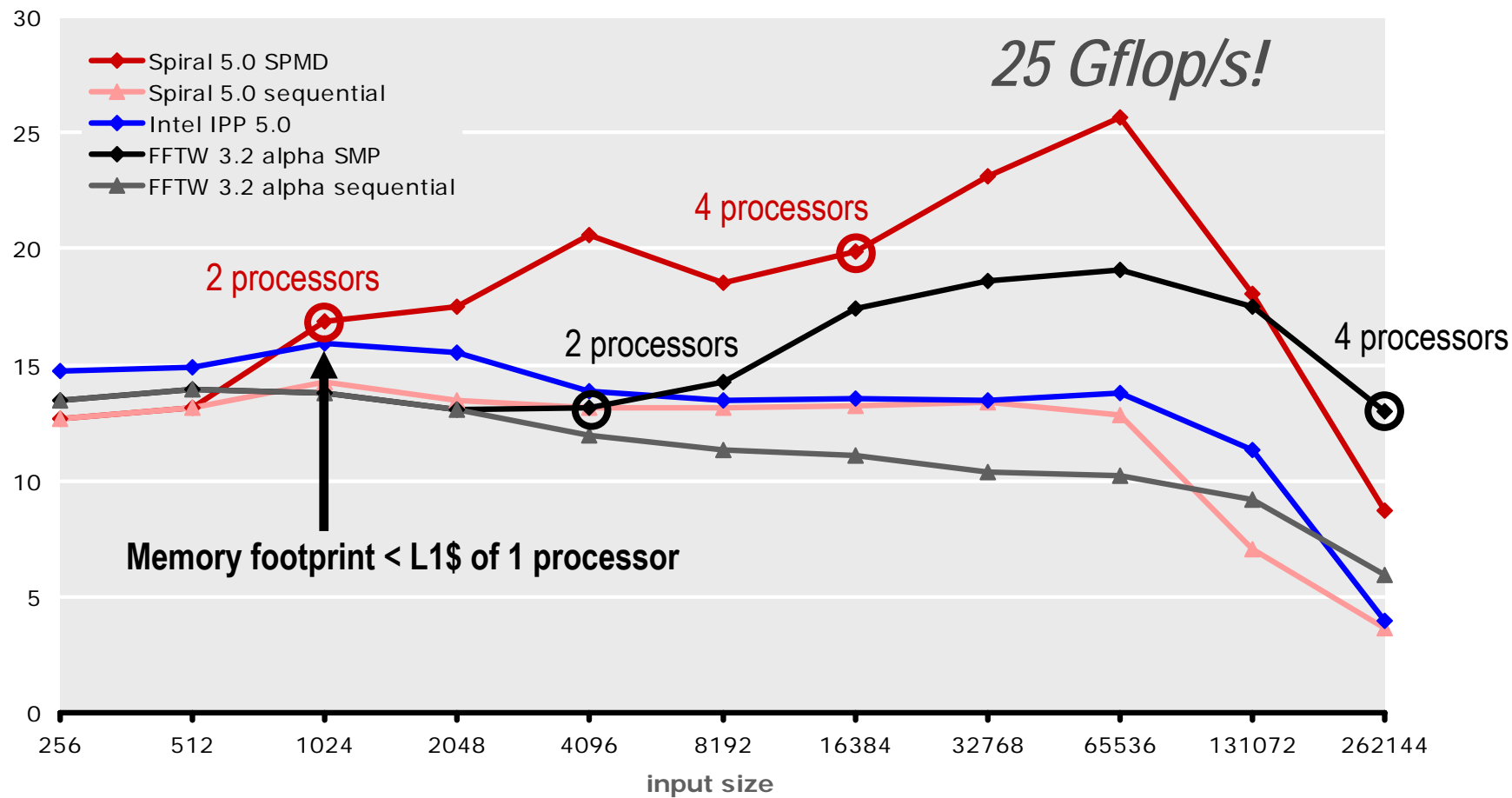
- Spiral overview
- Parallelization in Spiral
- **Results**
- Concluding remarks

