

Automatic Performance Tuning

Spiral, FFTW, ATLAS & Friends

Franz Franchetti

ECE, Carnegie Mellon University

www.ece.cmu.edu/~franzf

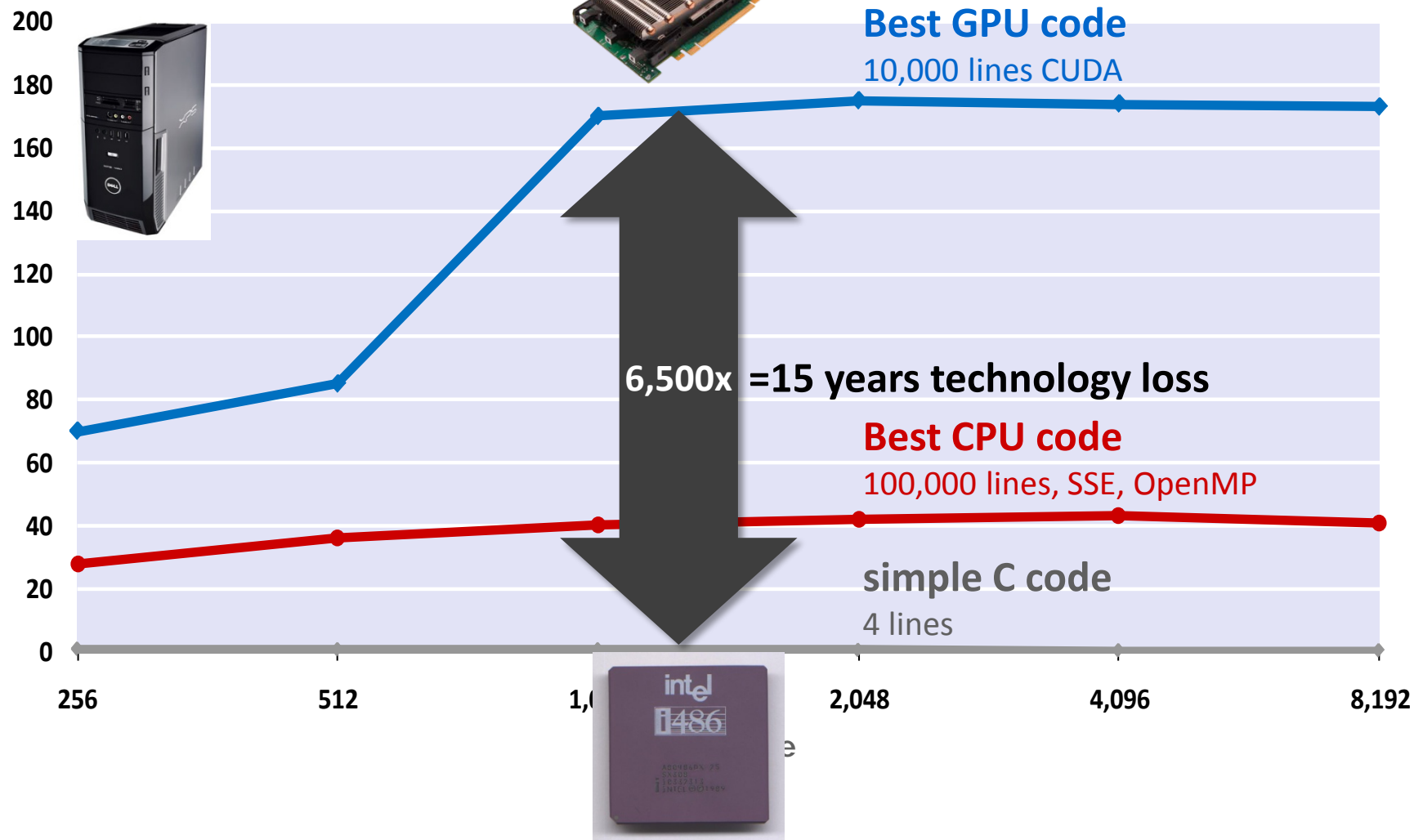
Co-Founder, SpiralGen

www.spiralgen.com

The work was sponsored by Defense Advanced Research Projects Agency (DARPA) under agreement No. HR0011-13-2-0007. The content, views and conclusions presented in this document do not necessarily reflect the position or the policy of DARPA or the U.S. Government. No official endorsement should be inferred.

