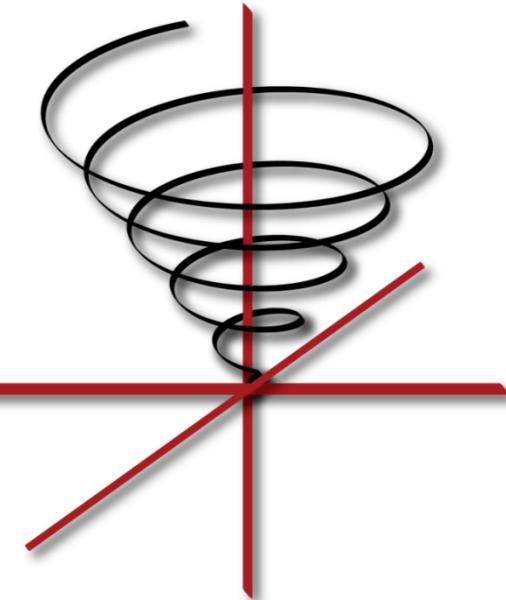


SPIRAL

FFT Library Generation and Autotuning

Thom Popovici, Franz Franchetti

Carnegie Mellon University



www.spiral.net

This work was supported by DARPA, ONR, DOE NSF, Intel, Mercury

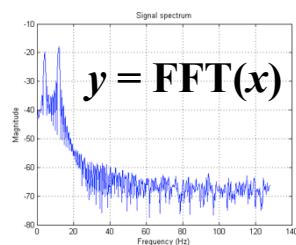
Idea: Go from Mathematics to Software

Given:

- Mathematical problem specification
does not change
- Computer platform
changes often

Wanted:

- Very good implementation of specification on platform
- Proof of correctness



automatic

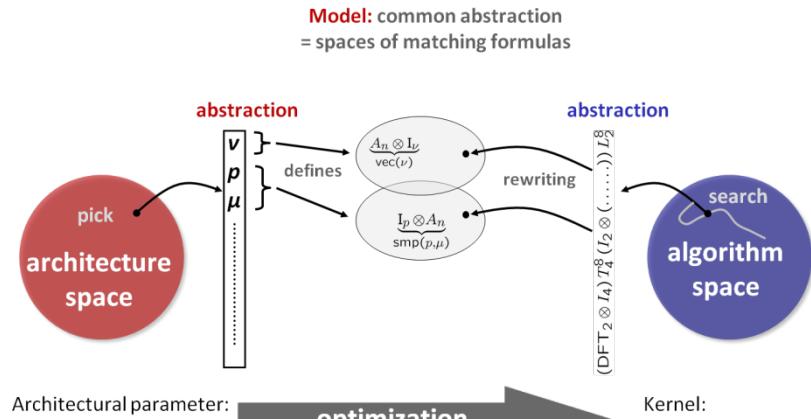
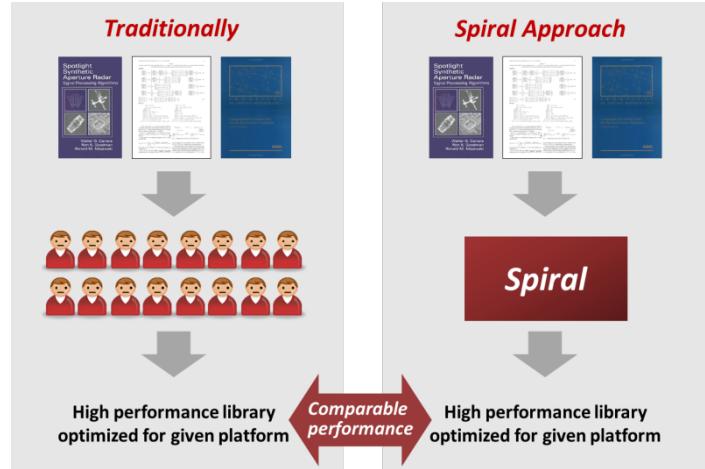
```
void fft64(double *Y, double *X) {
    ...
    s5674 = _mm256_permute2f128_pd(s5672, s5673, (0) | ((2) << 4));
    s5675 = _mm256_permute2f128_pd(s5672, s5673, (1) | ((3) << 4));
    s5676 = _mm256_unpacklo_pd(s5674, s5675);
    s5677 = _mm256_unpackhi_pd(s5674, s5675);
    s5678 = *((a3738 + 16));
    s5679 = *((a3738 + 17));
    s5680 = _mm256_permute2f128_pd(s5678, s5679, (0) | ((2) << 4));
    s5681 = _mm256_permute2f128_pd(s5678, s5679, (1) | ((3) << 4));
    s5682 = _mm256_unpacklo_pd(s5680, s5681);
    s5683 = _mm256_unpackhi_pd(s5680, s5681);
    t5735 = _mm256_add_pd(s5676, s5682);
    t5736 = _mm256_add_pd(s5677, s5683);
    t5737 = _mm256_add_pd(s5670, t5735);
    t5738 = _mm256_add_pd(s5671, t5736);
    t5739 = _mm256_sub_pd(s5670, _mm256_mul_pd(_mm_vbroadcast_sd(&(C22)), t5735));
    t5740 = _mm256_sub_pd(s5671, _mm256_mul_pd(_mm_vbroadcast_sd(&(C22)), t5736));
    t5741 = _mm256_mul_pd(_mm_vbroadcast_sd(&(C23)), _mm256_sub_pd(s5677, s5683));
    t5742 = _mm256_mul_pd(_mm_vbroadcast_sd(&(C23)), _mm256_sub_pd(s5676, s5682));
    ...
}
```

