

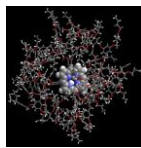
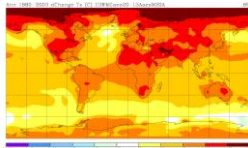
Program Generation for Performance Spiral

Markus Püschel

Computer Science
ETH zürich

SPiRAL
www.spiral.net

Mathematical Computing



Science simulations

Audio, image, Video processing



Signal processing, communication,
control

Security

Machine learning, data analytics



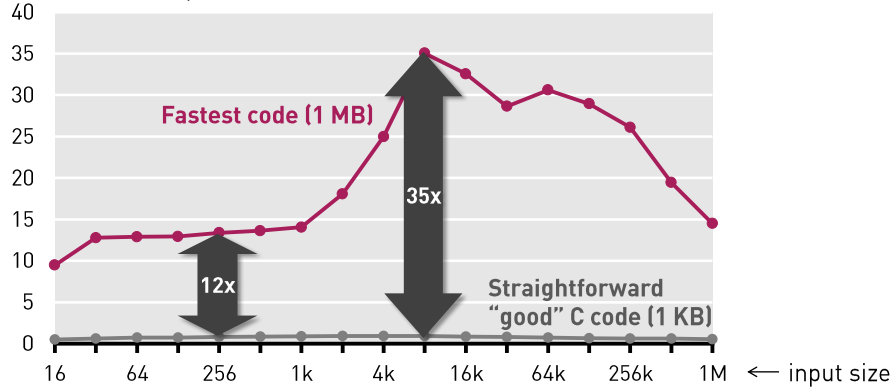
Optimization

***Highest performance
is often crucial***

Example: Discrete Fourier Transform

DFT (single precision) on Intel Core i7 (4 cores, 2.66 GHz)

Performance [Gflop/s]

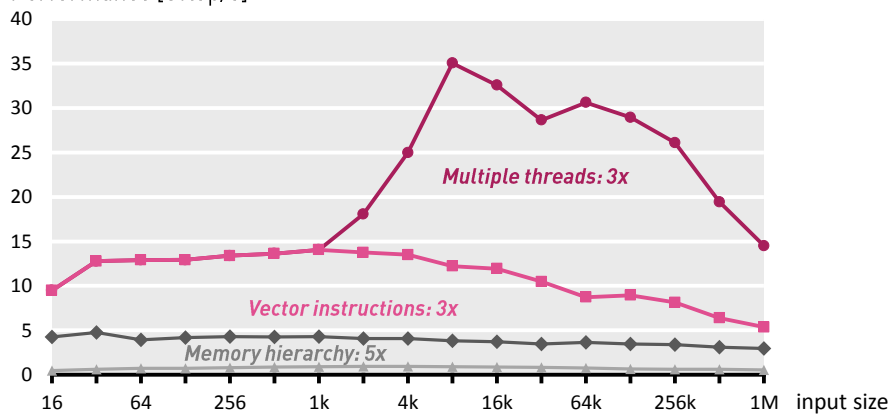


Vendor compiler, best flags

Roughly same operations count

DFT (single precision) on Intel Core i7 (4 cores, 2.66 GHz)

Performance [Gflop/s]



Compiler doesn't do the job

Doing by hand = restructure algorithm for locality & parallelism,
handle choices, choose proper code style, use vector intrinsics,
= nightmare

Model predictive control	Singular-value decomposition
Eigenvalues	Mean shift algorithm for segmentation
LU factorization	Stencil computations
Optimal binary search organization	Displacement based algorithms
Image color conversions	Motion estimation
Image geometry transformations	Multiresolution classifier
Enclosing ball of points	Kalman filter
Metropolis algorithm, Monte Carlo	Object detection
Seam carving	IIR filters
SURF feature detection	Arithmetic for large numbers
Submodular function optimization	Optimal binary search organization
Graph cuts, Edmond-Karps Algorithm	Software defined radio
Gaussian filter	Shortest path problem
Black Scholes option pricing	Feature set for biomedical imaging
Disparity map refinement	Biometrics identification

Same for most computational problems:
Straightforward code is highly suboptimal

Optimization: Register Locality and ILP

```
// straightforward code
for(i = 0; i < N; i += 1)
  for(j = 0; j < N; j += 1)
    for(k = 0; k < N; k += 1)
      c[i][j] += a[i][k]*b[k][j];
```

Concise and slow



Removes aliasing
 Enables register allocation and instruction scheduling

Compiler does not do well:

- often illegal
- many choices

```
// unrolling + scalar replacement
for(i = 0; i < N; i += MU) {
  for(j = 0; j < N; j += NU) {
    for(k = 0; k < N; k += KU) {
      t1 = A[i*N + k];
      t2 = A[i*N + k + 1];
      t3 = A[i*N + k + 2];
      t4 = A[i*N + k + 3];
      t5 = A[(i + 1)*N + k];
      <more copies>

      t10 = t1 * t9;
      t17 = t17 + t10;
      t21 = t1 * t8;
      t18 = t18 + t21;
      t12 = t5 * t9;
      t19 = t19 + t12;
      t13 = t5 * t8;
      t20 = t20 + t13;
      <more ops>

      C[i*N + j] = t17;
      C[i*N + j + 1] = t18;
      C[(i+1)*N + j] = t19;
      C[(i+1)*N + j + 1] = t20;
    }
  }
}
```

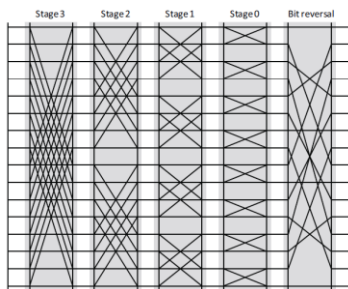
load

compute

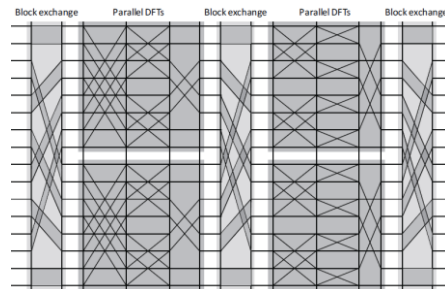
store

Ugly and fast

Optimization for Parallelism (Threads)



Parallelism is present, but is not in the “right shape”

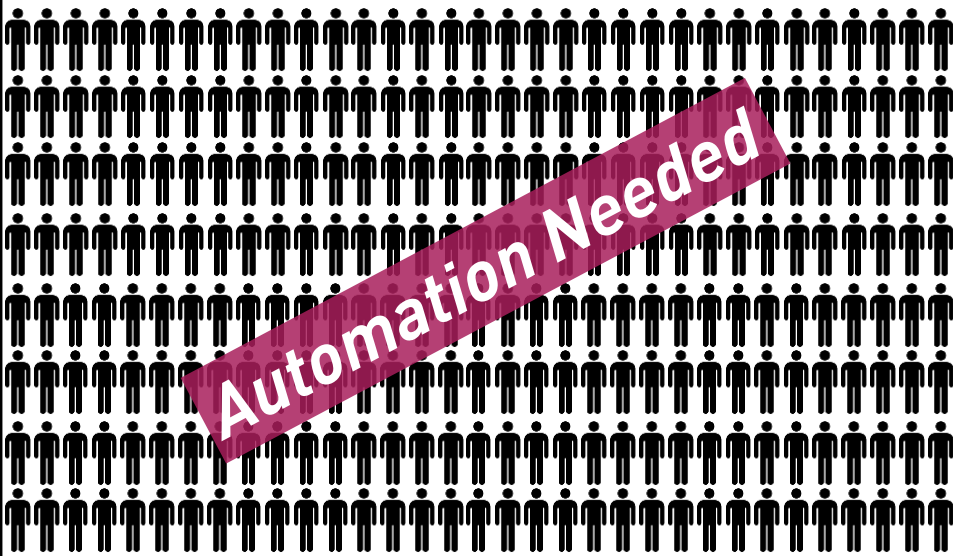


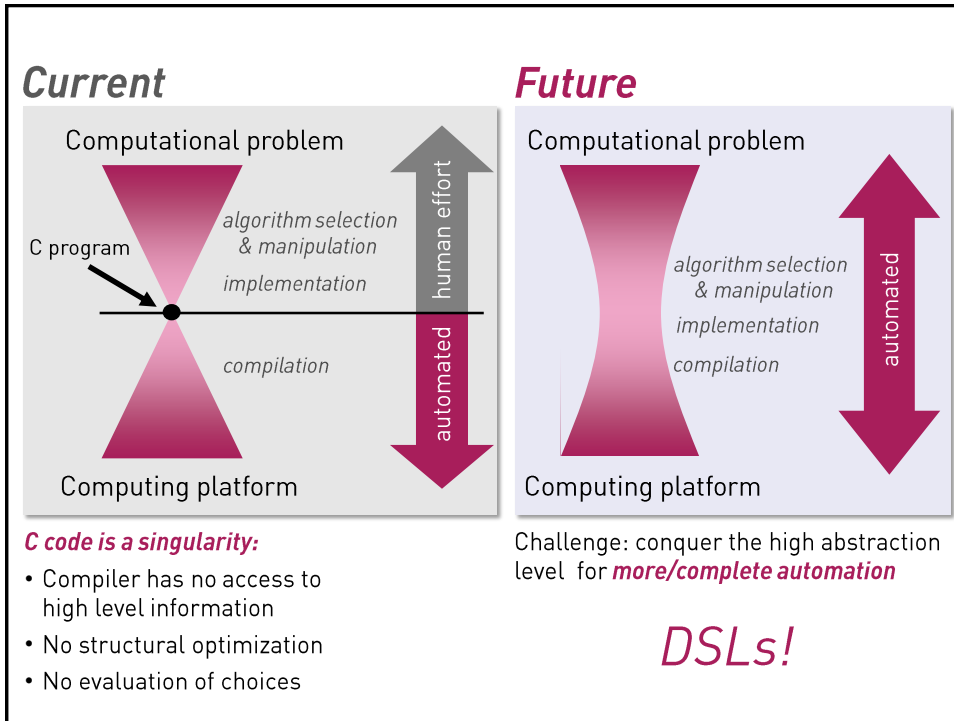
Restructured for locality and parallelism (shared memory, 2 cores, 2 elements per cache line)

