

Can We Teach Computers To Write Fast Libraries?

Markus Püschel

Collaborators:

Franz Franchetti

Yevgen Voronenko

Electrical and
Computer Engineering
Carnegie Mellon University

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



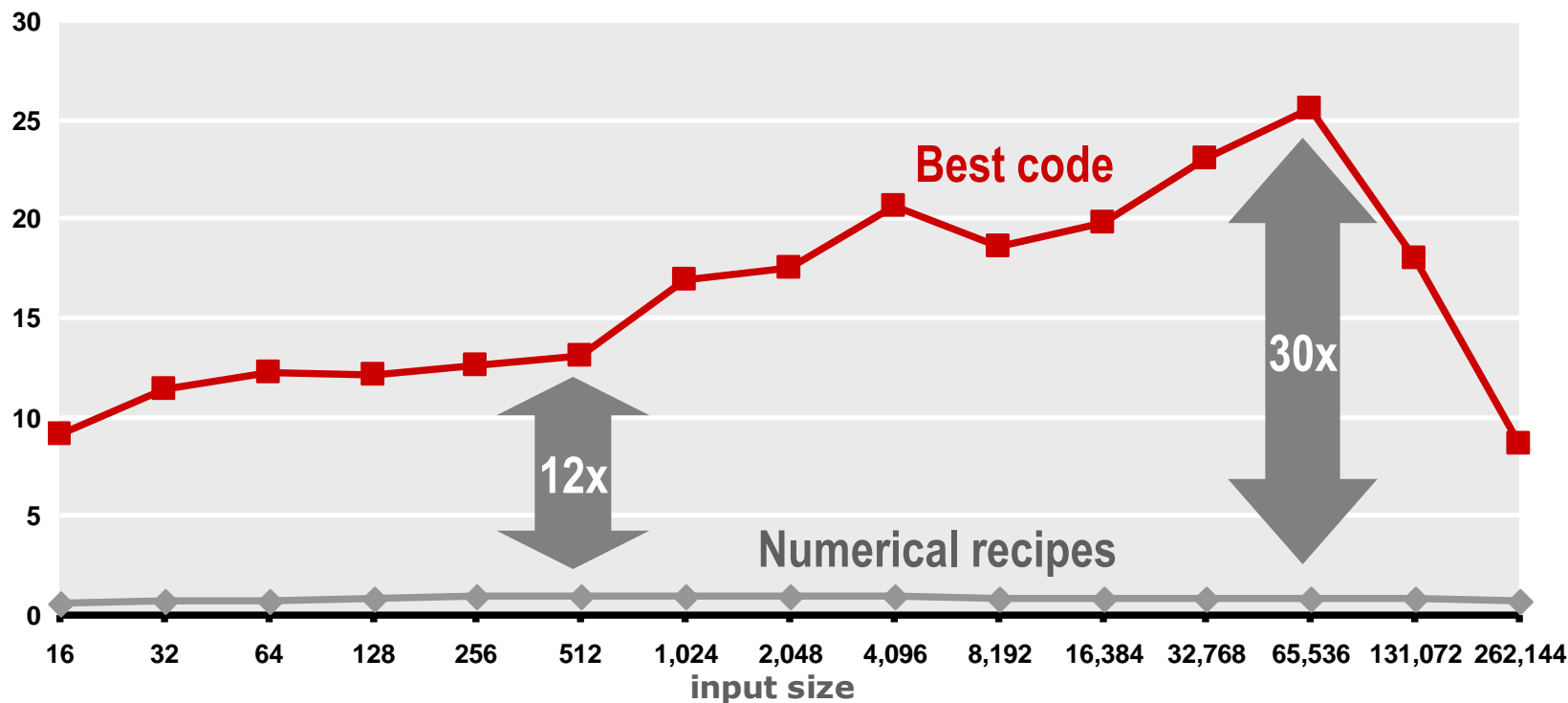
This work was supported by
DARPA DESA program, NSF-NGS/ITR, NSF-ACR, NSF-CPA, and Intel



Hiroshige, Temple in park near Osaka, ca. 1853-59

The Problem: Example DFT

Discrete Fourier Transform (DFT) on 2 x Core 2 Duo 3 GHz (single precision)
Performance [Gflop/s]

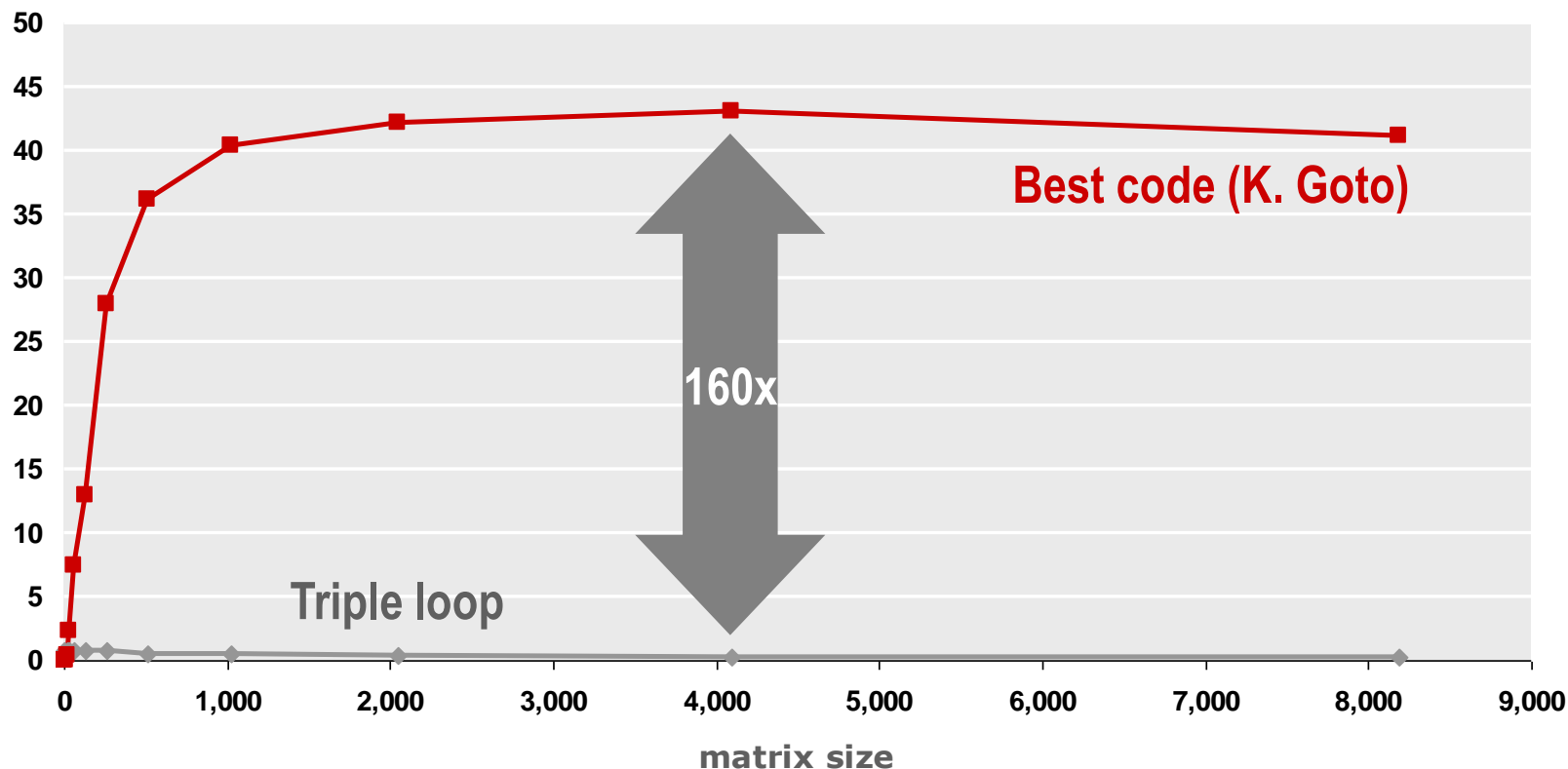


- Standard desktop computer, vendor compiler, using optimization flags
- All implementations have roughly the same operations count ($\sim 4n\log_2(n)$)
- *Maybe the DFT is just difficult?*

The Problem: Example MMM

Matrix-Matrix Multiplication (MMM) on 2 x Core 2 Duo 3 GHz (double precision)