The Cost Of Portability and Maintainability

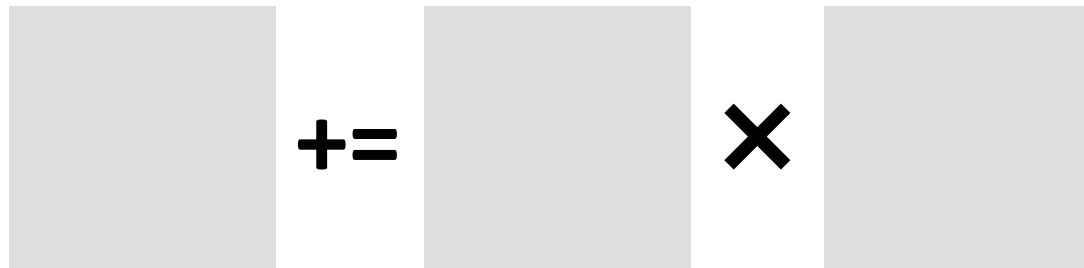
Matrix-Matrix Multiplication Performance [Gflop/s]



Compilers: The Fundamental Problem

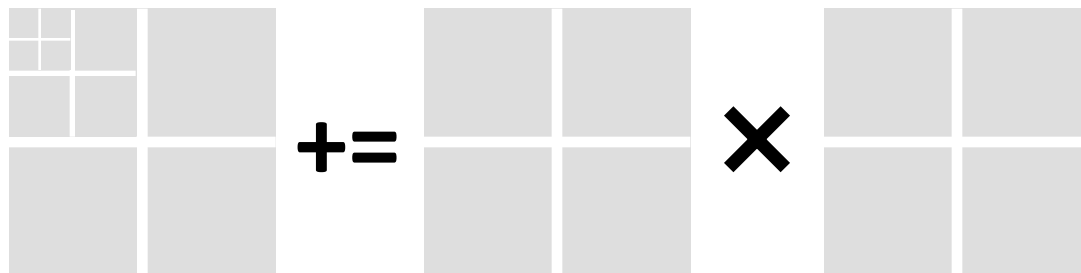
Matrix-matrix multiplication

```
for i=1:N
  for j=1:N
    for k=1:N
      C[i,j] += A[i,k]*B[k,j]
```



Tiled matrix-matrix multiplication

```
for i=1:NB:N
  for j=1:NB:N
    for k=1:NB:N
      for i0=i:NU:i+NB
        for j0=j:NU:j+NB
          for k0=k:NU:k+NB
            for k00=k0:1:k0+NU
              for j00=j0:1:j0+NU
                for i00=i0:1:i0+NU
                  C[i00,j00] +=
                    A[i00,k00]*B[k00,j00]
```

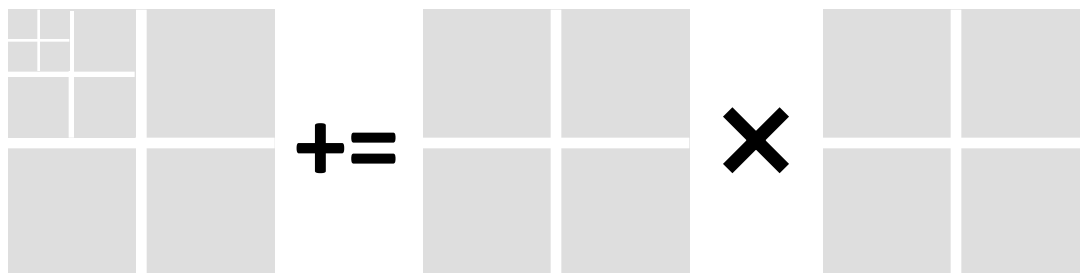


Problem: which transformation order? What parameter values?