performance
+
PROOF
QED.

Spiral Technology in a Nutshell

Library Generator

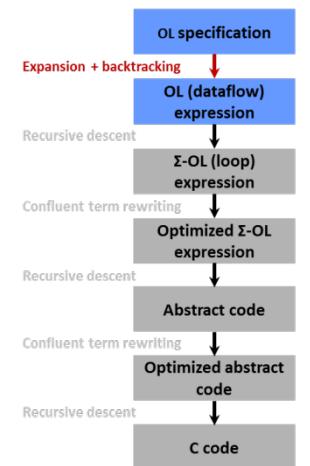
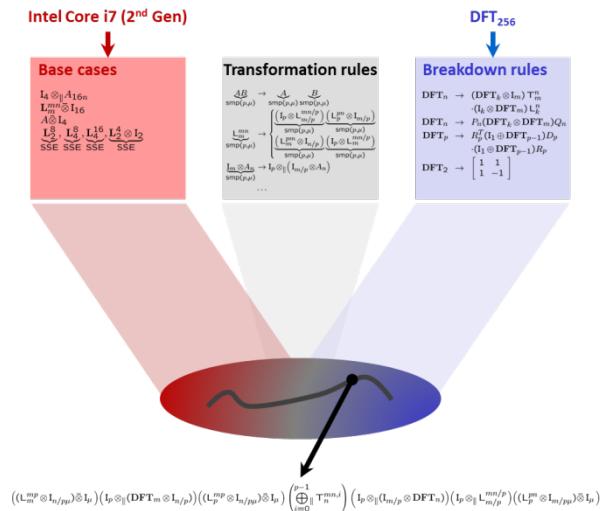
Mathematical Foundation