Benchmarks



Benchmark: Vector and SMP

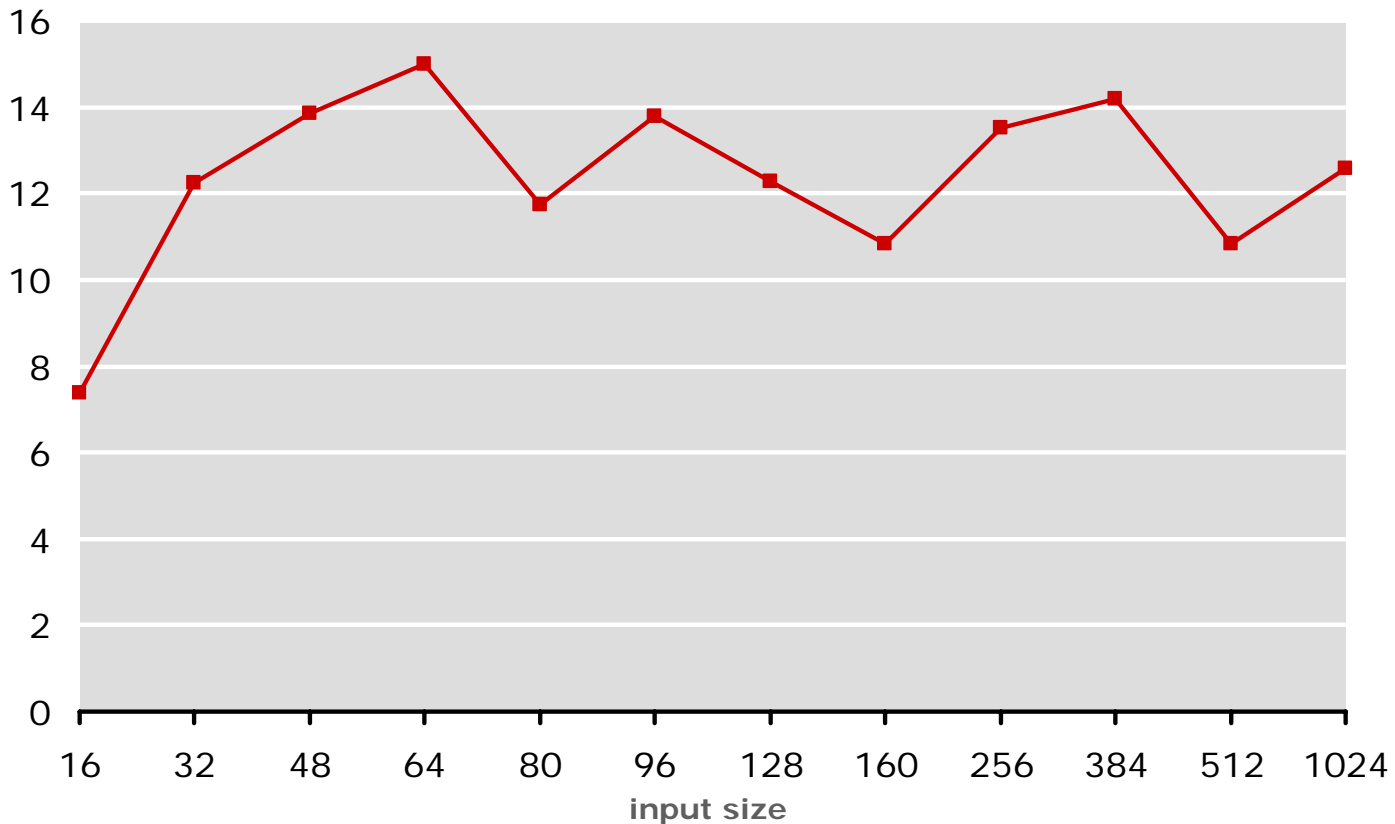
DFT (single precision): on 3 GHz 2 x Core 2 Extreme
 performance [Gflop/s]



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

Benchmark: Cell (1 processor = SPE)

DFT (single precision) on 3.2 GHz Cell BE (Single SPE)
performance [Gflop/s]