How to Overcome Compiler Limitations?

■ Autotuning

Write parameterized program and empirically find good parameters



■ Program generation

Automatically generate big basic blocks and special case code

Use special instructions (SSE/AVX/NEON/Altivec/QPX,...)

```
for (i = 0; i < 4; i++)  
  A[i] = B[i] * C[i];
```

```
A[0] = B[0] * C[0];  
A[1] = B[1] * C[1];  
A[2] = B[2] * C[2];  
A[3] = B[3] * C[3];
```

Two red arrows point from the loop body in the first code block to the four individual assignment statements in the second code block, illustrating the expansion of a loop into a basic block.

When done: 30MB of source code for DGEMM

Autotuning: From Simple to “Really Hard”

- **Level 0: simple C program**
implements the algorithm cleanly
- **Level 1: C macros plus search script**
use C preprocessor for meta-programming
- **Level 2: scripting for code specialization**
text-based program generation, e.g., ATLAS
- **Level 3: add compiler technology**
internal code representation, e.g., FFTW’s genfft
- **Level 4: synthesize the program from scratch**
high level representation, e.g., TCE and Spiral



Performance Portability Landscape

Synthesis from Domain Math

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

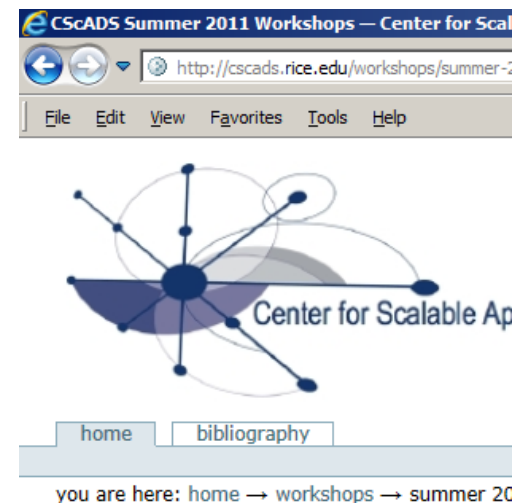
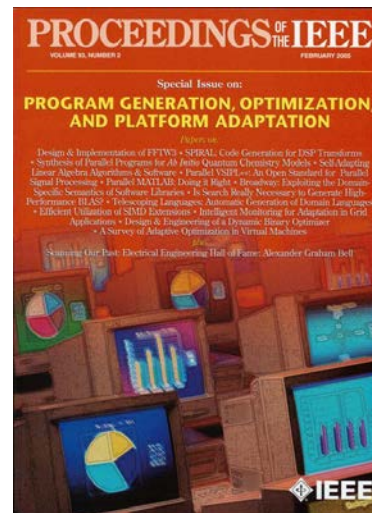
Numerical Libraries