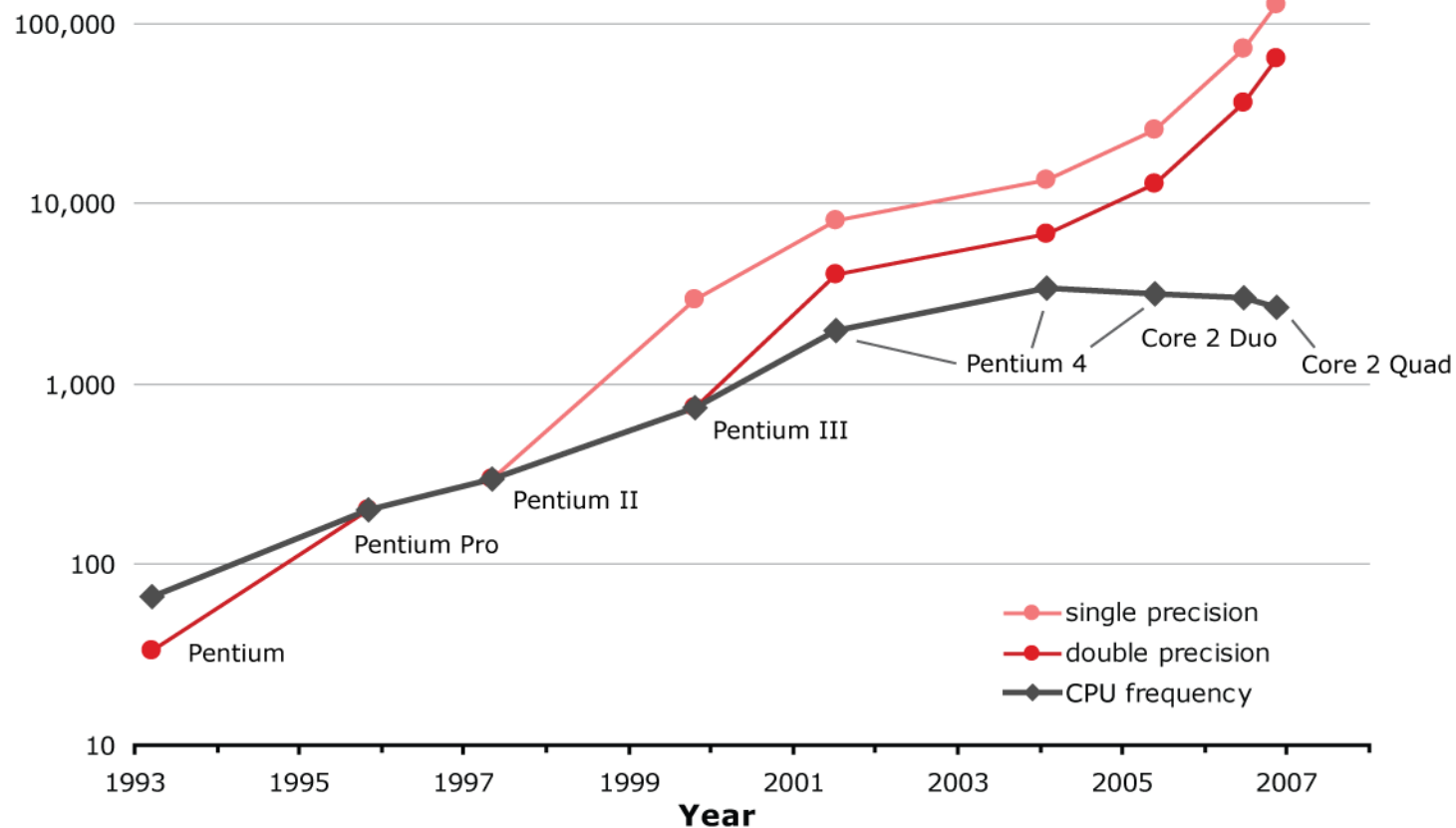
Performance [Gflop/s]



- Similar plots can be shown for all numerical kernels in linear algebra, signal processing, coding, crypto, ...
- *What's going on?*

Evolution of Processors (Intel)

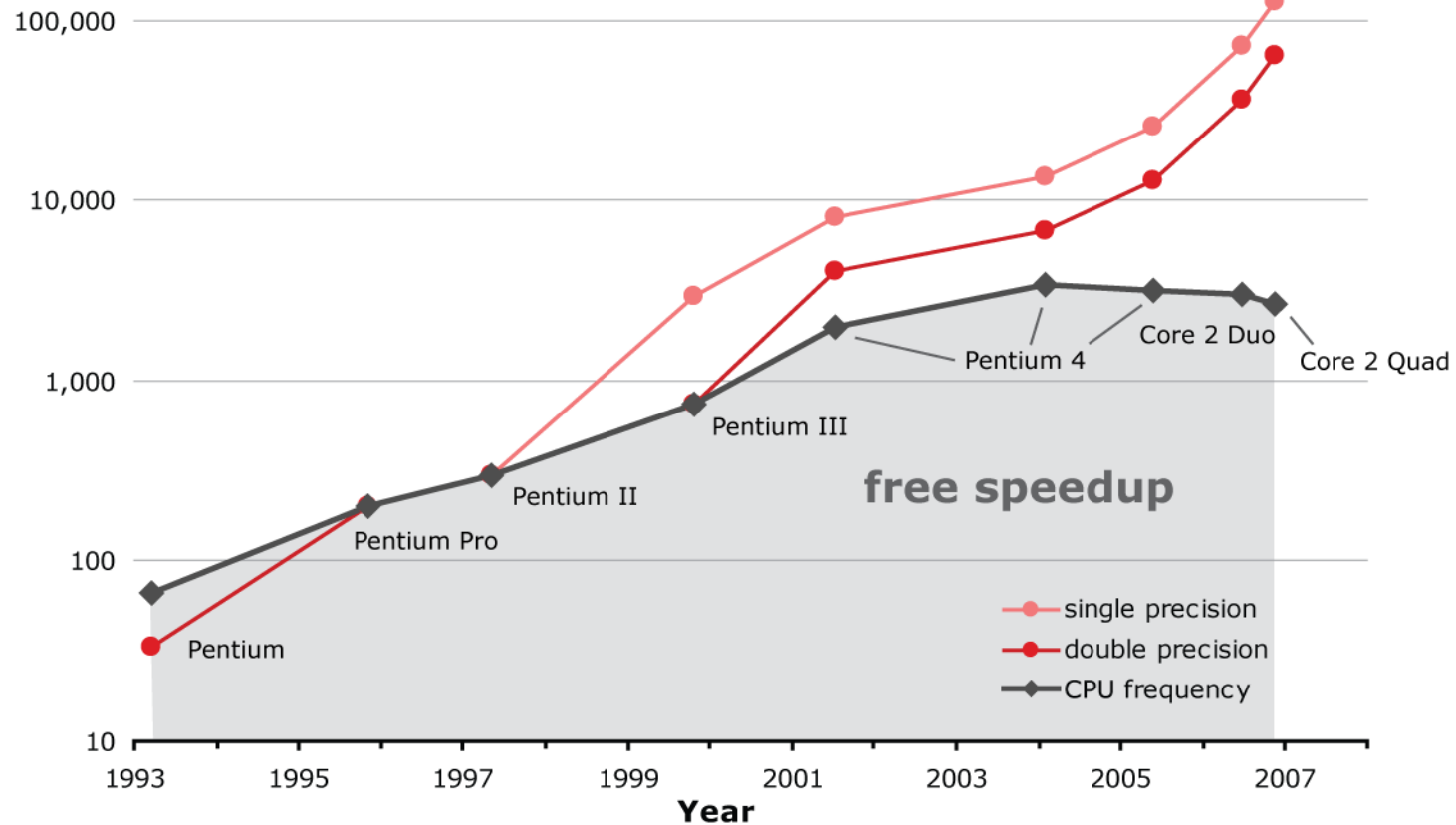
Floating point peak performance [Mflop/s]
CPU frequency [MHz]



data: www.sandpile.org

Evolution of Processors (Intel)

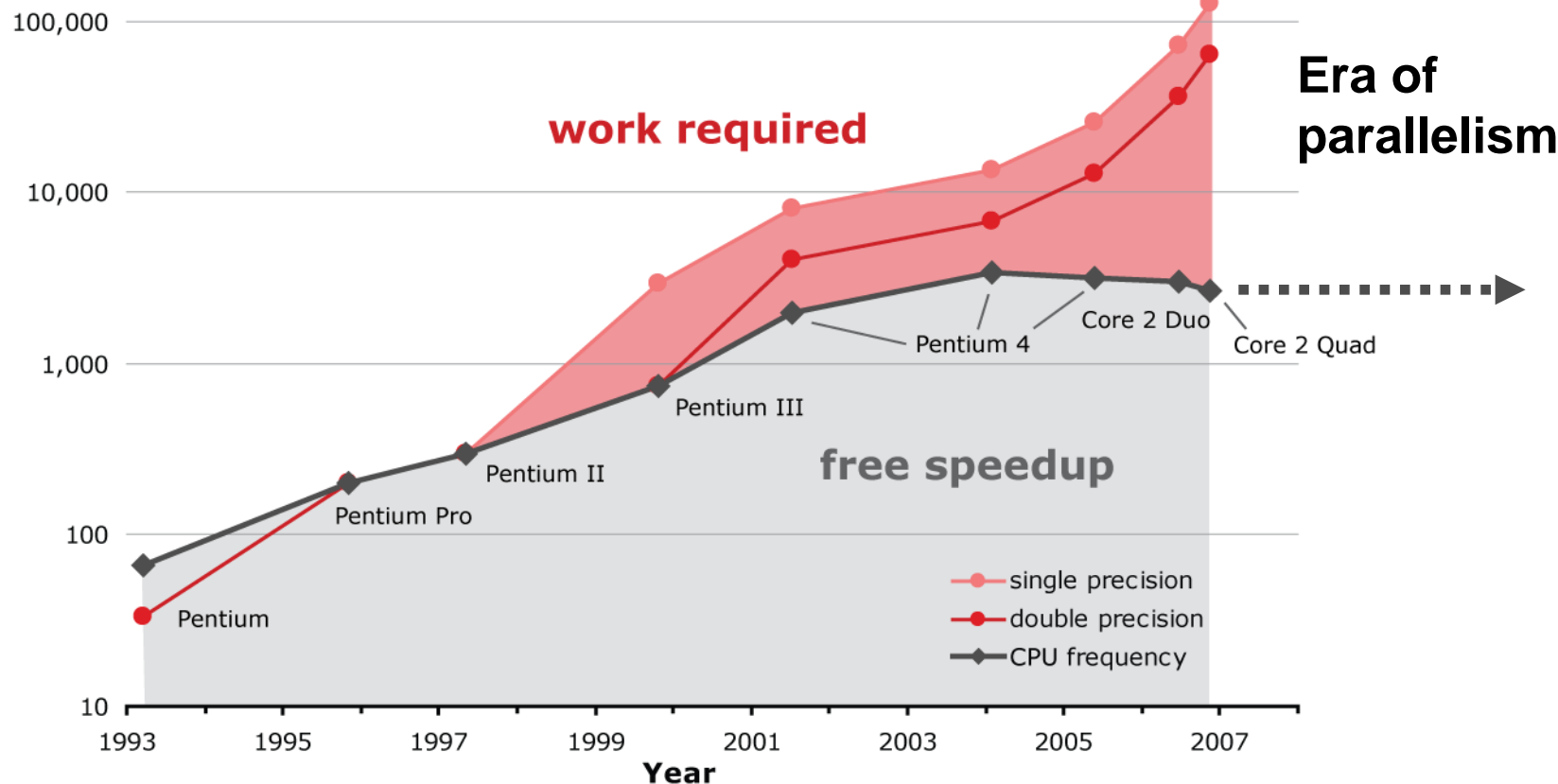
Floating point peak performance [Mflop/s]
CPU frequency [MHz]



data: www.sandpile.org

Evolution of Processors (Intel)