Extensible Algorithm Library

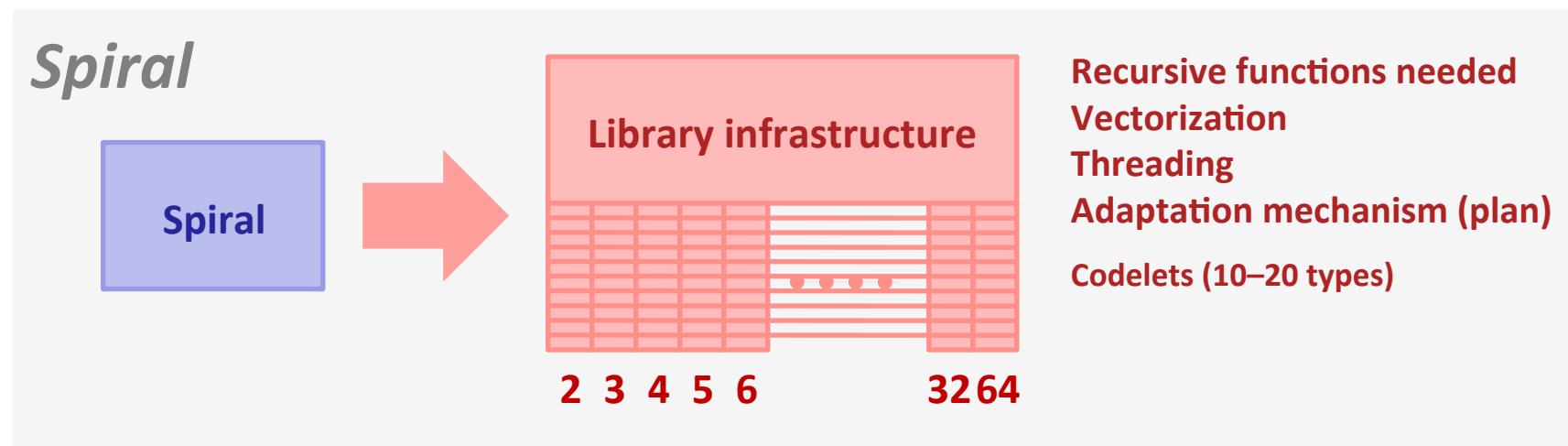
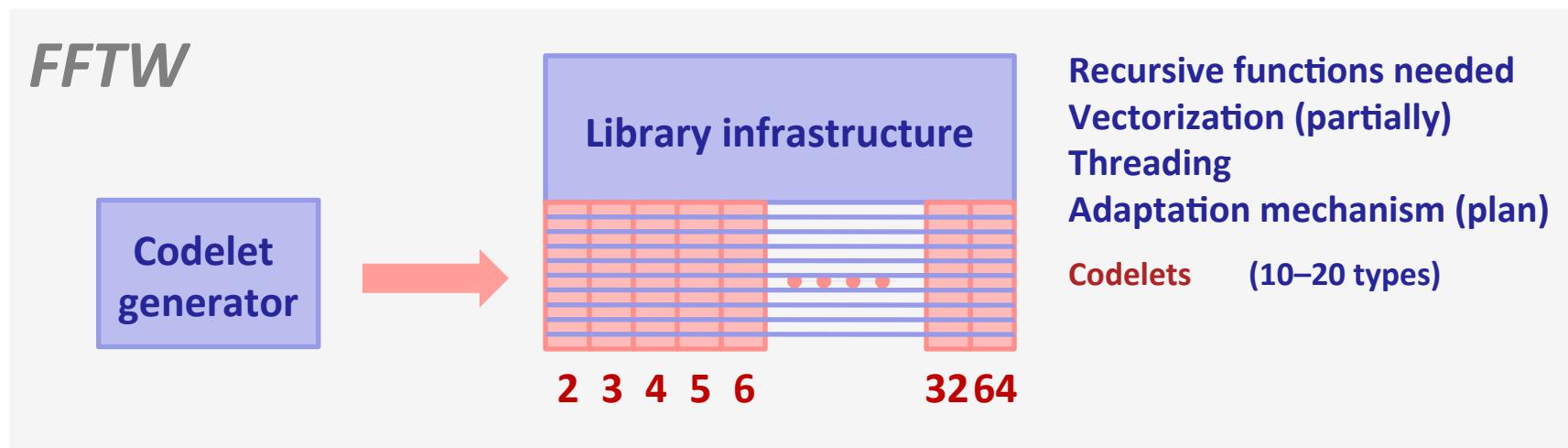
$$\begin{aligned}
 DFT_n &\rightarrow (DFT_k \otimes I_m) T_m^n (I_k \otimes DFT_m) L_k^n, \quad n = km \\
 DFT_n &\rightarrow P_n (DFT_k \otimes DFT_m) Q_n, \quad n = km, \quad \gcd(k, m) = 1 \\
 DFT_p &\rightarrow R_p^T (I_1 \oplus DFT_{p-1}) D_p (I_1 \oplus DFT_{p-1}) R_p, \quad p \text{ prime} \\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