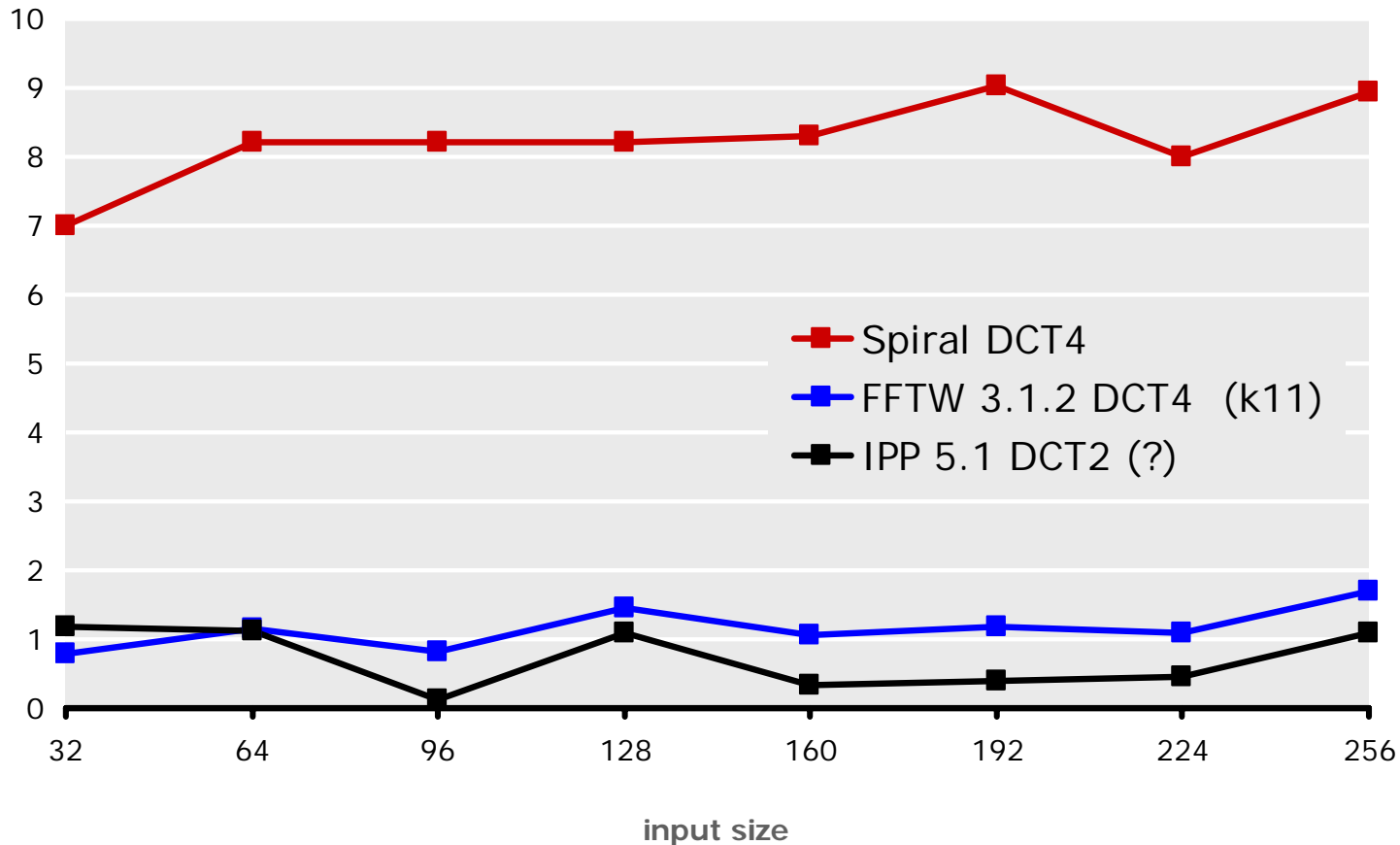


Generated using the simulator; run at Mercury (thanks to Robert Cooper)

Joint work with Th. Peter (ETH Zurich), S. Chellappa, M. Telgarsky, J. Moura (CMU)

DCT4, Multiples of 32: 4-way Vectorized

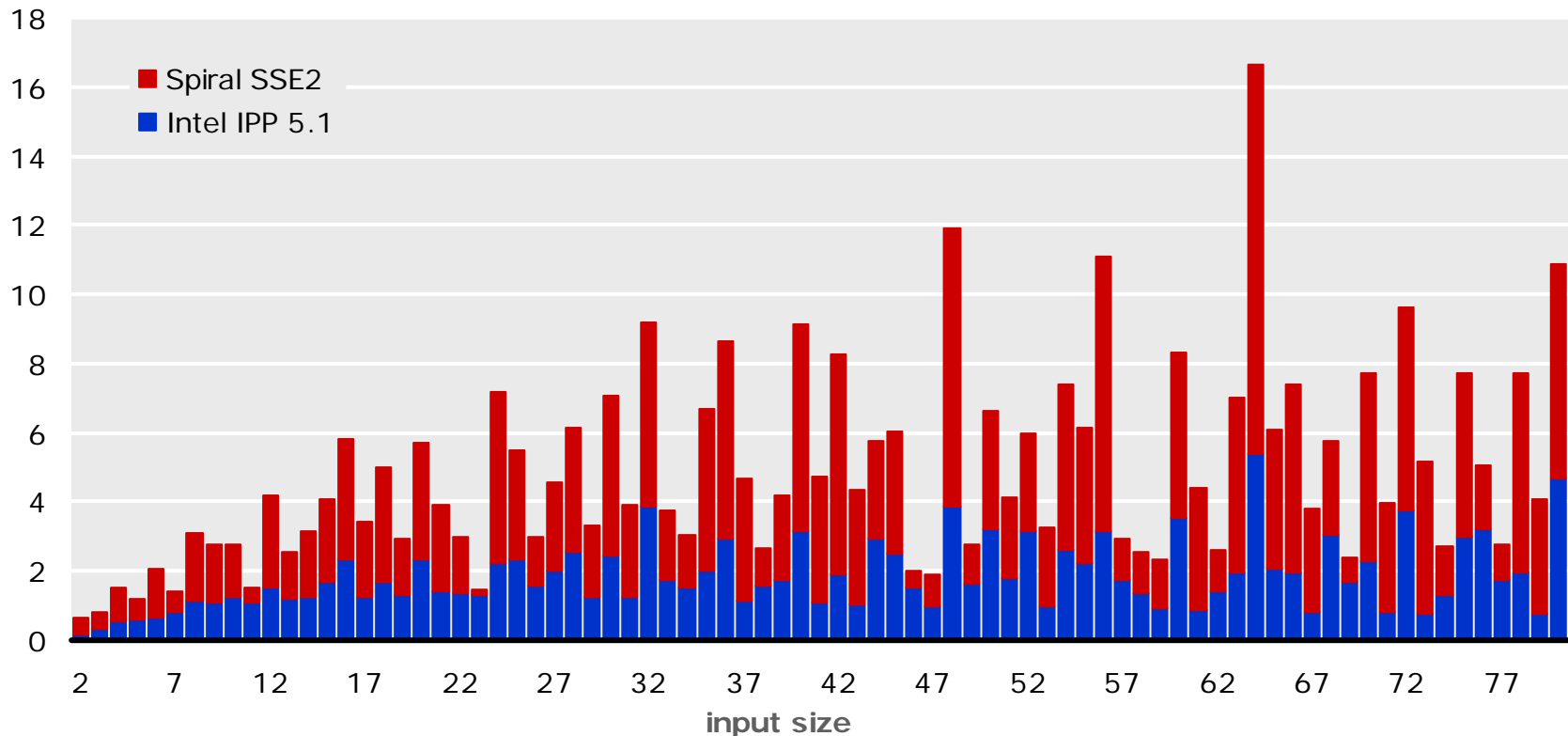
DCT (single precision) 2.66 GHz Core2 (4-way 32-bit SSE)
 performance [Gflop/s]