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

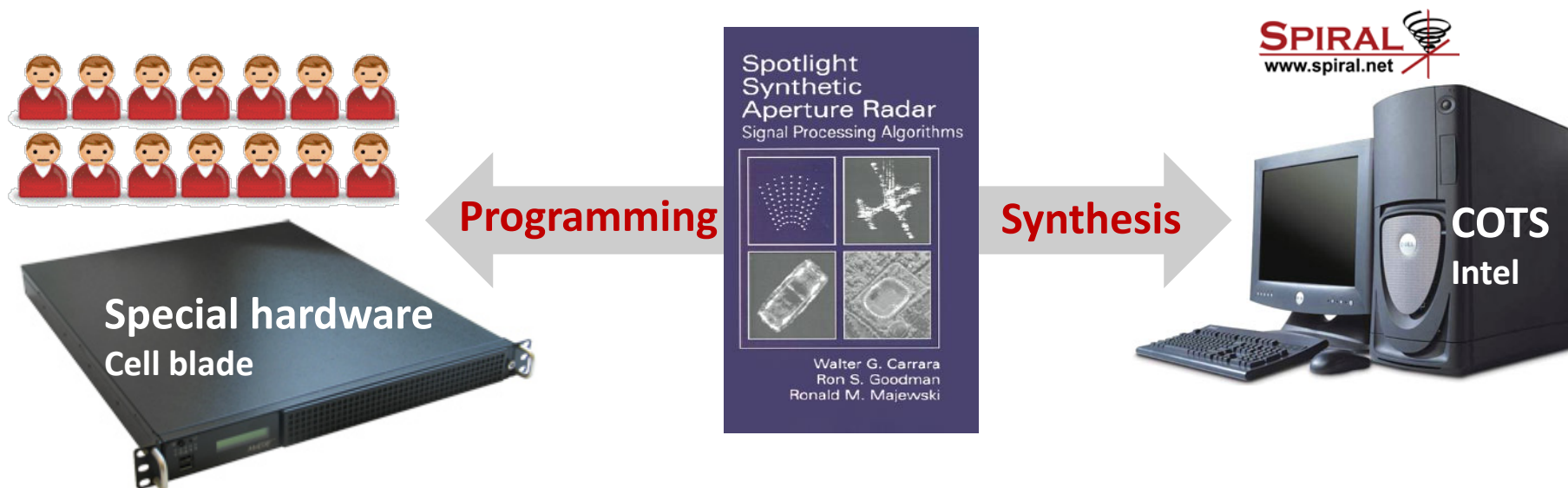
Compiler Synthesis/Autotuning

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

Autotuning Primer



Spiral: Computer Writes Best SAR Code



Result

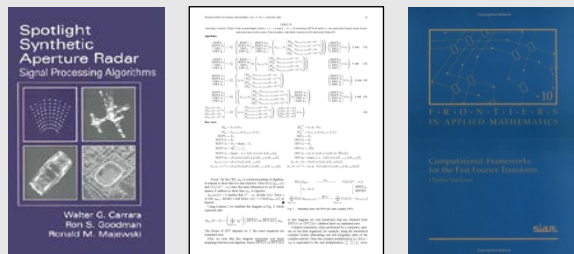
Same performance, 1/10th human effort, non-expert user

Key ideas

restrict domain, use mathematics, performance portability

What is Spiral?

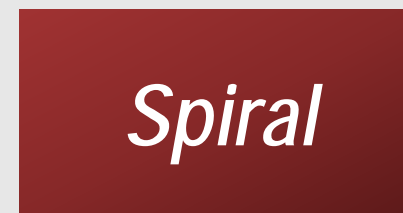
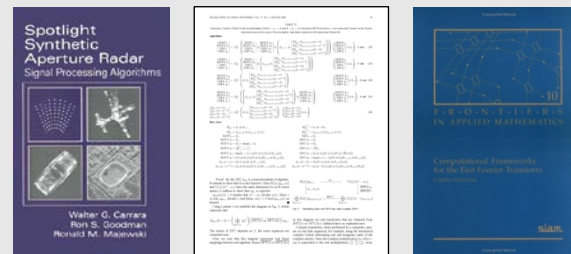
Traditionally



High performance library
optimized for given platform

*Comparable
performance*

Spiral Approach



High performance library
optimized for given platform

Algorithms: Rules in Domain Specific Language

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



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

Viterbi Decoding



$$\begin{aligned}
 \underline{\text{Vit}} &\rightarrow \underbrace{\left(\prod (L \times I) \circ (I \otimes C) \right)}_{\text{vec}(v)} \circ \text{Id} \\
 \text{vec}(v) &\rightarrow \left(\prod \underbrace{(L \times I) \circ (I \otimes C)}_{\text{vec}(v)} \right) \circ \text{Id} \\
 \times &\rightarrow \left(\prod (L \otimes \text{I}_v \times I) \circ (I \otimes C \otimes \text{I}_v) \circ (\vec{L} \times I) \right) \circ \text{Id} \\
 &\rightarrow \prod (L \otimes \text{I}_v \times I) \circ (I \otimes (B \otimes \text{I}_v)) \circ (\vec{L} \times I)
 \end{aligned}$$

