

Autotuning

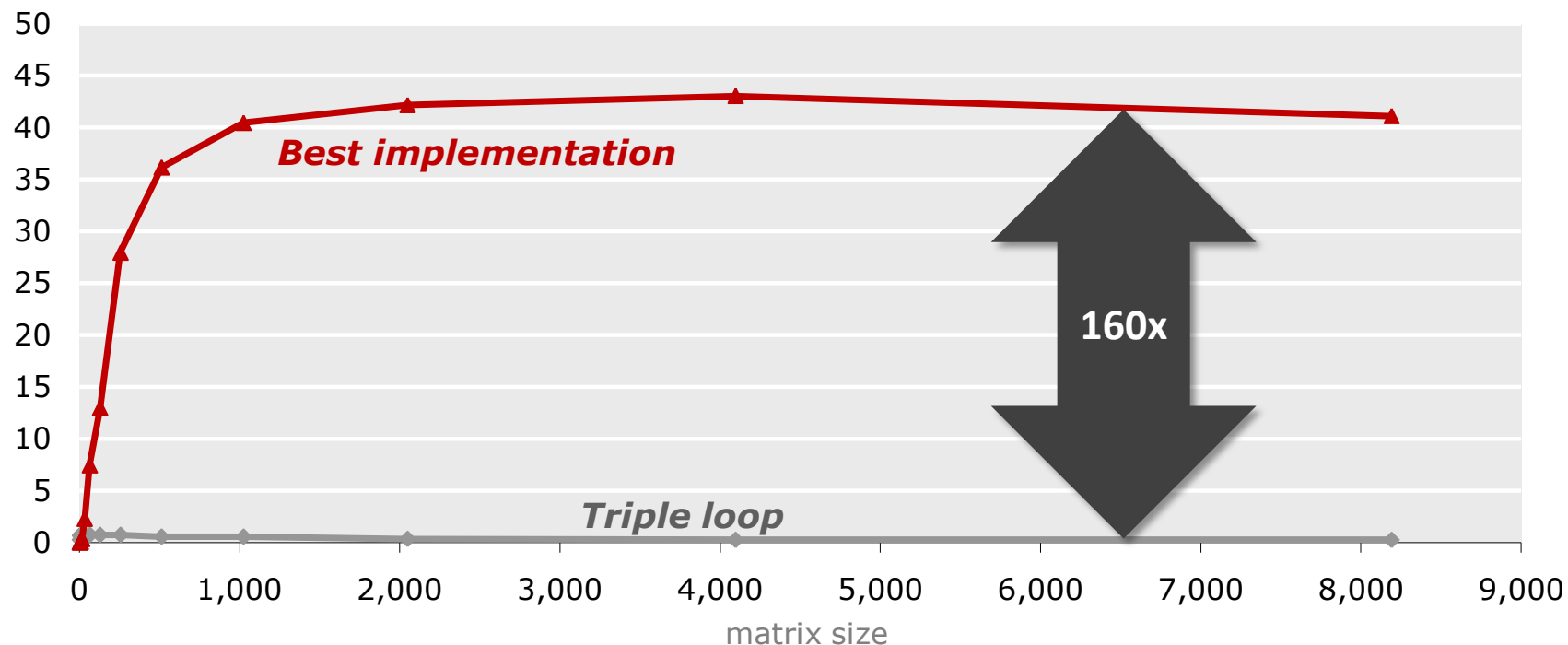
Markus Püschel

Electrical and Computer Engineering
Carnegie Mellon University

Why Autotuning?

Matrix-Matrix Multiplication (MMM) on quadcore Intel platform

Performance [Gflop/s]