Compiler usually does not do

- analysis may be unfeasible
- may require algorithm changes
- may require domain knowledge
- may require processor parameters

Current practice: Thousands of programmers re-implement and re-optimize the same functionality for every new processor and for every new processor generation





Goal:

Computer generation of high performance code for ubiquitous performance-critical components

Generate Code



Select convolutional code
Select a preset code or customize parameters

- custom
- Voyager
- NASA-DSN
- CCSDS/NASA-GSFC
- WiMax
- CDMA IS-95A
- LTE (3GPP - Long Term Evolution)
- UWB (802.15)
- CDMA 2000
- Cassini
- Mars Pathfinder & Stereo

rate: 1 / 2

K: 7

polynomials: 109

79

code rate (2)

constraint length (2)

polynomials for the code in decimal notation (2)

Select implementation options

frame length: 2048

unpadded frame length

Vectorization level: scalar C

type of code (2)

DFT IP Cores

parameter	value	range	explanation
Problem specification			
transform size	64	4-32768	Number of samples (2)
direction	forward		forward or inverse DFT (2)
data type	fixed point		fixed or floating point (2)
	16 bits	4-32 bits	fixed point precision (2)
	unscaled		scaling mode (2)
Parameters controlling implementation			
architecture	fully streaming		iterative or fully streaming (2)
radix	2	2, 4, 8, 16, 32, 64	size of DFT basic block (2)
streaming width	2	2-64	number of complex words per cycle (2)
data ordering	natural in / natural out		natural or digit-reversed data order (2)
BRAM budget	1000		maximum # of BRAMs to utilize (-1 for no limit) (2)

Viterbi Decoder

@ www.spiral.net

Possible Approach:

- Capturing algorithm knowledge:
Domain-specific languages (DSLs)
- Structural optimization:
Rewriting systems
- High performance code style:
Compiler
- Decision making for choices:
Machine learning

Spiral: Program Generation for Performance (www.spiral.net)



Franz Franchetti
 Yevgen Voronenko
 Jianxin Xiong
 Bryan Singer
 Srinivas Chellappa
 Frédéric de Mesmay
 Peter Milder
 José Moura
 David Padua
 Jeremy Johnson
 James Hoe
 <many more>

funding: DARPA, NSF, ONR, Intel

Linear Transforms

$$\begin{pmatrix} y_0 \\ y_1 \\ \vdots \\ y_{n-1} \end{pmatrix} = y = Tx \leftarrow \boxed{T} \leftarrow x = \begin{pmatrix} x_0 \\ x_1 \\ \vdots \\ x_{n-1} \end{pmatrix}$$

Output *Input*

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

Algorithms: Example FFT, n = 4

Fast Fourier transform (FFT):

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

12 adds, 4 mults

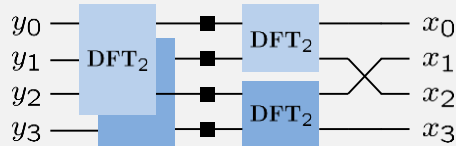
4 adds

1 mult

4 adds

0 adds/mults

Data flow graph



Description with matrix algebra (SPL)

$$\text{DFT}_4 = (\text{DFT}_2 \otimes I_2) \tau_2^4 (I_2 \otimes \text{DFT}_2) \iota_2^4$$

Decomposition Rules (>200 for >40 Transforms)

$$\begin{aligned} \text{DFT}_n &\rightarrow P_{k/2,2m}^\top (\text{DFT}_{2m} \oplus (J_{k/2-1} \otimes_i C_{2m} \text{rDFT}_{2m}(i/k))) (\text{RDFT}'_k \otimes I_m), \quad k \text{ even,} \\ \begin{bmatrix} \text{RDFT}'_n \\ \text{RDFT}''_n \\ \text{DHT}'_n \\ \text{DHT}''_n \end{bmatrix} &\rightarrow (P_{k/2,m}^\top \otimes I_2) \left(\begin{bmatrix} \text{RDFT}'_{2m} \\ \text{RDFT}''_{2m} \\ \text{DHT}'_{2m} \\ \text{DHT}''_{2m} \end{bmatrix} \oplus \left(I_{k/2-1} \otimes_i D_{2m} \begin{bmatrix} \text{rDFT}'_{2m}(i/k) \\ \text{rDFT}''_{2m}(i/k) \\ \text{rDHT}'_{2m}(i/k) \\ \text{rDHT}''_{2m}(i/k) \end{bmatrix} \right) \right) \left(\begin{bmatrix} \text{RDFT}'_k \\ \text{RDFT}''_k \\ \text{DHT}'_k \\ \text{DHT}''_k \end{bmatrix} \otimes I_m \right), \quad k \text{ even,} \\ \begin{bmatrix} \text{rDFT}'_{2n}(u) \\ \text{rDHT}'_{2n}(u) \end{bmatrix} &\rightarrow I_m^2 \left(I_k \otimes_i \begin{bmatrix} \text{rDFT}'_{2m}((i+u)/k) \\ \text{rDHT}'_{2m}((i+u)/k) \end{bmatrix} \right) \left(\begin{bmatrix} \text{rDFT}'_{2k}(u) \\ \text{rDHT}'_{2k}(u) \end{bmatrix} \otimes I_m \right), \\ \text{RDFT}'_{3n} &\rightarrow (Q_{k/2,m}^\top \otimes I_2) (I_k \otimes_i \text{rDFT}'_{2m}(i+1/2/k)) (\text{RDFT}'_{3k} \otimes I_m), \quad k \text{ even,} \\ \text{DCT}'_{2n} &\rightarrow P_{1,2,2m}^\top (\text{DCT}'_{2,2m} \oplus (I_{1,2,2} \otimes N_{2m} \text{RDFT}'_{3,2m})) B_n (I_{1,2}^{n/2} \otimes I_2) (I_m \otimes \text{RDFT}'_{Q_{m/2,k}}) \\ \text{DCT}'_{3n} &\rightarrow P_{1,2,2m}^\top (\text{DCT}'_{2,2m} \oplus (I_{1,2,2} \otimes N_{2m} \text{RDFT}'_{3,2m})) B_n (I_{1,2}^{n/2} \otimes I_2) (I_m \otimes \text{RDFT}'_{Q_{m/2,k}}) \\ \text{DCT}'_{4n} &\rightarrow P_{1,2,2m}^\top (\text{DCT}'_{2,2m} \oplus (I_{1,2,2} \otimes N_{2m} \text{RDFT}'_{3,2m})) B_n (I_{1,2}^{n/2} \otimes I_2) (I_m \otimes \text{RDFT}'_{Q_{m/2,k}}) \\ \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^\top (I_1 \oplus \text{DFT}'_{p-1}) D_p (I_1 \oplus \text{DFT}'_{p-1}) R_p, \quad p \text{ prime} \\ \text{DCT}'_{3n} &\rightarrow (I_m \oplus J_m) L_m (\text{DCT}'_{3m}(1/4) \oplus \text{DCT}'_{3m}(3/4)) \\ &\quad \cdot (F_2 \otimes I_m) \begin{bmatrix} I_m & 0 \oplus -J_{m-1} \\ \frac{1}{\sqrt{2}} (I_1 \oplus 2I_m) \end{bmatrix}, \quad n = 2m \\ \text{DCT}'_{4n} &\rightarrow S_n \text{DCT}'_{2n} \text{diag}_{0 \leq k < n} (1/(2 \cos((2k+1)\pi/4n))) \\ \text{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} \text{DCT}'_{4,2m} \\ \text{WHT}'_{2k} &\rightarrow \prod_{i=1}^k (I_{2^{k_i+1}} \otimes \text{WHT}'_{2^{k_i}} \otimes I_{2^{k_1+\dots+k_{i-1}}}), \quad k = k_1 + \dots + k_t \\ \text{DFT}'_2 &\rightarrow F_2 \\ \text{DCT}'_{2,2} &\rightarrow \text{diag}(1, 1/\sqrt{2}) F_2 \\ \text{DCT}'_{4,2} &\rightarrow J_2 R_{13\pi/8} \end{aligned}$$