Code Synthesis and Autotuning



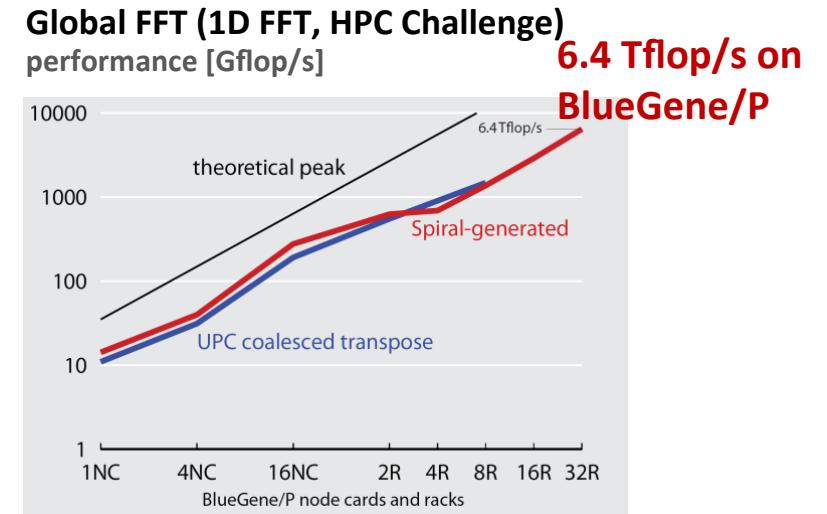
SPIRAL: Meta-Tool to FFTW

blue = hand-written, red = generated



SPIRAL: Success in HPC/Supercomputing

- **NCSA Blue Waters**
PAID Program, FFTs for Blue Waters
- **RIKEN K computer**
FFTs for the HPC-ACE ISA
- **LANL RoadRunner**
FFTs for the Cell processor
- **PSC/XSEDE Bridges**
Large size FFTs
- **LLNL BlueGene/L and P**
FFTW for BlueGene/L's Double FPU
- **ANL BlueGene/Q Mira**
Early Science Program, FFTW for BGQ QPX



