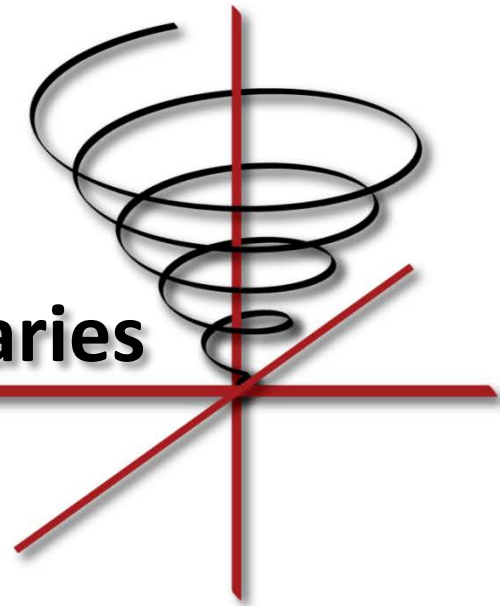


Spiral:

Automatic Generation of Industry Strength Performance Libraries



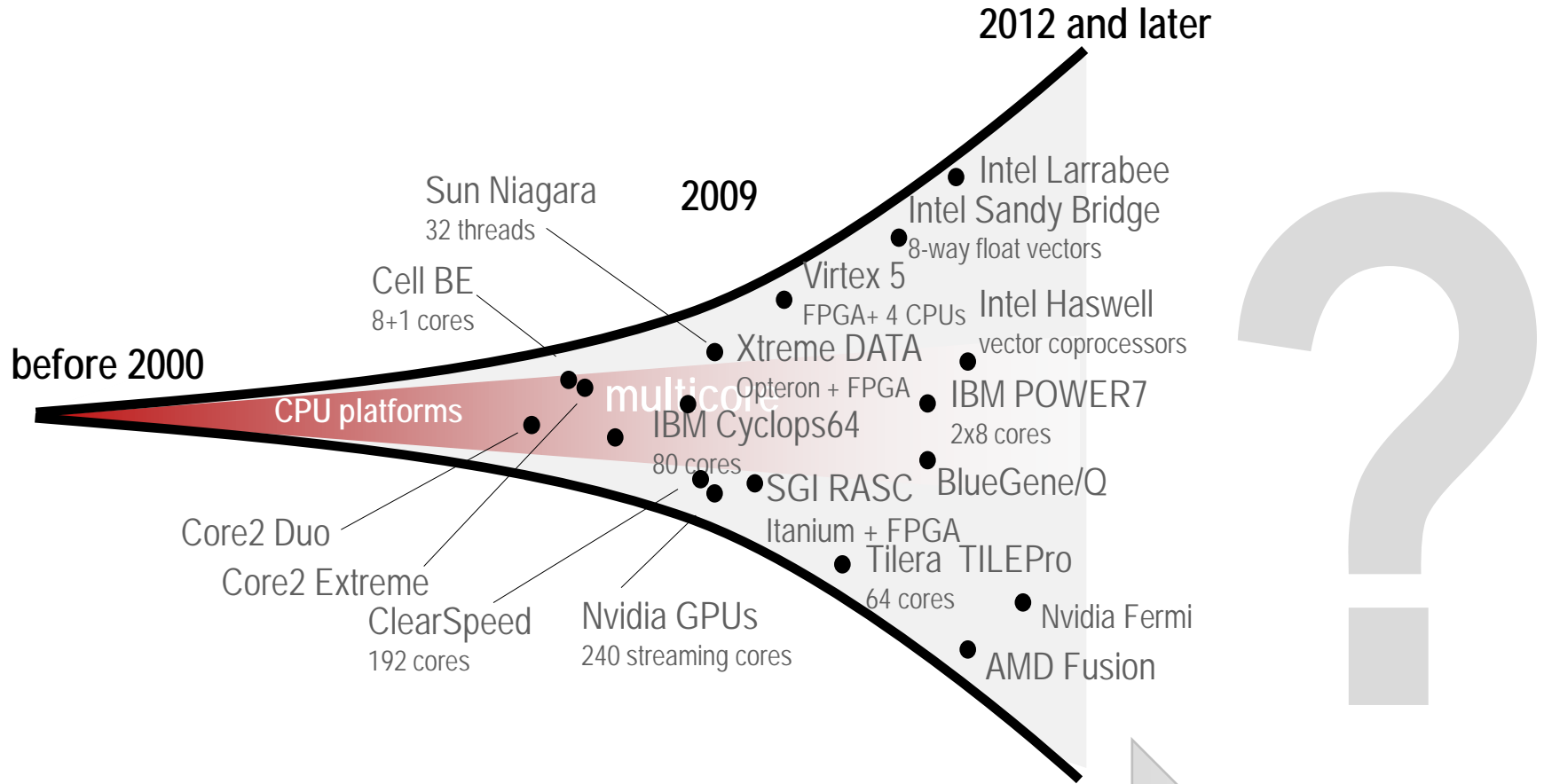
Franz Franchetti

Carnegie Mellon University
www.ece.cmu.edu/~franzf

CTO and Co-Founder, SpiralGen
www.spiralgen.com

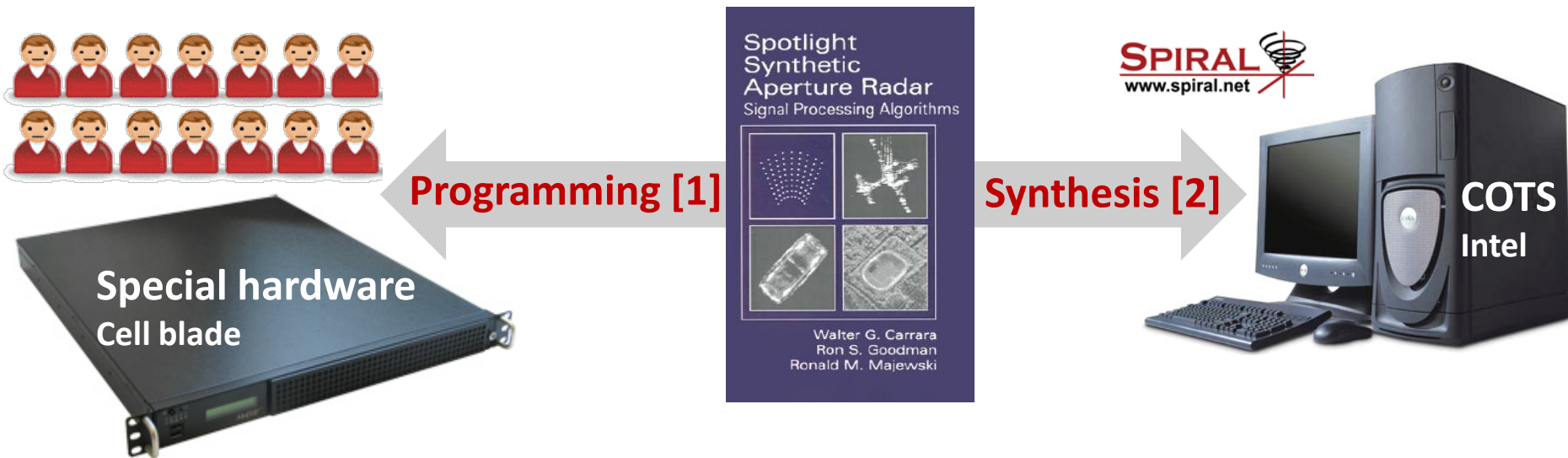
This work was supported by
DARPA DESA program, NSF, ONR, Mercury Inc., Intel, and Nvidia