Rules = algorithm knowledge

(≈100 journal papers)

SPL to Code

SPL S Pseudo code for $y = Sx$

$A_n B_n$ <code for: $t = Bx$ >
<code for: $y = At$ >

$I_m \otimes A_n$ for (i=0; i<m; i++)
<code for:
 $y[i*n:i*n+n-1] = A(x[i*n:i*n+n-1])$ >

$$I_m \otimes A_n = \begin{bmatrix} A_n & & \\ & \dots & \\ & & A_n \end{bmatrix}$$

$A_m \otimes I_n$ for (i=0; i<n; i++)
<code for:
 $y[i:n:i+m*n-n] = A(x[i:n:i+m*n-n])$ >

D_n for (i=0; i<n; i++)
 $y[i] = D[i]*x[i];$

$L_k^{k,m}$ for (i=0; i<k; i++)
 for (j=0; j<m; j++)
 $y[i*m+j] = x[j*k+i];$

F_2 $y[0] = x[0] + x[1];$
 $y[1] = x[0] - x[1];$

Gives reasonable, straightforward code

Program Generation in Spiral

Transform

DFT_8

Decomposition rules (algorithm knowledge)

Algorithm
(SPL)

$(\text{DFT}_2 \otimes I_4) T_4^8 (I_2 \otimes ((\text{DFT}_2 \otimes I_2)$
 $T_2^4 (I_2 \otimes \text{DFT}_2) L_2^4)) L_2^8$

parallelization
vectorization

Algorithm
(-SPL)

$\sum (S_j \text{DFT}_2 G_j) \sum (\sum (S_{k,l} \text{diag}(t_{k,l}) \text{DFT}_2 G_l)$
 $\sum (S_m \text{diag}(t_m) \text{DFT}_2 G_{k,m}))$

locality
optimization

C Program

```
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;
< more lines >
```

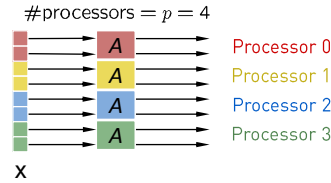
code style
code level
optimization

+ Search or
Learning for
Choices

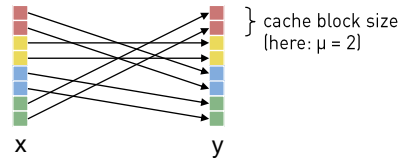
SPL to Shared Memory Code: Basic Idea

“Good” SPL structures

$$y = (I_p \otimes A)x$$



$$y = (P \otimes I_\mu)x$$

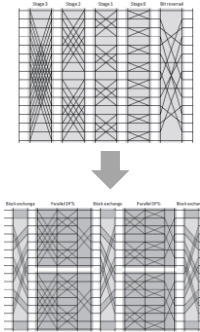


Rewriting: Bad structures \rightarrow good structures

Example: SMP Parallelization

Franchetti, Voronenko & P, SC 2006

$$\begin{aligned} \underbrace{\text{DFT}_{mn}}_{\text{smp}(p,\mu)} &\rightarrow \underbrace{((\text{DFT}_m \otimes I_n) \tau_n^{mn} (I_m \otimes \text{DFT}_n) L_m^{mn})}_{\text{smp}(p,\mu)} \\ &\dots \\ &\rightarrow \underbrace{(\text{DFT}_m \otimes I_n)}_{\text{smp}(p,\mu)} \underbrace{\tau_n^{mn}}_{\text{smp}(p,\mu)} \underbrace{(I_m \otimes \text{DFT}_n)}_{\text{smp}(p,\mu)} \underbrace{L_m^{mn}}_{\text{smp}(p,\mu)} \\ &\dots \\ &\rightarrow \underbrace{((L_m^{mp} \otimes I_{n/p\mu}) \otimes I_\mu)}_{\text{smp}(p,\mu)} \underbrace{(I_p \otimes (\text{DFT}_m \otimes I_{n/p}))}_{\text{smp}(p,\mu)} \underbrace{((L_p^{mp} \otimes I_{n/p\mu}) \otimes I_\mu)}_{\text{smp}(p,\mu)} \\ &\quad \underbrace{\left(\bigoplus_{i=0}^{p-1} \tau_n^{mn,i}\right)}_{\text{smp}(p,\mu)} \underbrace{(I_p \otimes (I_{m/p} \otimes \text{DFT}_n))}_{\text{smp}(p,\mu)} \underbrace{(I_p \otimes L_{m/p}^{mn/p})}_{\text{smp}(p,\mu)} \underbrace{((L_p^{pn} \otimes I_{m/p\mu}) \otimes I_\mu)}_{\text{smp}(p,\mu)} \end{aligned}$$



load-balanced, no false sharing

**One rewriting system for every platform paradigm:
SIMD, distributed memory parallelism, FPGA, ...**

Same Approach for Different Paradigms

Threading:

$$\begin{aligned} \frac{\text{DFT}_{mn}}{\text{smp}(\mu,\mu)} &\rightarrow \frac{(\text{DFT}_m \otimes I_n) \text{T}_{mn}^{\text{inv}} (I_m \otimes \text{DFT}_n) \text{L}_{mn}^{\text{inv}}}{\text{smp}(\mu,\mu)} \\ &\dots \\ &\rightarrow \frac{(\text{DFT}_m \otimes I_n)}{\text{smp}(\mu,\mu)} \frac{\text{T}_{mn}^{\text{inv}}}{\text{smp}(\mu,\mu)} \frac{(I_m \otimes \text{DFT}_n)}{\text{smp}(\mu,\mu)} \frac{\text{L}_{mn}^{\text{inv}}}{\text{smp}(\mu,\mu)} \\ &\dots \\ &\rightarrow \left((L_{mn/p}^{\text{inv}} \otimes I_{n/p}) \otimes_p I_p \right) \left(I_p \otimes_p (\text{DFT}_m \otimes I_{n/p}) \right) \left((L_{m/p}^{\text{inv}} \otimes I_{n/p}) \otimes_p I_p \right) \\ &\quad \left(\bigoplus_{i=0}^{p-1} \text{T}_{mn,i}^{\text{inv}} \right) \left(I_p \otimes_p (I_{m/p} \otimes \text{DFT}_n) \right) \left(I_p \otimes_p (L_{m/p}^{\text{inv}}) \right) \left((L_{m/p}^{\text{inv}} \otimes I_{n/p}) \otimes_p I_p \right) \end{aligned}$$

Vectorization:

$$\begin{aligned} \frac{(\text{DFT}_{mn})}{\text{vec}(r)} &\rightarrow \frac{(\text{DFT}_m \otimes I_n) \text{T}_{mn}^{\text{inv}} (I_m \otimes \text{DFT}_n) \text{L}_{mn}^{\text{inv}}}{\text{vec}(r)} \\ &\dots \\ &\rightarrow \frac{(\text{DFT}_m \otimes I_n)}{\text{vec}(r)} \left(\frac{\text{T}_{mn}^{\text{inv}}}{\text{vec}(r)} \right) \frac{(I_m \otimes \text{DFT}_n) \text{L}_{mn}^{\text{inv}}}{\text{vec}(r)} \\ &\dots \\ &\rightarrow \left((L_{mn/p} \otimes I_{n/p}^2) \otimes_{\text{SSE}} (\text{DFT}_m \otimes I_{n/p} \otimes I_p) \right) \left(\frac{\text{T}_{mn}^{\text{inv}}}{\text{SSE}} \right) \\ &\quad \left(I_{m/p} \otimes (I_p^2 \otimes I_p) \right) \left(I_{n/p} \otimes (L_{n/p}^2 \otimes I_p) \right) \left(I_2 \otimes \frac{\text{L}_{m/p}^{\text{inv}}}{\text{SSE}} \right) \left((L_{m/p}^{\text{inv}} \otimes I_p) \right) \left(\text{DFT}_n \otimes I_p \right) \\ &\quad \left((L_{m/p}^{\text{inv}} \otimes I_2) \otimes I_p \right) \left(I_{m/p} \otimes \frac{\text{L}_{n/p}^2}{\text{SSE}} \right) \end{aligned}$$

GPUs:

$$\begin{aligned} \frac{(\text{DFT}_{p,k})}{\text{gpu}(r,c)} &\rightarrow \frac{\left(\prod_{i=0}^{k-1} L_r^{\text{inv}} (I_{k-i} \otimes \text{DFT}_r) \right) \left(L_{k-i-1} (I_r \otimes \text{T}_{k-i-1}^{\text{inv}}) L_{r+i}^{\text{inv}} \right) R_r^{\text{inv}}}{\text{gpu}(r,c)} \\ &\dots \\ &\rightarrow \frac{\left(\prod_{i=0}^{k-1} (L_r^{\text{inv}2} \otimes I_2) \right) (I_{k-1/2} \otimes (\text{DFT}_r \otimes I_2) L_r^{\text{inv}2}) \text{T}_i}{\text{shd}(r,c)} \\ &\quad \left((L_r^{\text{inv}2} \otimes I_2) (I_{k-1/2} \otimes \frac{\text{L}_r^{\text{inv}2}}{\text{shd}(r,c)}) (R_r^{\text{inv}2} \otimes I_r) \right) \end{aligned}$$

