

# A Rewriting System for the Vectorization of Signal Transforms

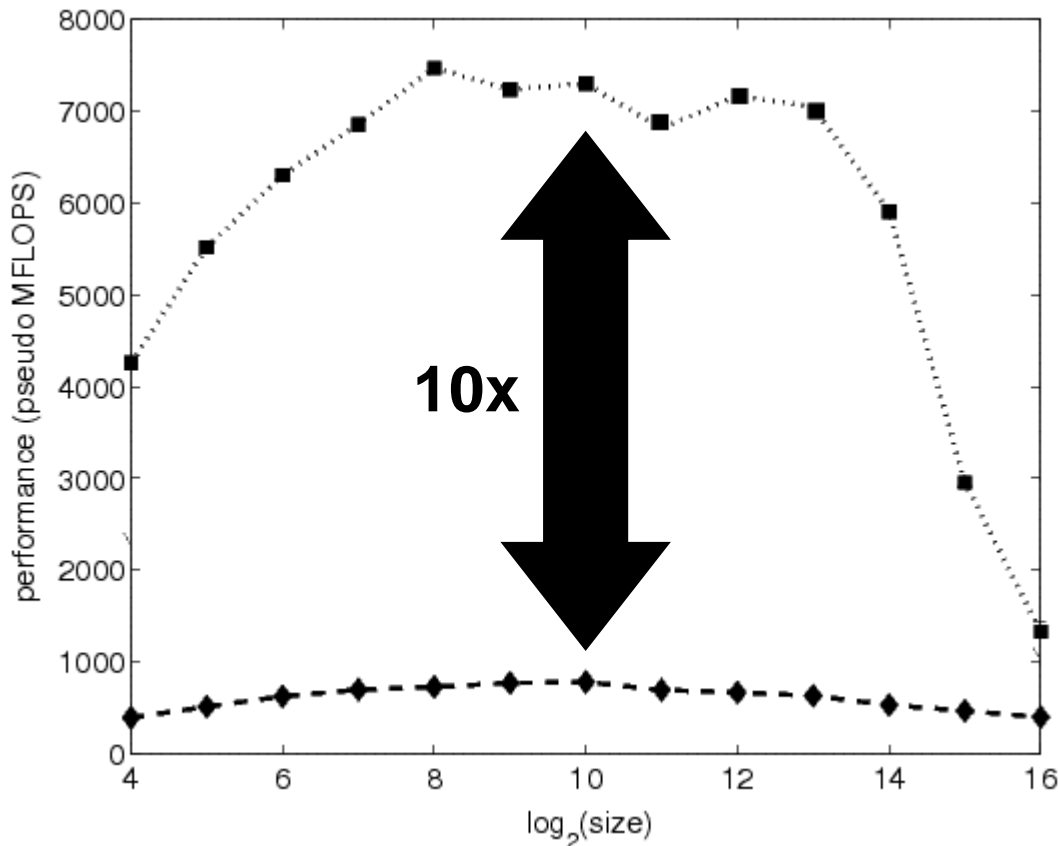
Franz Franchetti  
Yevgen Voronenko  
Markus Püschel