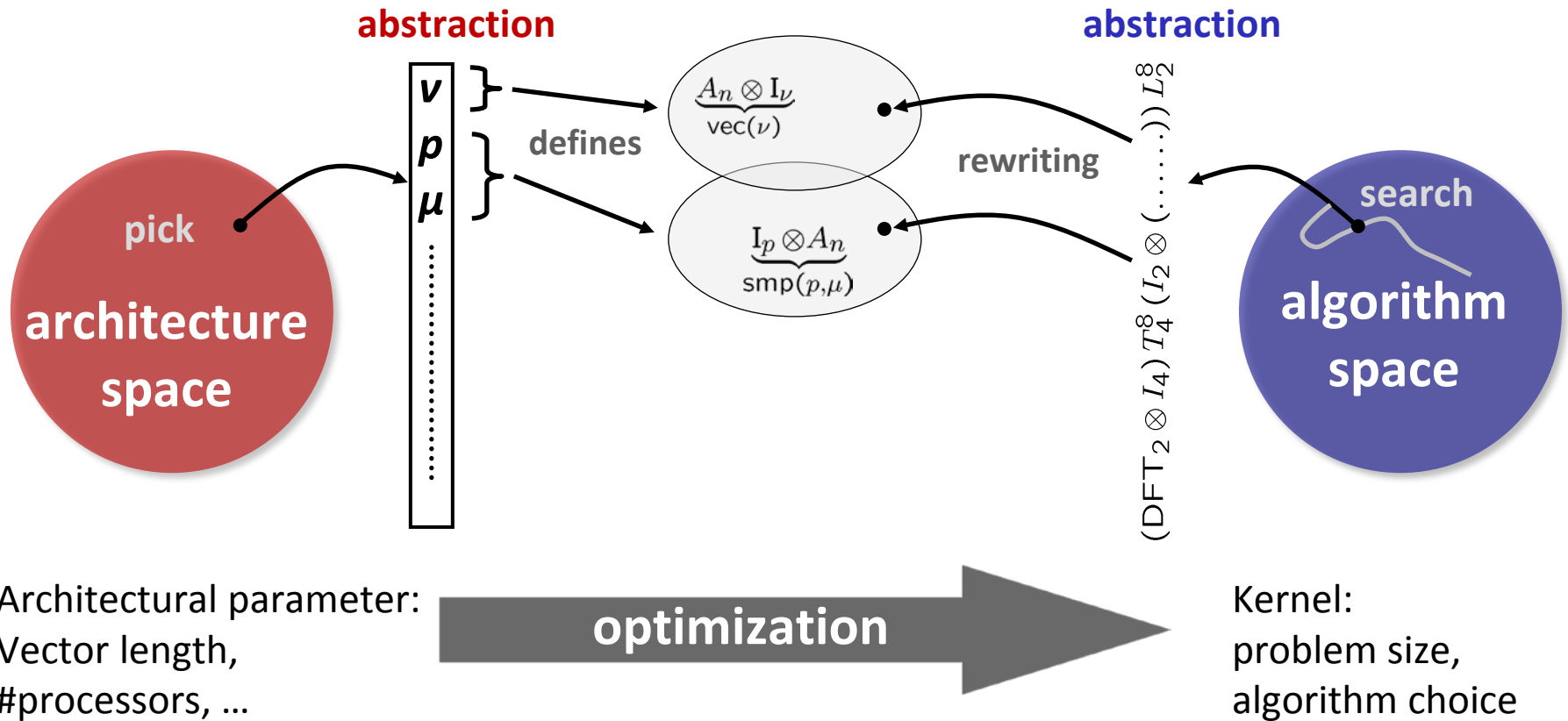
Synthetic Aperture Radar (SAR)