Verilog for FPGAs:

$$\begin{aligned} \frac{(\text{DFT}_{p,k})}{\text{stream}(r^*)} &\rightarrow \frac{\left[\prod_{i=0}^{k-1} L_r^{\text{inv}} (I_{k-i} \otimes \text{DFT}_r) \right] \left(L_{k-i-1} (I_r \otimes \text{T}_{k-i-1}^{\text{inv}}) L_{r+i}^{\text{inv}} \right) R_r^{\text{inv}}}{\text{stream}(r^*)} \\ &\dots \\ &\rightarrow \frac{\left[\prod_{i=0}^{k-1} \frac{\text{L}_r^{\text{inv}}}{\text{stream}(r^*)} \left(I_{k-i} \otimes \text{DFT}_r \right) \right] \left(L_{k-i-1} (I_r \otimes \text{T}_{k-i-1}^{\text{inv}}) L_{r+i}^{\text{inv}} \right) R_r^{\text{inv}}}{\text{stream}(r^*)} \\ &\dots \\ &\rightarrow \frac{\left[\prod_{i=0}^{k-1} \frac{\text{L}_r^{\text{inv}}}{\text{stream}(r^*)} \left(I_{k-i-1} \otimes (I_{r+1} \otimes \text{DFT}_r) \right) \right] \text{T}_i}{\text{stream}(r^*)} R_r^{\text{inv}} \end{aligned}$$

- Rigorous, correct by construction
- Overcomes compiler limitations

Computer generated Functions for Intel IPP



3984 C functions
1M lines of code

Transforms: DFT (fwd+inv), RDFT (fwd+inv), DCT2, DCT3, DCT4, DHT, WHT
Sizes: 2-64 (DFT, RDFT, DHT); 2-powers (DCTs, WHT)
Precision: single, double
Data type: scalar, SSE, AVX (DFT, DCT), LRB (DFT)

Computer generated

Results: SpiralGen Inc.

Challenge: General Size Libraries

So far:

Code specialized to fixed input size

```
DFT_384(x, y) {
  ...
  for(i = ...) {
    t[2i] = x[2i] + x[2i+1]
    t[2i+1] = x[2i] - x[2i+1]
  }
  ...
}
```

- Algorithm fixed
- Nonrecursive code

Challenge:

Library for general input size

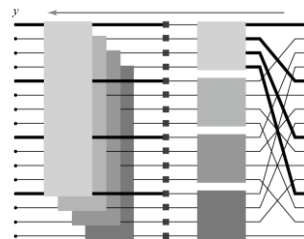
```
DFT(n, x, y) {
  ...
  for(i = ...) {
    DFT_strided(m, x+mi, y+i, 1, k)
  }
  ...
}
```

- Algorithm cannot be fixed
- Recursive code
- Creates many challenges

Challenge: Recursive Steps Needed

Cooley-Tukey FFT

$$y = (\text{DFT}_k \otimes I_m) T_m^{km} (I_k \otimes \text{DFT}_m) L_k^{km} x$$



Implementation that increases locality (e.g., FFTW 2.x)

```
void DFT(int n, cpx *y, cpx *x) {
  int k = choose_dft_radix(n);

  for (int i=0; i < k; ++i)
    DFTrec(m, y + m*i, x + i, k, 1);
  for (int j=0; j < m; ++j)
    DFTscaled(k, y + j, t[j], m);
}
```

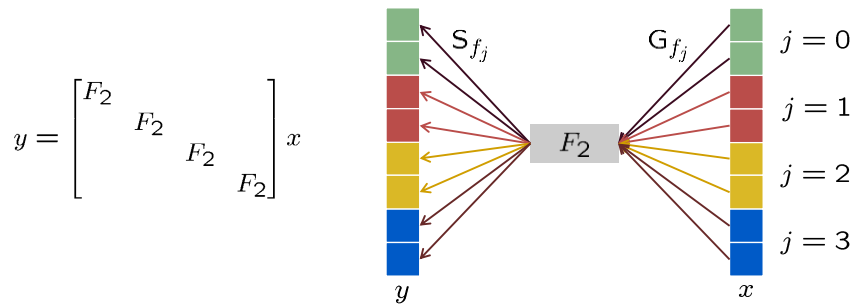
Σ-SPL : Basic Idea

Four additional matrix constructs: Σ , G , S , Perm

- Σ (sum) matrix sum (explicit loop)
- G_f (gather) load data with index mapping f
- S_f (scatter) store data with index mapping f
- Perm_f permute data with the index mapping f

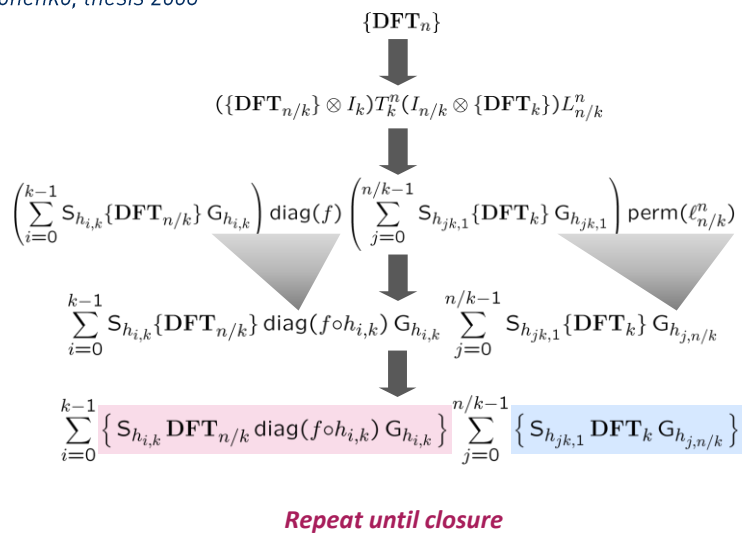
Σ-SPL formulas = matrix factorizations

Example: $y = (I_4 \otimes F_2)x \rightarrow y = \sum_{j=0}^3 S_{f_j} F_2 G_{f_j} x$



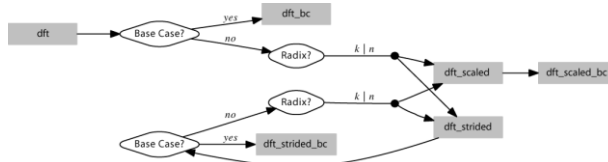
Find Recursion Step Closure

Voronenko, thesis 2008

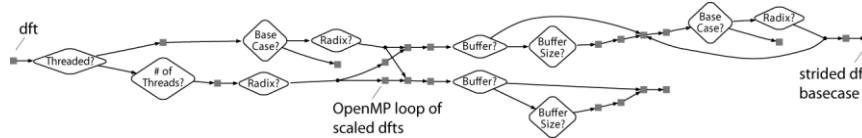


Recursion Step Closure: Examples

DFT: scalar code



DFT: full-fledged (vectorized and parallel code)



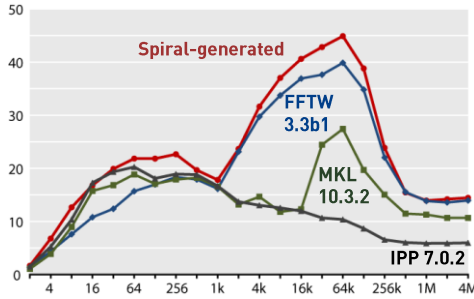
Generating Dozens of “FFTWs”