The Future is Parallel and Heterogeneous



Programmability?
Performance portability?
Rapid prototyping?

Spiral: Computer Writes Best SAR Code



Result

Same performance, $1/10^{\text{th}}$ human effort, non-expert user

Key ideas

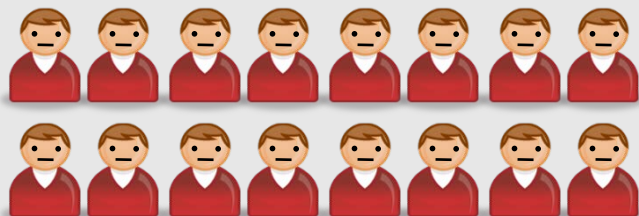
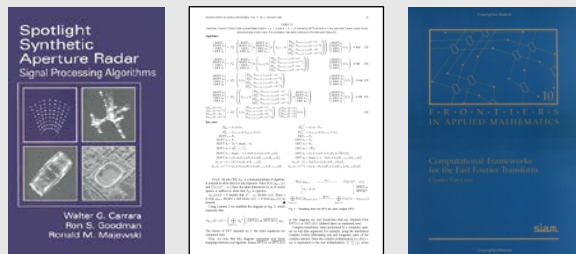
restrict domain, use mathematics, program synthesis

[1] Rudin, J., **Implementation of Polar Format SAR Image Formation on the IBM Cell Broadband Engine**, in Proceedings High Performance Embedded Computing (HPEC), 2007. *Best Paper Award*.

[2] D. McFarlin, F. Franchetti, M. Püschel, and J. M. F. Moura: **High Performance Synthetic Aperture Radar Image Formation On Commodity Multicore Architectures**. in Proceedings SPIE, 2009.

What is Spiral?

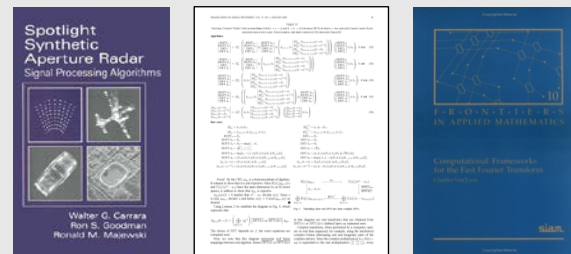
Traditionally



High performance library
optimized for given platform

*Comparable
performance*

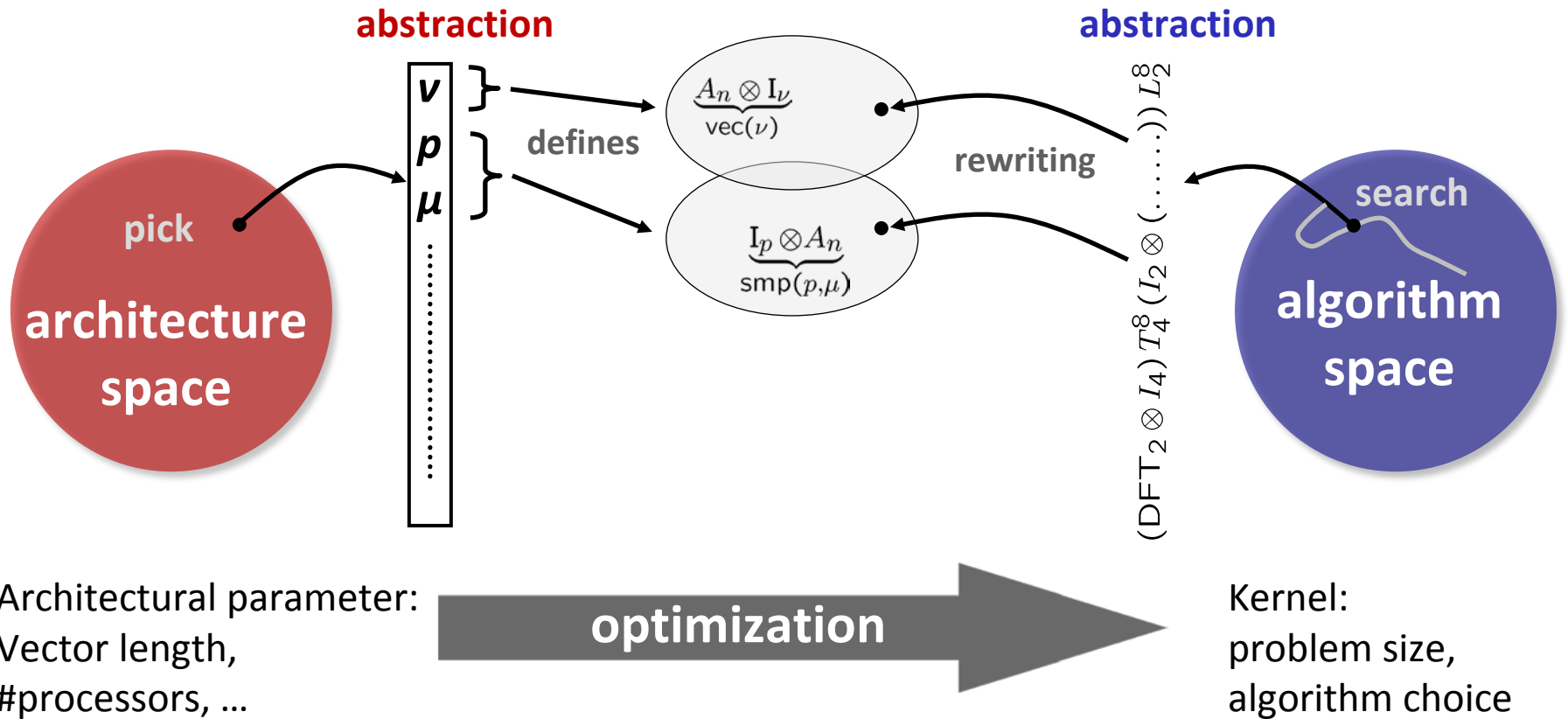
Spiral Approach



High performance library
optimized for given platform

Spiral's Domain-Specific Program Synthesis

Model: common abstraction
 = spaces of matching formulas



Related Work

Synthesis from Domain Math

- **Spiral**
Signal and image processing, SDR
- **Tensor Contraction Engine**
Quantum Chemistry Code Synthesizer
- **FLAME**
Numerical linear algebra (LAPACK)

