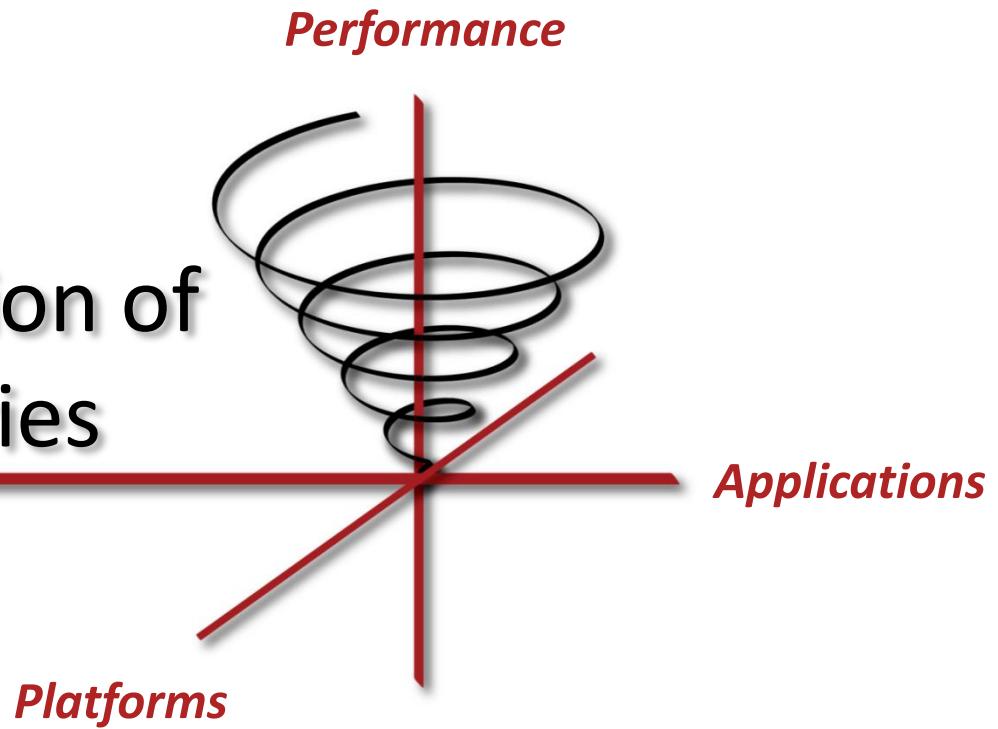


Spiral

Computer Generation of
Performance Libraries

José M. F. Moura
Markus Püschel
Franz Franchetti
& the Spiral Team

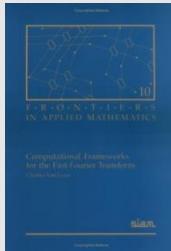
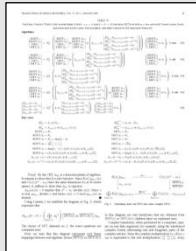
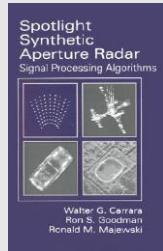


Carnegie
Mellon



What is Spiral?

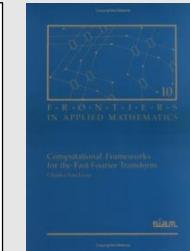
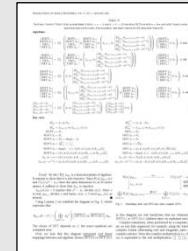
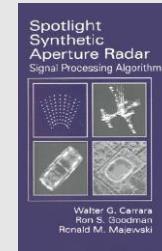
Traditionally