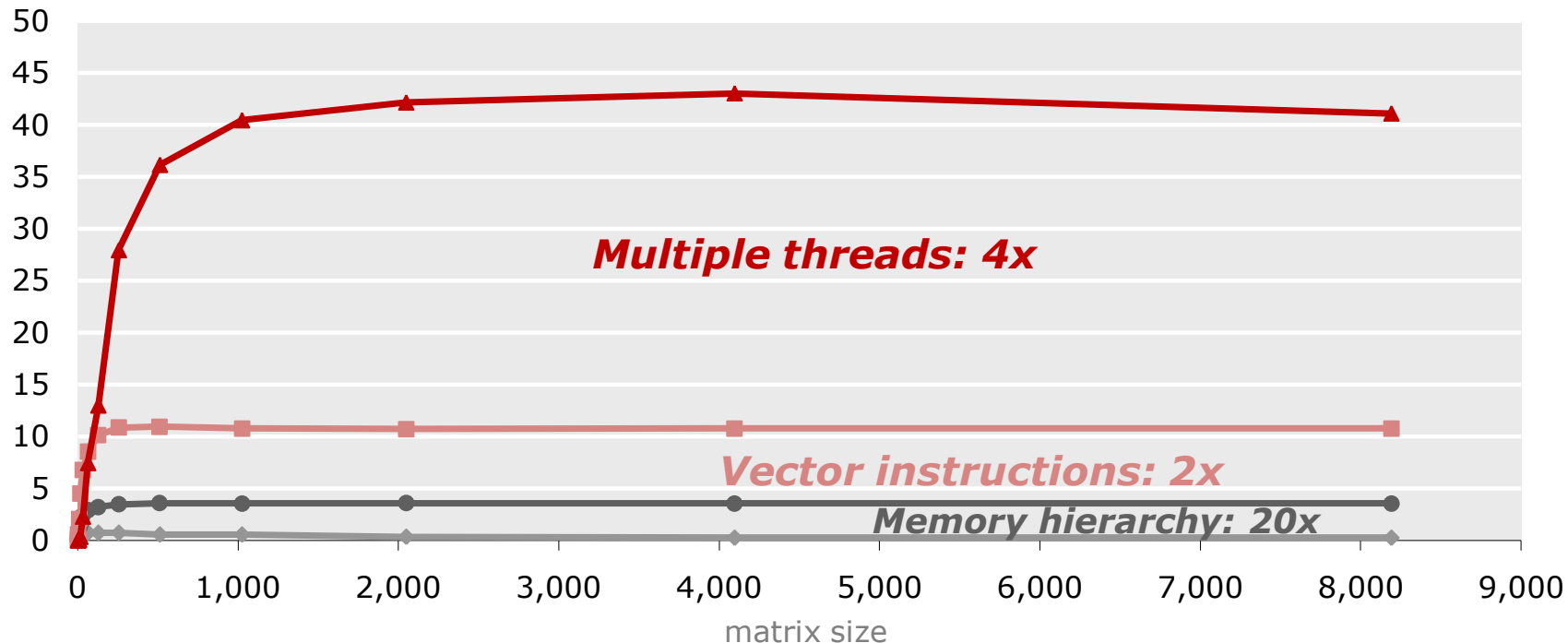


- Same (mathematical) operation count ($2n^3$)
- Compiler underperforms by 160x

Why Autotuning?

Matrix-Matrix Multiplication (MMM) on 2 x Core 2 Extreme 3 GHz

Performance [Gflop/s]

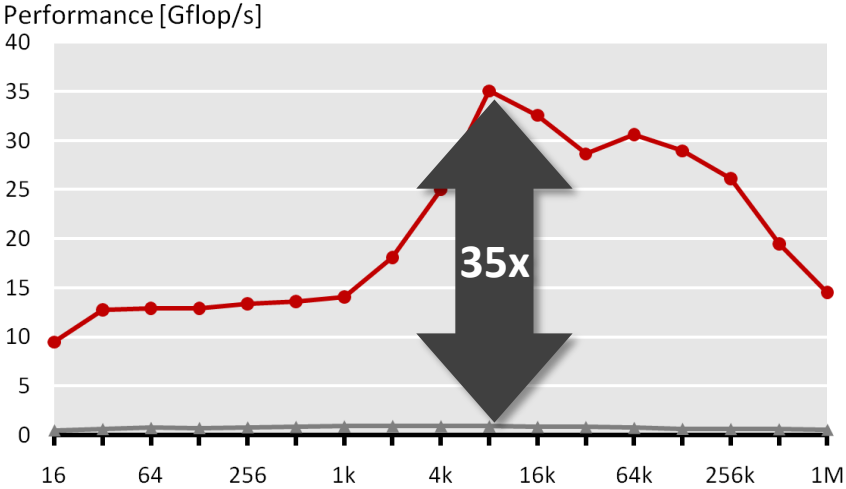


■ Code size:

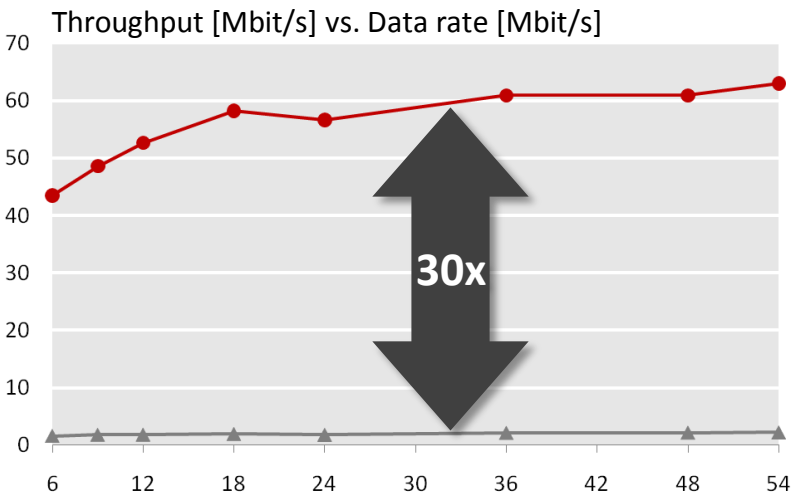
- Triple loop: < 1 KB
- Best code: about 100 KB