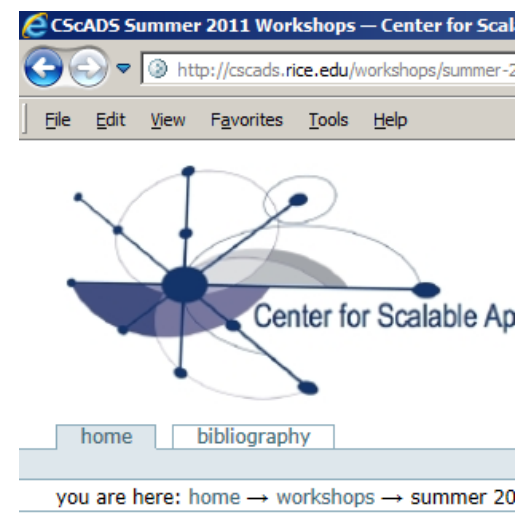
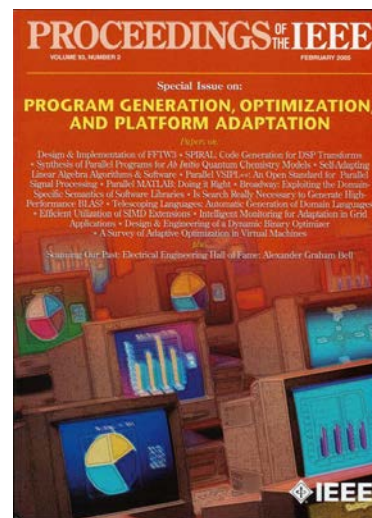
Autotuning Numerical Libraries

- **ATLAS**
BLAS generator
- **FFTW**
kernel generator
- **Vendor math libraries**
Code generation scripts

Compiler-Based Autotuning

- **Polyhedral framework**
IBM XL, Pluto, CHiLL
- **Transformation prescription**
CHiLL, POET
- **Profile guided optimization**
Intel C, IBM XL

Autotuning Primer



Organization

- **Spiral overview**
- Validation and Verification
- Results
- Concluding remarks

M. Püschel, F. Franchetti, Y. Voronenko: **Spiral**. Encyclopedia of Parallel Computing, D. A. Padua (Editor), 2011.

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. Special issue, Proceedings of the IEEE 93(2), 2005.

Spiral's Origin: Linear Transforms

- **Transform = Matrix-vector multiplication**

Example: Discrete Fourier transform (DFT)

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

- **Fast algorithm = sparse matrix factorization = SPL formula**

Example: Cooley-Tukey FFT algorithm

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

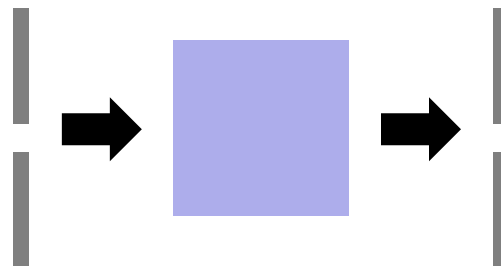
$$\text{DFT}_4 = (\text{DFT}_2 \otimes \text{I}_2) \text{T}_2^4 (\text{I}_2 \otimes \text{DFT}_2) \text{L}_2^4$$

Beyond Transforms: General Operators

- 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

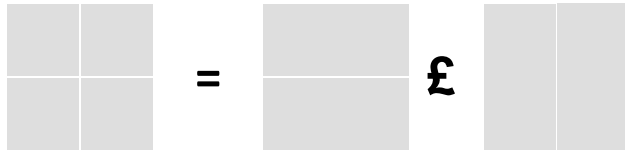


Expressing Kernels as Operator Formulas

Linear Transforms

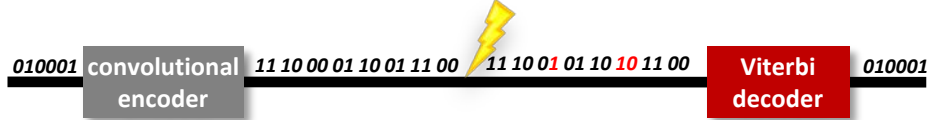
$$\begin{aligned} \text{DFT}_n &\rightarrow (\text{DFT}_k \otimes \text{I}_m) \text{T}_m^n (\text{I}_k \otimes \text{DFT}_m) \text{L}_k^n, \quad n = km \\ \text{DFT}_n &\rightarrow P_n (\text{DFT}_k \otimes \text{DFT}_m) Q_n, \quad n = km, \text{ gcd}(k, m) = 1 \\ \text{DFT}_p &\rightarrow R_p^T (\text{I}_1 \oplus \text{DFT}_{p-1}) D_p (\text{I}_1 \oplus \text{DFT}_{p-1}) R_p, \quad p \text{ prime} \\ \text{DCT-3}_n &\rightarrow (\text{I}_m \oplus \text{J}_m) \text{L}_m^n (\text{DCT-3}_m(1/4) \oplus \text{DCT-3}_m(3/4)) \\ &\quad \cdot (\text{F}_2 \otimes \text{I}_m) \begin{bmatrix} \text{I}_m & 0 \oplus -\text{J}_{m-1} \\ \frac{1}{\sqrt{2}}(\text{I}_1 \oplus 2\text{I}_m) \end{bmatrix}, \quad n = 2m \\ \text{DCT-4}_n &\rightarrow S_n \text{DCT-2}_n \text{diag}_{0 \leq k < n} (1/(2 \cos((2k+1)\pi/4n))) \\ \text{IMDCT}_{2m} &\rightarrow (\text{J}_m \oplus \text{I}_m \oplus \text{I}_m \oplus \text{J}_m) \left(\left(\begin{bmatrix} 1 \\ -1 \end{bmatrix} \otimes \text{I}_m \right) \oplus \left(\begin{bmatrix} -1 \\ -1 \end{bmatrix} \otimes \text{I}_m \right) \right) \text{J}_{2m} \text{DCT-4}_{2m} \\ \text{WHT}_{2^k} &\rightarrow \prod_{i=1}^t (\text{I}_{2^{k_1+\dots+k_{i-1}}} \otimes \text{WHT}_{2^{k_i}} \otimes \text{I}_{2^{k_{i+1}+\dots+k_t}}), \quad k = k_1 + \dots + k_t \\ \text{DFT}_2 &\rightarrow \text{F}_2 \\ \text{DCT-2}_2 &\rightarrow \text{diag}(1, 1/\sqrt{2}) \text{F}_2 \\ \text{DCT-4}_2 &\rightarrow \text{J}_2 \text{R}_{13\pi/8} \end{aligned}$$