High performance library
optimized for given platform

*Comparable
performance*

Spiral Approach



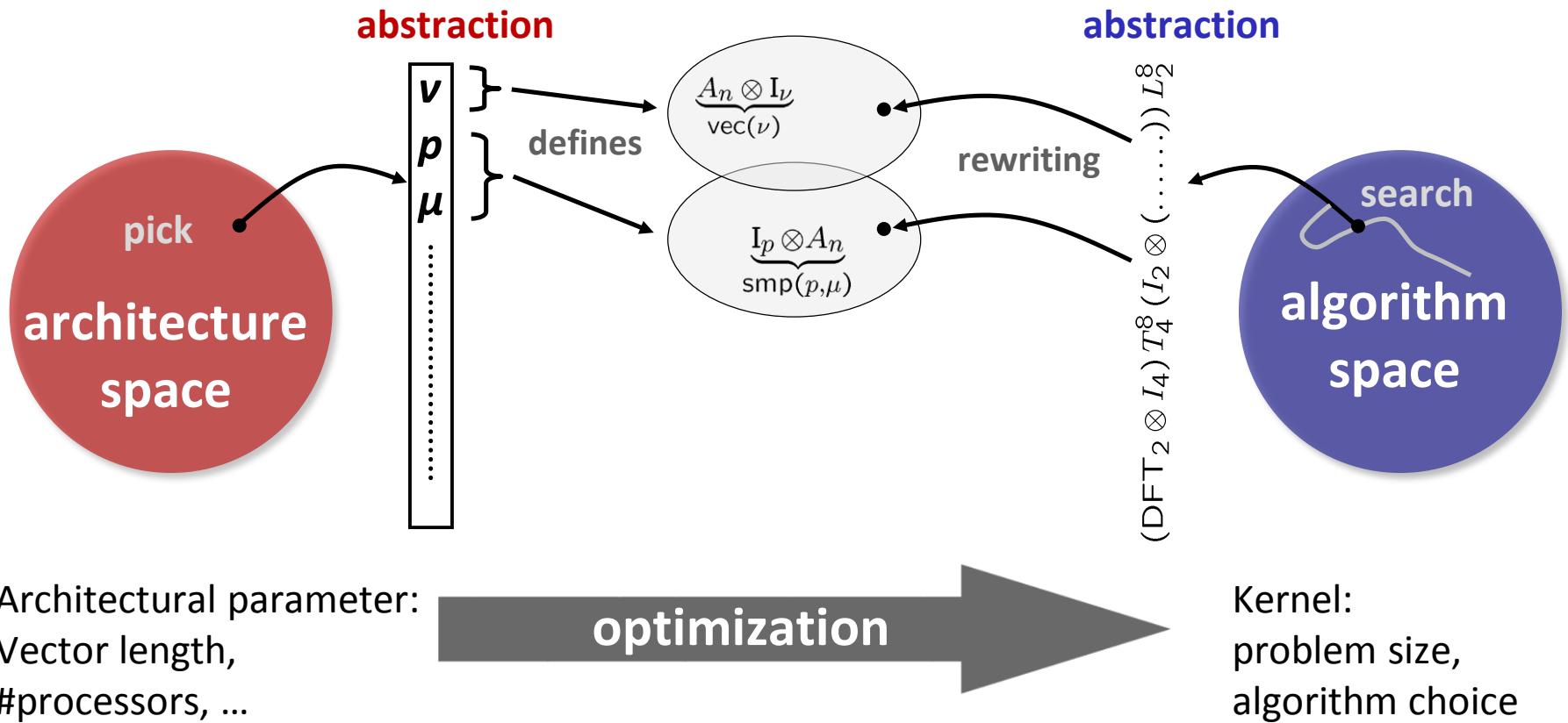
Spiral



High performance library
optimized for given platform

Main Idea: Program Generation

Model: common abstraction
= spaces of matching formulas

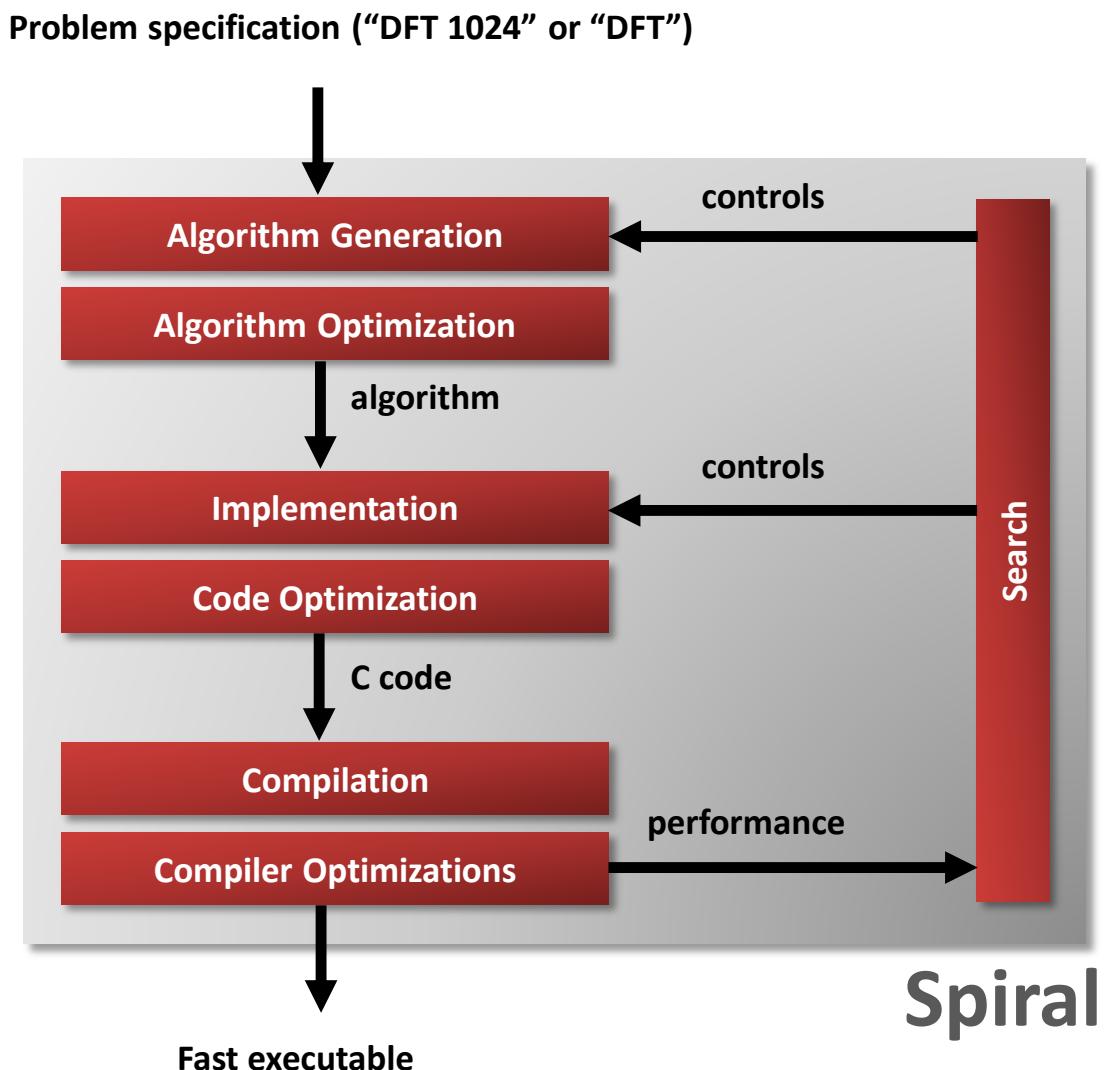


How Spiral Works

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

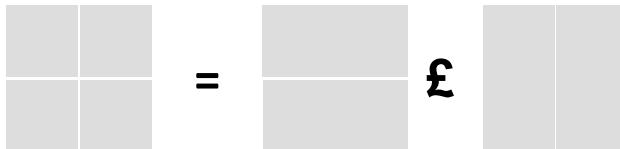


Algorithms: Rules in Domain Specific Language

Linear Transforms

$$\begin{aligned}
DFT_n &\rightarrow (DFT_k \otimes I_m) T_m^n (I_k \otimes DFT_m) L_k^n, \quad n = km \\
DFT_n &\rightarrow P_n(DFT_k \otimes DFT_m) Q_n, \quad n = km, \quad \gcd(k, m) = 1 \\
DFT_p &\rightarrow R_p^T (I_1 \oplus DFT_{p-1}) D_p (I_1 \oplus DFT_{p-1}) R_p, \quad p \text{ prime} \\
DCT-3_n &\rightarrow (I_m \oplus J_m) L_m^n (DCT-3_m(1/4) \oplus DCT-3_m(3/4)) \\
&\quad \cdot (F_2 \otimes I_m) \begin{bmatrix} I_m & 0 \oplus -J_{m-1} \\ 0 & \frac{1}{\sqrt{2}}(I_1 \oplus 2I_m) \end{bmatrix}, \quad n = 2m \\
DCT-4_n &\rightarrow S_n DCT-2_n \text{diag}_{0 \leq k < n} (1/(2 \cos((2k+1)\pi/4n))) \\