$$\begin{aligned}
 \text{DFT}_n &= P_{(2,2)}^{(2,2)} \left[\text{DFT}_{2m} \oplus (k_2 \oplus 1) \otimes \text{DFT}_{2m}(i/k) \right] \left[\text{RDFT}_k^T \otimes I_n \right], \quad k \text{ even.} \\
 \text{RDFT}_n &= \left(P_{(2,2)}^{(2,2)} \otimes I_2 \right) \left[\text{RDFT}_{2m}^T \oplus \left(\begin{array}{c} \text{RDFT}_{2m}^T(i/k) \\ \text{RDFT}_{2m}^T(i/k) \end{array} \right) \right] \left(\text{RDFT}_k^T \otimes I_n \right), \quad k \text{ even.} \\
 \text{DHT}_n &= \left(P_{(2,2)}^{(2,2)} \otimes I_2 \right) \left[\text{DHT}_{2m}^T \oplus \left(\begin{array}{c} \text{DHT}_{2m}^T(i/k) \\ \text{DHT}_{2m}^T(i/k) \end{array} \right) \right] \left(\text{DHT}_k^T \otimes I_n \right), \quad k \text{ even.} \\
 \text{DHT}_n &= \left(P_{(2,2)}^{(2,2)} \otimes I_2 \right) \left[\text{DHT}_{2m}^T \oplus \left(\begin{array}{c} \text{DHT}_{2m}^T(i/k) \\ \text{DHT}_{2m}^T(i/k) \end{array} \right) \right] \left(\text{DHT}_k^T \otimes I_n \right), \quad k \text{ even.} \\
 \text{RDFT}_n &= \left(P_{(2,2)}^{(2,2)} \otimes I_2 \right) \left[\text{RDFT}_{2m} \oplus (i+1/2)/k \right] \left[\text{RDFT}_k \otimes I_n \right], \quad k \text{ even.} \\
 \text{DCT-2}_n &= \left(P_{(2,2)}^{(2,2)} \otimes I_2 \right) \left[\text{DCT-2}_{2m} \oplus (i+1/2)/k \right] \left[\text{RDFT}_k \otimes I_n \right], \quad k \text{ even.} \\
 \text{DCT-3}_n &= \left(P_{(2,2)}^{(2,2)} \otimes I_2 \right) \left[\text{DCT-3}_{2m} \oplus (i+1/2)/k \right] \left[\text{RDFT}_k \otimes I_n \right], \quad k \text{ even.} \\
 \text{DCT-4}_n &= \left(P_{(2,2)}^{(2,2)} \otimes I_2 \right) \left[\text{DCT-4}_{2m} \oplus (i+1/2)/k \right] \left[\text{RDFT}_k \otimes I_n \right], \quad k \text{ even.} \\
 \text{DFT}_n &= \left(\text{DFT}_n \otimes I_n \right) \left[\text{DFT}_n \oplus \text{DFT}_n \right], \quad n = 4m \\
 \text{DFT}_n &= \left(\text{DFT}_n \otimes I_n \right) \left[\text{DFT}_n \oplus \text{DFT}_n \right], \quad n = 4m, \text{ gcd}(k, m) = 1 \\
 \text{DFT}_n &= \left(\text{DFT}_n \otimes I_n \right) \left[\text{DFT}_n \oplus \text{DFT}_n \right], \quad n = 4m, \text{ gcd}(k, m) = 1 \\
 \text{DCT-3}_n &= \left(\text{DCT-3}_n \otimes I_n \right) \left[\text{DCT-3}_n \oplus \text{DCT-3}_n \right], \quad n = 2m \\
 \text{DCT-4}_n &= \left(\text{DCT-4}_n \otimes I_n \right) \left[\text{DCT-4}_n \oplus \text{DCT-4}_n \right], \quad n = 2m \\
 \text{DCT-4}_n &= \left(\text{DCT-4}_n \otimes I_n \right) \left[\text{DCT-4}_n \oplus \text{DCT-4}_n \right], \quad n = 2m \\
 \text{IMDCT}_{2m} &= \left(\text{DCT-4}_{2m} \otimes I_n \right) \left[\text{DCT-4}_{2m} \oplus \text{DCT-4}_{2m} \right], \quad n = 2m \\
 \text{WHT}_p &= \left[\left(\text{DFT}_p \otimes I_n \right) \oplus \text{WHT}_{2p} \otimes I_{2p} \right] \left[\text{DFT}_p \oplus \text{DFT}_p \right], \quad k = k_1 + \dots + k_p \\
 \text{DFT}_2 &= \text{DFT}_2 \\
 \text{DCT-2}_2 &= \text{DCT-2}_2 \\
 \text{DCT-4}_2 &= \text{DCT-4}_2
 \end{aligned}$$


Transform	Code size	
	non-parallelized	parallelized
<i>no vectorization</i>		
DFT	13.1 KLOC / 0.59 MB	10.3 KLOC / 0.45 MB
RDFT	8.5 KLOC / 0.36 MB	8.8 KLOC / 0.39 MB
DHT	9.1 KLOC / 0.40 MB	9.4 KLOC / 0.39 MB
DCT-2	12.0 KLOC / 0.55 MB	12.4 KLOC / 0.57 MB
DCT-3	12.0 KLOC / 0.56 MB	12.3 KLOC / 0.59 MB
DCT-4	6.8 KLOC / 0.33 MB	7.1 KLOC / 0.35 MB
WHT	5.6 KLOC / 0.21 MB	—
<i>2-way vectorization</i>		
DFT	14.8 KLOC / 0.73 MB	15.0 KLOC / 0.74 MB
RDFT	13.6 KLOC / 0.76 MB	16.0 KLOC / 0.81 MB
scaled RDFT	16.0 KLOC / 0.78 MB	—
DHT	16.9 KLOC / 0.83 MB	17.2 KLOC / 0.87 MB
DCT-2	20.7 KLOC / 1.10 MB	21.0 KLOC / 1.09 MB
DCT-3	27.9 KLOC / 1.56 MB	28.2 KLOC / 1.59 MB
DCT-4	7.8 KLOC / 0.47 MB	8.1 KLOC / 0.50 MB
WHT	6.9 KLOC / 0.32 MB	5.8 KLOC / 0.26 MB
FIR Filter	167 KLOC / 7.75 MB	120 KLOC / 5.12 MB
Downsampled FIR Filter	100 KLOC / 4.2 MB	68 KLOC / 2.76 MB
<i>4-way vectorization</i>		
DFT	17.9 KLOC / 1.09 MB	18.2 KLOC / 1.11 MB
RDFT	16.2 KLOC / 0.86 MB	16.5 KLOC / 0.91 MB
scaled RDFT	16.5 KLOC / 0.88 MB	—
DHT	17.9 KLOC / 1.02 MB	18.3 KLOC / 1.04 MB
DCT-2	23.3 KLOC / 1.50 MB	23.6 KLOC / 1.53 MB
DCT-3	32.0 KLOC / 2.17 MB	32.3 KLOC / 2.20 MB
DCT-4	8.3 KLOC / 0.63 MB	8.6 KLOC / 0.66 MB
WHT	8.5 KLOC / 0.53 MB	6.9 KLOC / 0.4 MB
2D DFT	20.6 KLOC / 1.32 MB	20.8 KLOC / 1.33 MB
2D DCT-2	27.0 KLOC / 2.1 MB	27.2 KLOC / 2.11 MB
FIR Filter	109 KLOC / 5.69 MB	74 KLOC / 3.4 MB
Downsampled FIR Filter	151 KLOC / 7.7 MB	92 KLOC / 4.61 MB

It Really Works

DFT on Sandybridge (3.3 GHz, 4 Cores, AVX)
Performance [Gflop/s]

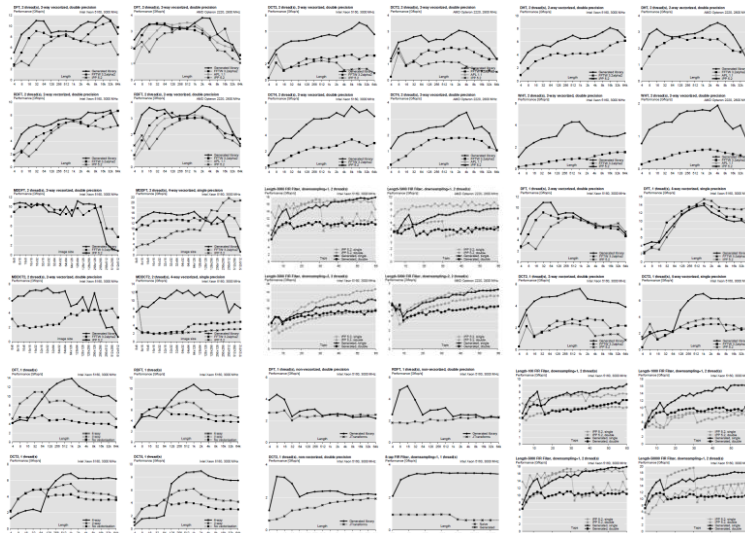


$$\begin{aligned}
 \text{DFT}_N &\rightarrow (\text{DFT}_k \otimes I_m) Y_m^T (I_k \otimes \text{DFT}_m) L_k^T \\
 \text{DFT}_N &\rightarrow P_{k/2,2m}^T (\text{DFT}_{2m} \otimes (I_{k/2-1} \otimes C_{2m} \text{rDFT}_{2m}(l/k))) (\text{RDFT}_k \otimes I_m) \\
 \text{RDFT}_N &\rightarrow (P_{k/2,2m}^T \otimes I_2) (\text{RDFT}_{2m} \otimes (I_{k/2-1} \otimes D_{2m} \text{rDFT}_{2m}(l/k))) (\text{RDFT}_k \otimes I_m) \\
 \text{rDFT}_{2m}(u) &\rightarrow L_m^T (I_k \otimes \text{rDFT}_{2m}(l+u/k)) (\text{rDFT}_{2k}(u) \otimes I_m)
 \end{aligned}$$

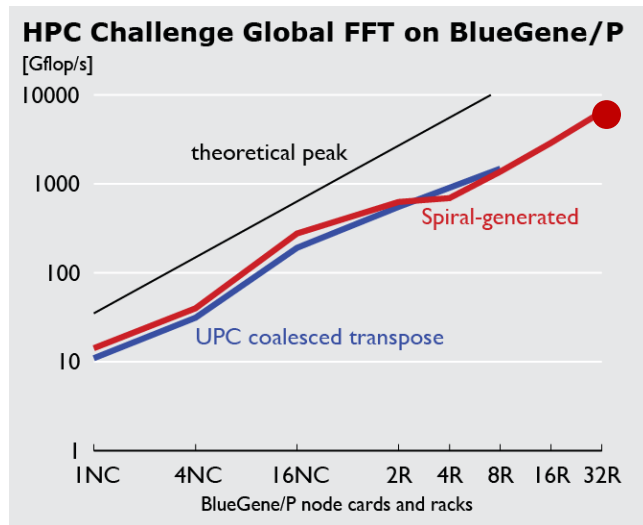


vectorized, threaded,
platform-tuned, adaptive library
(5 MB source code)

Generating Dozens of “FFTWs”



Very Large Scale: BG/P



6.4 Tflop/s

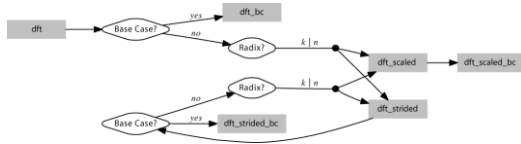
32 racks
= 32K node cards
= 128K cores

Is everything automated?