Matrix-Matrix Multiplication



$$\begin{aligned} \text{MMM}_{1,1,1} &\rightarrow (\cdot)_1 \\ \text{MMM}_{m,n,k} &\rightarrow (\otimes)_{m/m_b \times 1} \otimes \text{MMM}_{m_b,n,k} \\ \text{MMM}_{m,n,k} &\rightarrow \text{MMM}_{m,n_b,k} \otimes (\otimes)_{1 \times n/n_b} \\ \text{MMM}_{m,n,k} &\rightarrow ((\Sigma_{k/k_b} \circ (\cdot)_{k/k_b}) \otimes \text{MMM}_{m,n,k_b}) \circ \\ &\quad ((L_{k/k_b}^{nk/k_b} \otimes \text{I}_{k_b}) \times \text{I}_{kn}) \\ \text{MMM}_{m,n,k} &\rightarrow (L_m^{mn/n_b} \otimes \text{I}_{n_b}) \circ \\ &\quad ((\otimes)_{1 \times n/n_b} \otimes \text{MMM}_{m,n_b,k}) \circ \\ &\quad (\text{I}_{km} \times (L_{n/n_b}^{kn/n_b} \otimes \text{I}_{n_b})) \end{aligned}$$

Software Defined Radio

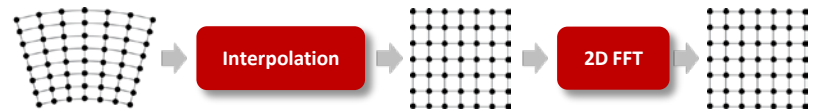


$$\underline{\mathbf{F}}_{K,F} \rightarrow \prod_{i=1}^F \left((\text{I}_{2^{K-2}} \otimes_j B_{F-i,j}) L_{2^{K-2}}^{2^{K-1}} \right)$$

$$\underline{\mathbf{F}}_{K,F} \nu \rightarrow \prod_{i=1}^F \left((\text{I}_{2^{K-2/\nu}} \otimes_{j_1} \text{L}_{\nu}^{-2\nu} \text{B}_{F-i,j_1}^{\nu}) (L_{2^{K-2/\nu}}^{2^{K-1/\nu}} \otimes \text{I}_{\nu}) \right)$$

$$B_{i,j} : \begin{cases} \pi_U = \min_{d_U} (\pi_A + \beta_{A \rightarrow U}, \pi_B + \beta_{B \rightarrow U}) \\ \pi_V = \min_{d_V} (\pi_A + \beta_{A \rightarrow V}, \pi_B + \beta_{B \rightarrow V}) \end{cases}$$

Synthetic Aperture Radar (SAR)



$$\begin{aligned} \text{SAR}_{k \times m \rightarrow n \times n} &\rightarrow \text{DFT}_{n \times n} \circ \text{Interp}_{k \times m \rightarrow n \times n} \\ \text{DFT}_{n \times n} &\rightarrow (\text{DFT}_n \otimes \text{I}_n) \circ (\text{I}_n \otimes \text{DFT}_n) \\ \text{Interp}_{k \times m \rightarrow n \times n} &\rightarrow (\text{Interp}_{k \rightarrow n} \otimes_i \text{I}_n) \circ (\text{I}_k \otimes_i \text{Interp}_{m \rightarrow n}) \\ \text{Interp}_{r \rightarrow s} &\rightarrow \left(\bigoplus_{i=0}^{n-2} \text{InterpSeg}_k \right) \oplus \text{InterpSegPruned}_{k,l} \\ \text{InterpSeg}_k &\rightarrow G_f^{u \cdot n \rightarrow k} \circ \text{iPrunedDFT}_{n \rightarrow u \cdot n} \circ \left(\frac{1}{n} \right) \circ \text{DFT}_n \end{aligned}$$

One Approach for all Types of Parallelism

- Multithreading (Multicore)

$$I_p \otimes_{\parallel} A_{\mu n}, \quad L_m^{mn} \bar{\otimes} I_{\mu}$$

- Vector SIMD (SSE, VMX/AltiVec,...)

$$A \hat{\otimes} I_{\nu} \quad \underbrace{L_2^{2\nu}}_{isa}, \quad \underbrace{L_{\nu}^{2\nu}}_{isa}, \quad \underbrace{L_{\nu}^{\nu^2}}_{isa}$$

- Message Passing (Clusters, MPP)

$$I_p \otimes_{\parallel} A_n, \quad \underbrace{L_p^{p^2} \bar{\otimes} I_{n/p^2}}_{\text{all-to-all}}$$

- Streaming/multibuffering (Cell)

$$I_n \otimes_2 A_{\mu n}, \quad L_m^{mn} \bar{\otimes} I_{\mu}$$

- Graphics Processors (GPUs)

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

- Gate-level parallelism (FPGA)

$$\prod_{i=0}^{n-1} A_i, \quad I_s \tilde{\otimes} A, \quad \underbrace{L_n^m}_{\text{bram}}$$

- HW/SW partitioning (CPU + FPGA)