IMDCT_{2m} &\rightarrow (J_m \oplus I_m \oplus I_m \oplus J_m) \left(\left(\begin{bmatrix} 1 \\ -1 \end{bmatrix} \otimes I_m \right) \oplus \left(\begin{bmatrix} -1 \\ -1 \end{bmatrix} \otimes I_m \right) \right) J_{2m} DCT-4_{2m} \\
WHT_{2^k} &\rightarrow \prod_{i=1}^t (I_{2^{k_1+\dots+k_{i-1}}} \otimes WHT_{2^{k_i}} \otimes I_{2^{k_{i+1}+\dots+k_t}}), \quad k = k_1 + \dots + k_t \\
DFT_2 &\rightarrow F_2 \\
DCT-2_2 &\rightarrow \text{diag}(1, 1/\sqrt{2}) F_2 \\
DCT-4_2 &\rightarrow J_2 R_{13\pi/8}
\end{aligned}$$

Matrix-Matrix Multiplication



$$MMM_{1,1,1} \rightarrow (\cdot)_1$$

$$MMM_{m,n,k} \rightarrow (\otimes)_{m/m_b \times 1} \otimes MMM_{m_b, n, k}$$

$$MMM_{m,n,k} \rightarrow MMM_{m,nb,k} \otimes (\otimes)_{1 \times n/nb}$$

$$MMM_{m,n,k} \rightarrow ((\Sigma_{k/k_b} \circ (\cdot)_{k/k_b}) \otimes MMM_{m,n,k_b}) \circ ((L_{k/k_b}^{mk/k_b} \otimes I_{k_b}) \times I_{kn})$$

$$MMM_{m,n,k} \rightarrow (L_m^{mn/n_b} \otimes I_{n_b}) \circ ((\otimes)_{1 \times n/n_b} \otimes MMM_{m,n_b,k}) \circ (I_{km} \times (L_{n/n_b}^{kn/n_b} \otimes I_{n_b}))$$

Viterbi Decoding



$$\begin{aligned}
\text{Vit} &\rightarrow \underbrace{\left(\prod_{v \in v} (L \times I) \circ (I \otimes C) \right)}_{\text{vec}(v)} \circ Id \\