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



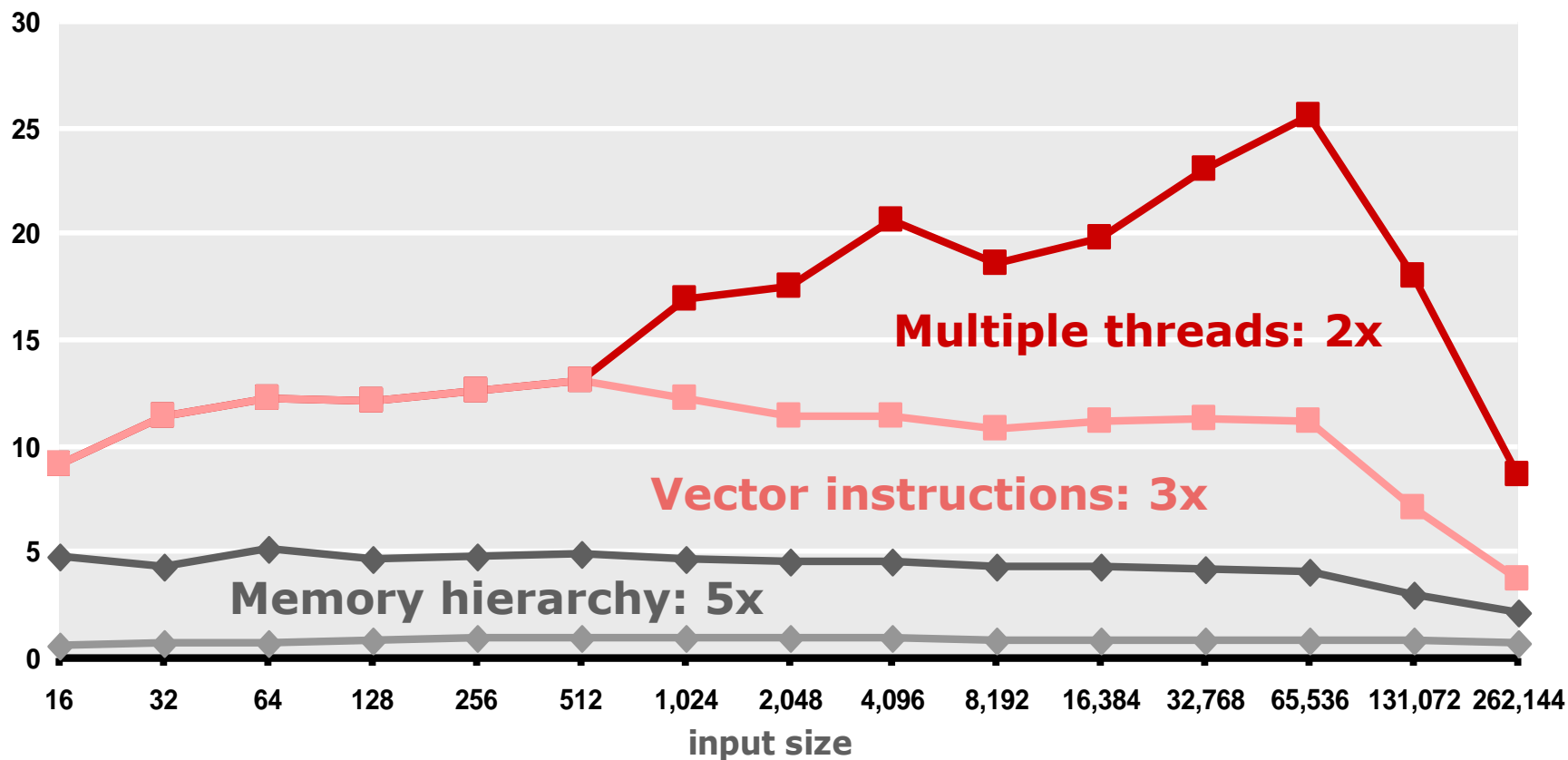
data: www.sandpile.org

High performance library development becomes increasingly difficult

DFT Plot: Analysis

Discrete Fourier Transform (DFT) on 2 x Core 2 Duo 3 GHz

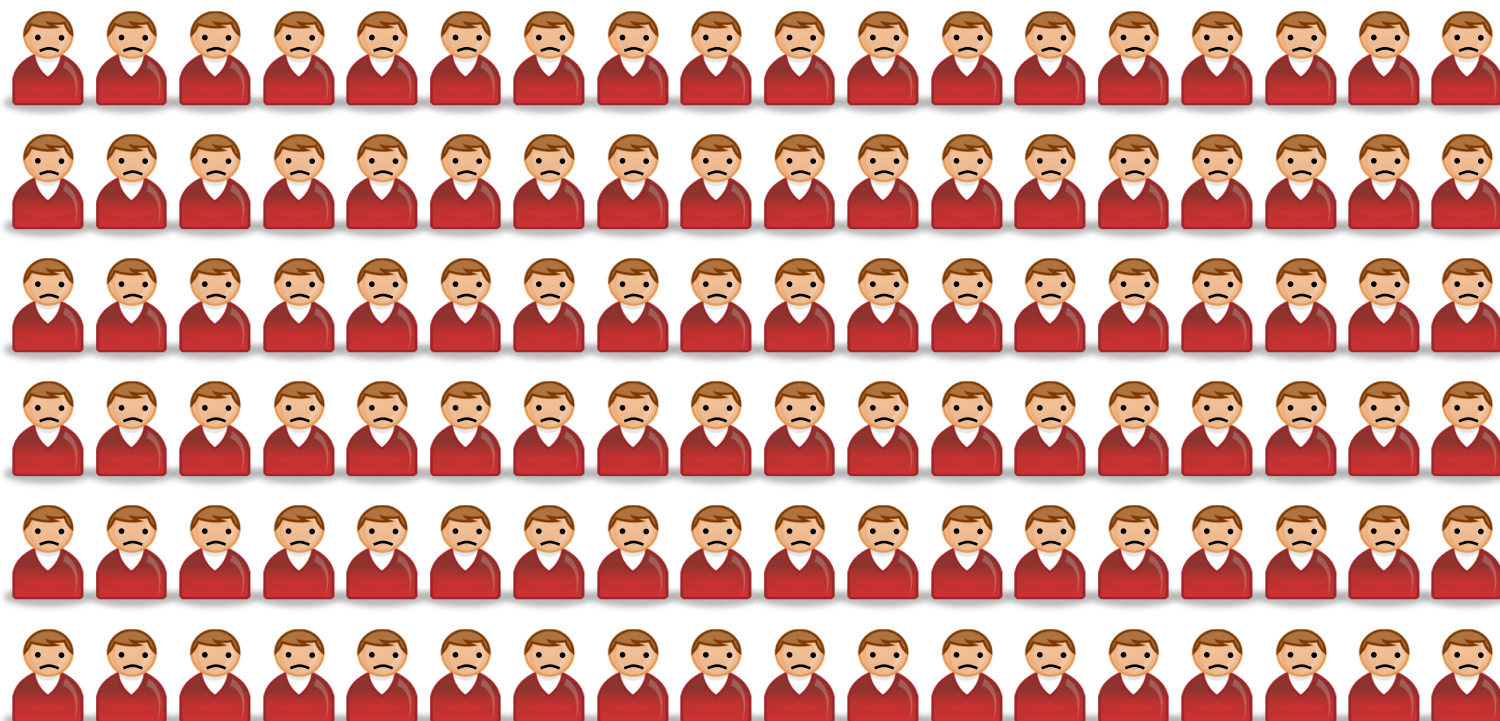
Gflop/s



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: Automatic Performance Tuning

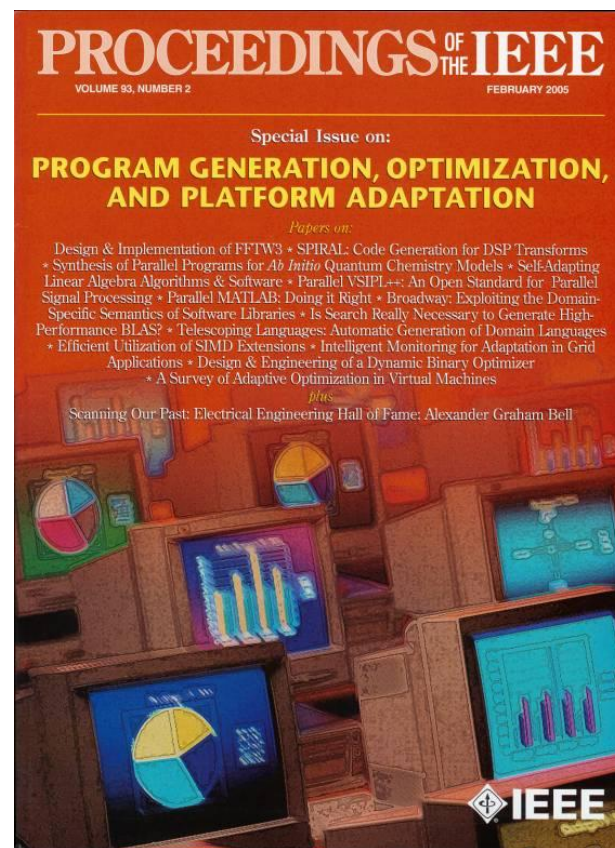
- Automate (parts of) the implementation or optimization



- Research efforts

- Linear algebra: Phipac/ATLAS (UTK), Sparsity/Bebop (Berkeley), Flame (UT Austin)
- Tensor computations (Ohio State)
- PDE/finite elements: Fenics
- Adaptive sorting (UIUC)
- Fourier transform: FFTW (MIT)
- Linear transforms: **Spiral**
- ...others
- New compiler techniques