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

Autotuning in Constraint Solution Space

Intel MIC

DFT₂₅₆

Base cases

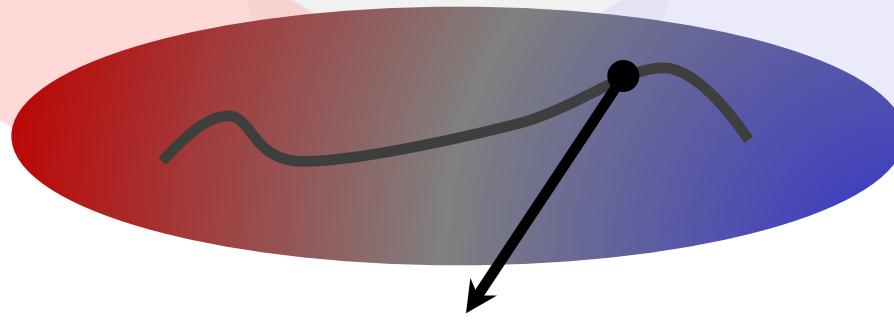
$$\begin{aligned}
 & I_4 \otimes_{\parallel} A_{16n} \\
 & L_m^{mn} \otimes I_{16} \\
 & A \otimes I_4 \\
 & \underbrace{L_2^8}_{\text{SSE}}, \underbrace{L_4^8}_{\text{SSE}}, \underbrace{L_4^{16}}_{\text{SSE}}, \underbrace{L_2^4 \otimes I_2}_{\text{SSE}}
 \end{aligned}$$

Transformation rules

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

Breakdown rules

$$\begin{aligned}
 \text{DFT}_n & \rightarrow (\text{DFT}_k \otimes I_m) T_m^n \cdot (I_k \otimes \text{DFT}_m) L_k^n \\
 \text{DFT}_n & \rightarrow P_n (\text{DFT}_k \otimes \text{DFT}_m) Q_n \\
 \text{DFT}_p & \rightarrow R_p^T (I_1 \oplus \text{DFT}_{p-1}) D_p \cdot (I_1 \oplus \text{DFT}_{p-1}) R_p \\
 \text{DFT}_2 & \rightarrow \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}
 \end{aligned}$$



$$\left((L_m^{mp} \otimes I_{n/p\mu}) \otimes I_\mu \right) \left(I_p \otimes_{\parallel} (\text{DFT}_m \otimes I_{n/p}) \right) \left((L_p^{mp} \otimes I_{n/p\mu}) \otimes I_\mu \right) \left(\bigoplus_{i=0}^{p-1} T_n^{mn,i} \right) \left(I_p \otimes_{\parallel} (I_{m/p} \otimes \text{DFT}_n) \right) \left(I_p \otimes_{\parallel} L_{m/p}^{mn/p} \right) \left((L_p^{pn} \otimes I_{m/p\mu}) \otimes I_\mu \right)$$

Translating a Formula into Code

Constraint Solver Input:

$\underbrace{\text{DFT}_8}_{\text{double}}$



Output =

OL Formula: $(\text{DFT}_2 \otimes I_4) T_4^8 \left(I_2 \otimes \left((\text{DFT}_2 \otimes I_2) T_2^4 (I_2 \otimes \text{DFT}_2) L_2^4 \right) \right) L_2^8$



Σ -OL: $\sum_{j=0}^3 \left(S_j \text{DFT}_2 G_j \right) \sum_{k=0}^1 \left(\sum_{l=0}^1 \left(S_{k,l} \text{diag}(t_{k,l}) \text{DFT}_2 G_l \right) \sum_{m=0}^1 \left(S_m \text{diag}(t_m) \text{DFT}_2 G_{k,m} \right) \right)$



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;
}
```

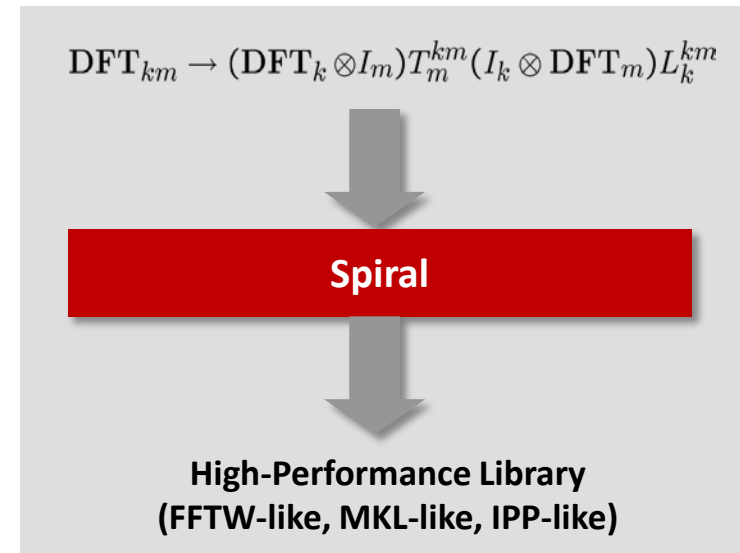
Auto-Generation of Performance Library

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



Core Idea: Recursion Step Closure

- **Input:** transform T and a breakdown rules
- **Output:** problem specifications for recursive function and codelets

- **Algorithm:**

1. Apply the breakdown rule

$$\begin{array}{c} \{\text{DFT}_n\} \\ \downarrow \\ (\{\text{DFT}_{n/k}\} \otimes I_k) T_k^n (I_{n/k} \otimes \{\text{DFT}_k\}) L_{n/k}^n \end{array}$$

2. Convert to Σ -SPL

$$\left(\sum_{i=0}^{k-1} S_{h_{i,k}} \{\text{DFT}_{n/k}\} G_{h_{i,k}} \right) \text{diag}(f) \left(\sum_{j=0}^{n/k-1} S_{h_{jk,1}} \{\text{DFT}_k\} G_{h_{jk,1}} \right) \text{perm}(\ell_{n/k}^n)$$

3. Apply loop merging + index simplification rules.

$$\sum_{i=0}^{k-1} S_{h_{i,k}} \{\text{DFT}_{n/k}\} \text{diag}(f \circ h_{i,k}) G_{h_{i,k}} \quad \sum_{j=0}^{n/k-1} S_{h_{jk,1}} \{\text{DFT}_k\} G_{h_{j,n/k}}$$

4. Extract recursion steps

$$\sum_{i=0}^{k-1} \left\{ S_{h_{i,k}} \text{DFT}_{n/k} \text{diag}(f \circ h_{i,k}) G_{h_{i,k}} \right\} \quad \sum_{j=0}^{n/k-1} \left\{ S_{h_{jk,1}} \text{DFT}_k G_{h_{j,n/k}} \right\}$$

5. Repeat until closure is reached

Spiral-Generated Code (Intel MIC/LRBni)