Same for All Critical Compute Functions

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



WiFi Receiver (Physical layer) on one Intel Core

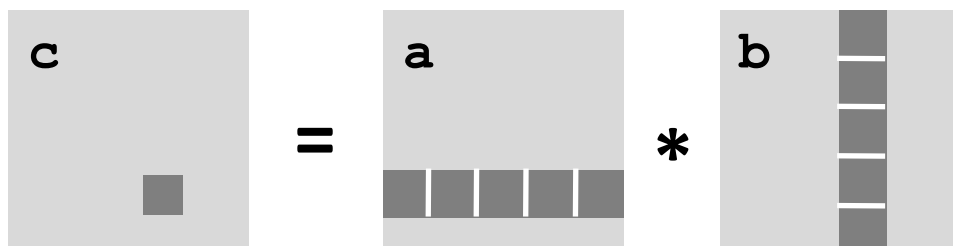


Summary of the Situation: *Huge Problem*

- For computing functions compilers underperform by 10–100x
- Reason: the following are unsolved problems
 - Memory hierarchy optimizations
 - Vectorization
 - Parallelization
- Optimization by hand requires highly skilled programmers
- Performance does not port well
- *Solution (autotuning): Automating performance optimization with tools that complement/aid the compiler or programmer*

PhiPac/ATLAS: MMM Generator

Whaley, Bilmes, Demmel, Dongarra, ...



Blocking improves locality

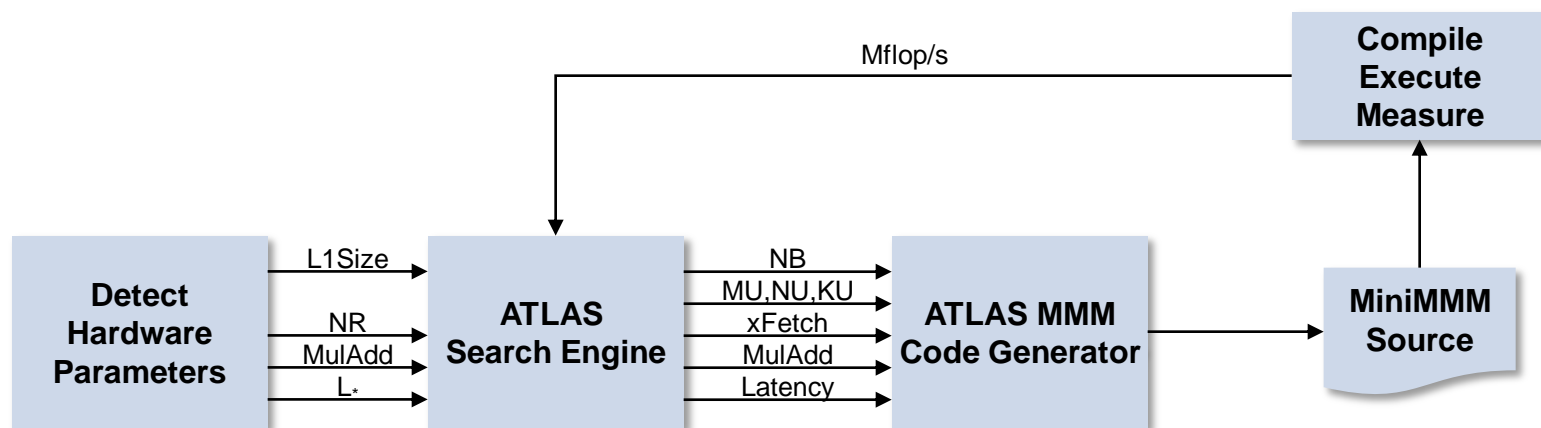
```

c = (double *) calloc(sizeof(double), n*n);

/* Multiply n x n matrices a and b */
void mmm(double *a, double *b, double *c, int n) {
    int i, j, k;
    for (i = 0; i < n; i+=B)
        for (j = 0; j < n; j+=B)
            for (k = 0; k < n; k+=B)
                /* B x B mini matrix multiplications */
                for (i1 = i; i1 < i+B; i++)
                    for (j1 = j; j1 < j+B; j++)
                        for (k1 = k; k1 < k+B; k++)
                            c[i1*n+j1] += a[i1*n + k1]*b[k1*n + j1];
}