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

Spiral's Domain-Specific Program Synthesis

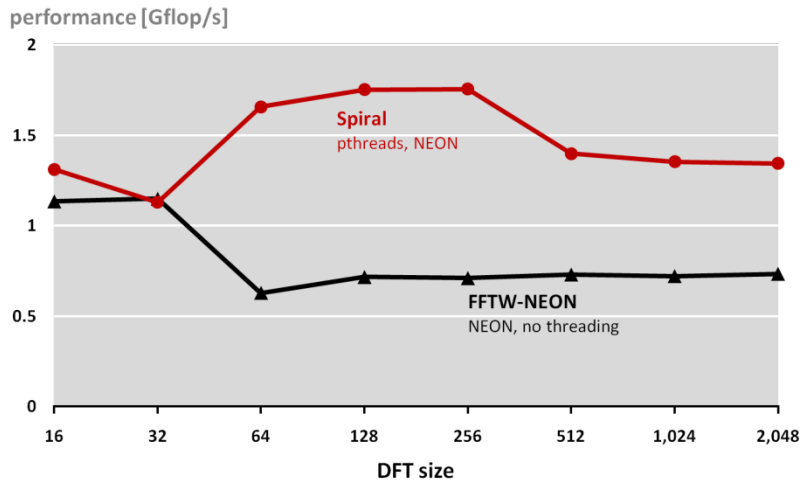
Model: common abstraction
= spaces of matching formulas



From Cell Phone To Supercomputer

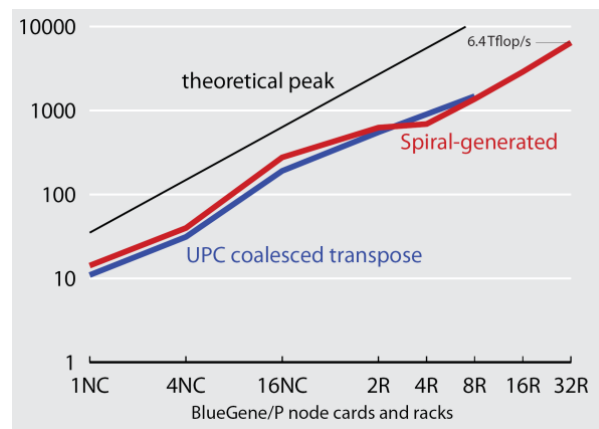
DFT on Samsung Galaxy S II

Dual-core 1.2 GHz Cortex-A9 with NEON ISA



Global FFT (1D FFT, HPC Challenge)

performance [Gflop/s]



6.4 Tflop/s on BlueGene/P

Samsung i9100 Galaxy S II

Dual-core ARM at 1.2GHz with NEON ISA

SIMD vectorization + multi-threading



BlueGene/P at Argonne National Laboratory

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

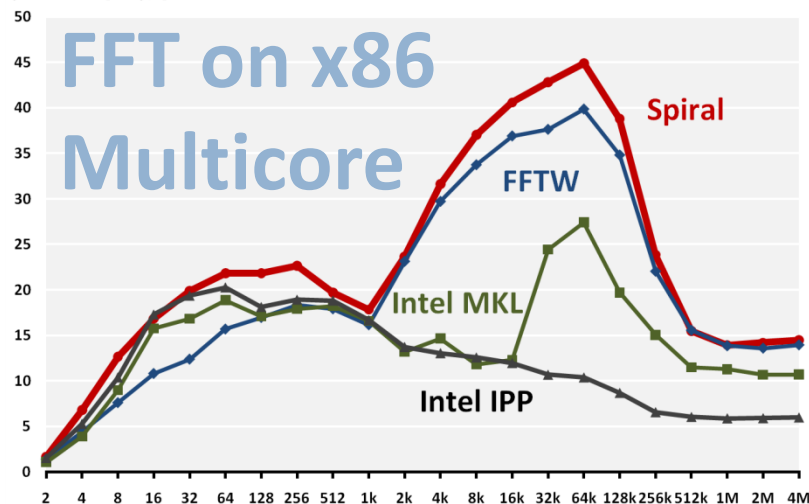
SIMD vectorization + multi-threading + MPI