**Promising new area but more work needed
In particular for parallelism ...**



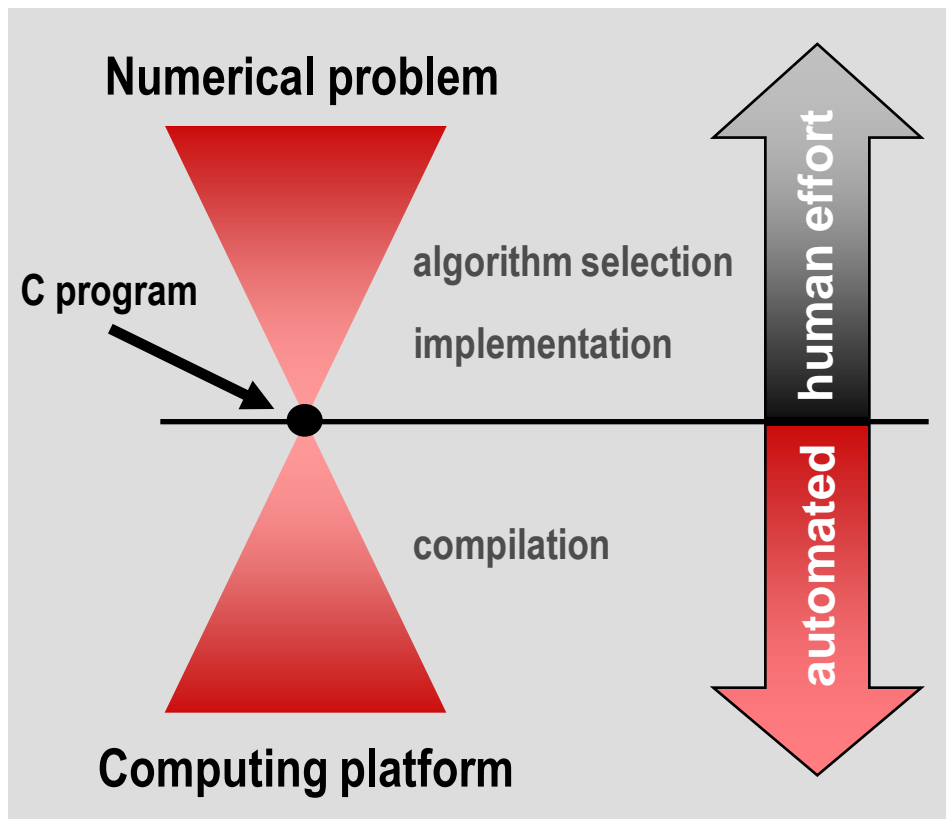
Organization

- **Spiral: Brief overview**
- Parallelization in Spiral
- Results
- Conclusion

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

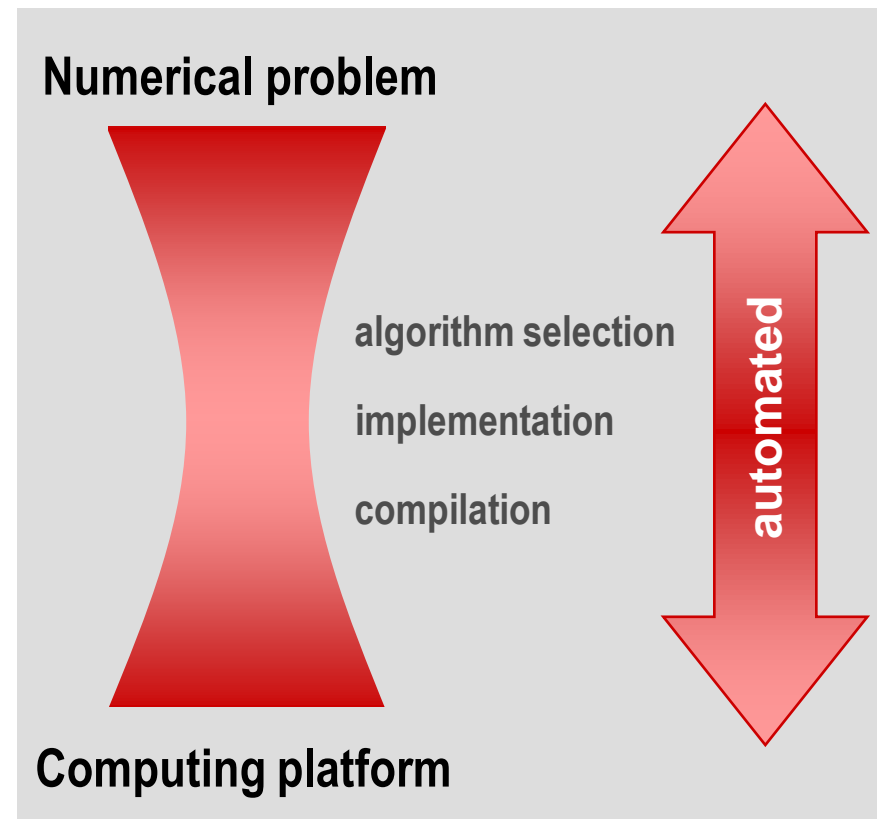
Vision Behind Spiral

Current



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

Future



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

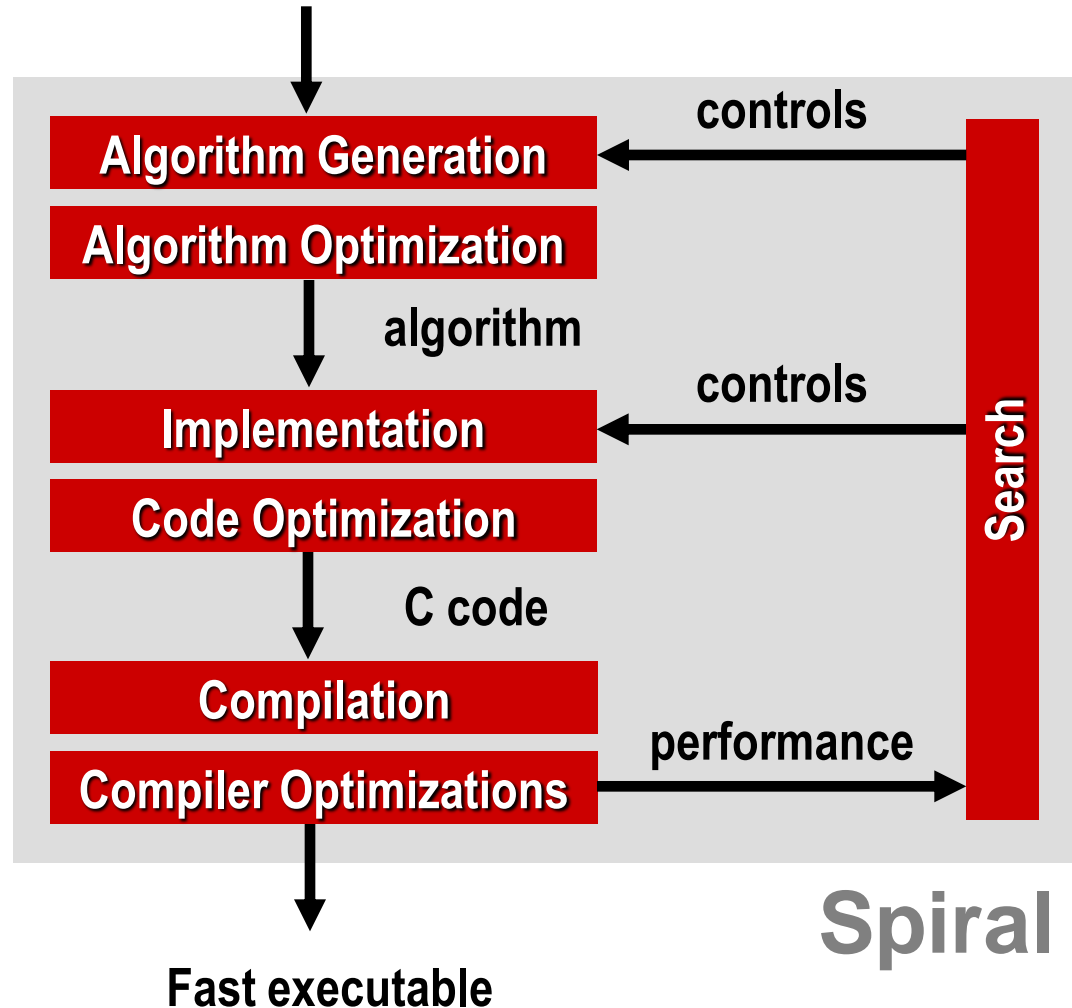
Spiral

- Library generator for linear transforms (DFT, DCT, DWT, 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

Problem specification (transform)



Spiral:

Complete automation of the implementation and optimization task

Basic ideas:

Declarative representation of algorithms

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

Spiral

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

Browse Archive expand

Filter by Platform: All Platforms Selected

Filter by Transform: All Transforms Selected

Filter by Size: All Sizes Selected

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)

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

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

Examples: Transforms (currently 55)

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

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

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

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

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

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

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

Examples: Breakdown Rules (currently ≈ 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)