```

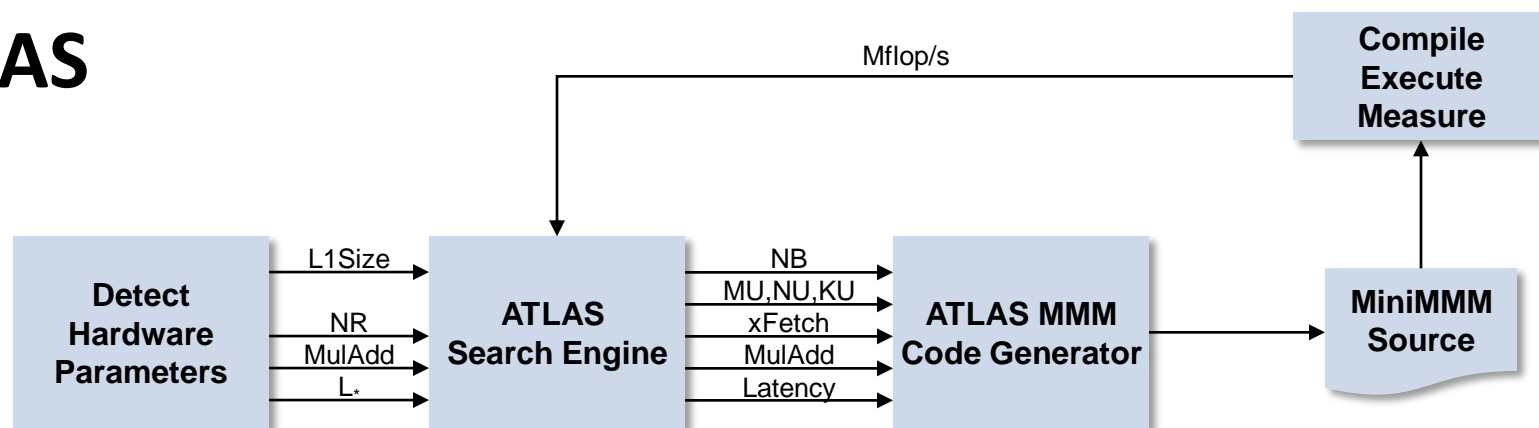
PhiPac/ATLAS: MMM Generator



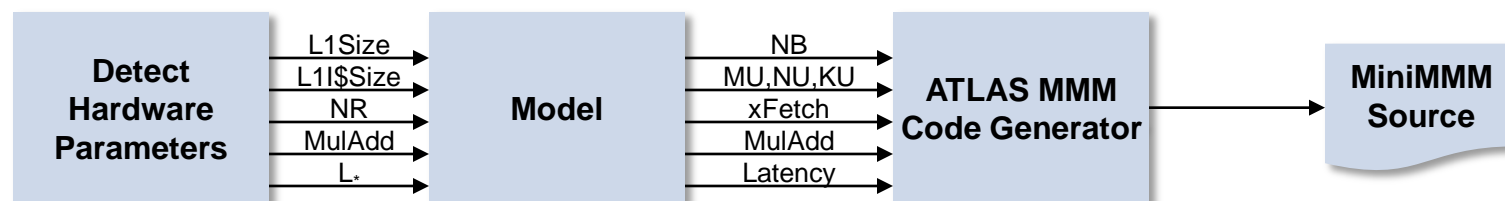
■ *Techniques:*

- Program generation (here: template-based)
- Feedback-driven search over a set of parameters

ATLAS



Model-Based ATLAS (Yotov et al.)



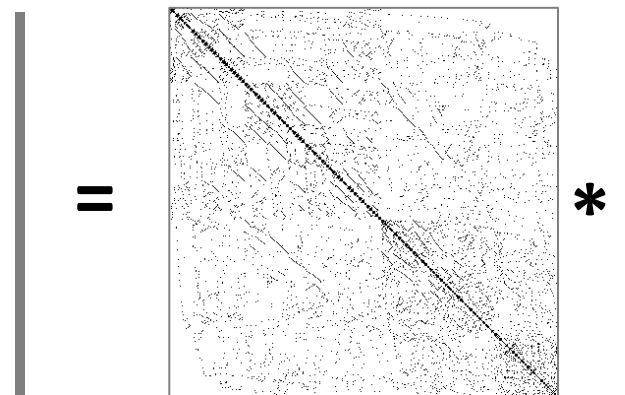
$$\left\lceil \frac{N_B^2}{B_1} \right\rceil + 3 \left\lceil \frac{N_B \times N_U}{B_1} \right\rceil + \left\lceil \frac{M_U}{B_1} \right\rceil \times N_U \leq \frac{C_1}{B_1}$$

■ *Techniques:*

- Hardware parameter based model

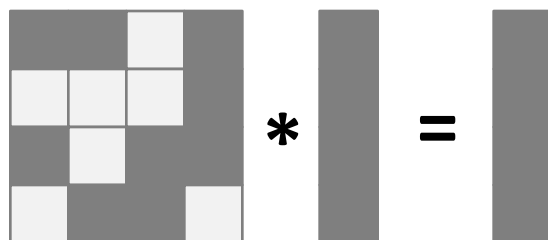
OSKI: Sparse Matrix-Vector Multiplication