Results: Spiral Outperforms Humans

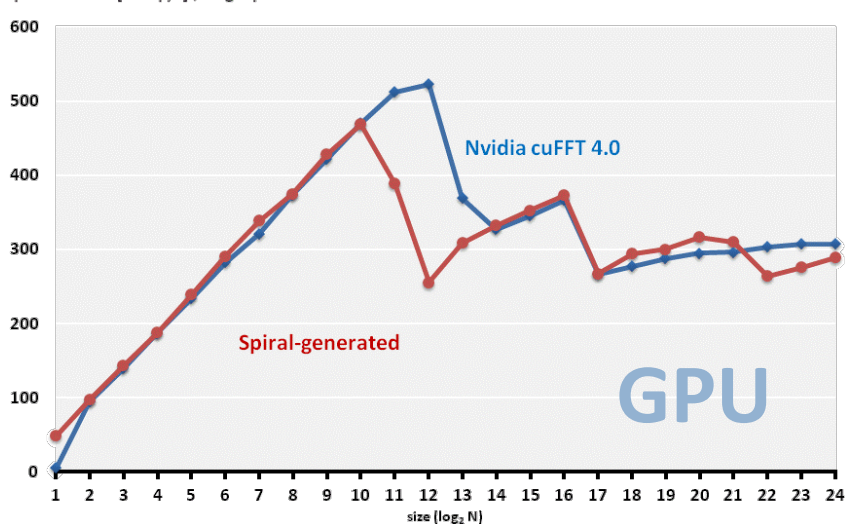
1D DFT on 3.3 GHz Sandy Bridge (4 Cores, AVX)

performance [Gflop/s]



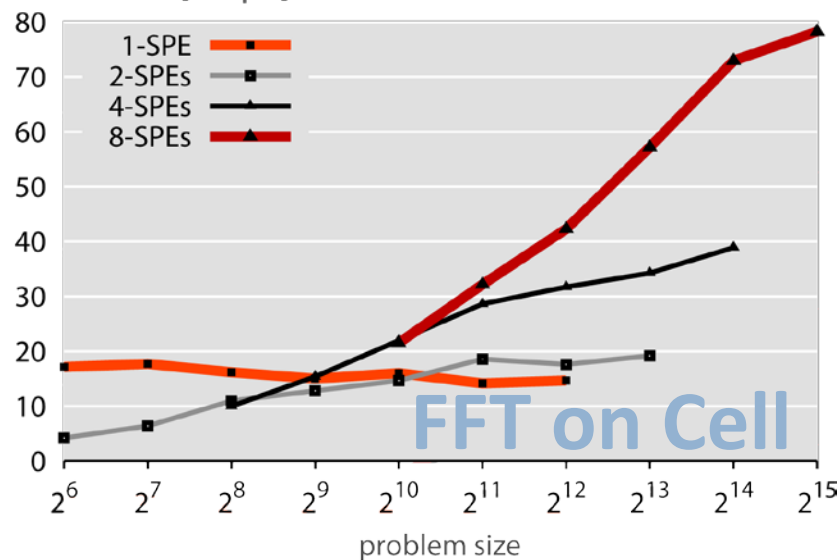
1D Batch DFT (Nvidia GTX 480)

performance [GFlop/s], single precision



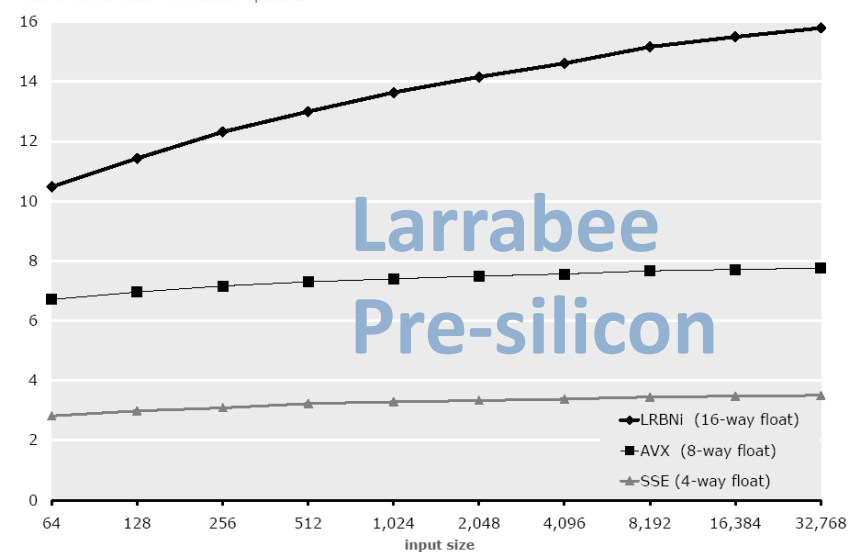
DFT on multiple SPEs (LS-LS, block-cyclic)