■ *novel algorithm (algebraic algorithm theory)*

DFT, 8-way Vectorized: All Sizes Up To 80

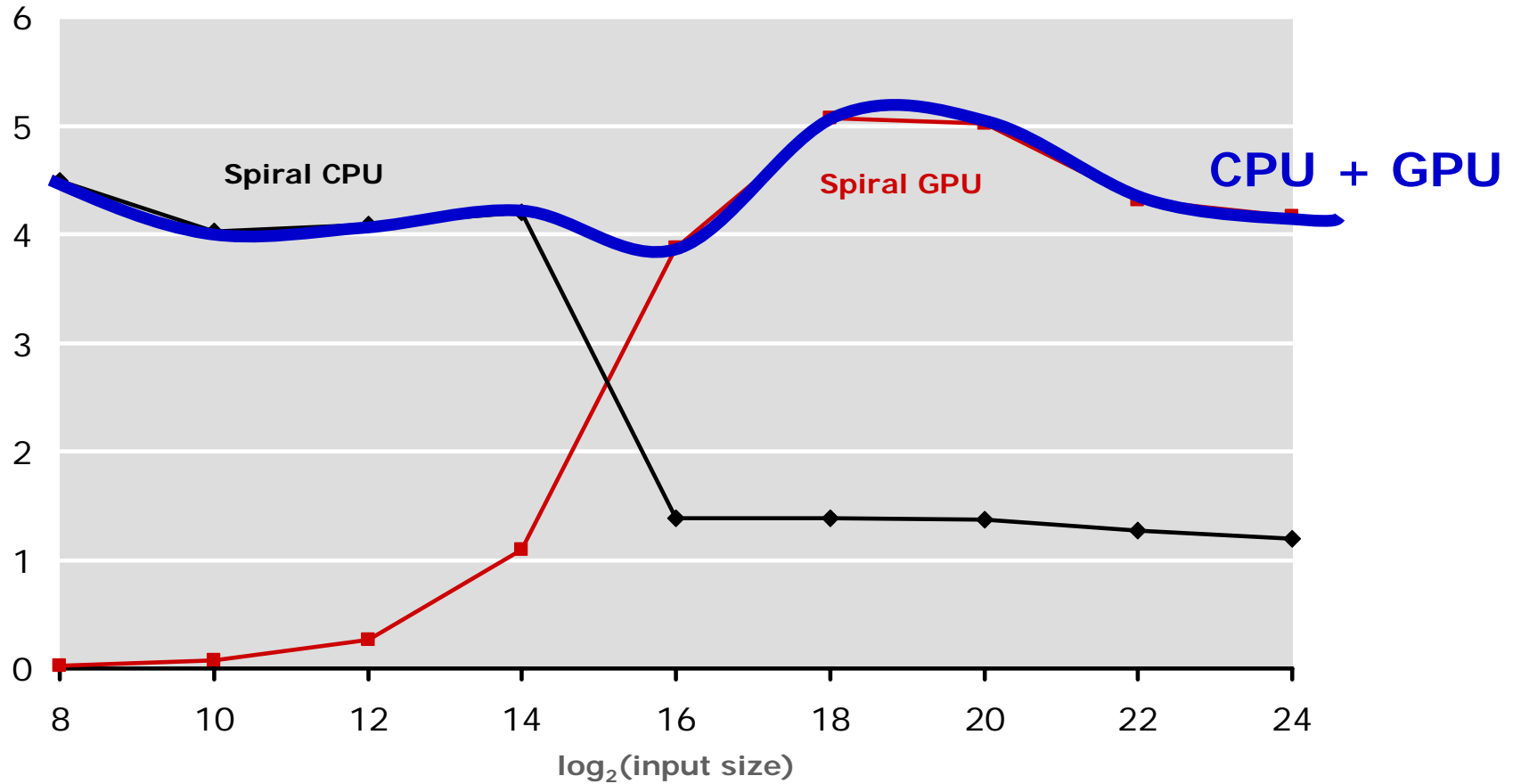
DFT (16-bit integer) on 2.66 GHz Core2 Duo (8-way SSE2)
 performance [Gflop/s]



- *first 8-way DFTs for all sizes*
- *arbitrary vector length /arbitrary DFT sizes in principle solved*