Vuduc, Im, Yelick, Demmel



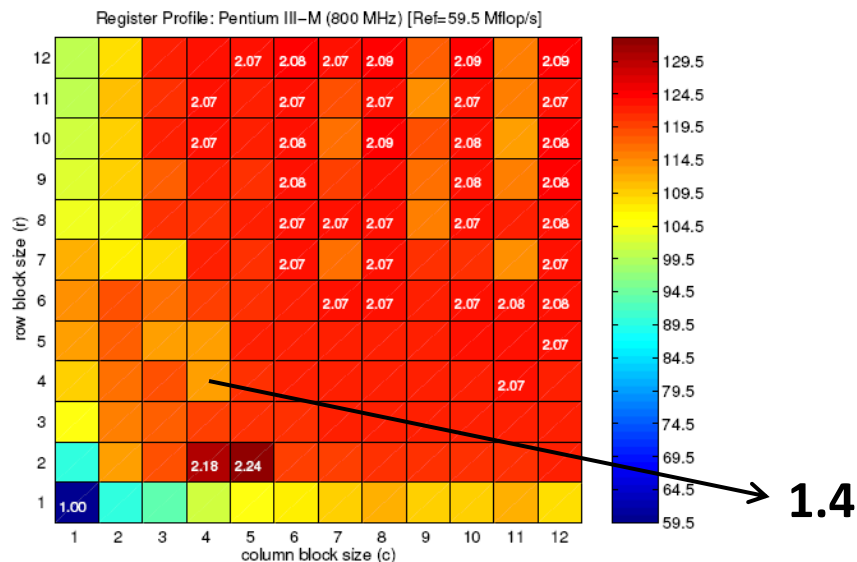
■ Blocking for registers:

- Improves locality (reuse of input vector)
- But creates overhead (zeros in block)

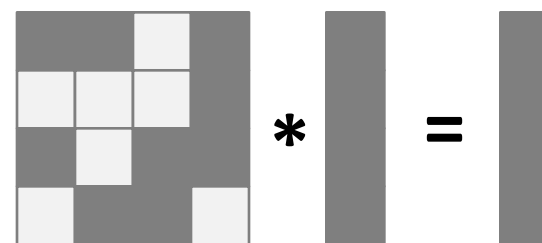


OSKI: Sparse Matrix-Vector Multiplication

Gain by blocking (dense MVM)