Yes, except efficient tuning

32

Simple Generated Library (~FFTW 2.x)



```

void dft(int n, cpx *y, cpx *x) {
  if (use_dft_base_case(n))
    dft_bc(n, y, x);
  else {
    int k = choose_dft_radix(n);
    for (int i=0; i < k; ++i)
      dft_strided(m, k, t + m*i, x + m*i);
    for (int i=0; i < m; ++i)
      dft_scaled(k, m, precomp_d[i], y + i, t + i);
  }
}

```

Choices used for tuning/adaptation

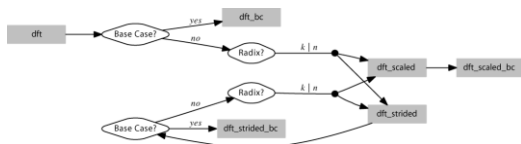
Recursion

```

void dft_strided(int n, int istr, cpx *y, cpx *x) { ... }
void dft_scaled(int n, int str, cpx *d, cpx *y, cpx *x) { ... }

```

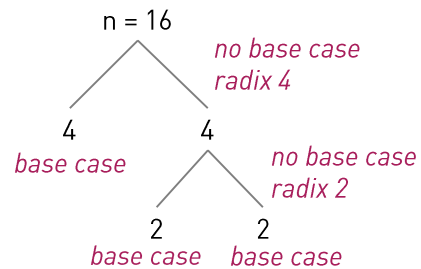
FFT Search Space (Simple Library)



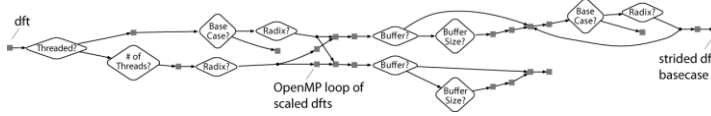
Recursive choice:

$$n = 2^k \begin{matrix} \text{base case?} \\ \text{radix?} \end{matrix}$$

Example selections for $n = 16$:



FFT Search Space (The Real Thing)

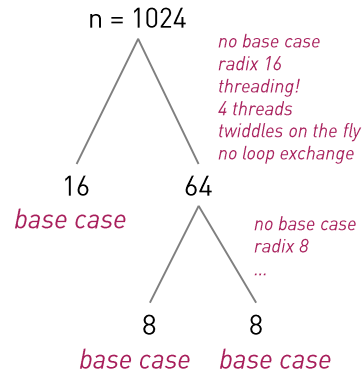


Recursive choice:

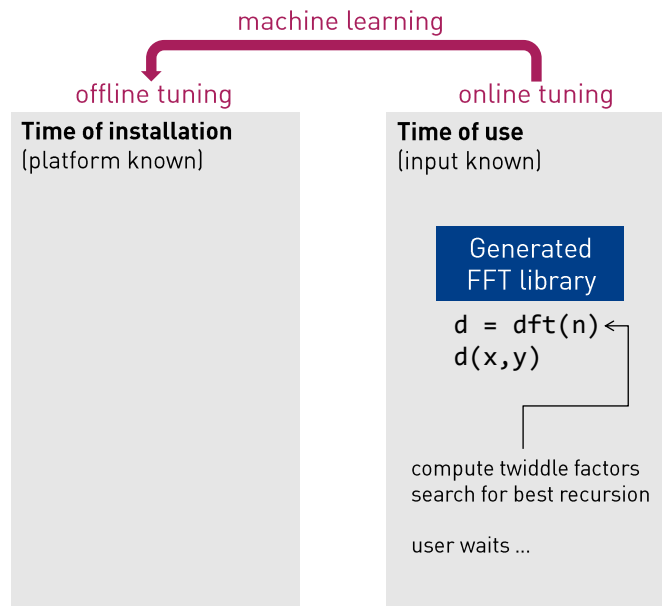
$$n = 2^k$$

*base case?
radix?
threading?
#threads?
twiddles?
loop exchange?*

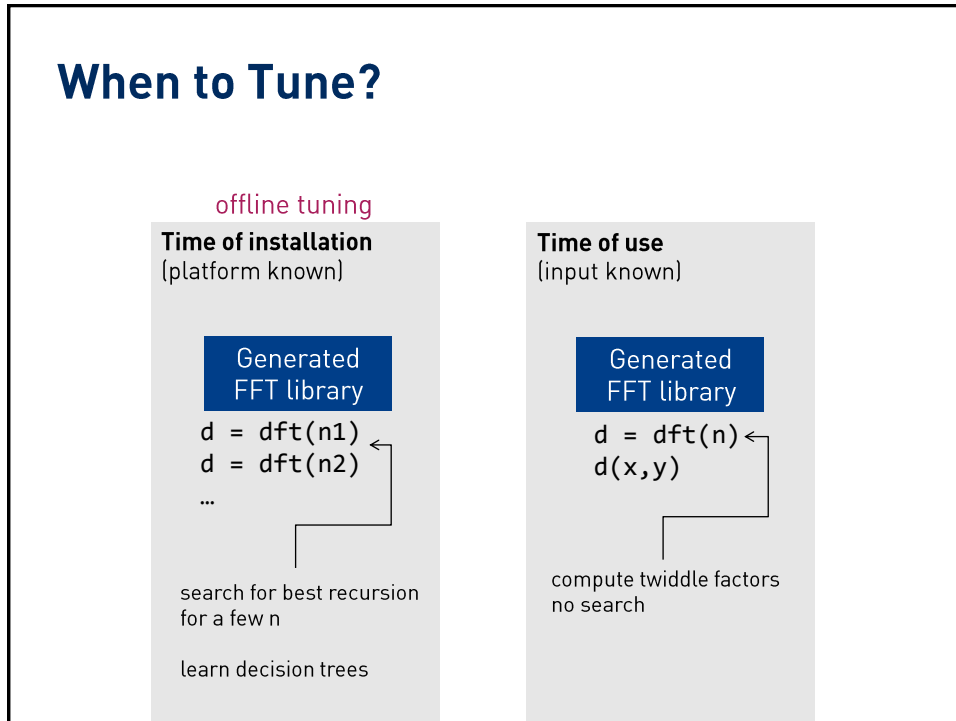
Example selections for $n = 1024$:



When to Tune?



When to Tune?



Online Tuning → Offline Tuning

At installation time, run search for a few n

Learn decision trees

Insert into library

```
void dft(int n, cpx *y, cpx *x) {
  if (use_dft_base_case(n))
    dft_bc(n, y, x);
  else {
    int k = choose_dft_radix(n);
    for (int i=0; i < k; ++i)
      dft_strided(m, k, t + m*i, x + m*i);
    for (int i=0; i < m; ++i)
      dft_scaled(k, m, precomp_d[i], y + i, t + i);
  }
}
void dft_strided(int n, int istr, cpx *y, cpx *x) { ... }
void dft_scaled(int n, int str, cpx *d, cpx *y, cpx *x) { ... }
```

Choices used for tuning/adaptation

Online Tuning → Offline Tuning

At installation time, run search for a few n

Learn decision trees

Insert into library

```
void dft(int n, cpx *y, cpx *x) {
    if (use_dft_base_case(n))
        dft_bc(n, y, x);
    else {
        int k = choose_dft_radix(n);
        for (int i=0; i < k; ++i)
            dft_strided(m, k, t + m*i, x + m*i);
        for (int i=0; i < m; ++i)
            dft_scaled(k, m, precomp_d[i], y + i, t + i);
    }
}

void dft_strided(int n, int istr, cpx *y, cpx *x) { ... }
void dft_scaled(int n, int str, cpx *d, cpx *y, cpx *x) { ... }
```

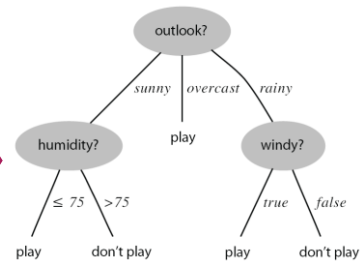
Example decision tree

```
if ( n <= 65536 ) {
    if ( n <= 32 ) {
        if ( n <= 4 ) {return 2;}
        else {return 4;}
    }
    else {
        if ( n <= 1024 ) {
            if ( n <= 256 ) {return 8;}
            else {return 32;}
        }
        else {
            .....
        }
    }
}
```

Decision Tree Generation: C4.5

Features (events)

Outlook	Temperature	Humidity	Windy	Decision
sunny	85	85	false	don't play
sunny	80	90	true	don't play
overcast	83	78	false	play
rain	70	96	false	play
rain	68	80	false	play
rain	65	70	true	don't play
overcast	64	65	true	play
sunny	72	95	false	don't play
sunny	69	70	false	play
rain	75	80	false	play
sunny	75	70	true	play
overcast	72	90	true	play
overcast	81	75	false	play
rain	71	80	true	don't play



- P(play|windy=false) = 6/8
- P(don't play|windy=false) = 2/8
- P(play|windy=true) = 1/2
- P(don't play|windy=true) = 1/2



Entropy of Features
 H(windy) = 0.89
H(outlook) = 0.69
 H(humidity) = ...