Department of Electrical & Computer Engineering  
Carnegie Mellon University

<http://www.spiral.net>

**Supported by:** NSF ACR-0234293, ITR/NGS-0325687,  
DARPA NBCH-105000, Intel, Austrian FWF

# The Problem (Example FFT Performance)



**best available  
implementation**  
(FFTW, Intel IPP, Spiral)



**roughly the same  
operations count**



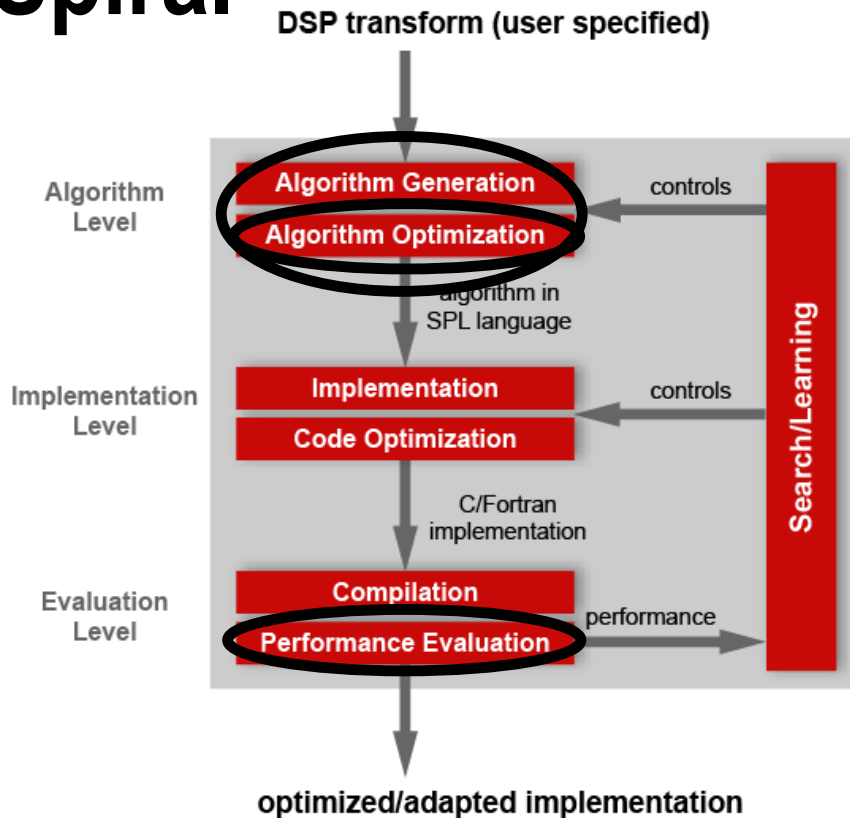
**reasonable  
implementation**  
(Numerical recipes.  
GNU scientific library)

**Solution: program generators like Atlas and Spiral,  
adaptive libraries like FFTW**

# Organization

- **Spiral overview**
- **SIMD vector instructions**
- **Vectorization by rewriting**
- **Extension to SMP and Multicore**
- **Experimental results**
- **Summary**

# Spiral



- **Program generation from a problem specification** for linear digital signal processing (DSP) transforms (DFT, DCT, DWT, filters, ....)
- **Goal 1:** A flexible push-button program generation framework for an entire domain of algorithms
- **Goal 2:** With new architectures, update the tool rather than the individual programs in the library

**Spiral:** Principle 2: Optimization at a high level of abstraction and memory, multicore, distributed memory, FPGAs, embedded CPUs

# What is a DSP Transform?

- **Mathematically: Matrix-vector multiplication**

$$x \mapsto y = T \cdot x$$

input vector (signal)  $\xrightarrow{\quad}$  output vector (signal)

$T$   $\uparrow$   
**transform = matrix**

- **Example: Discrete Fourier transform (DFT)**

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

# DSP Algorithms: Example 4-point DFT

- Algorithm = **sparse matrix factorization**
- Reduce computation cost from  $O(n^2)$  to  $O(n \log n)$
- For every transform there are **many** fast algorithms

$$\begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & i & -1 & -i \\ 1 & -1 & 1 & -1 \\ 1 & -i & -1 & i \end{pmatrix} x \rightarrow \begin{pmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & -1 & 0 \\ 0 & 1 & 0 & -1 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & i \end{pmatrix} \begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & -1 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} x$$

12 adds  
4 mults

4 adds

1 mult

(when multiplied with input vector x)

4 adds

$$\text{DFT}_4 \rightarrow (\text{DFT}_2 \otimes \text{I}_2) D (\text{I}_2 \otimes \text{DFT}_2) P$$

- SPIRAL generates the space of algorithms using **breakdown rules** in the domain-specific **Signal Processing Language (SPL)**

$$\text{DFT}_{mn} \rightarrow (\text{DFT}_m \otimes \text{I}_n) D (\text{I}_m \otimes \text{DFT}_n) P$$

# Some Transforms

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

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

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

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

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

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

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

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

***Spiral currently contains 45 transforms***

# Some Breakdown Rules

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

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

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

$$\begin{aligned} \mathbf{DCT-3}_n \rightarrow & (\mathbf{I}_m \oplus \mathbf{J}_m) \mathbf{L}_m^n (\mathbf{DCT-3}_m(1/4) \oplus \mathbf{DCT-3}_m(3/4)) \\ & \cdot (\mathbf{F}_2 \otimes \mathbf{I}_m) \begin{bmatrix} \mathbf{I}_m & 0 \oplus -\mathbf{J}_{m-1} \\ \frac{1}{\sqrt{2}}(\mathbf{I}_1 \oplus 2\mathbf{I}_m) \end{bmatrix}, \quad n = 2m \end{aligned}$$

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

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

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

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

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

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

**Base case rules**

**Spiral currently contains 165 rules**



# SPL (Signal Processing Language)

- SPL expresses transform algorithms as structured sparse matrix factorization

- Examples:

$$F_2 = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$$
$$A \otimes B = [a_{k,l} B]_{k,l}$$
$$A \oplus B = \begin{bmatrix} A & \\ & B \end{bmatrix}$$
$$I_n \otimes B = \begin{bmatrix} B & & \\ & \dots & \\ & & B \end{bmatrix}$$