Overhead by blocking



$$16/9 = 1.77$$



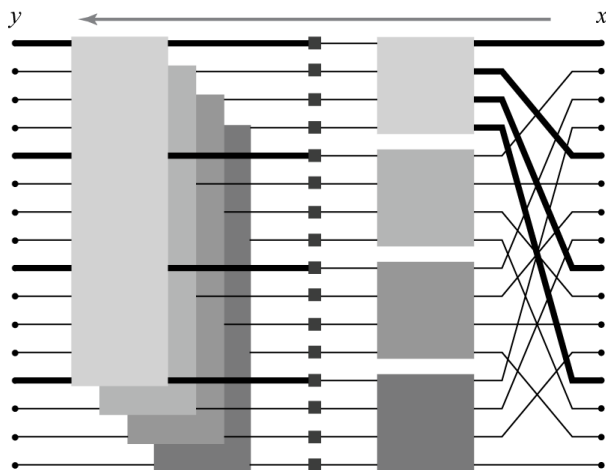
$$1.4/1.77 = 0.79 \text{ (no gain)}$$

■ *Techniques:*

- Measurement-based model
- Data structure adaptation

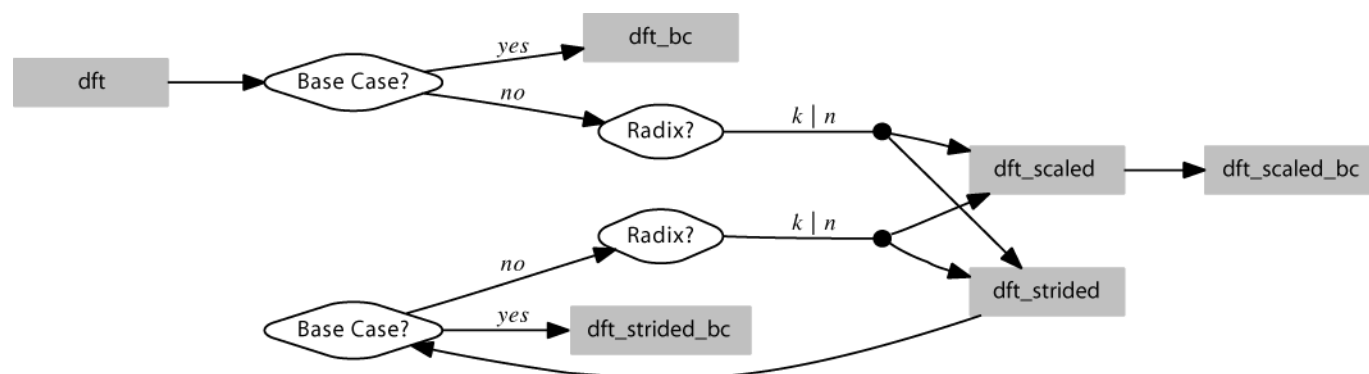
FFTW: Discrete Fourier Transform

Frigo, Johnson



```
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 adaptation



Vectorization, threading, etc. add more choices

FFTW: Discrete Fourier Transform

Installation

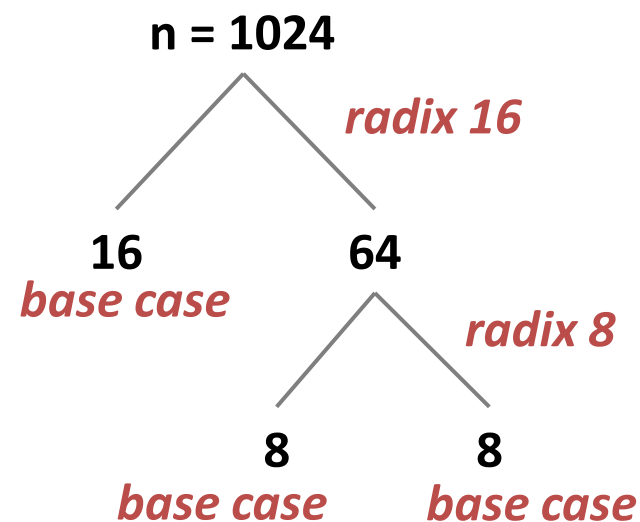
configure/make

Usage

$d = \text{dft}(n)$
 $d(x, y)$

Twiddles

Search for fastest
computation strategy



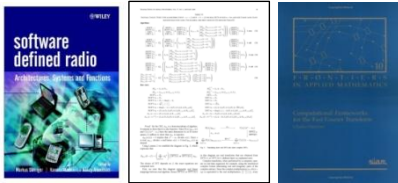
■ *Techniques:*

- (Online) Adaptive library
- Dynamic programming search
- Not explained: Program generator for basic blocks

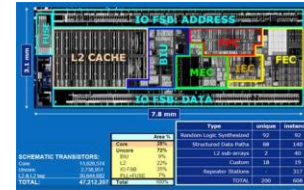
Spiral: Linear Transforms & More

Franchetti, Voronenko, Püsichel, Xiong, Singer, Moura, Johnson, Padua, ...

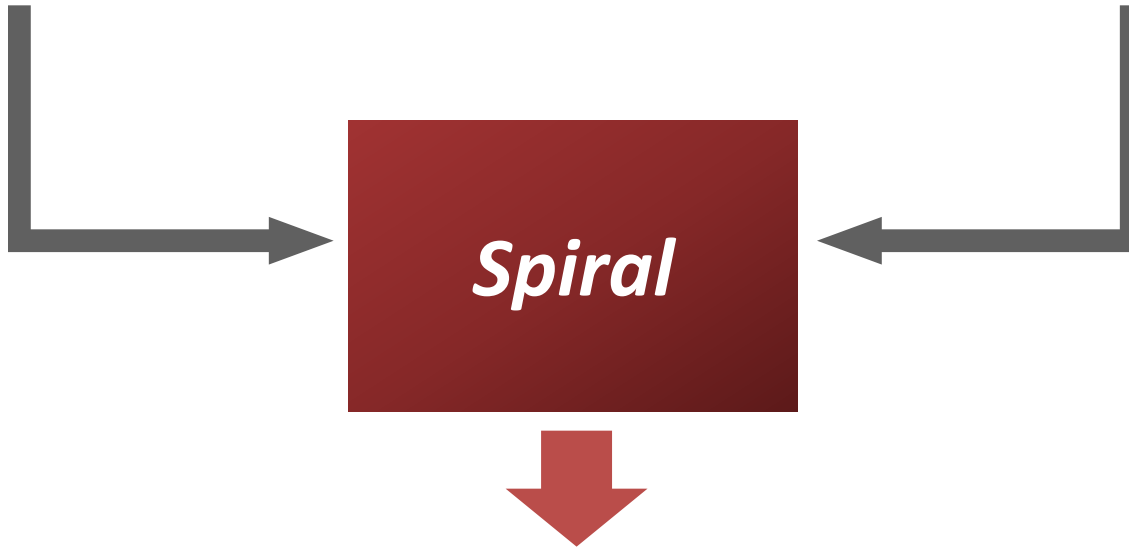
Algorithm knowledge



Platform description