BlueGene/P at Argonne National Laboratory
128k cores (quad-core CPUs) at 850 MHz

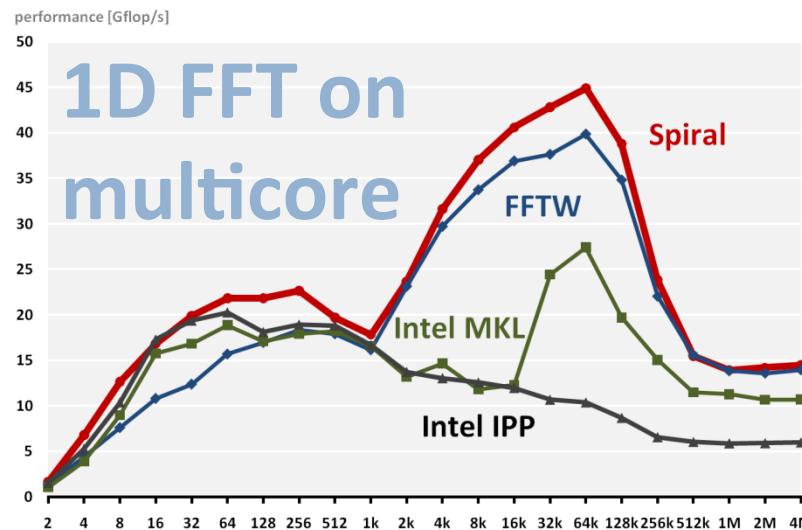


2006 Gordon Bell Prize (Peak Performance Award) with LLNL and IBM

2010 HPC Challenge Class II Award (Most Productive System) with ANL and IBM

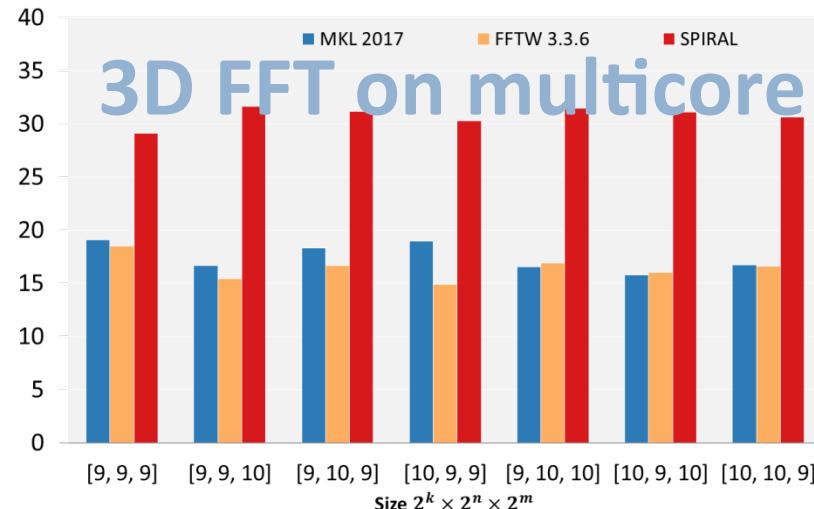
SPIRAL: Performance Across Platforms

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

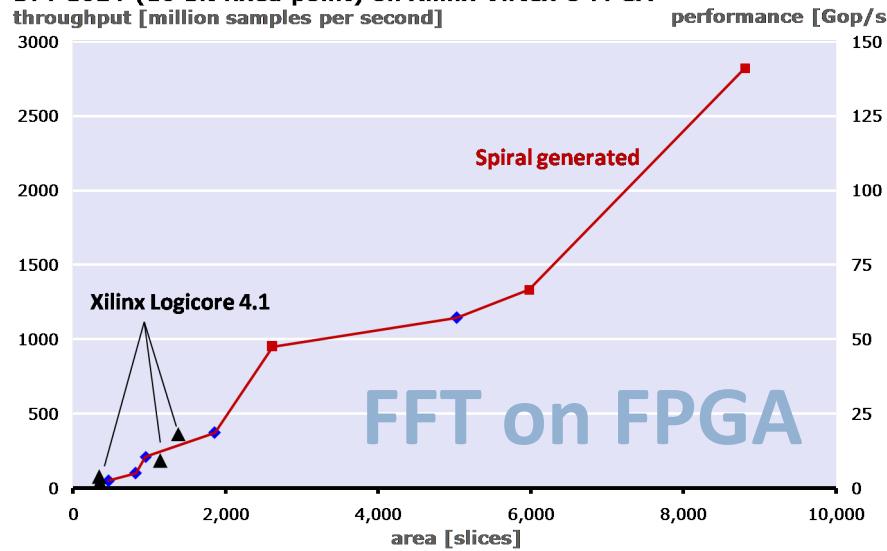


3D FFT Performance on Intel Haswell 4770K

3.5 GHz, 4/8 cores/threads, double-precision, AVX
Performance [Gflop/s]

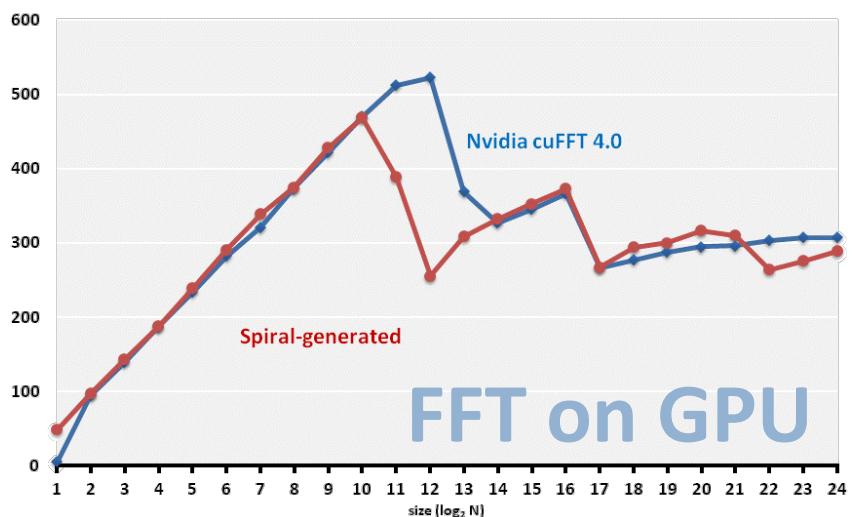


DFT 1024 (16 bit fixed point) on Xilinx Virtex-5 FPGA



1D Batch DFT (Nvidia GTX 480)

performance [GFlop/s], single precision



Thrust 1: Repeat The Success of LAPACK for FFTs

Numerical Linear Algebra

LAPACK

ScaLAPACK

LU factorization

Eigensolves

SVD

BLAS, BLACS

BLAS-1

BLAS-2

BLAS-3

Spectral Algorithms

Convolution

Correlation

Upsampling

Poisson solver



FFTW

DFT, RDFT

1D, 2D, 3D,...

batch

No LAPACK equivalent for spectral methods

- **Medium size 1D FFT (1k–10k data points) is most common library call**
applications break down 3D problems themselves and then call the 1D FFT library
- **Higher level FFT calls rarely used**
FFTW *guru* interface is powerful but hard to use, leading to performance loss
- **Low arithmetic intensity and variation of FFT use make library approach hard**
Algorithm specific decompositions and FFT calls intertwined with non-FFT code

Thrust 2: Updated FFTW Interface for Node FFTs

Front End: application facing

- **Backwards compatible to FFTW 2.X and 3.X**
old code runs unmodified and gains substantially but not fully
- **Delayed execution/futures**
asynchronous execution, offloading, triggering of execution
- **Data location and layout control**
specify location of data (host/device, NUMA, NVRAM,...)
- **Performance metrics**
user control of resource/performance/energy requirements and adaptation

Backend: developer and hardware facing

- **Extensible to novel and yet-unknown hardware and system architectures**
future-proof platform extension API: transparently enable accelerators, etc.
- **Plug-in mechanism for special case acceleration packs**
transparent access to application/architecture specialized FFT kernels

More Information:

www.spiral.net

www.spiralgen.com