Application to Libraries

Features = arguments of functions (except variable pointers)

```
dft(int n, cpx *y, cpx *x)
dft_strided(int n, int istr, cpx *y, cpx *x)
dft_scaled(int n, int str, cpx *d, cpx *y, cpx *x)
etc.
```

Experimental Setup

3GHz Intel Xeon 5160 (2 Core 2 Duos = 4 cores)

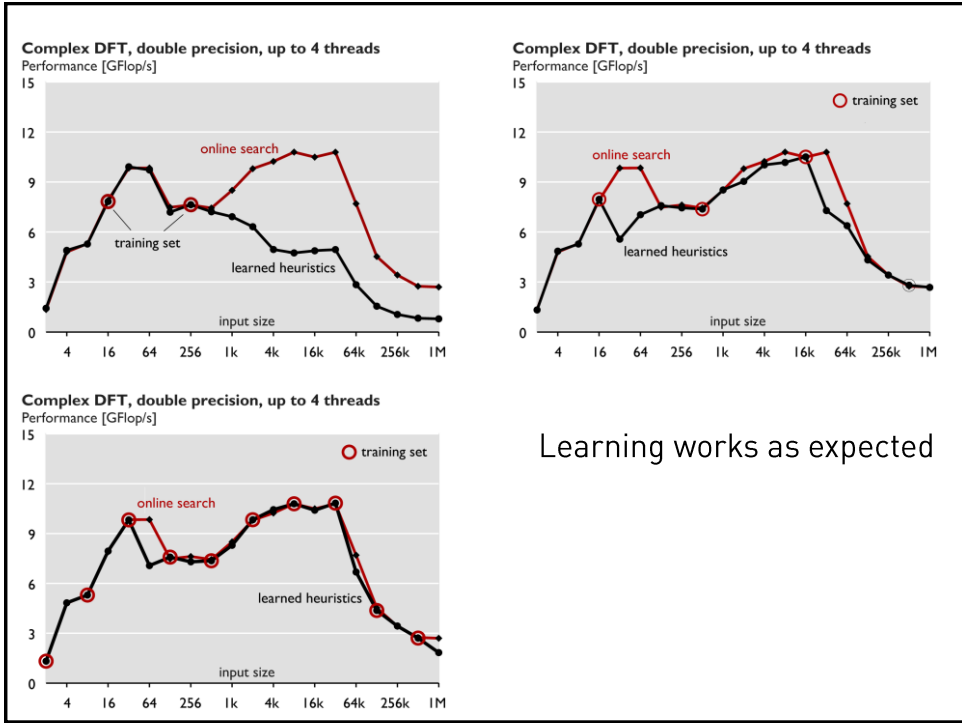
icc 10.1

Spiral-generated “FTW” [Voronenko et al., CGO, 2009]:

Recursive choice:

$n = 2^k$

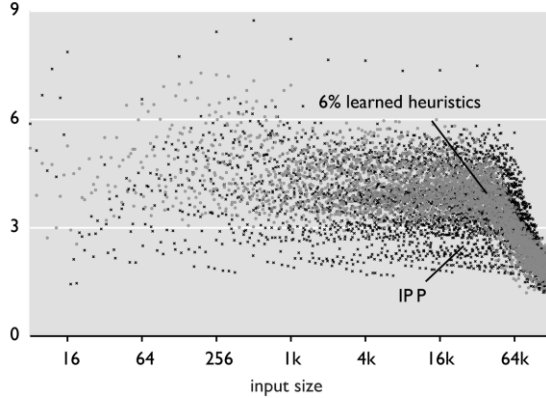
- base case?*
- radix?*
- threading?*
- #threads?*
- twiddles?*
- loop exchange?*



“All” Sizes

Complex DFT, double precision, mixed sizes

Performance [GFlop/s]

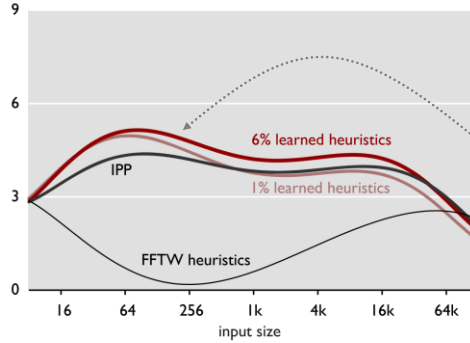


All sizes $n \leq 2^{18}$, with prime factors ≤ 19

“All” Sizes

$$\begin{aligned}
 \text{DFT}_n &= (\text{DFT}_k \otimes I_m) T_m^T (I_k \otimes \text{DFT}_m) L_k^T \\
 \text{DFT}_n &= P_{1/2,2m}^T (\text{DFT}_{2m} \oplus (I_{k/2-1} \otimes C_{2m} \text{rDFT}_{2m}(i/k))) (\text{RDFT}_k \otimes I_m) \\
 \text{RDFT}_n &= (P_{1/2,m}^T \otimes I_2) (\text{RDFT}_{2m} \oplus (I_{k/2-1} \otimes D_{2m} \text{rDFT}_{2m}(i/k))) (\text{RDFT}_k \otimes I_m) \\
 \text{rDFT}_{2m}(u) &= L_{2m}^T (I_k \otimes \text{rDFT}_{2m}((i+u)/k)) (\text{rDFT}_{2k}(u) \otimes I_m)
 \end{aligned}$$

Complex DFT, double precision, mixed sizes
Performance [GFlop/s]



↓ Spiral
vectorized, threaded,
platform-tuned, **adaptive** library
(5 MB source code)

↓ Learning
vectorized, threaded,
platform-tuned, **adapted** library
(5 MB source code)

All sizes $n \leq 2^{18}$, with prime factors ≤ 19

Higher order fit of all sizes

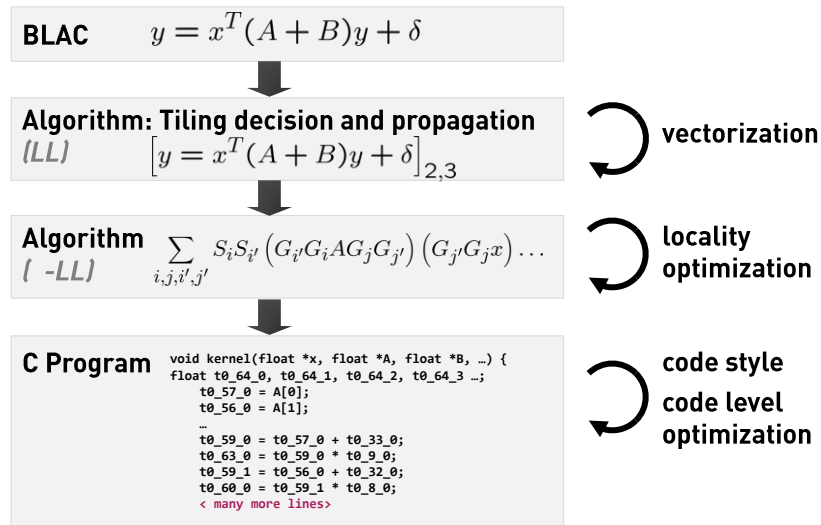
FFT & Co. are solved, what now?