```

void dft64(float *Y, float *X) {
    __m512 U912, U913, U914, U915, U916, U917, U918, U919, U920, U921, U922, U923, U924, U925,...;
    a2153 = ((__m512 *) X);  s1107 = *(a2153);
    s1108 = *((a2153 + 4));  t1323 = _mm512_add_ps(s1107,s1108);
    ...
    U926 = _mm512_swizupconv_r32(_mm512_set_1to16_ps(0.70710678118654757),_MM_SWIZ_REG_CDAB);
    s1121 = _mm512_madd231_ps(_mm512_mul_ps(_mm512_mask_or_pi(
        _mm512_set_1to16_ps(0.70710678118654757),0xAAAA,a2154,U926),t1341),
        _mm512_mask_sub_ps(_mm512_set_1to16_ps(0.70710678118654757),0x5555,a2154,U926),
        _mm512_swizupconv_r32(t1341,_MM_SWIZ_REG_CDAB));
    U927 = _mm512_swizupconv_r32(_mm512_set_16to16_ps(0.70710678118654757, (-0.70710678118654757),
        0.70710678118654757, (-0.70710678118654757), 0.70710678118654757, (-0.70710678118654757),
        0.70710678118654757, (-0.70710678118654757), 0.70710678118654757, (-0.70710678118654757),
        0.70710678118654757, (-0.70710678118654757), 0.70710678118654757, (-0.70710678118654757),
        0.70710678118654757, (-0.70710678118654757)),_MM_SWIZ_REG_CDAB);
    ...
    s1166 = _mm512_madd231_ps(_mm512_mul_ps(_mm512_mask_or_pi(_mm512_set_16to16_ps(
        0.70710678118654757, (-0.70710678118654757), 0.70710678118654757, (-0.70710678118654757),
        0.70710678118654757, (-0.70710678118654757), 0.70710678118654757, (-0.70710678118654757),
        0.70710678118654757, (-0.70710678118654757), 0.70710678118654757, (-0.70710678118654757),
        0.70710678118654757, (-0.70710678118654757), 0.70710678118654757, (-0.70710678118654757)),
        0xAAAA,a2154,U951),t1362),
        _mm512_mask_sub_ps(_mm512_set_16to16_ps(0.70710678118654757,
        (-0.70710678118654757), 0.70710678118654757, (-0.70710678118654757), 0.70710678118654757,
        (-0.70710678118654757), 0.70710678118654757, (-0.70710678118654757), 0.70710678118654757,
        (-0.70710678118654757), 0.70710678118654757, (-0.70710678118654757), 0.70710678118654757,
        (-0.70710678118654757), 0.70710678118654757, (-0.70710678118654757)),0x5555,a2154,U951),
        _mm512_swizupconv_r32(t1362,_MM_SWIZ_REG_CDAB));
    ...
}

```


Support For Library-Specific Interfaces

Complex FFT

```
IPPAPI(IppStatus, name data type, size scaling  
(const Ipp32fc *pSrc, Ipp32fc *pDst, int length, int flag) )
```

Walsh-Hadamard Transform

```
IPPAPI(IppStatus, ippgWHT_32f, log(size) scaling memory  
(const Ipp32f *pSrc, Ipp32f *pDst, int order, int flag, Ipp8u *pBuf)
```