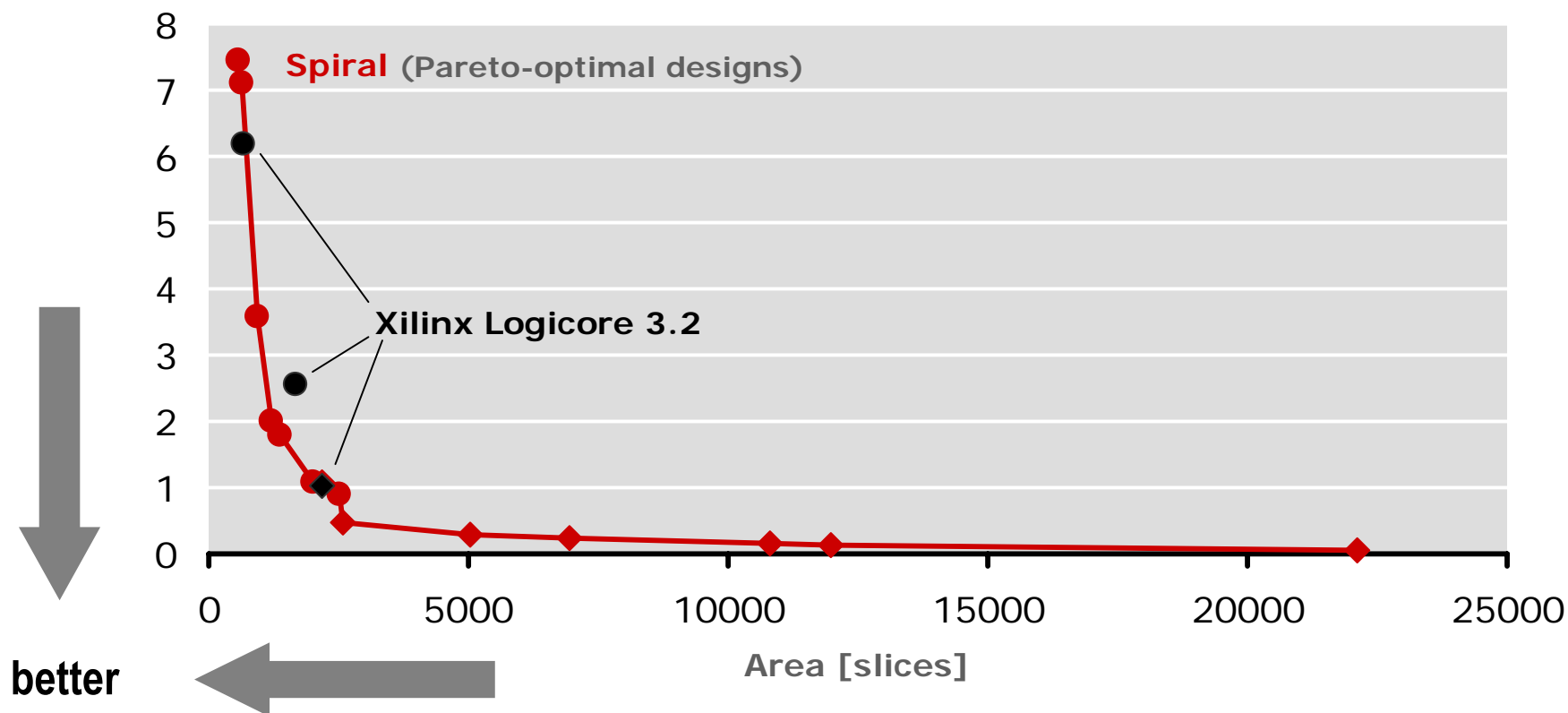
Benchmark: GPU

WHT (single precision) on 3.6 GHz Pentium 4 with Nvidia 7900 GTX
performance [Gflops/s]



Benchmark: FPGA

DFT 256 on Xilinx Virtex 2 Pro FPGA
inverse throughput (gap) [us]

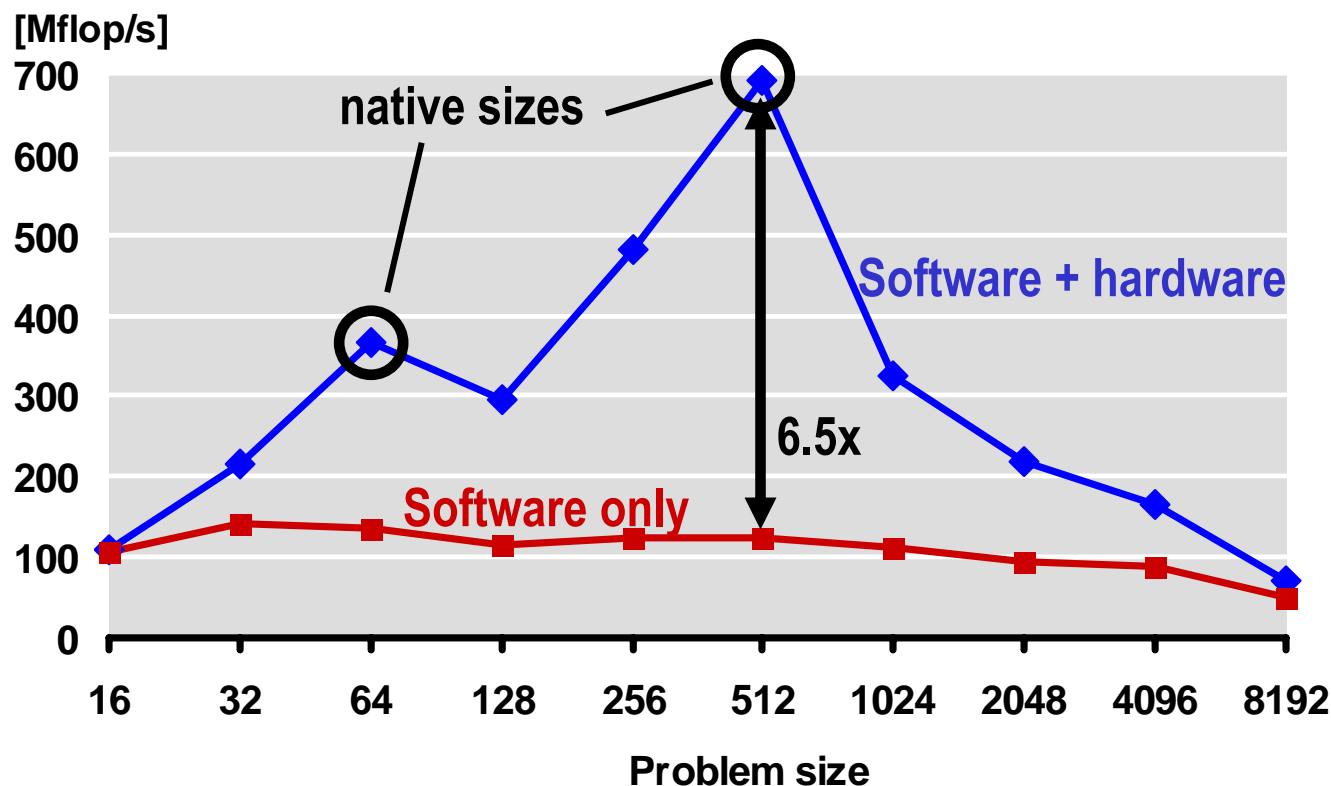


- competitive with professional designs
- much larger set of performance/area trade-offs