- Kronecker product = loop (parallel, vector)

$$y = (I_n \otimes B_{m \times m})x$$



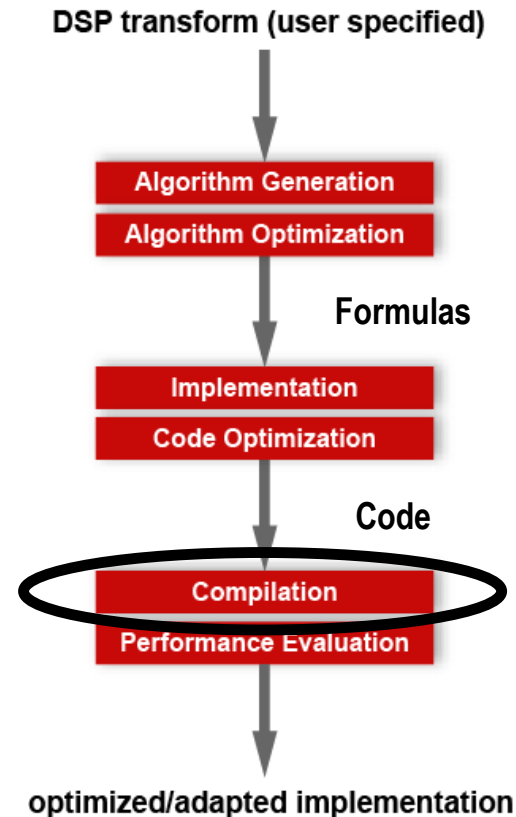
```
for i = 0:n-1
    y[im:im+m-1] = B · x[im:im+m-1]
endfor
```

# Formula Level Optimization: Idea

**Move optimizations to higher abstraction level:**  
Domain knowledge overcomes compiler limitations

**Traditionally optimizations by C/Fortran compilers**

**Formula level optimizations in Spiral:**  
Implemented through rewriting systems

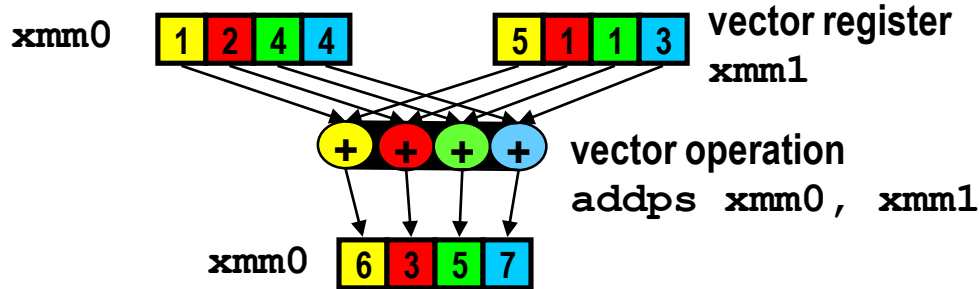


- Loop merging
- **Vectorization**
- Parallelization

# SIMD (Signal Instruction Multiple Data) Vector Instructions in a Nutshell

## ■ What are these instructions?

- Extension of the ISA. Data types and instructions for parallel computation on short (**2-way–16-way**) **vectors** of integers and floats



## ■ Problems:

- Not standardized
- Compiler vectorization limited
- Low-level issues (data alignment,...)
- Reordering data kills runtime

- Intel MMX
- AMD 3DNow!
- Intel SSE
- AMD Enhanced 3DNow!
- Motorola AltiVec
- AMD 3DNow! Professional
- Itanium
- Intel XScale
- Intel SSE2
- AMD-64
- IBM BlueGene/L PPC440FP2
- Intel Wireless MMX
- Intel SSE3
- ...