```
IPPAPI(IppStatus, ippgWHTGetBufferSize_32f, memory  
(int order, Ipp32u *pBufferSize) )
```


Organization

- Spiral overview
- **Validation and Verification**
- Results
- Concluding remarks

Symbolic Verification

- Transform = Matrix-vector multiplication
matrix fully defines the operation

$$\text{DFT}_4 = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & j & -1 & -j \\ 1 & -1 & 1 & -1 \\ 1 & -j & -1 & j \end{bmatrix}$$



= ?

- Algorithm = Formula
represents a matrix expression, can be evaluated to a matrix

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

Empirical Verification

- Run program on all basis vectors, compare to columns of transform matrix

$$\begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & j & -1 & -j \\ 1 & -1 & 1 & -1 \\ 1 & -j & -1 & j \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad \text{_____} \quad \text{= ?}$$

`DFT4([0,1,0,0])` _____

- Compare program output on random vectors to output of a random implementation of same kernel

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

`DFT4_rnd([0.1, 1.77, 2.28, -55.3])` _____

Verification of the Generator

- Rule replaces left-hand side by right-hand side when preconditions match

$$I_m \otimes A_n \rightarrow L_m^{mn} (A_n \otimes I_m) L_n^{mn}$$

- Test rule by evaluating expressions before and after rule application and compare result

$$I_2 \otimes \text{DFT}_2 = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & -1 \end{bmatrix}$$

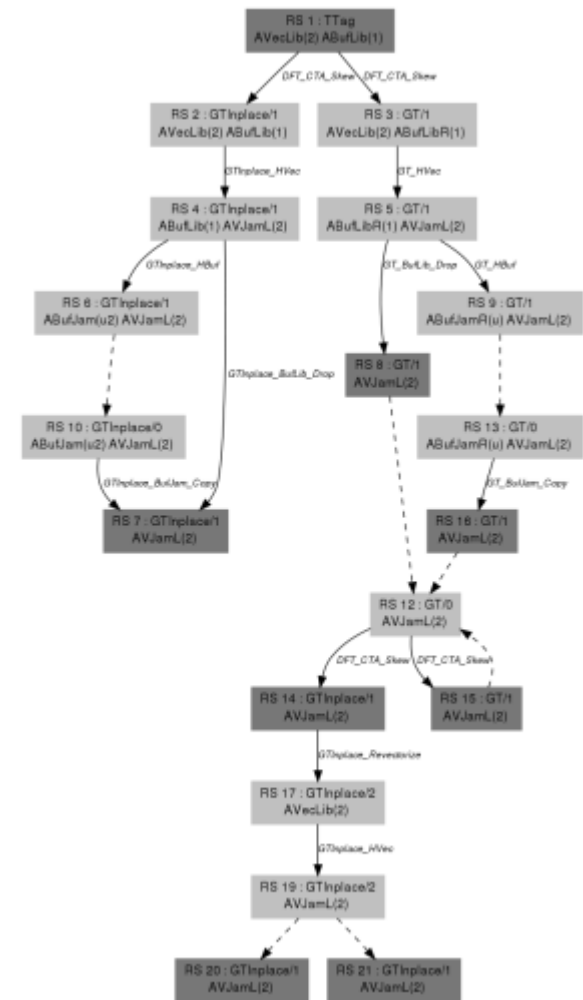
= ?

$$L_2^4(\text{DFT}_2 \otimes I_2) L_2^4 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & -1 & 0 \\ 0 & 1 & 0 & -1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

Verification of Autotuning Libraries

Auto-generated FFTW-like library

- Need verifier for each function
- Auto-generated from specification
- Auto-generate test harness
- Drop-in replacement into existing infrastructure

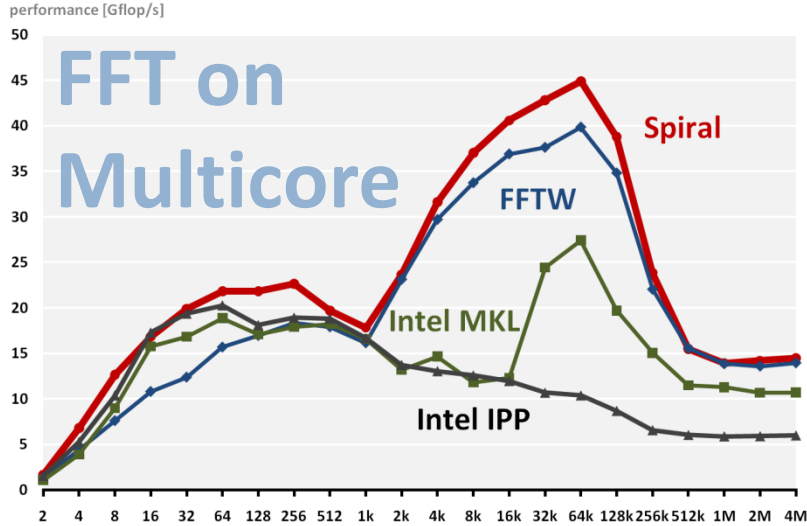


Organization

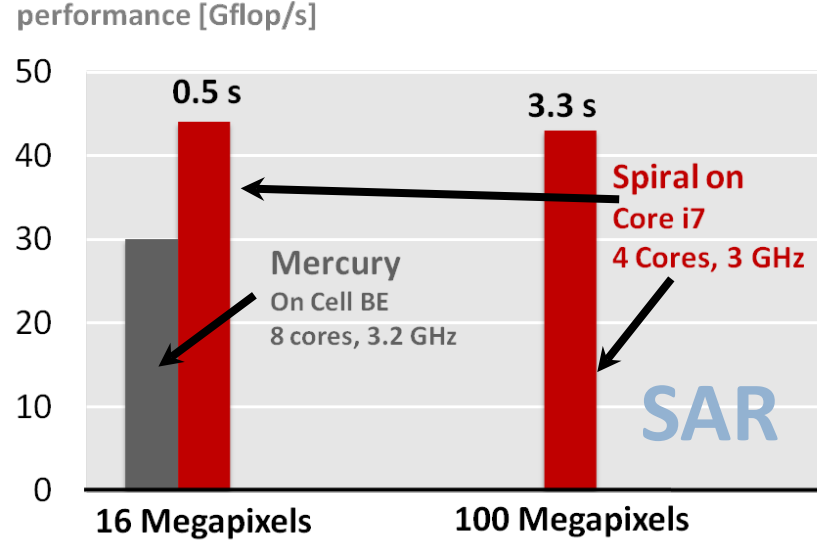
- Spiral overview
- Validation and Verification
- **Results**
- Concluding remarks

Results: Spiral Outperforms Humans

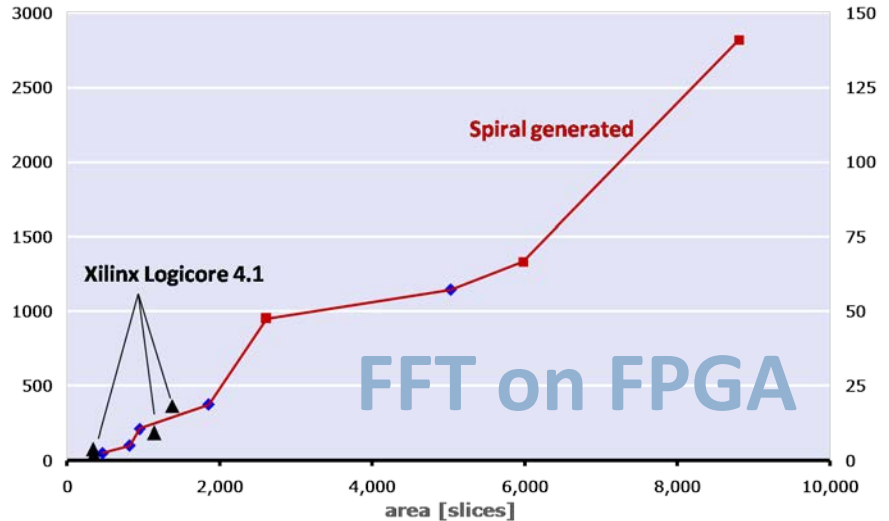