```
_mm_set1_epi8(x) = ...  
_mm_xor_si128(x,y) = ...  
_mm_avg_epu8(x,y) = ...  
_mm_cmpeq_epi8(x,y) = ...  
_mm_unpacklo_epi8(x,y) = ...  
...
```



Optimized implementation
(regenerated for every new platform)

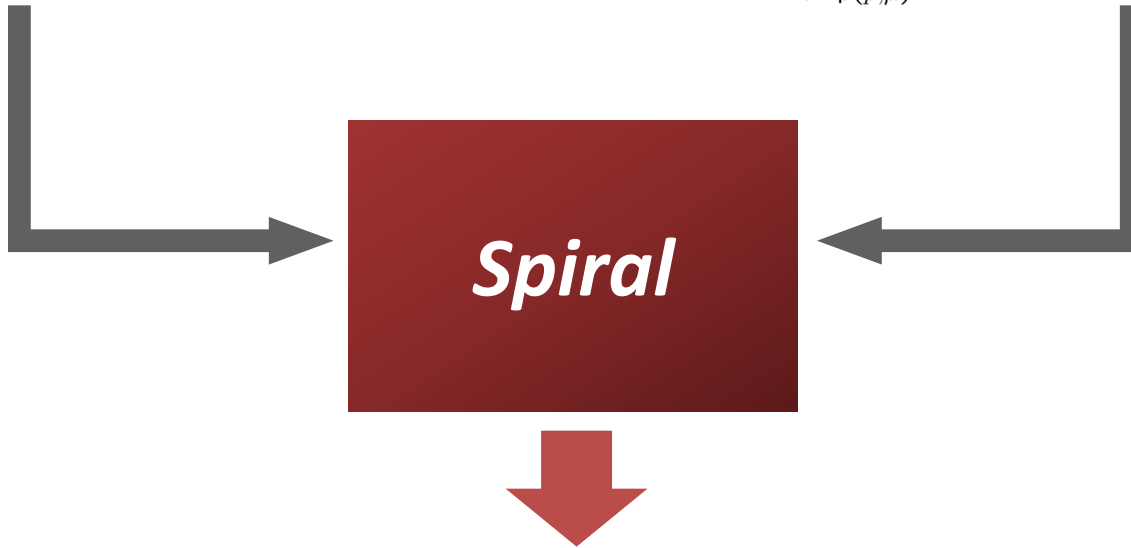
Spiral: Linear Transforms & More

Algorithm knowledge

$$\begin{aligned} \text{DFT}_n &\rightarrow P_{k/2,2m}^\top \left(\text{DFT}_{2m} \oplus \left(I_{k/2-1} \otimes_i C_{2m} \text{rDFT}_{2m}(i/k) \right) \right) \left(\text{RDFT}'_k \otimes I_m \right) \\ \begin{vmatrix} \text{rDFT}_{2n}(u) \\ \text{rDHT}_{2n}(u) \end{vmatrix} &\rightarrow L_m^{2n} \left(I_k \otimes_i \begin{vmatrix} \text{rDFT}_{2m}((i+u)/k) \\ \text{rDHT}_{2m}((i+u)/k) \end{vmatrix} \right) \left(\begin{vmatrix} \text{rDFT}_{2k}(u) \\ \text{rDHT}_{2k}(u) \end{vmatrix} \otimes I_m \right) \\ \text{RDFT-3}_n &\rightarrow (Q_{k/2,m}^\top \otimes I_2) (I_k \otimes_i \text{rDFT}_{2m})(i+1/2/k) (\text{RDFT-3}_k \otimes I_m) \end{aligned}$$

Platform description

$$\begin{aligned} \underbrace{A_m \otimes I_n}_{\text{smp}(p,\mu)} &\rightarrow \underbrace{\left(L_m^{mp} \otimes I_{n/p} \right) \left(I_p \otimes (A_m \otimes I_{n/p}) \right) \left(L_p^{mp} \otimes I_{n/p} \right)}_{\text{smp}(p,\mu)} \\ \underbrace{I_m \otimes A_n}_{\text{smp}(p,\mu)} &\rightarrow I_p \otimes_{||} \left(I_{m/p} \otimes A_n \right) \\ \underbrace{(P \otimes I_n)}_{\text{smp}(p,\mu)} &\rightarrow (P \otimes I_{n/\mu}) \overline{\otimes} I_\mu \end{aligned}$$