&\rightarrow \left(\prod_{v \in v} (L \times I) \circ (I \otimes C) \right) \circ Id \\
\mathfrak{L} &\rightarrow \left(\prod_{v \in v} (L \otimes I_v \times I) \circ (I \otimes C \otimes I_v) \circ (\bar{L} \times I) \right) \circ Id \\
&\rightarrow \prod_{v \in v} (L \otimes I_v \times I) \circ (I \otimes (B \otimes I_v)) \circ (\bar{L} \times I)
\end{aligned}$$

Synthetic Aperture Radar (SAR)



$$\begin{aligned}
SAR_{k \times m \rightarrow n \times n} &\rightarrow DFT_{n \times n} \circ \text{Interp}_{k \times m \rightarrow n \times n} \\
DFT_{n \times n} &\rightarrow (DFT_n \otimes I_n) \circ (I_n \otimes DFT_n) \\
\text{Interp}_{k \times m \rightarrow n \times n} &\rightarrow (\text{Interp}_{k \rightarrow n} \otimes_i I_n) \circ (I_k \otimes_i \text{Interp}_{m \rightarrow n}) \\
\text{Interp}_{r \rightarrow s} &\rightarrow \left(\bigoplus_{i=0}^{n-2} \text{InterpSeg}_k \right) \oplus \text{InterpSegPruned}_{k,\ell} \\
\text{InterpSeg}_k &\rightarrow G_f^{u \cdot n \rightarrow k} \circ \text{iPrunedDFT}_{n \rightarrow u \cdot n} \circ \left(\frac{1}{n} \right) \circ DFT_n
\end{aligned}$$