1D DFT on 3.3 GHz Sandy Bridge (4 Cores, AVX)



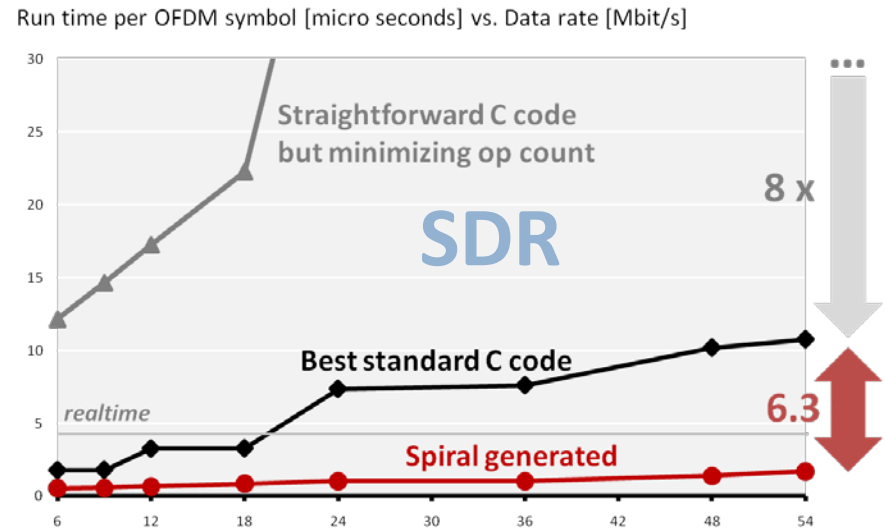
SAR Image Formation on Intel platforms



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



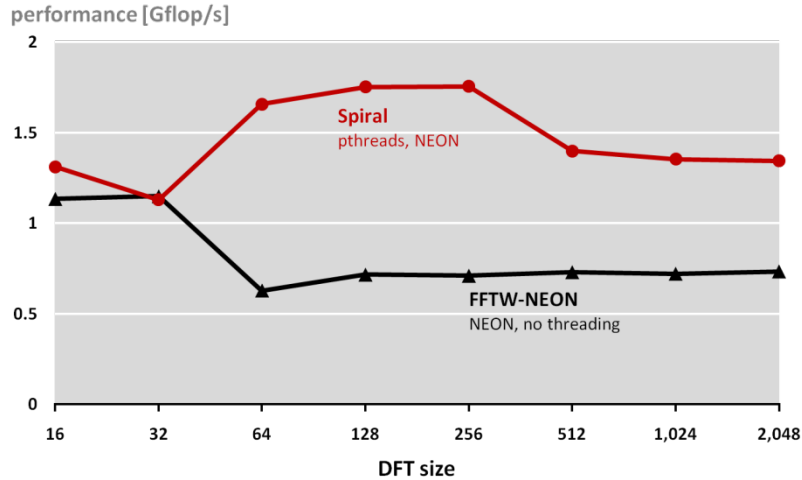
WiFi transmitter on Dualcore Intel Atom



From 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

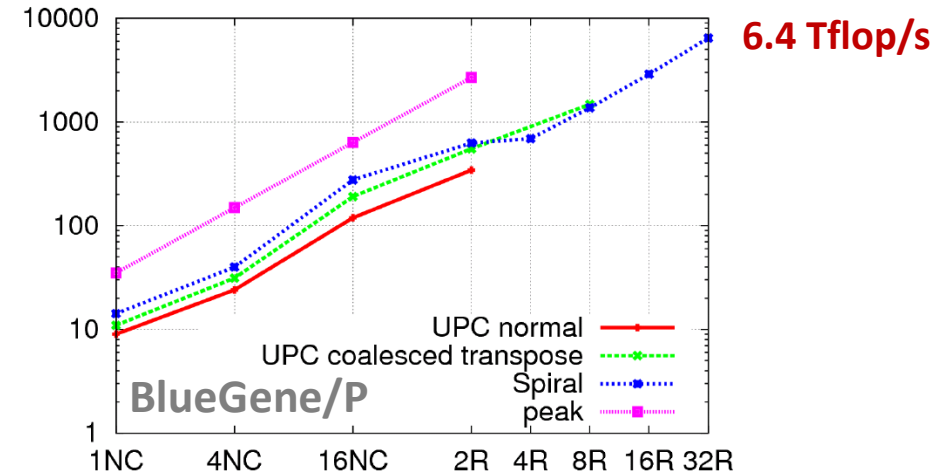
Dual-core ARM at 1.2GHz with NEON ISA

SIMD vectorization + multi-threading



Global FFT (1D FFT, HPC Challenge)

performance [Gflop/s]



BlueGene/P at Argonne National Laboratory

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

SIMD vectorization + multi-threading + MPI



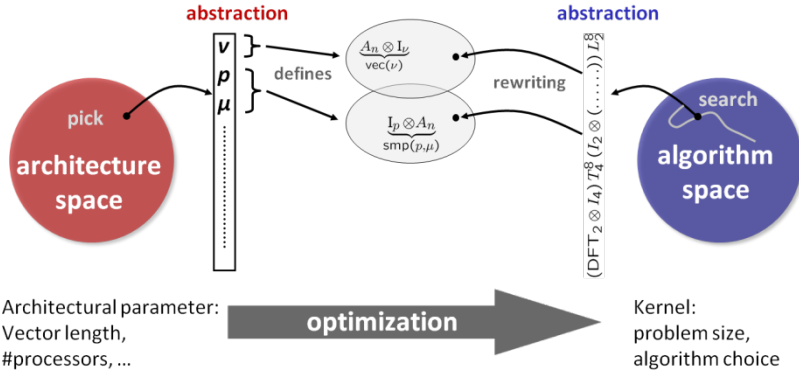
Organization

- Spiral overview
- Validation and Verification
- Results
- **Concluding remarks**

Summary: Spiral in a Nutshell

Joint Abstraction

Model: common abstraction
= spaces of matching formulas

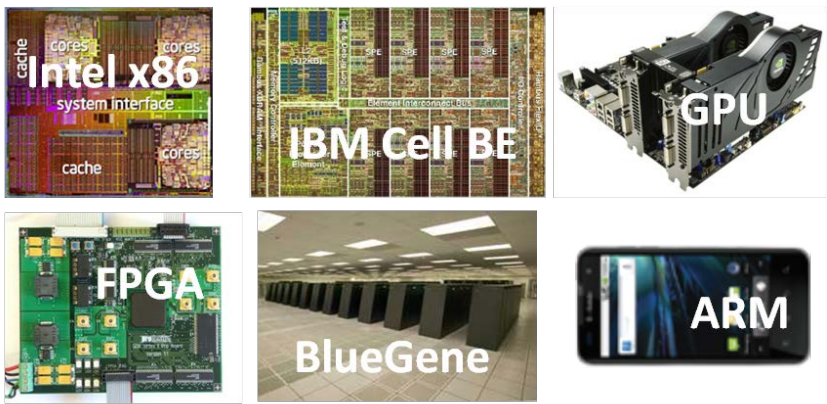


Verification

$$\begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & j & -1 & -j \\ 1 & -1 & 1 & -1 \\ 1 & -j & -1 & j \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} = ?$$

$DFT_4([0, 1, 0, 0])$

Target Machines



Application Domains

Signal Processing
 $DFT_n \rightarrow (DFT_k \otimes I_m) T_m^n (I_k \otimes DFT_m) L_k^n$

Matrix Algorithms

$$\begin{bmatrix} \square & \square \\ \square & \square \end{bmatrix} = \begin{bmatrix} \square & \square \\ \square & \square \end{bmatrix} \times \begin{bmatrix} \square \\ \square \end{bmatrix}$$

Software Defined Radio

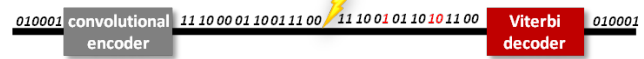


Image Formation (SAR)



Acknowledgement

James C. Hoe

Jeremy Johnson

José M. F. Moura

David Padua

Markus Püschel

Volodymyr Arbatov

Paolo D'Alberto

Peter A. Milder

Yevgen Voronenko

Qian Yu

Berkin Akin

Christos Angelopoulos

Srinivas Chellappa

Frédéric de Mesmay

Daniel S. McFarlin

Marek R. Telgarsky



Special thanks to:

Randi Rost, Scott Buck (Intel), Jon Greene (Mercury Inc.), Yuanwei Jin (UMES)

Gheorghe Almasi, Jose E. Moreira, Jim Sexton (IBM), Saeed Maleki (UIUC)

Francois Gygi (LLNL, UC Davis), Kim Yates (LLNL), Kalyan Kumaran (ANL)

More Information:

www.spiral.net

www.spiralgen.com