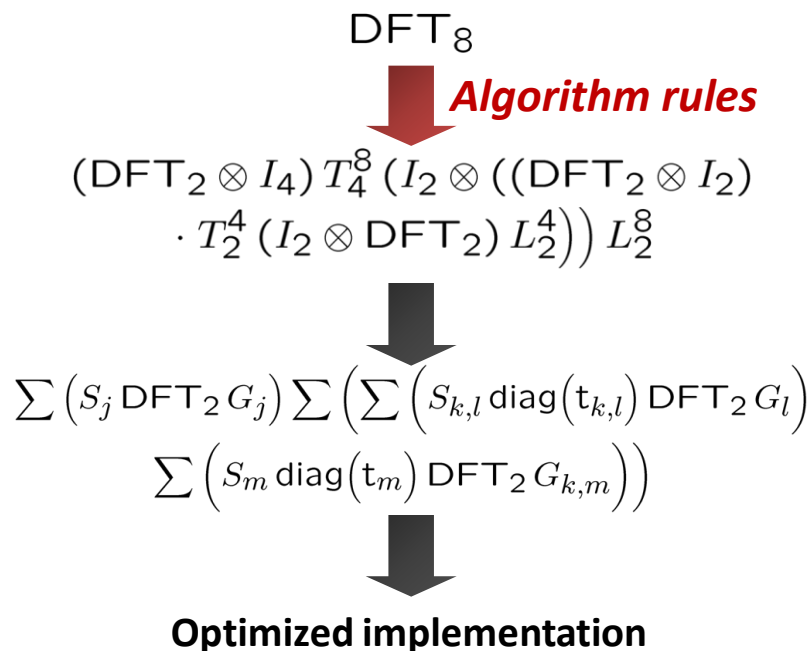
Optimized implementation
(regenerated for every new platform)

Program Generation in Spiral (Sketched)

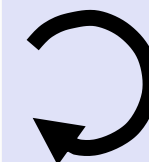
Transform
user specified

Fast algorithm
in SPL
many choices

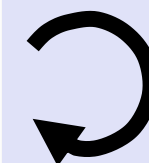
Σ -SPL



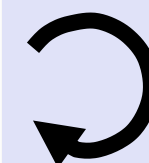
*Optimization at all
abstraction levels*



parallelization
vectorization



loop
optimizations



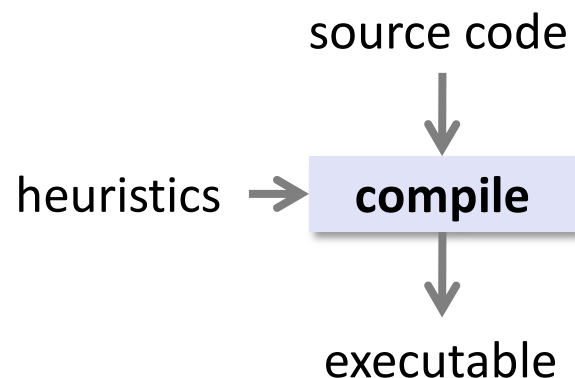
constant folding
scheduling
.....

■ *Techniques:*

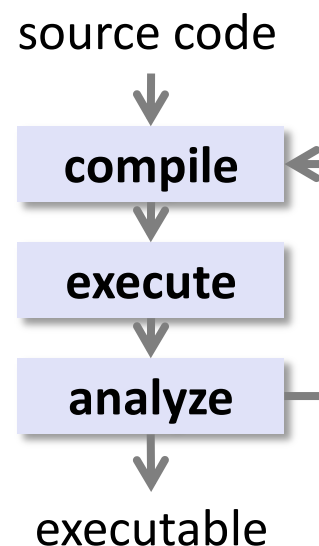
- Domain-specific language (declarative, mathematical, point-free)
- Rewriting for optimization
- Search techniques
- ...

Adaptive Compilation

Traditional:



Iterative:



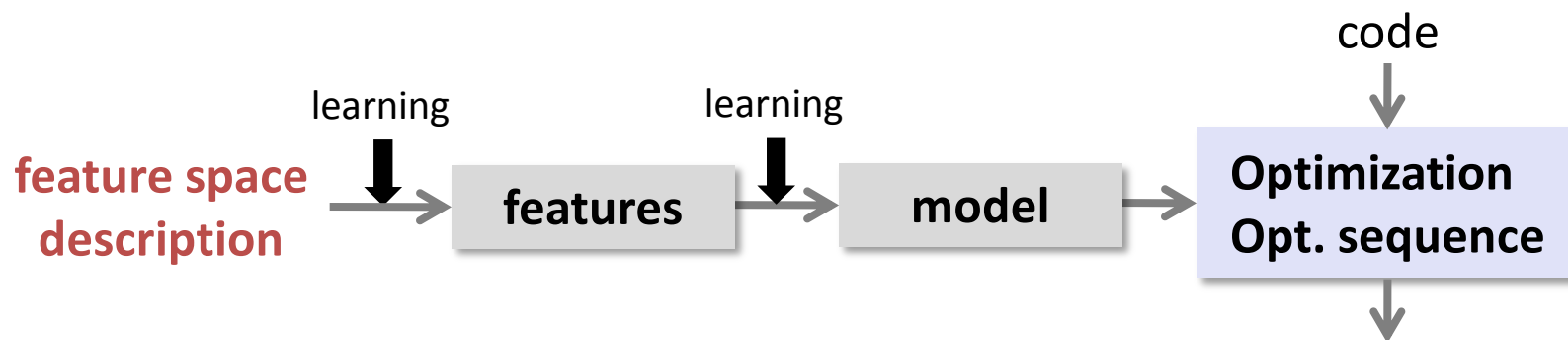
- Optimization ordering: Cooper et al., Kulkarni et al., Triantifyllis et al.,
- Parameter selection: Kisuki et al., Stephenson et al.

■ *Techniques:*

- Feedback-driven search based on measurements

Adaptive Compilation

Moss, Boyle, Cavazos, Stephenson, Beaty, ...

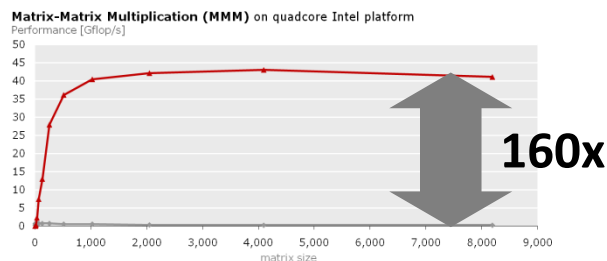


■ *Techniques:*

- Machine learning

Autotuning: Conclusions

- Is necessary



- Search over space of parameters

Need for a broader set of techniques:

- Program generation
- Domain-specific languages (declarative)
- Models
- Machine learning
- Data structures