One Approach for all Types of Parallelism

- **Multithreading (Multicore)**
- **Vector SIMD (SSE, VMX/Altivec,...)**
- **Message Passing (Clusters, MPP)**
- **Streaming/multibuffering (Cell)**
- **Graphics Processors (GPUs)**
- **Gate-level parallelism (FPGA)**
- **HW/SW partitioning (CPU + FPGA)**

$$\mathbf{I}_p \otimes_{\parallel} A_{\mu n}, \quad \mathbf{L}_m^{mn} \bar{\otimes} \mathbf{I}_{\mu}$$

$$A \hat{\otimes} \mathbf{I}_{\nu} \quad \underbrace{\mathbf{L}_2^{2\nu}}_{\text{isa}}, \quad \underbrace{\mathbf{L}_{\nu}^{2\nu}}_{\text{isa}}, \quad \underbrace{\mathbf{L}_{\nu}^{\nu^2}}_{\text{isa}}$$

$$\mathbf{I}_p \otimes_{\parallel} A_n, \quad \underbrace{\mathbf{L}_p^{p^2} \bar{\otimes} \mathbf{I}_{n/p^2}}_{\text{all-to-all}}$$

$$\mathbf{I}_n \otimes_2 A_{\mu n}, \quad \mathbf{L}_m^{mn} \bar{\otimes} \mathbf{I}_{\mu}$$

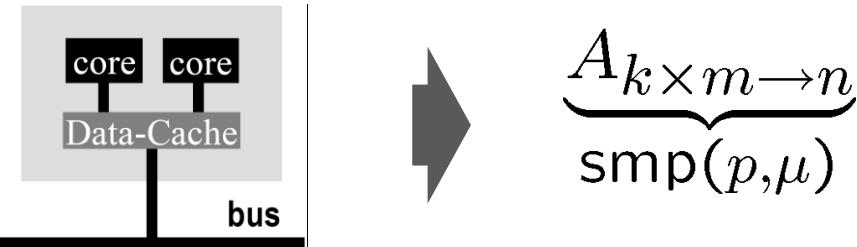
$$\prod_{i=0}^{n-1} A_i, \quad A_n \hat{\otimes} \mathbf{I}_w, \quad P_n \otimes Q_w$$

$$\prod_{i=0}^{n-1} \overset{\text{ir}}{A}, \quad \mathbf{I}_s \tilde{\otimes} A, \quad \underbrace{\mathbf{L}_n^m}_{\text{bram}}$$

$$\underbrace{A_1}_{\text{fpga}}, \quad \underbrace{A_2}_{\text{fpga}}, \quad \underbrace{A_3}_{\text{fpga}}, \quad \underbrace{A_4}_{\text{fpga}}$$

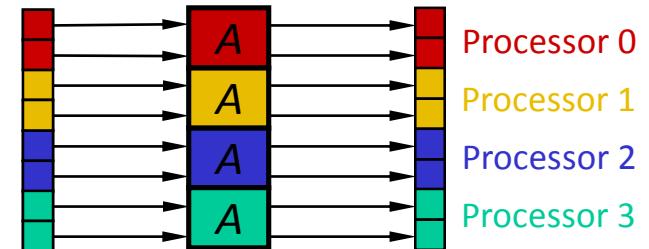
Example: Code Generation for Multicore CPUs

- Hardware abstraction: shared cache with cache lines



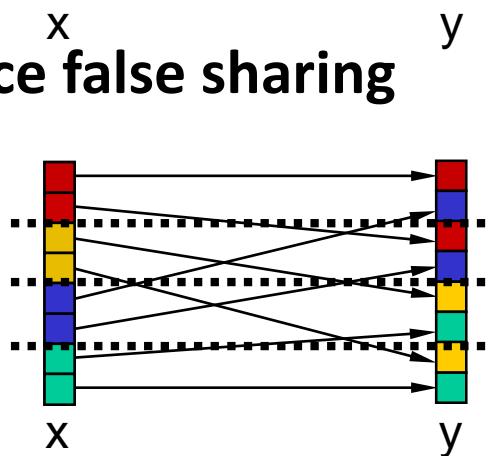
- Tensor product: embarrassingly parallel operator