BSSkewDFT 4				Filt 10		PDHT2 2	
BSkewDFT_Base2				RuleFilt_Base		PDHT2_Base2	
BSkewDFT_Base4				RuleFilt_Nest		PDHT2_CT	
BSkewDFT_Fact				RuleFilt_OverlapSave		PDHT3 4	PolyDFT 4
BSkewDFT_Decom				RuleFilt_OverlapSaveFreq		PDHT3_Base2	PolyDFT_ToNormal
BSSkewPRDFT 1				RuleFilt_OverlapAdd		PDHT3_CT	PolyDFT_Base2
BSkewPRDFT_Decom				RuleFilt_OverlapAdd2		PDHT3_Trig	PolyDFT_SkewBase2
BSSkewRDFT 1				RuleFilt_KaratsubaSimple		PDHT3_CT_Radix2	PolyDFT_SkewDST3_CT
BSkewPRDFT_Decom				RuleFilt_Circulant		PDHT4 3	RDFT 4
Circulant 10				RuleFilt_Blocking		PDHT4_Base2	RDFT_Base
RuleCirculant_Base				RuleFilt_KaratsubaFast		PDHT4_CT	RDFT_Trigonometric
RuleCirculant_toFilt				IPRDFT 5		PDHT4_Trig	RDFT_CT_Radix2
RuleCirculant_Blocking				IPRDFT1_Base1		PDST 2	RDFT_toDCT2
RuleCirculant_BlockingDense				IPRDFT1_Base2		PDST4 2	RHT 2
RuleCirculant_BlockingDense				IPRDFT1_CT		PDST4_Base2	RHT_Base
RuleCirculant_DiagonalizeStep				IPRDFT1_Complex		PDST4_CT	RHT_CooleyTukey
RuleCirculant_DFT				IPRDFT1_CT_Radix2		PRDFT 10	SinDFT 2
RuleCirculant_RDFT				IPRDFT2 4		PRDFT1_Base1	SinDFT_Base
RuleCirculant_PRDFT				IPRDFT2_Base1		PRDFT1_Base2	SinDFT_Trigonometric
RuleCirculant_PRDFT1				IPRDFT2_Base2		PRDFT1_CT	SkewDFT 2
RuleCirculant_DHT				IPRDFT2_CT		PRDFT1_Complex	SkewDFT_Base2
CosDFT 2				IPRDFT2_CT_Radix2		PRDFT1_Complex_T	SkewDFT_Fact
CosDFT_Base				IPRDFT3 3		PRDFT1_Trig	SkewDFT 3
CosDFT_Trigonometric				IPRDFT3_Base1		PRDFT1_PF	SkewDFT_Base2
DCT1 3				IPRDFT3_Base2		PRDFT1_CT_Radix2	SkewDCT3_VarSteidl
DCT1_Base2				IPRDFT3_CT			
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_Decom							
DCT5 1							
DCT5_Rader							
DFT 22							
DFT_Base							
DFT_Canonicalize							
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 1							
DST1_Base2							
DST1_DST3and1							
DST2 3							
DST2_Base2							
DST2_toDCT2							
DST2_DST2and4							
DST4 2							
DST4_Base							
DST4_toDCT4							
DTT 15							
DTT_C3_Base2							
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_Dimless							
MDDFT_tSPL_RowCol							
PDCT 4 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_Decom							
WHT 8							
WHT_Base							
WHT_GeneralSplit							
WHT_BinSplit							
WHT_DDL							
WHT_Dirsum							
WHT_tSPL_BinSplit							
WHT_tSPL_STerm							
WHT_tSPL_Pease							

Cooley-Tukey FFT

Breakdown rules in Spiral:

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

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

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

Correct code: easy

fast code: very difficult

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:
[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

.....

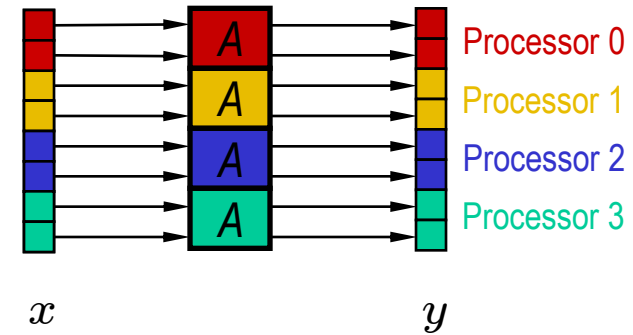
Organization

- Spiral: Brief overview
- **Parallelization in Spiral**
- Results
- Conclusion

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

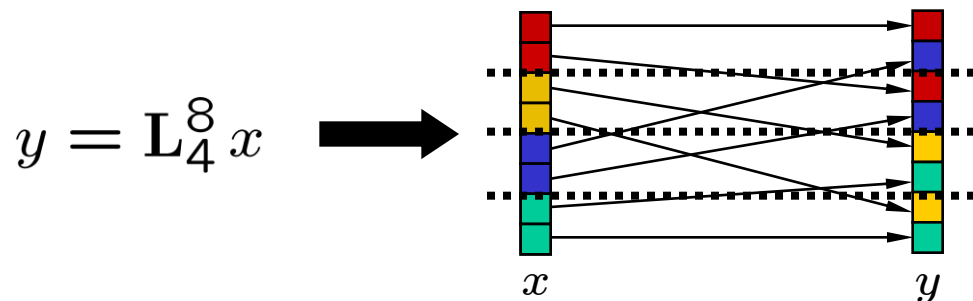
- Governing construct: tensor product

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



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: 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** for (p, μ) 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

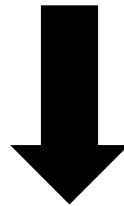
■ Goal: Transform formulas into fully optimized formulas

- Formulas rewritten, tags propagated
- There may be choices

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

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

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

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

Same Approach for Other Parallel Paradigms

Message Passing: [ISPA 06]

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

With Bonelli, Lorenz, Ueberhuber, TU Vienna

Cg/OpenGL for GPUs:

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

With Shen, TU Denmark

Vectorization: [IPDPS 02, VecPar 06]

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

With Milder, Hoe, CMU

Verilog for FPGAs: [DAC 05]

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

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

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: $MMM \rightarrow I \otimes MMM$

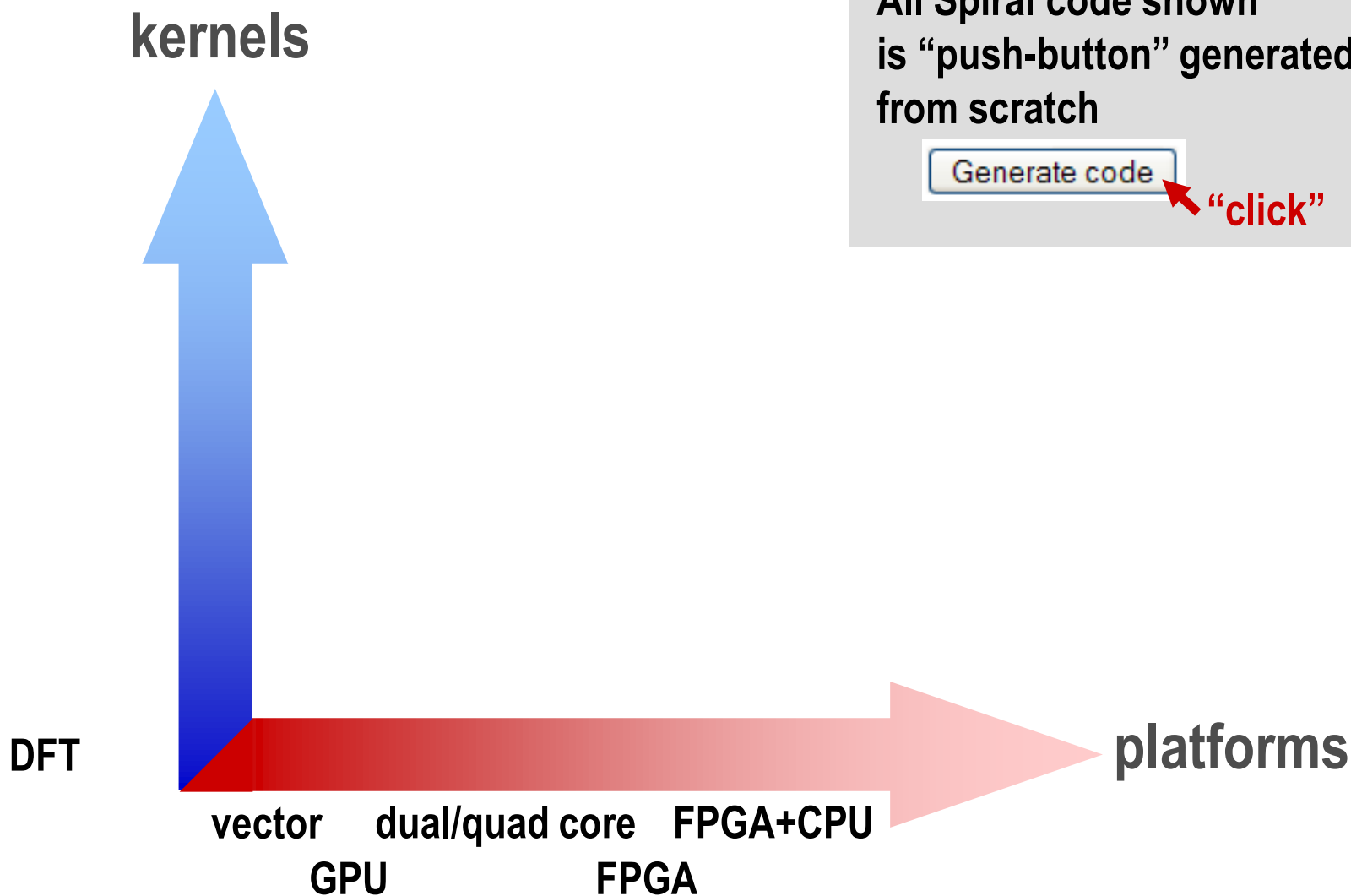
$MMM \rightarrow (I \otimes L) \circ (MMM \otimes I) \circ (I \times (I \otimes L))$

- We have first results beyond transforms!