Joint work with P. Milder, J. Hoe (CMU)

Benchmark: Hardware Accelerator (FPGA)

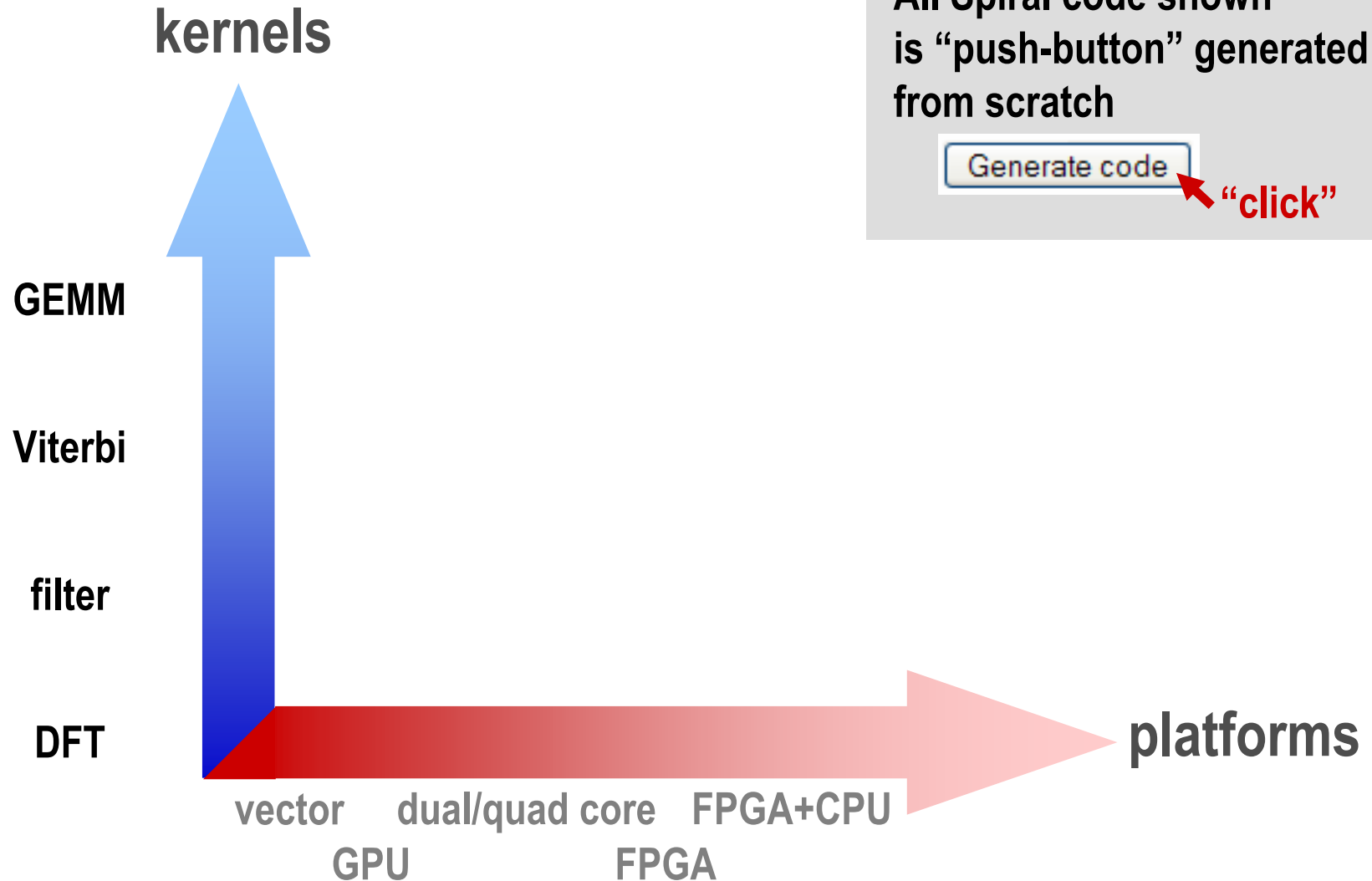
Xilinx Virtex 2 Pro FPGA: 1M gates @ 100 MHz + 2 PowerPC 405 @ 300 MHz



Fixed set of accelerators speed up a whole library

Joint work with P. D'Alberto (Yahoo), A. Sandryhaila, P. Milder, J. Hoe,
J. M. F. Moura (CMU), J. Johnson (Drexel)