Performance [Gflop/s]



Vectorization Efficiency

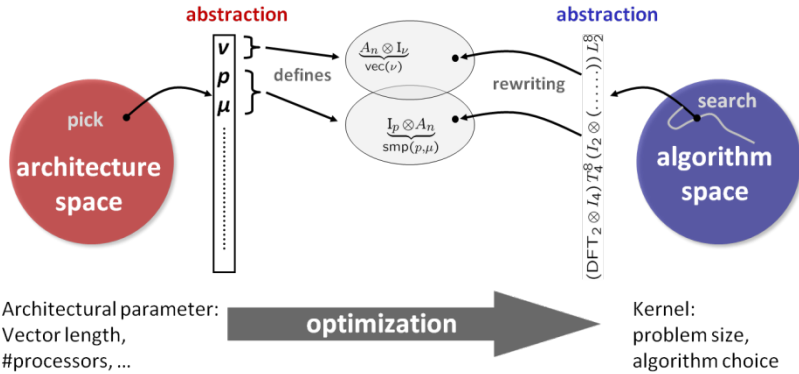
Ratio of x87 to Vector Architecture opcounts



Summary: Spiral in a Nutshell

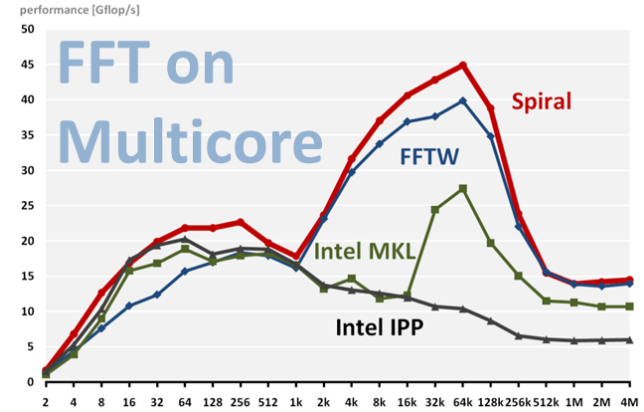
Joint Abstraction

Model: common abstraction
= spaces of matching formulas

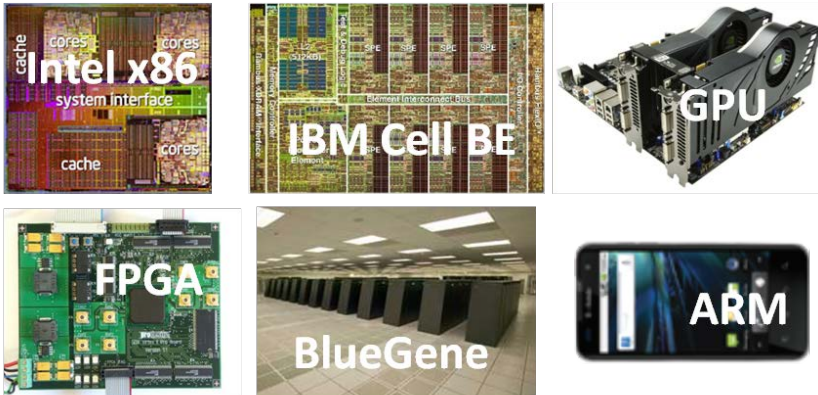


Outperforms Human Experts

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



Performance Portability



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

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

Software Defined Radio



Image Formation (SAR)