***One can easily slow down a program by vectorizing it***

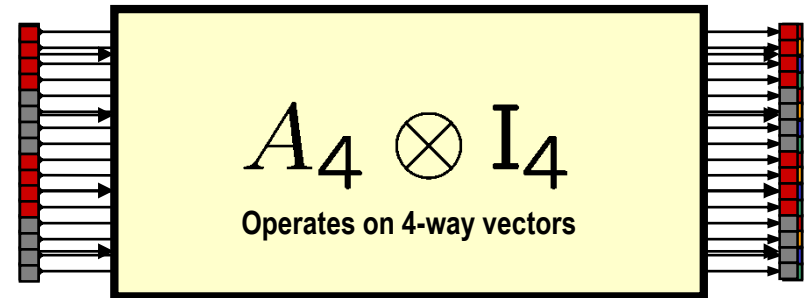
# Vectorization of Formulas by Rewriting

- Naturally vectorizable construct

Franchetti and Püschel (IPDPS 2002/2003)

$$y = (A \otimes I_\nu) x$$

vector length (any two-power)



- Rewriting rules to vectorize formulas

Introduces **data reorganization (permutations)**

$$I_n \otimes A^{k \times m} \rightarrow I_{n/\nu} \otimes \underline{L_\nu^{k\nu}} \left( \underline{A^{k \times m} \otimes I_\nu} \right) \underline{L_m^{m\nu}}$$

$$L_m^{m\nu} \rightarrow \left( \underline{I_{m/\nu} \otimes L_\nu^{\nu^2}} \right) \left( \underline{L_{m/\nu}^m \otimes I_\nu} \right)$$

vector construct  
further rewriting  
base case

**Definition:** *Vectorized formula* := **vector constructs** and **base cases**,  $A \not\subseteq B$ , and  $IA$  of vectorized formulas

# Example: DFT

$$\underbrace{(\overline{\text{DFT}}_{mn})}_{\text{vec}(\nu)} \rightarrow \underbrace{\left( (\text{DFT}_m \otimes \text{I}_n) \text{T}_n^{mn} (\text{I}_m \otimes \text{DFT}_n) \text{L}_m^{mn} \right)}_{\text{vec}(\nu)}$$

...

$$\rightarrow \underbrace{\left( \text{DFT}_m \otimes \text{I}_n \right)}_{\text{vec}(\nu)}^\nu \underbrace{\left( \text{T}_n^{mn} \right)}_{\text{vec}(\nu)}^\nu \underbrace{\left( \text{I}_m \otimes \text{DFT}_n \right)}_{\text{vec}(\nu)}^\nu \text{L}_m^{mn}$$

...

$$\rightarrow \left( \text{I}_{mn/\nu} \otimes \underbrace{\left( \text{L}_\nu^{2\nu} \right)}_{\text{vec}(\nu)} \right) \underbrace{\left( \text{DFT}_m \otimes \text{I}_{n/\nu} \otimes_\nu \text{I}_\nu \right)}_{\text{vec}(\nu)} \underbrace{\left( \text{T}_n^{mn} \right)}_{\text{vec}(\nu)}^\nu$$

$$\left( \text{I}_{m/\nu} \otimes \underbrace{\left( \text{L}_\nu^n \otimes_\nu \text{I}_\nu \right)}_{\text{vec}(\nu)} \right) \left( \text{I}_{n/\nu} \otimes \underbrace{\left( \text{L}_\nu^{2\nu} \otimes_\nu \text{I}_\nu \right)}_{\text{vec}(\nu)} \right) \left( \text{I}_2 \otimes \underbrace{\left( \text{L}_\nu^{\nu^2} \right)}_{\text{vec}(\nu)} \right) \underbrace{\left( \text{L}_\nu^{2\nu} \otimes_\nu \text{I}_\nu \right)}_{\text{vec}(\nu)} \underbrace{\left( \overline{\text{DFT}}_n \otimes_\nu \text{I}_\nu \right)}_{\text{vec}(\nu)}$$

$$\left( \underbrace{\left( \text{L}_m^{mn} \otimes \text{I}_2 \right)}_{\text{vec}(\nu)} \otimes_\nu \text{I}_\nu \right) \left( \text{I}_{mn/\nu} \otimes \underbrace{\left( \text{L}_\nu^{2\nu} \right)}_{\text{vec}(\nu)} \right)$$

vector constructs  
base cases

**Formula is vectorized w.r.t. Definition**

# Some Vectorization Rules

$$\underbrace{(\bar{A})}_{\text{vec}(\nu)} \rightarrow \overleftarrow{\underbrace{A}_{\text{vec}(\nu)}}^{\nu}$$

$$\overleftarrow{AB}^{\nu} \rightarrow \overleftarrow{A}^{\nu} \overleftarrow{B}^{\nu}$$

$$\overleftarrow{AB}^{\nu} \rightarrow \overleftarrow{A}^{\nu} \overleftarrow{B}^{\nu}$$

$$\underbrace{A \otimes I_m}_{\text{vec}(\nu)} \rightarrow (A \otimes I_{m/\nu}) \otimes_{\nu} I_{\nu}$$