Organization

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

Benchmarks



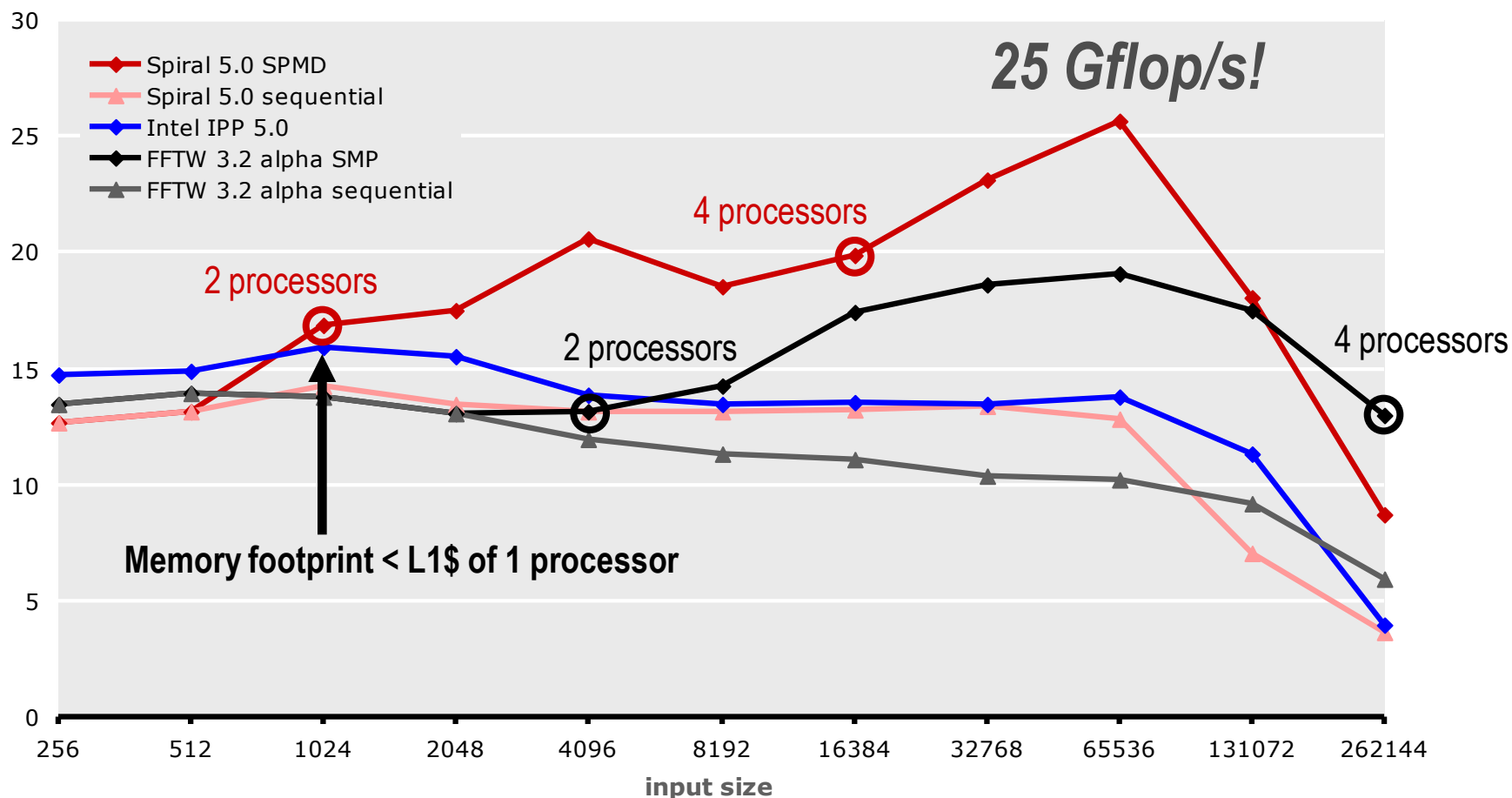
All Spiral code shown is “push-button” generated from scratch

Generate code

“click”

Benchmark: Vector and SMP

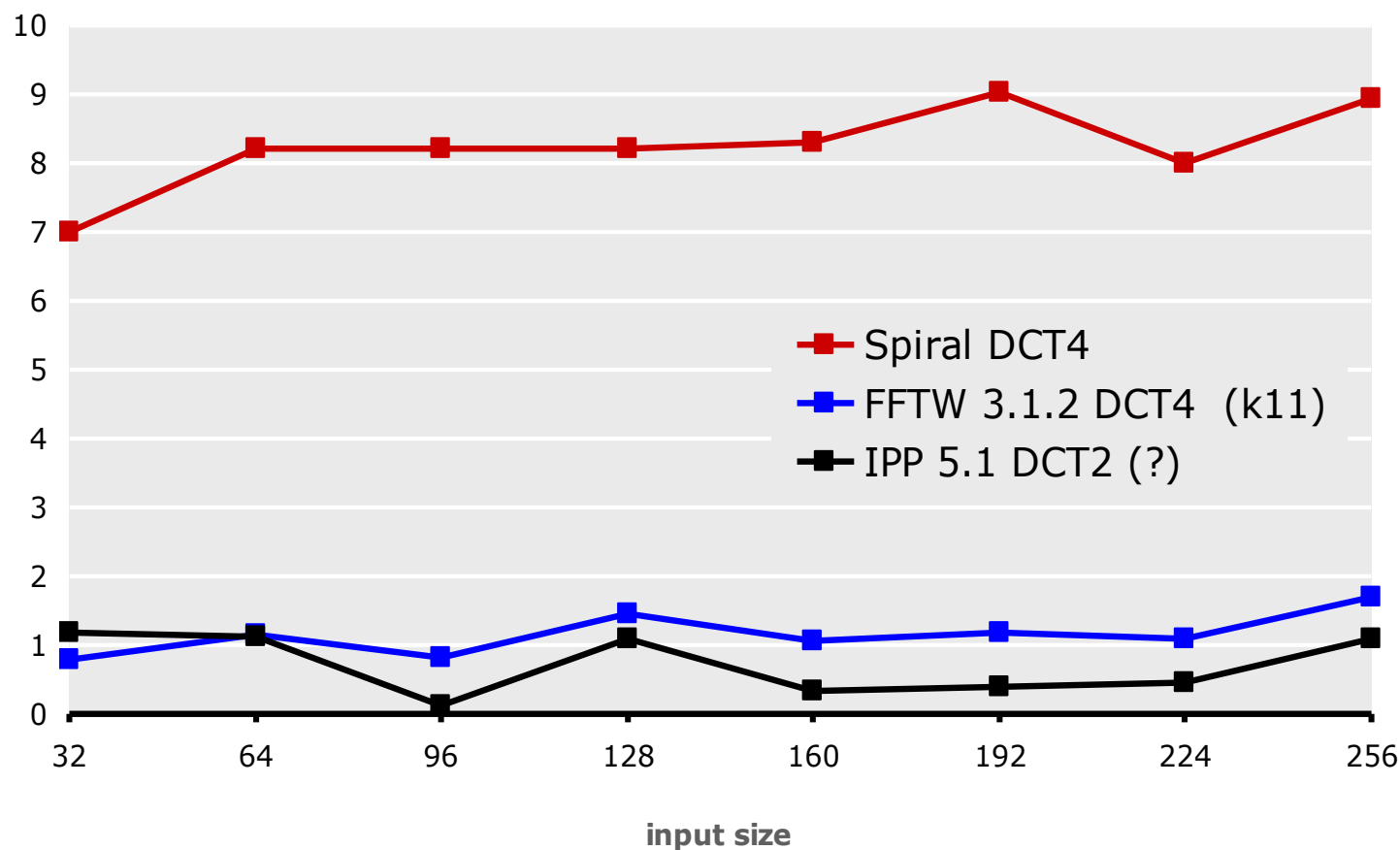
DFT (single precision): on 3 GHz 2 x Core 2 Extreme
performance [Gflop/s]



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

DCT4, Multiples of 32: 4-way Vectorized

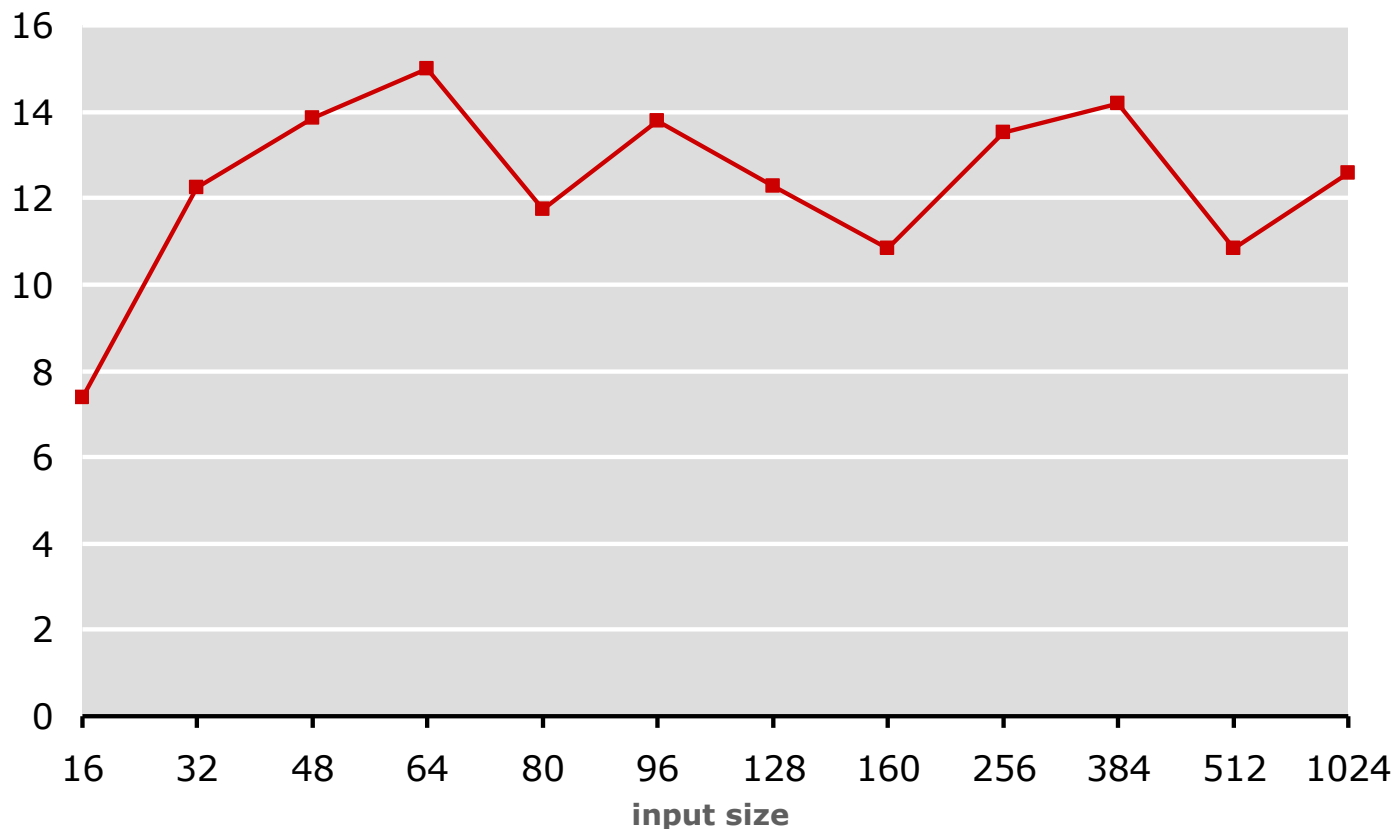
DCT (single precision) 2.66 GHz Core2 (4-way 32-bit SSE)
performance [Gflop/s]



Less popular transform usually means no good code

Benchmark: Cell (1 processor = SPE)