$$y = (\mathbf{I}_p \otimes A)(x)$$



- Permutation: problematic; may produce false sharing

$$y = L_4^8(x)$$



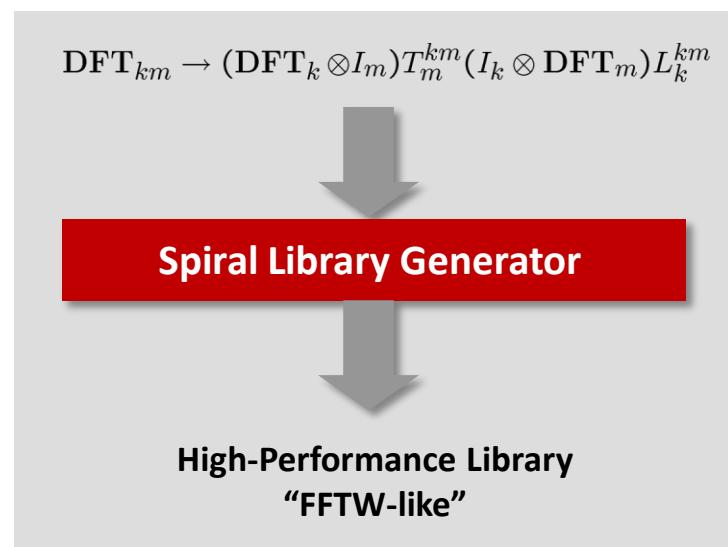
Spiral: Meta-Tool to Autotuning Libraries

Input:

- **Transform:** DFT_n
- **Algorithms:** $\text{DFT}_{km} \rightarrow (\text{DFT}_k \otimes I_m) T_m^{km} (I_k \otimes \text{DFT}_m) L_k^{km}$
 $\text{DFT}_2 \rightarrow \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$
- **Vectorization:** 2-way SSE
- **Threading:** Yes

Output:

- Optimized library (10,000 lines of C++)
- For general input size
(**not** collection of fixed sizes)
- Vectorized
- Multithreaded
- With runtime adaptation mechanism
- Performance competitive with hand-written code



Verification and Testing

■ Verify algorithms symbolically

$$\text{DFT}_4 = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & j & -1 & -j \\ 1 & -1 & 1 & -1 \\ 1 & -j & -1 & j \end{bmatrix} \quad = ? \quad (\text{DFT}_2 \otimes \text{I}_2) \text{T}_2^4 (\text{I}_2 \otimes \text{DFT}_2) \text{L}_2^4 = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & -1 & 0 \\ 0 & 1 & 0 & -1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & -1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & j \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

■ Check rules through verification of instances

$$\text{I}_m \otimes A_n \rightarrow \text{L}_m^{mn} (A_n \otimes \text{I}_m) \text{L}_n^{mn}$$

$$\text{I}_2 \otimes \text{DFT}_2 = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & -1 \end{bmatrix} \quad = ? \quad \text{L}_2^4 (\text{DFT}_2 \otimes \text{I}_2) \text{L}_2^4 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & -1 & 0 \\ 0 & 1 & 0 & -1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

■ Check code empirically

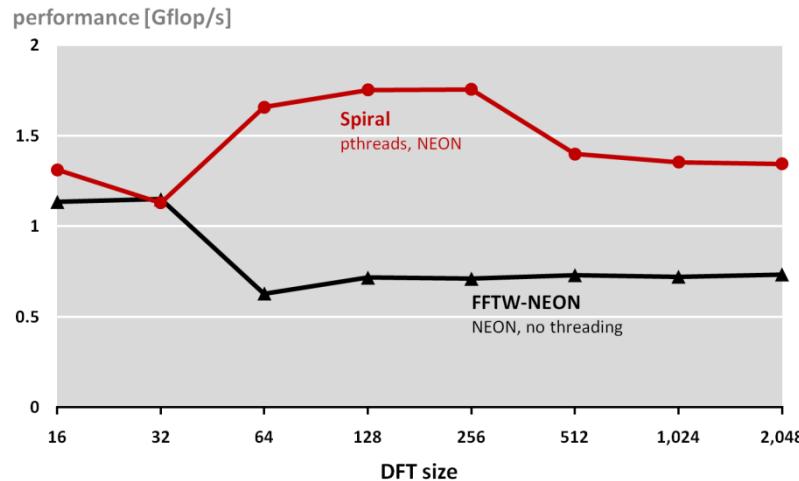
$$\begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & j & -1 & -j \\ 1 & -1 & 1 & -1 \\ 1 & -j & -1 & j \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad = ? \quad \text{DFT4}([0, 1, 0, 0])$$

$$\text{DFT4}([0.1, 1.77, 2.28, -55.3]) \quad = ? \quad \text{DFT4_rnd}([0.1, 1.77, 2.28, -55.3]))$$