Benchmarks



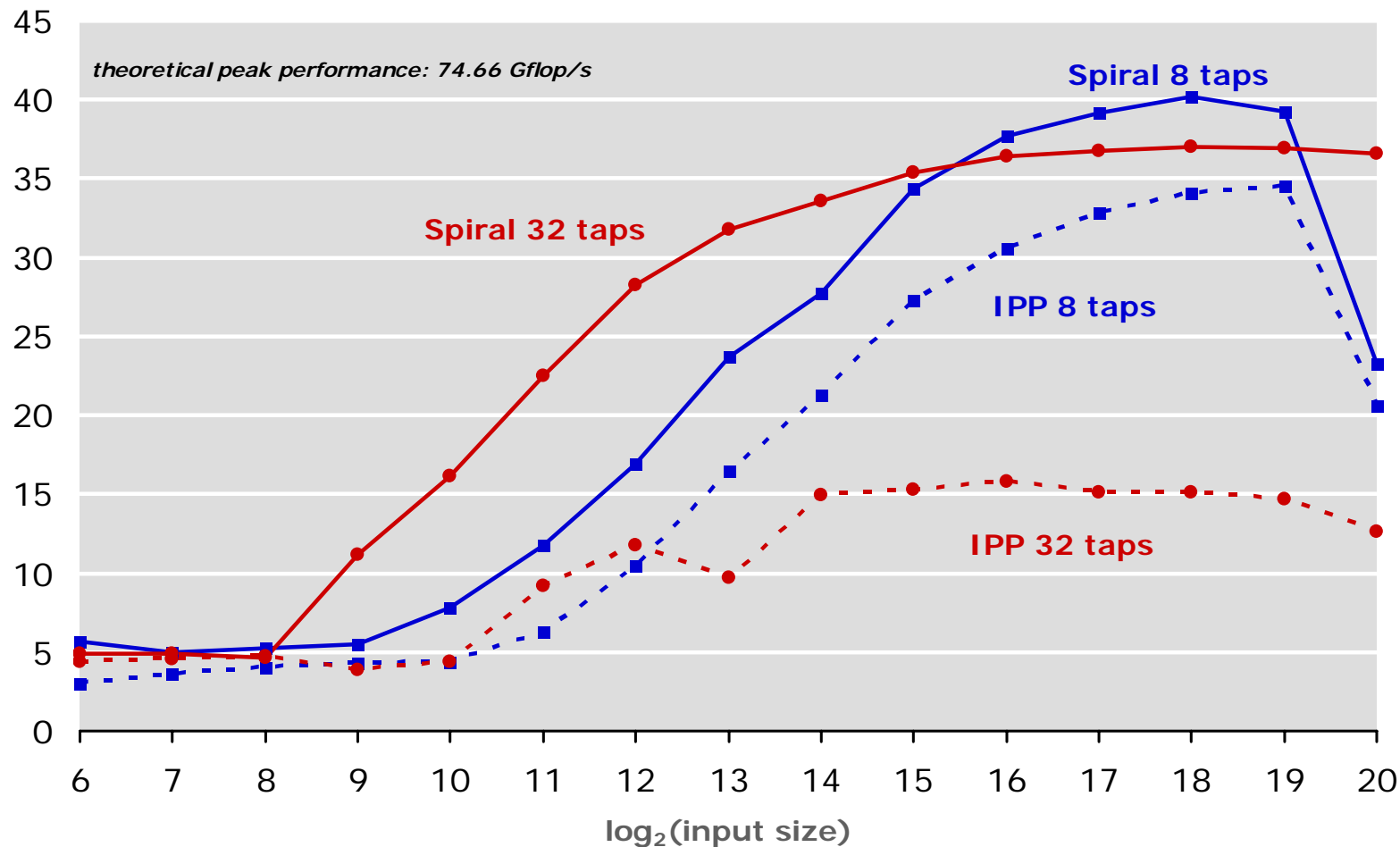
All Spiral code shown is “push-button” generated from scratch

Generate code

“click”

Benchmark: Finite Impulse Response Filter

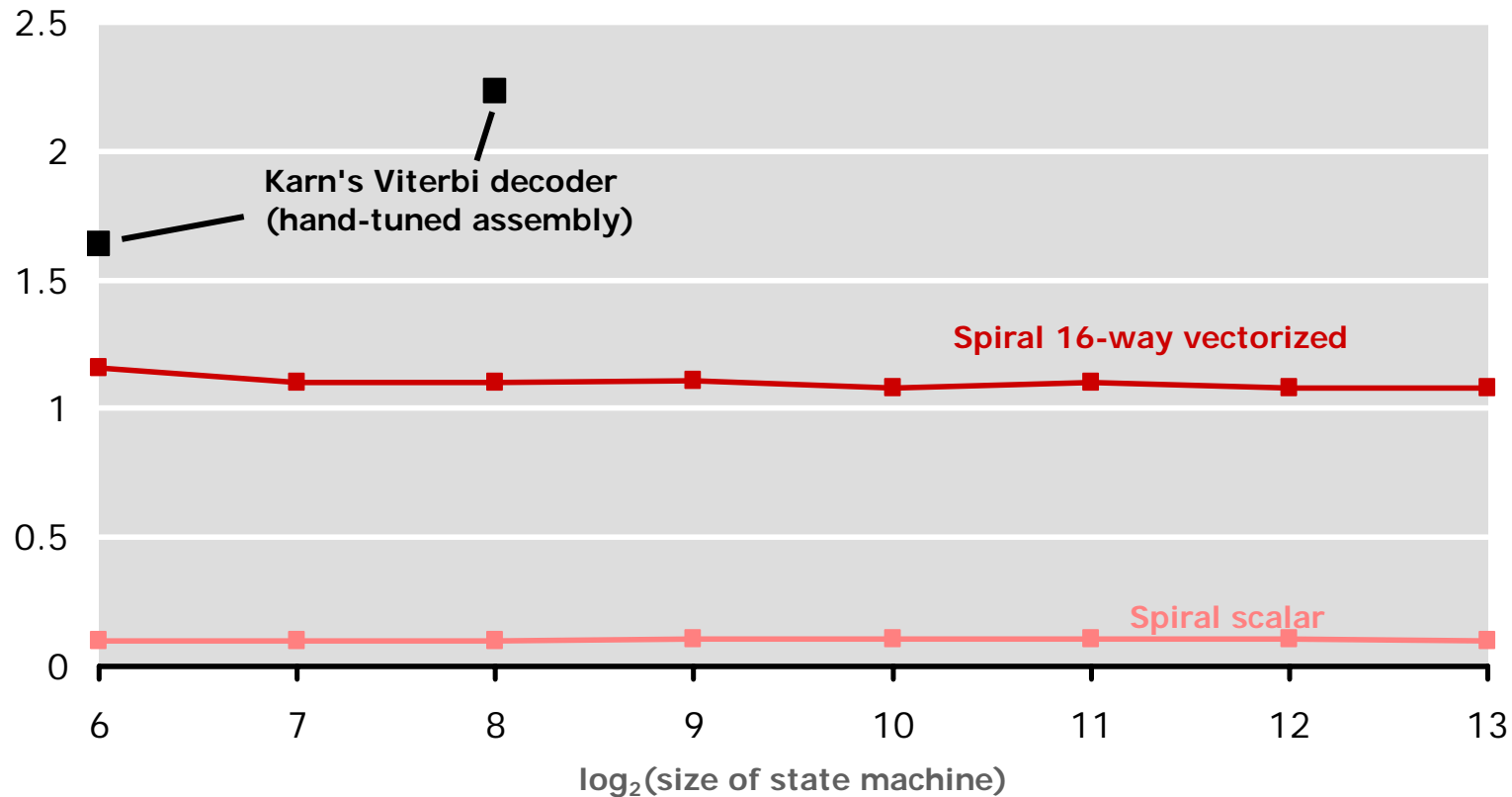
FIR filter (double precision) on 2.33 GHz 2x Core 2 Quad (8 threads)
 performance [Gflop/s]