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

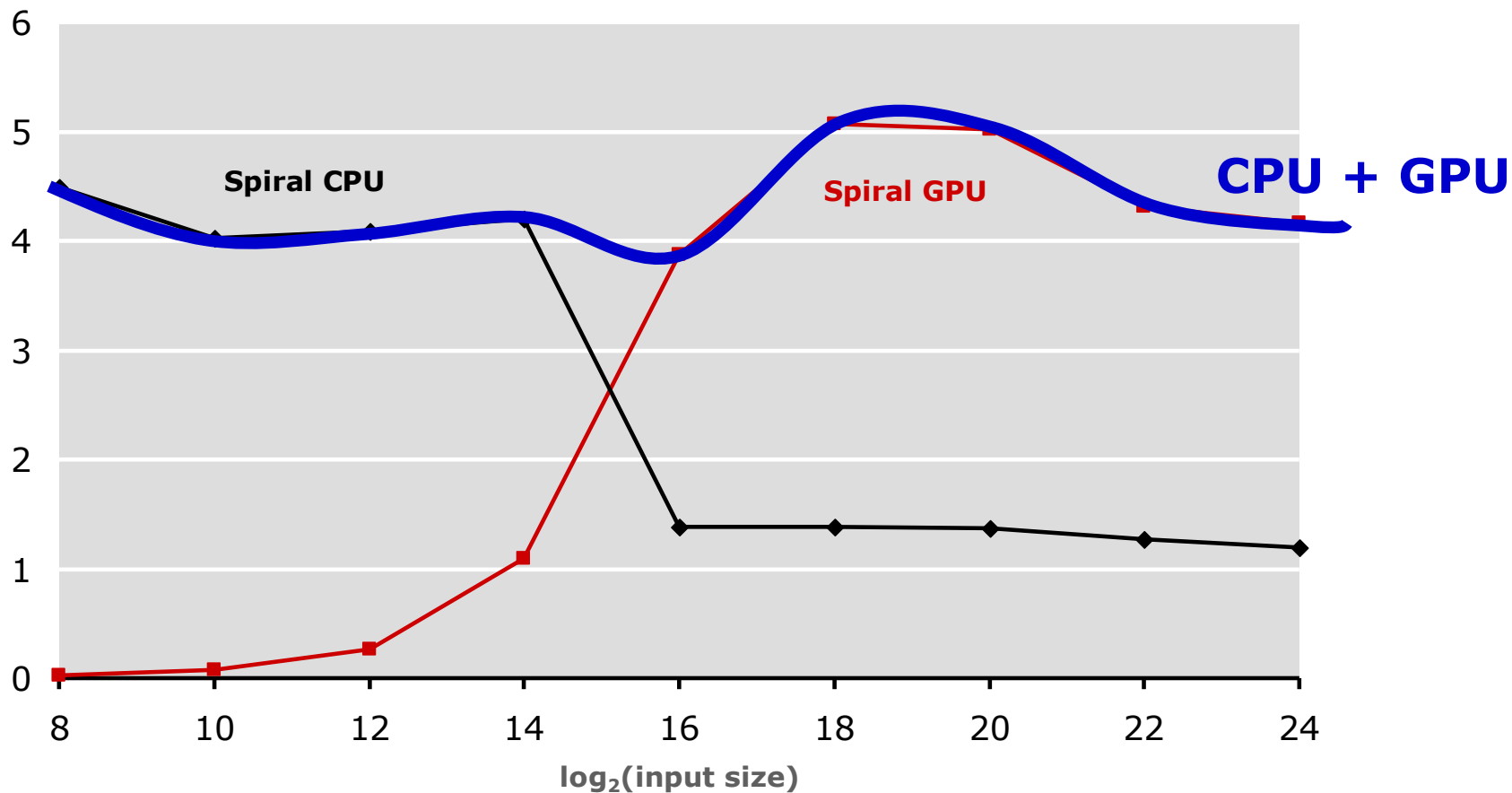


Generated using the simulator; run at Mercury (thanks to Robert Cooper)

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

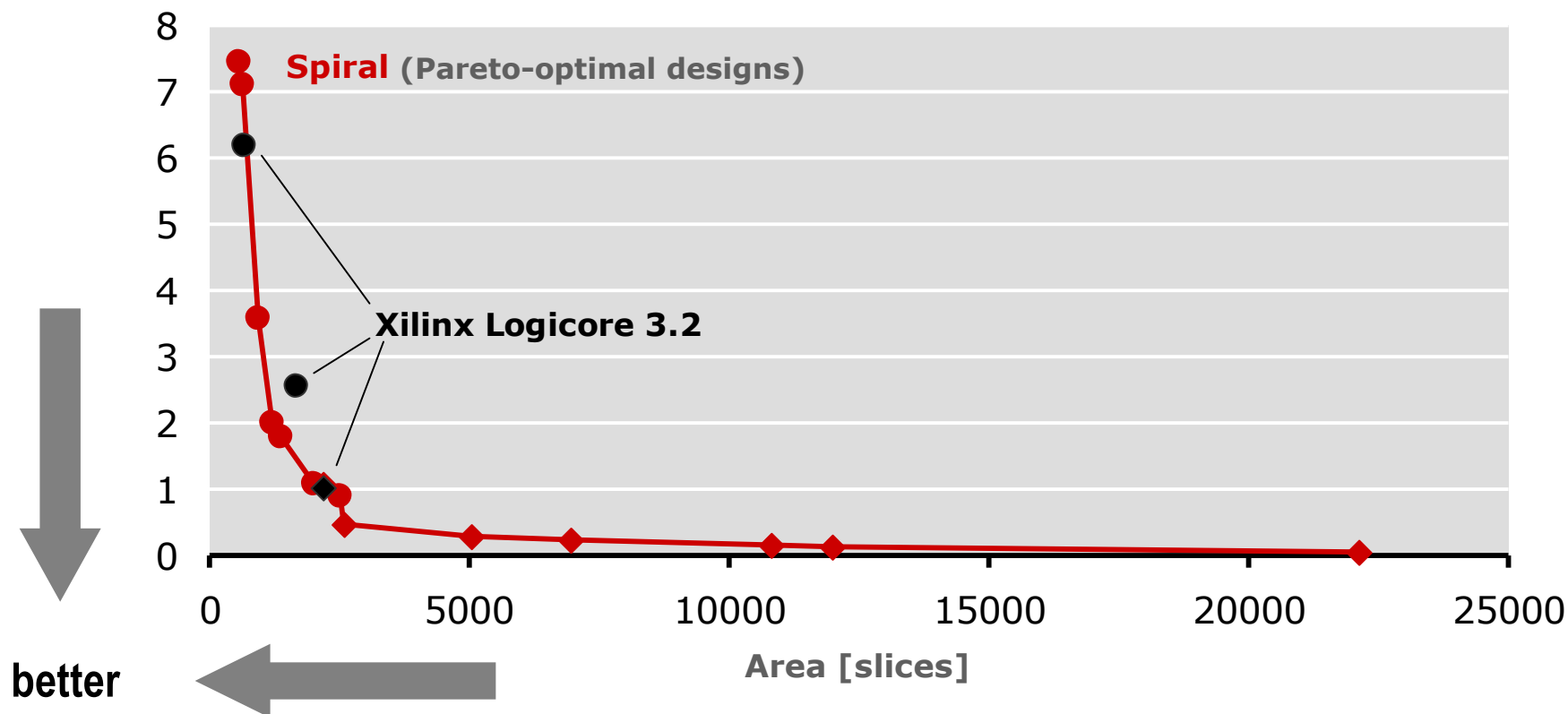
Benchmark: GPU

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



Benchmark: FPGA

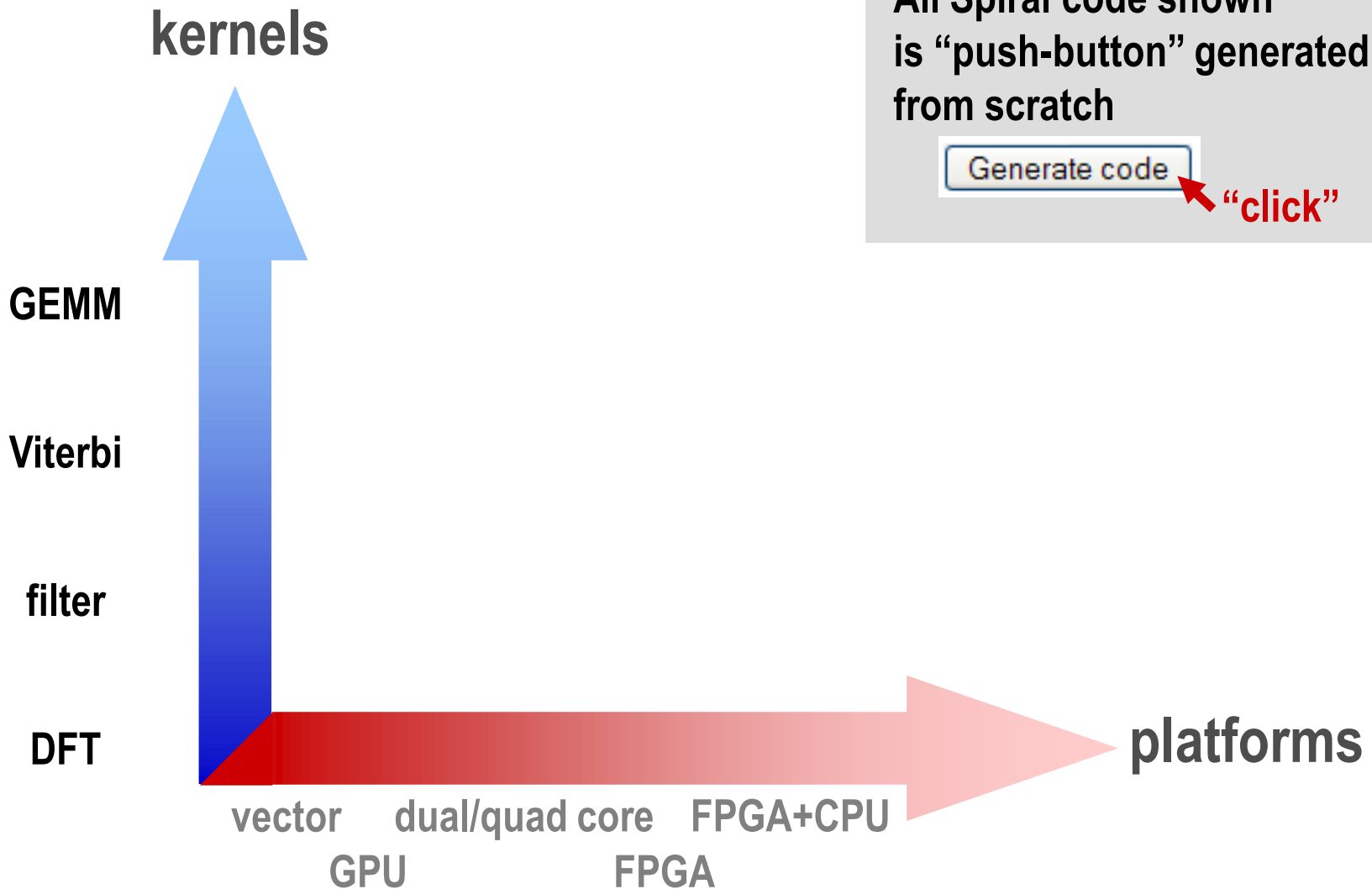
DFT 256 on Xilinx Virtex 2 Pro FPGA
inverse throughput (gap) [us]