Range: Cell Phone To Supercomputer

DFT on Samsung Galaxy S II

Dual-core 1.2 GHz Cortex-A9 with NEON ISA



Samsung i9100 Galaxy S II

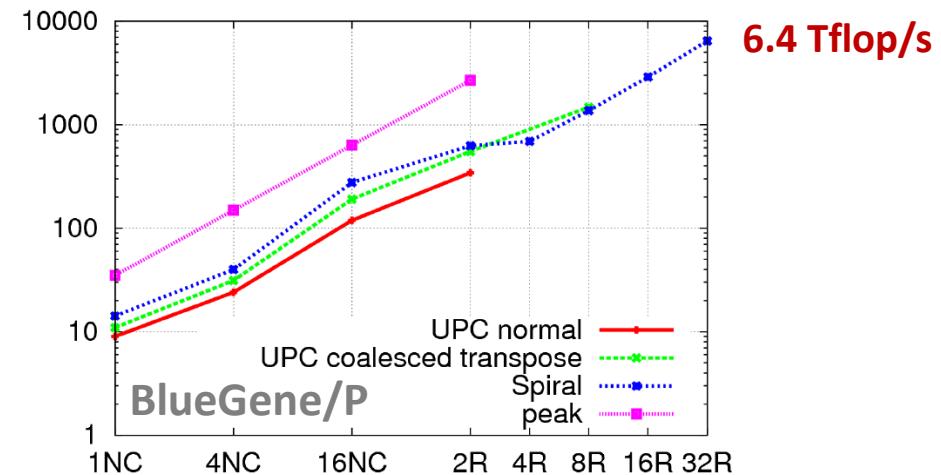
Dual-core ARM at 1.2GHz with NEON ISA

SIMD vectorization + multi-threading



Global FFT (1D FFT, HPC Challenge)

performance [Gflop/s]