Beyond Transforms : Viterbi Decoding

Viterbi decoding (8-bit) on 2.66 GHz Core 2 Duo
 performance [Gbutterflies/s]

1 butterfly
 = ~22 ops

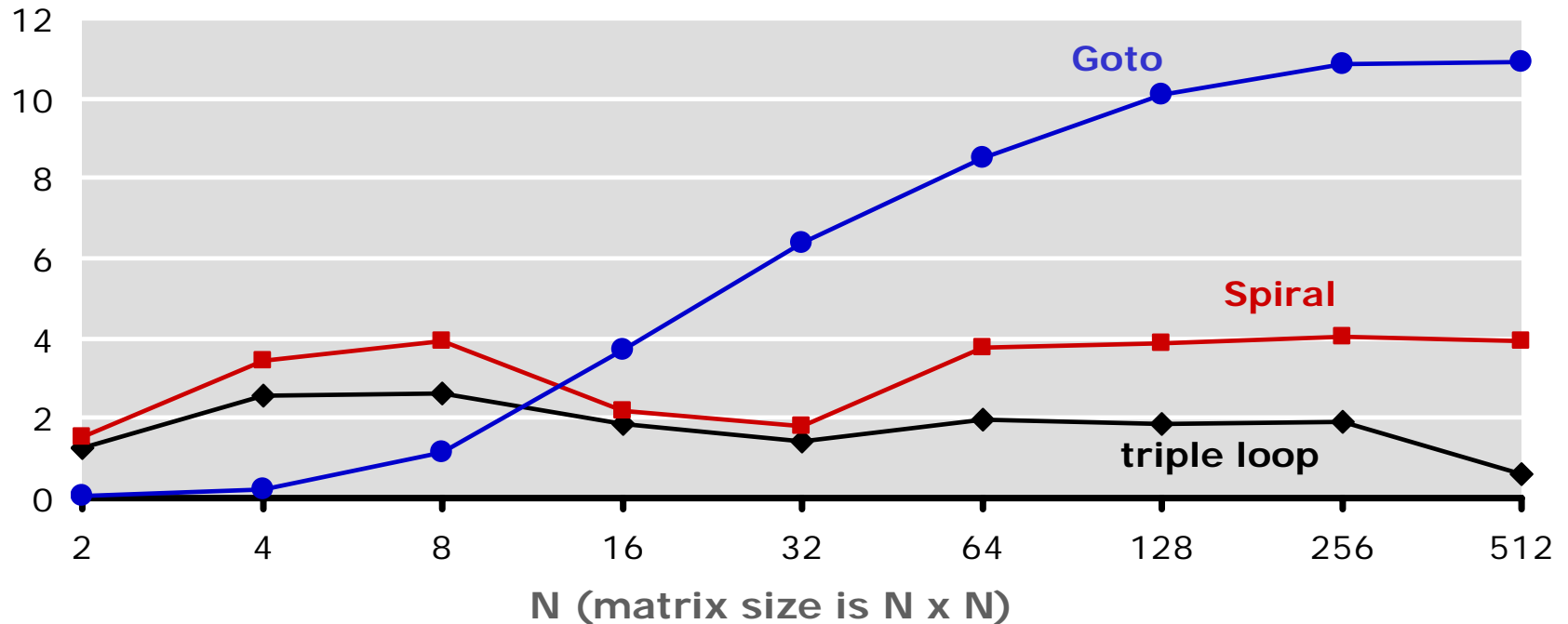


Vectorized using practically the same rules as for DFT

First Results: Matrix-Matrix-Multiply

DGEMM on 3 GHz Core 2 Duo (1 thread)

performance [Gflop/s]



Organization

- Spiral overview
- Parallelization in Spiral
- Results
- **Concluding remarks**

Conclusions

- **Automatic generation of very fast and fastest numerical kernels is possible and desirable**

- **High level language and approach**
 - Algorithm generation
 - Algorithm optimization

- **Same approach for loop optimization, different forms of parallelism, SW and HW implementations**

Spiral Web Interface @spiral.net (beta version)

1. Select platform

2. Select functionality

3. **Generate code** "click"

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