$$\underbrace{(I_m \otimes A) L_m^{mn}}_{\text{vec}(\nu)} \rightarrow \left( I_{m/\nu} \otimes \underbrace{L_{\nu}^{n\nu}}_{\text{vec}(\nu)} (A \otimes_{\nu} I_{\nu}) \right) (L_{m/\nu}^{mn/\nu} \otimes_{\nu} I_{\nu})$$

$$\underbrace{L_{\nu}^{n\nu}}_{\text{vec}(\nu)} \rightarrow (L_{\nu}^n \otimes_{\nu} I_{\nu}) \left( I_{n/\nu} \otimes \underbrace{L_{\nu}^{\nu^2}}_{\text{vec}(\nu)} \right)$$

$$\overleftarrow{A \otimes_{\nu} I_{\nu}}^{\nu} \rightarrow (I_{n/\nu} \otimes \underbrace{L_{\nu}^{2\nu}}_{\text{vec}(\nu)}) (\bar{A} \otimes_{\nu} I_{\nu})$$

$$\overline{\underbrace{L_{\nu}^{\nu^2}}_{\text{vec}(\nu)}}^{\nu} \rightarrow (L_{\nu}^{2\nu} \otimes_{\nu} I_{\nu}) (I_2 \otimes \underbrace{L_{\nu}^{\nu^2}}_{\text{vec}(\nu)}) (L_2^{2\nu} \otimes_{\nu} I_{\nu})$$

$$\overleftarrow{I_m \otimes A}^{\nu} \rightarrow I_m \otimes \overleftarrow{A}^{\nu}$$

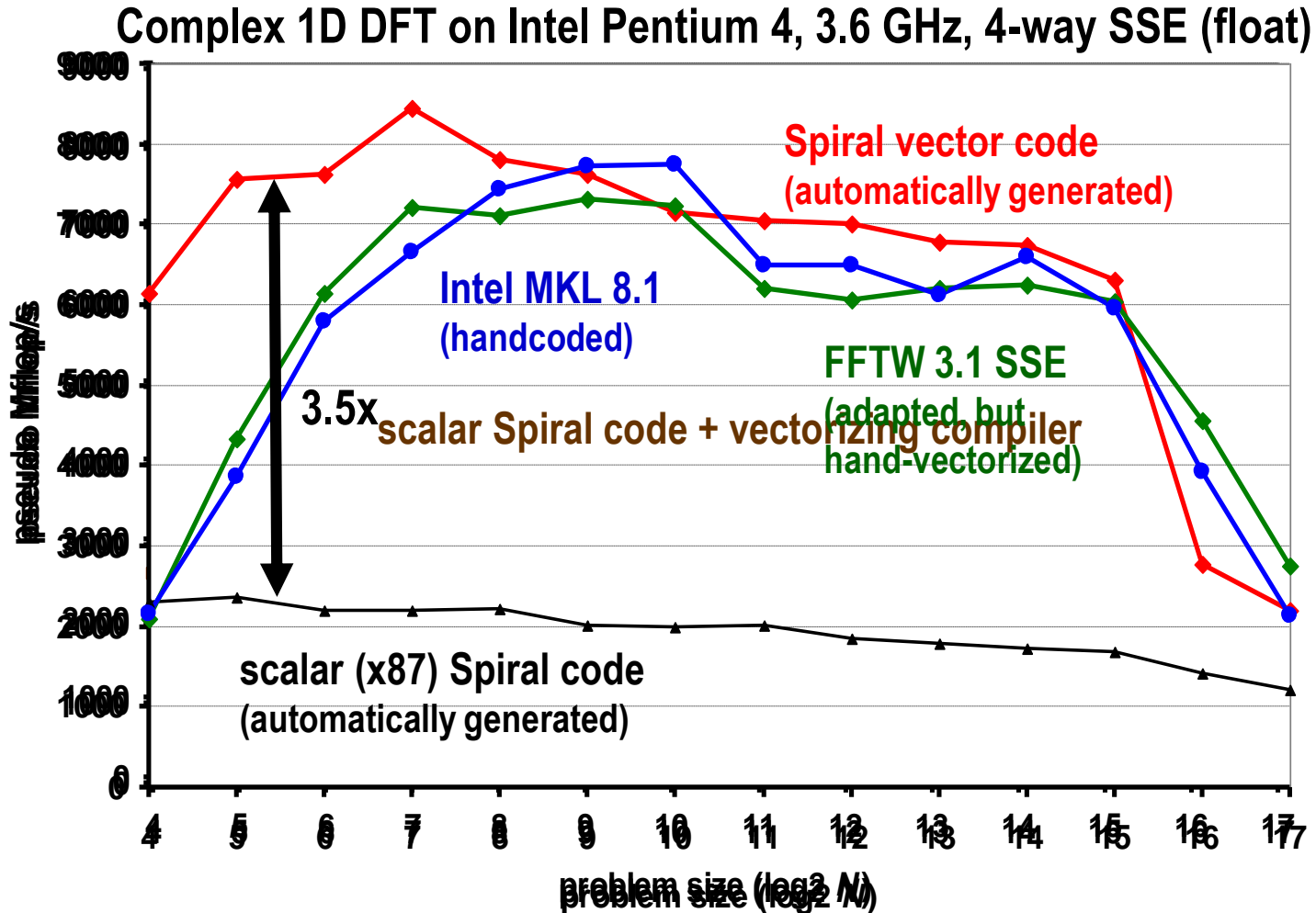
# Shared Memory Parallelization by Rewriting