6.4 Tflop/s

BlueGene/P at Argonne National Laboratory

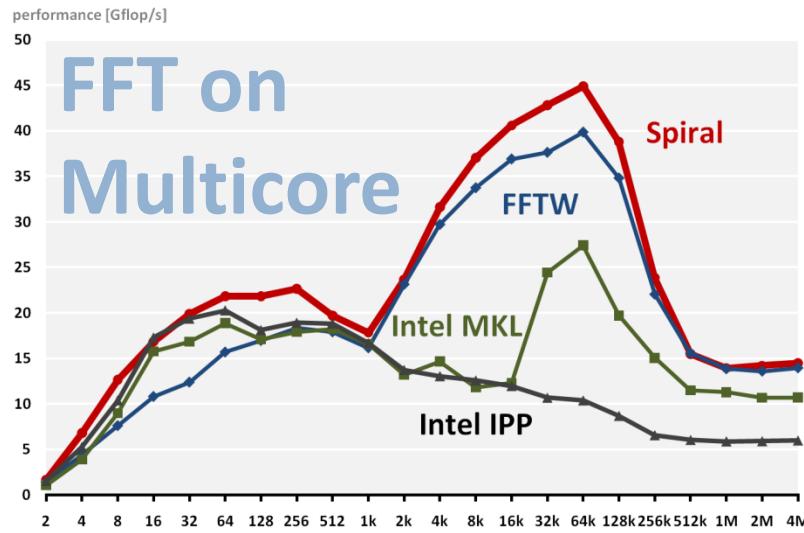
128k cores (quad-core CPUs) at 850 MHz

SIMD vectorization + multi-threading + MPI

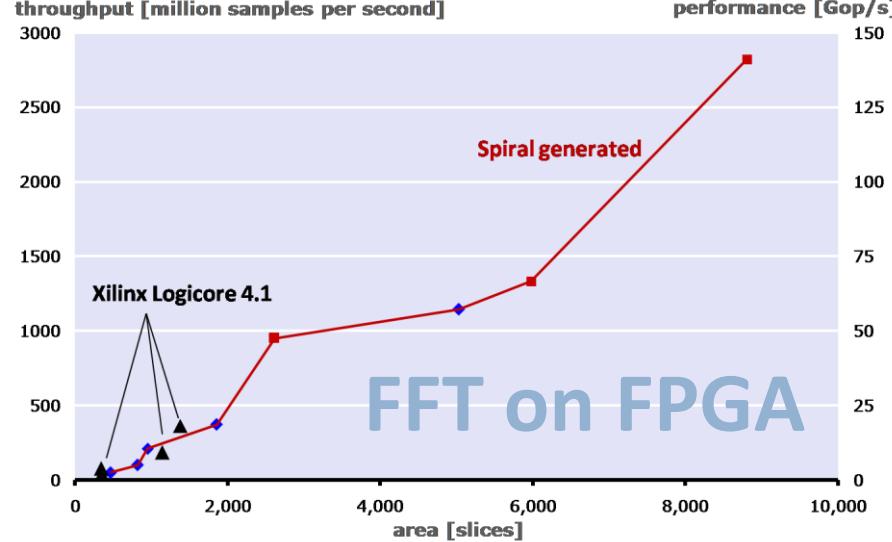


More Results: Spiral Outperforms Humans

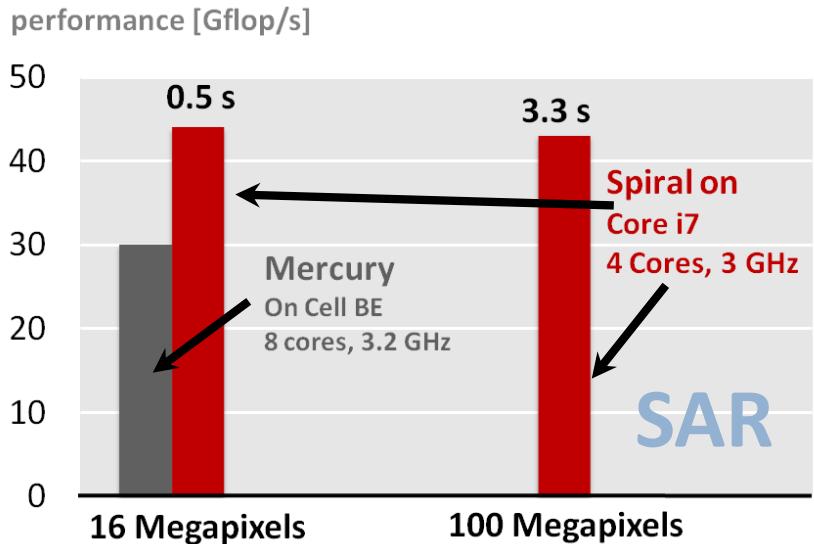
1D DFT on 3.3 GHz Sandy Bridge (4 Cores, AVX)



DFT 1024 (16 bit fixed point) on Xilinx Virtex-5 FPGA

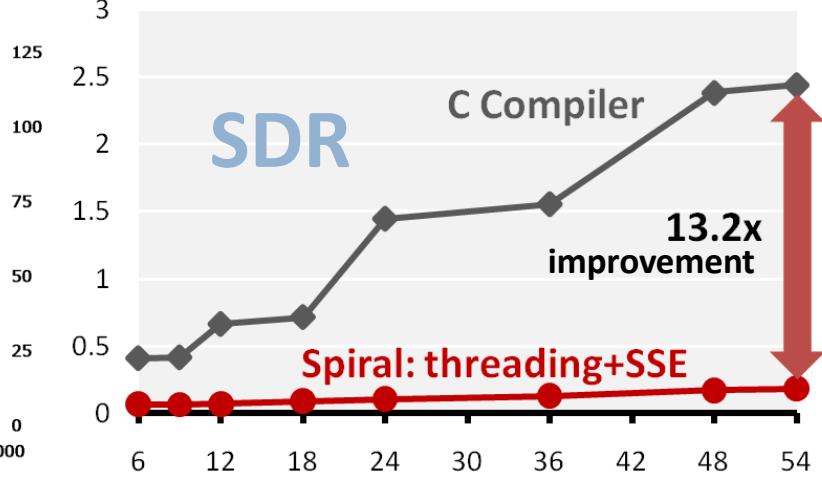


SAR Image Formation on Intel platforms



WiFi transmitter Core i7

Run time per symbol [micro seconds] vs. Data rate [Mbit/s]



More Information:

www.spiral.net

www.spiralgen.com