1. Generators for more domains
2. Support for building such generators



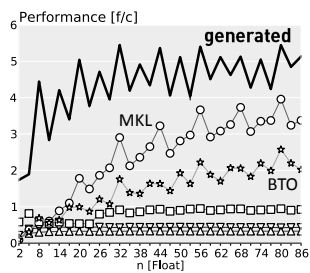
LGen: Generator for Basic Linear Algebra

Spampinato & P, CGO 2014



LGen: Sample Results

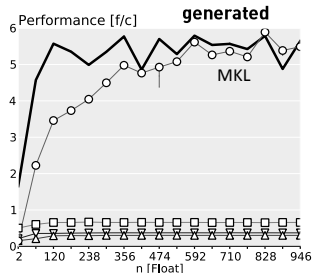
$$C = \alpha AB + \beta C$$



$$A \in \mathbb{R}^{n \times 4}$$

$$B \in \mathbb{R}^{4 \times n}$$

$$C = \alpha(A_0 + A_1)^T B + \beta C$$



$$A_0 \in \mathbb{R}^{4 \times 4}$$

$$B \in \mathbb{R}^{4 \times n}$$

- LGen
- ▽ Handwritten fixed size
- △ Handwritten gen size
- MKL 11.0
- Eigen 3.1.3
- * BTO 1.3
- ◇ IPP 7.1

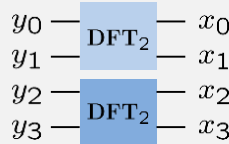
PL Support: Example Code Style

Ofenbeck, Rompf, Stojanov, Odersky & P, GPCE 2012



SPL $y = (I_2 \otimes \text{DFT}_2)x$

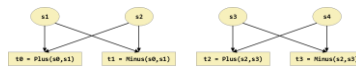
Data flow graph



Scala function

```
def f(x: Array[Double], y: Array[Double]) = {
  for (i <- 0 until 2) {
    y(2*i) = x(i*2) + x(i*2+1)
    y(2*i+1) = x(i*2) - x(i*2+1)
  }
}
```

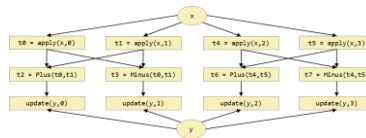
```
def f(x: Array[Rep[Double]],
      y: Array[Rep[Double]]) = {
  for (i <- 0 until 2) {
    y(2*i) = x(i*2) + x(i*2+1)
    y(2*i+1) = x(i*2) - x(i*2+1)
  }
}
```



scalarized

```
t0 = s0 + s1;
t1 = s0 - s1;
t2 = s2 + s3;
t2 = s2 - s3;
```

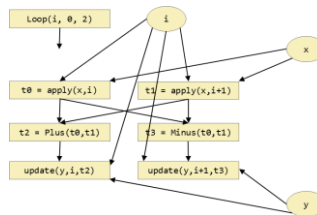
```
def f(x: Rep[Array[Double]],
      y: Rep[Array[Double]]) = {
  for (i <- 0 until 2) {
    y(2*i) = x(i*2) + x(i*2+1)
    y(2*i+1) = x(i*2) - x(i*2+1)
  }
}
```



unrolled, scalar repl.

```
t0 = x[0];
t1 = x[1];
t2 = t0 + t1;
t3 = t0 - t1;
t4 = x[0];
t5 = x[1];
t6 = t4 + x5;
y[0] = t6;
t7 = t4 - x5;
y[3] = t7;
```

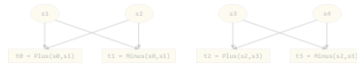
```
def f(x: Rep[Array[Double]],
      y: Rep[Array[Double]]) = {
  for (i <- 0 until 2: Rep[Range]) {
    y(2*i) = x(i*2) + x(i*2+1)
    y(2*i+1) = x(i*2) - x(i*2+1)
  }
}
```



looped, scalar repl.

```
for (int i=0; i < 2; i++)
{
  t0 = x[i];
  t1 = x[i+1];
  t2 = t0 + t1;
  y[i] = t2;
  t3 = t0 - t1;
  y[i+1] = t3;
}
```

```
def f(x: Array[Rep[Double]],
    y: Array[Rep[Double]]) = {
  for (i <- 0 until 2) {
    y(2*i) = x(i*2) + x(i*2+1)
    y(2*i+1) = x(i*2) - x(i*2+1)
  }
}
```



scalarized

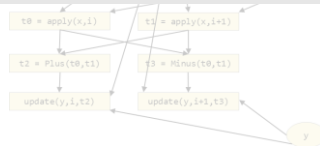
```
t0 = s0 + s1;
t1 = s0 - s1;
t2 = s2 + s3;
t2 = s2 - s3;
```

Staging enables program generation

**Abstracting over code style =
abstracting over staging decisions**

```
def f[L[_],A[_],T](looptype: L, x: A[Array[T]], y: A[Array[T]]) = {
  for (i <- 0 until 2: L[Range]) {
    y(2*i) = x(i*2) + x(i*2+1)
    y(2*i+1) = x(i*2) - x(i*2+1)
  }
}
```

```
y: Rep[Array[Double]] = {
  for (i <- 0 until 2: Rep[Range]) {
    y(2*i) = x(i*2) + x(i*2+1)
    y(2*i+1) = x(i*2) - x(i*2+1)
  }
}
```



```
t0 = x[i];
t1 = x[i+1];
t2 = t0 + t1;
y[i] = t2;
t3 = t0 - t1;
y[i+1] = t3;
```

Related Work

Program generators for performance

[FFTW codelet generator \(Frigo\)](#)

[Flame](#) (van de Geijn/Quinta, na-Orti, Bientinesi, ...)

[cvxgen](#) (Mattingley, Boyd)

[PetaBricks](#) (Ansel, Amarasinghe, ...)

[Spiral](#)

Autotuning

[ATLAS/PhiPAC](#) (Whaley, Bilmes, Demmel, Dongarra, ...)

[FFTW adaptive library](#) (Frigo, Johnson)

[OSKI](#) (Vuduc et al.)

[Adaptive sorting](#) (Li et al.)

Environments for DSLs and program generation

[Scala and lightweight modular staging \(LMS\)](#)

[More examples](#)

Automatically from Math to Fast Code

Principles

Generate Code

Capturing algorithm knowledge:
Mathematical DSLs



Structural optimization:
Rewriting

$$\begin{aligned} \text{DFT}_n &\rightarrow (\text{DFT}_k \otimes \text{I}_m) \text{T}_m (\text{I}_k \otimes \text{DFT}_m) \text{L}_k^T, \quad n = km \\ \text{DFT}_n &\rightarrow P_n (\text{DFT}_k \otimes \text{DFT}_m) Q_n, \quad n = km, \gcd(k, m) = 1 \\ \text{DCT-4}_n &\rightarrow S_n \text{DCT-2}_n \text{diag}_{0 \leq k < n} (1 / (2 \cos((2k + 1)\pi / 4n))) \\ \text{IMDCT}_{2m} &\rightarrow (J_m \oplus \text{I}_m \oplus J_m) \left(\begin{pmatrix} 1 \\ -1 \end{pmatrix} \otimes \text{I}_m \right) \oplus \left(\begin{pmatrix} -1 \\ 1 \end{pmatrix} \otimes \text{I}_m \right) J_{2m} \text{DCT-4}_{2m} \end{aligned}$$

Decision making:
Search and learning

$$\underbrace{A_m \otimes I_n}_{\text{smp}(p, \mu)} \rightarrow \underbrace{I_m^m}_{\text{smp}(p, \mu)} \left(I_p \otimes (I_{n/p} \otimes A_m) \right) \underbrace{I_n^m}_{\text{smp}(p, \mu)}$$

Key Challenges

New domains (linear algebra, filters, ...)

Programming language support (DSLs, staging)

More information: www.spiral.net

Further Reading

Spiral (www.spiral.net)

Markus Püschel, Franz Franchetti and Yevgen Voronenko

[Spiral](#)

in Encyclopedia of Parallel Computing, Eds. David Padua, Springer 2011

Georg Ofenbeck, Tiark Rompf, Alen Stojanov, Martin Odersky and Markus Püschel

[Spiral in Scala: Towards the Systematic Construction of Generators for Performance Libraries](#)

Proc. International Conference on Generative Programming: Concepts & Experiences (GPCE), pp. 125-134, 2013

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](#)

Proceedings of the IEEE special issue on "Program Generation, Optimization, and Adaptation," Vol. 93, No. 2, 2005, pp. 232-275

LGen

Daniele G. Spampinato and Markus Püschel

[A Basic Linear Algebra Compiler](#)

Proc. International Symposium on Code Generation and Optimization (CGO), pp. 23-32, 2014

Daniele G. Spampinato and Markus Püschel

[A basic linear algebra compiler for structured matrices](#)

Proc. International Symposium on Code Generation and Optimization (CGO), pp. 117-127, 2016

Tuning through machine learning

Frédéric de Mesmay, Yevgen Voronenko and Markus Püschel

[Offline Library Adaptation Using Automatically Generated Heuristics](#)

Proc. International Parallel and Distributed Processing Symposium (IPDPS), pp. 1-10, 2010

Frédéric de Mesmay, Arpad Rimmel, Yevgen Voronenko and Markus Püschel

[Bandit-Based Optimization on Graphs with Application to Library Performance Tuning](#)

Proc. International Conference on Machine Learning (ICML), pp. 729-736, 2009

Bryan Singer and Manuela Veloso

[Learning to Generate Fast Signal Processing Implementations](#)

Proc. International Conference on Machine Learning (ICML), pp. 529-536, 2001