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

**Load balanced, contiguous blocks**

**No false sharing (entire cache lines are swapped)**

# How Good is Our Generated Vector Code?

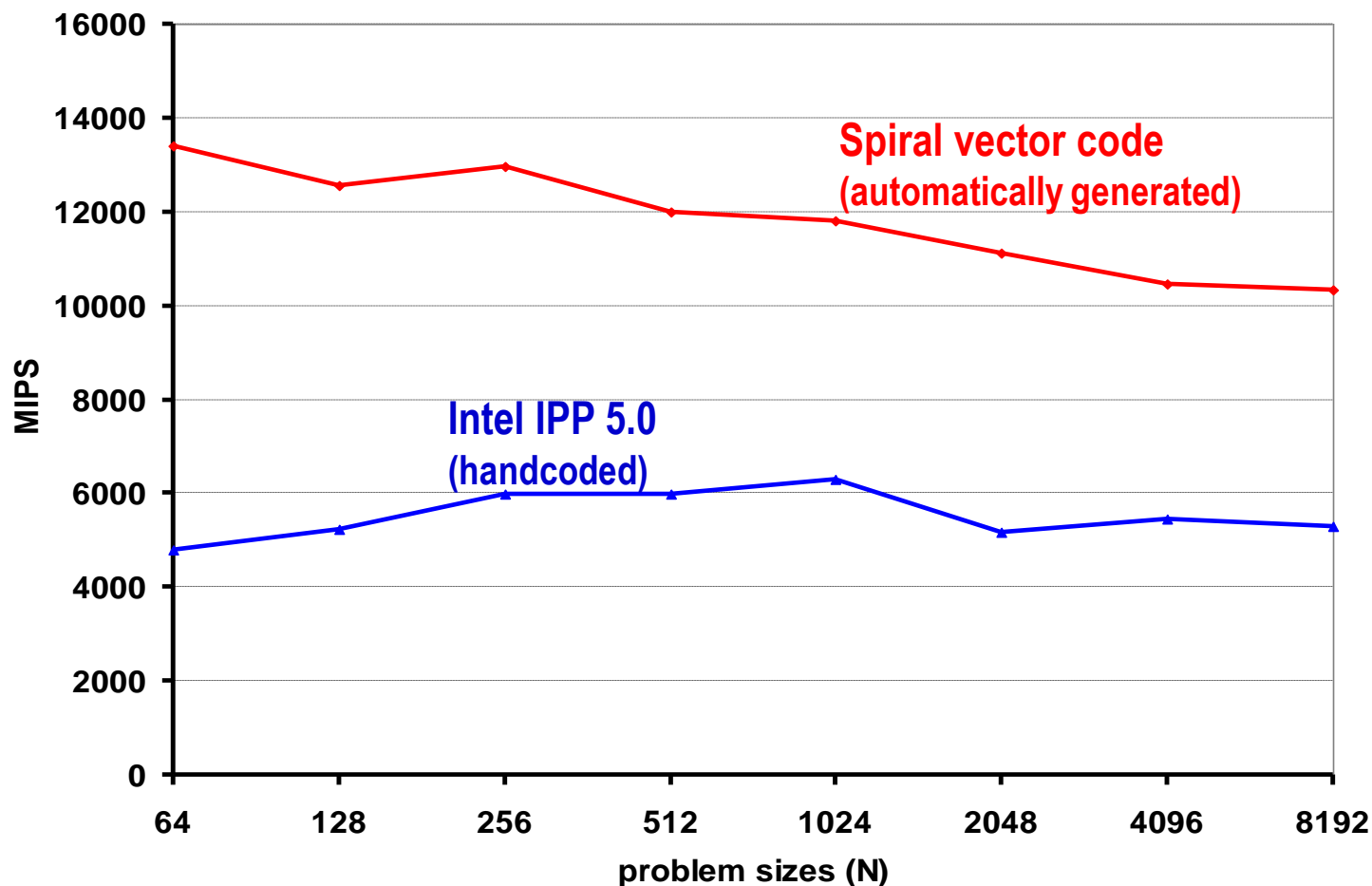