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

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

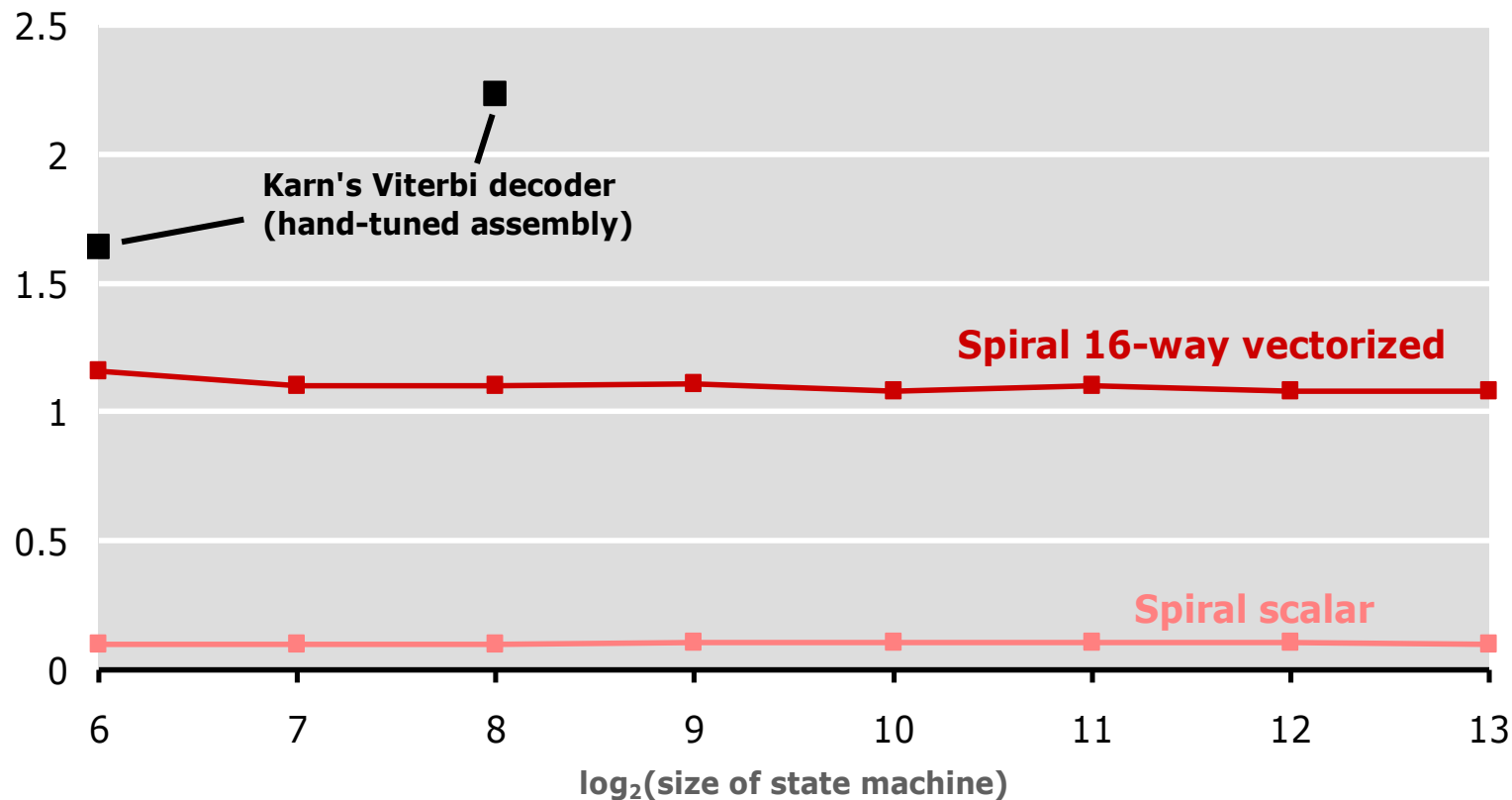
Benchmarks



Beyond Transforms : Viterbi Decoding

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

1 butterfly
= ~22 ops

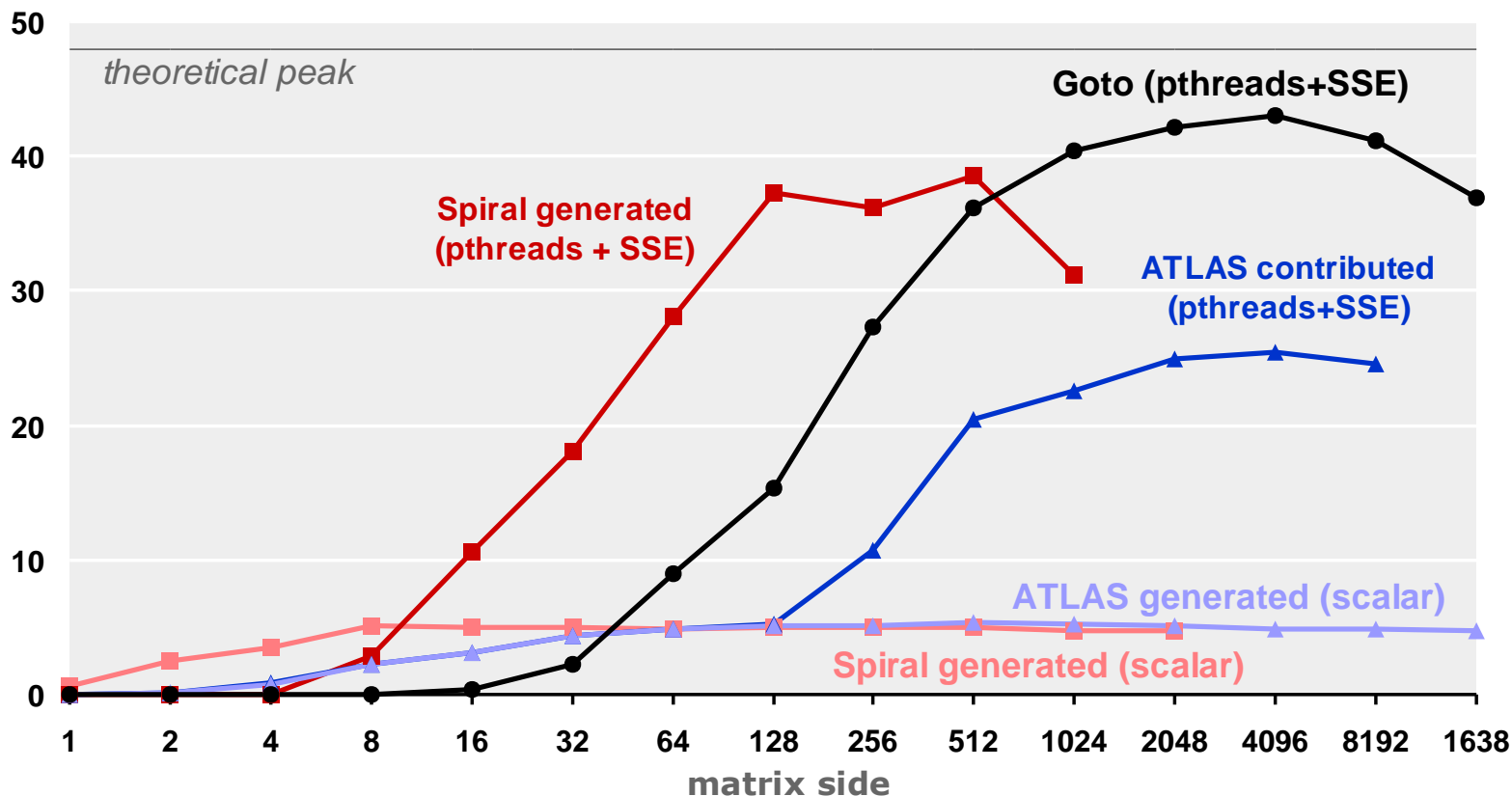


Vectorized using practically the same rules as for DFT

Beyond Transforms: Matrix-Matrix-Multiply

MMM(square matrix real double) 2 x Core2Duo 3Ghz

performance [Gflop/s]



All measurements with warm cache

ATLAS code is optimized for cold cache

Organization

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

Spiral Summary

■ Spiral: The computer implements and optimizes transforms

- Rigorous formal framework
- Support for different platform paradigms (vector, SMP, GPU, FPGA, ...)
- Optimization at high abstraction level through rewriting
+ empirical search over alternatives
- The generated code is often faster than human written code
- We have extended the framework beyond transforms

■ What we have learned

- Declarative representation of algorithms (mathematical DSL)
- 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

Attacking the Parallel Library Challenge

Automation in High Performance Library Development

Programming languages
Program generation

Symbolic Computation
Rewriting

*High-Performance
Parallel Library
Development*

Software
Scientific Computing

Algorithms
Mathematics

Compilers

We Need to Work Together