*Spiral generated code performs comparable to expertly hand-tuned code*



# What About 8-way Vector Code?

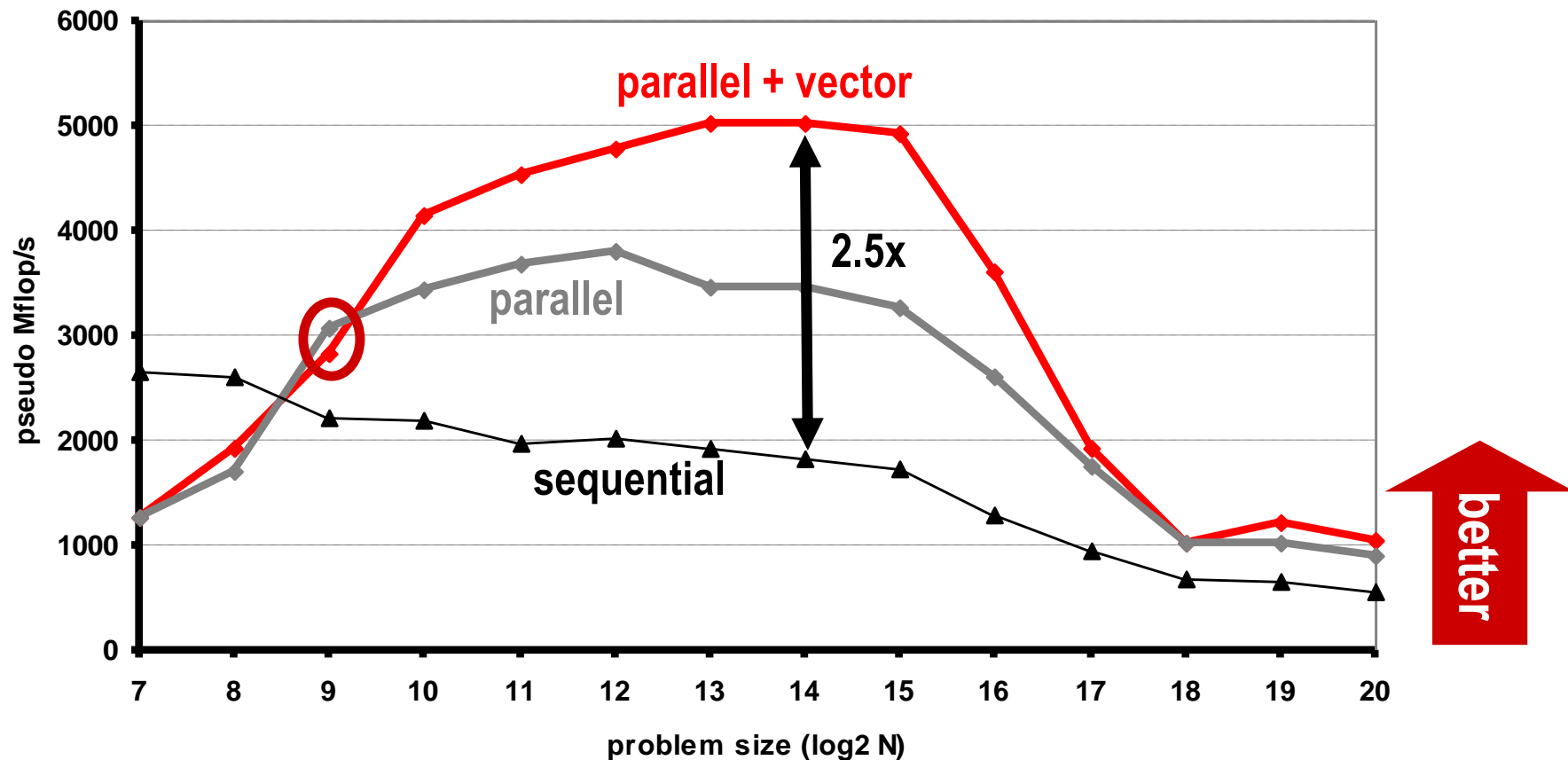
Complex 1D DFT on Intel Pentium 4, 3.6 GHz, 8-way SSE2 (16-bit int)



*Spiral generated code clearly outperforms expertly hand-tuned code*

# Combined Multicore and Vector Code

Pentium D 3.6 GHz (Dual Core, 2-way SIMD), double precision 1-D DFT

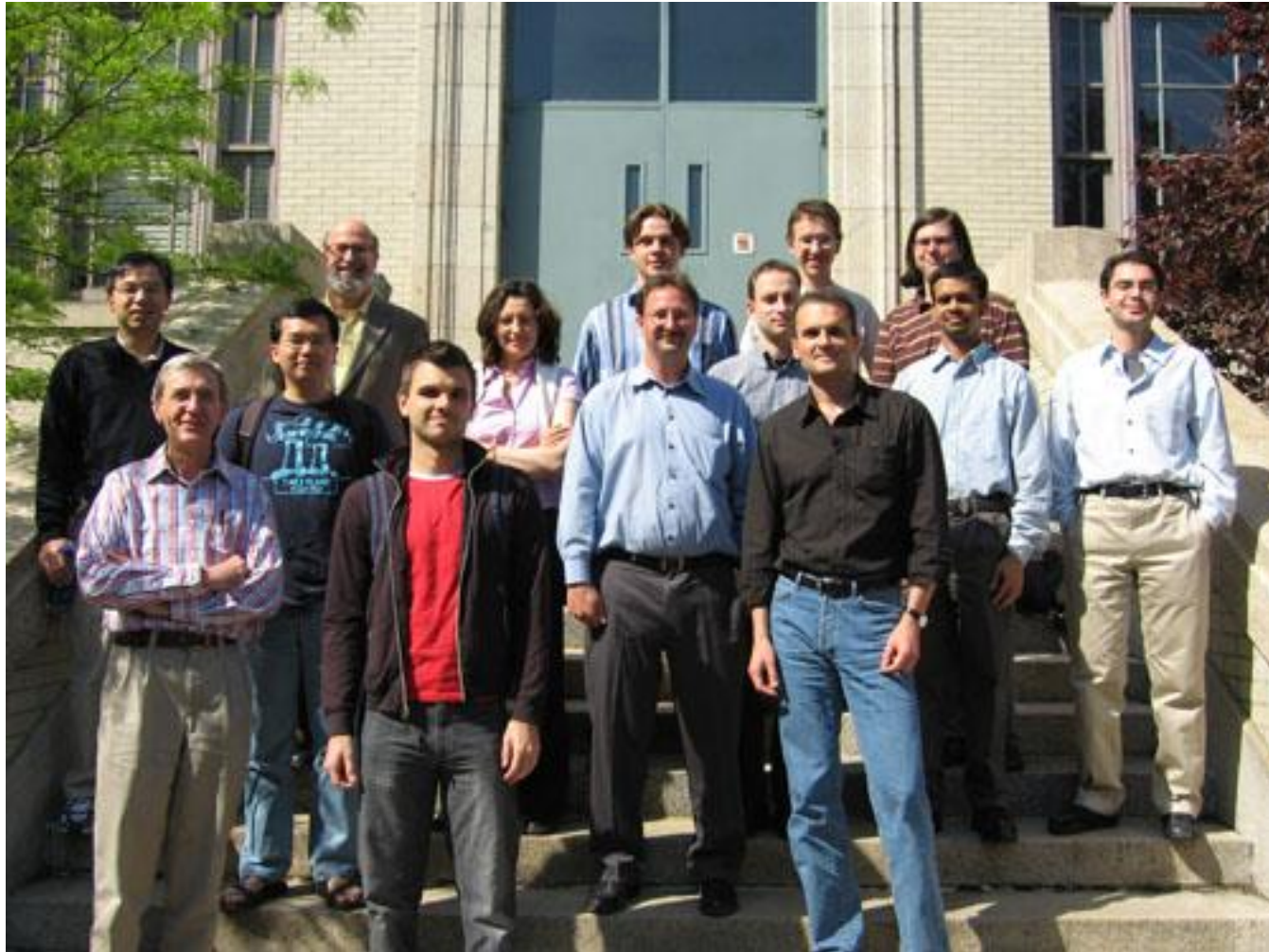


- *2.5x speed-up from parallel + vector*
- *Parallelization speed-up for small problems*

# Summary

- **Parallelization and vectorization in Spiral**
  - Entirely automatic
  - Principled approach
  - Rewriting system
  - Generated code is very fast
- **Works for other hardware as well**
  - **Distributed memory: MPI**  
with C.W. Ueberhuber, A. Bonelli, and J. Lorenz, Vienna University of Technology
  - **Hardware: FPGAs**  
with J.C. Hoe and Peter Milder, Carnegie Mellon University

# (Part of the) Spiral Team



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