



Autotuning (2/2): Specialized code generators

Prof. Richard Vuduc

Georgia Institute of Technology

CSE/CS 8803 PNA: Parallel Numerical Algorithms

[L.18] Thursday, March 6, 2008



Today's sources

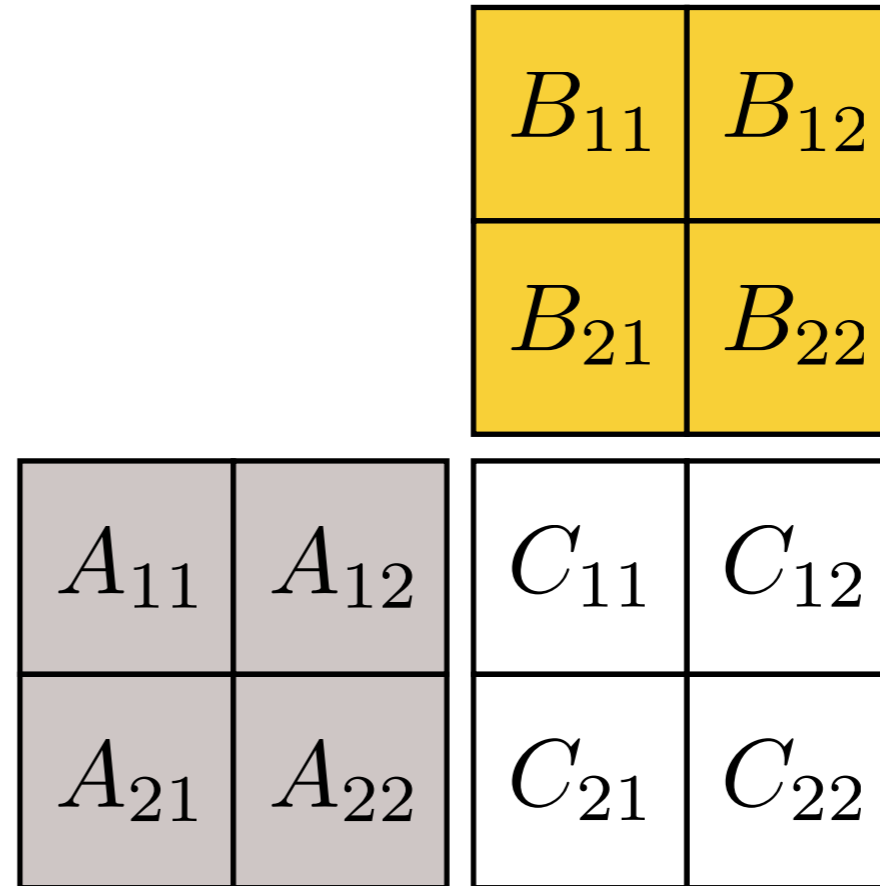
- CS 267 at UCB (Demmel & Yelick)
- Papers from various autotuning projects
 - PHiPAC, ATLAS, FFTW, SPIRAL, TCE
- See: Proc. IEEE 2005 special issue on Program Generation, Optimization, and Platform Adaptation
- Me (for once!)



Review: Cache-oblivious algorithms

A recursive algorithm for matrix-multiply

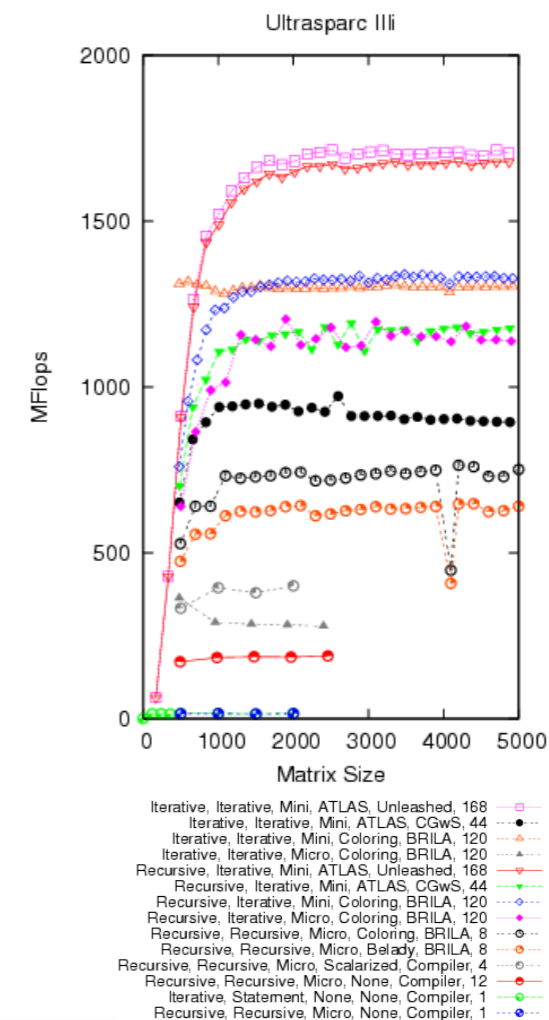
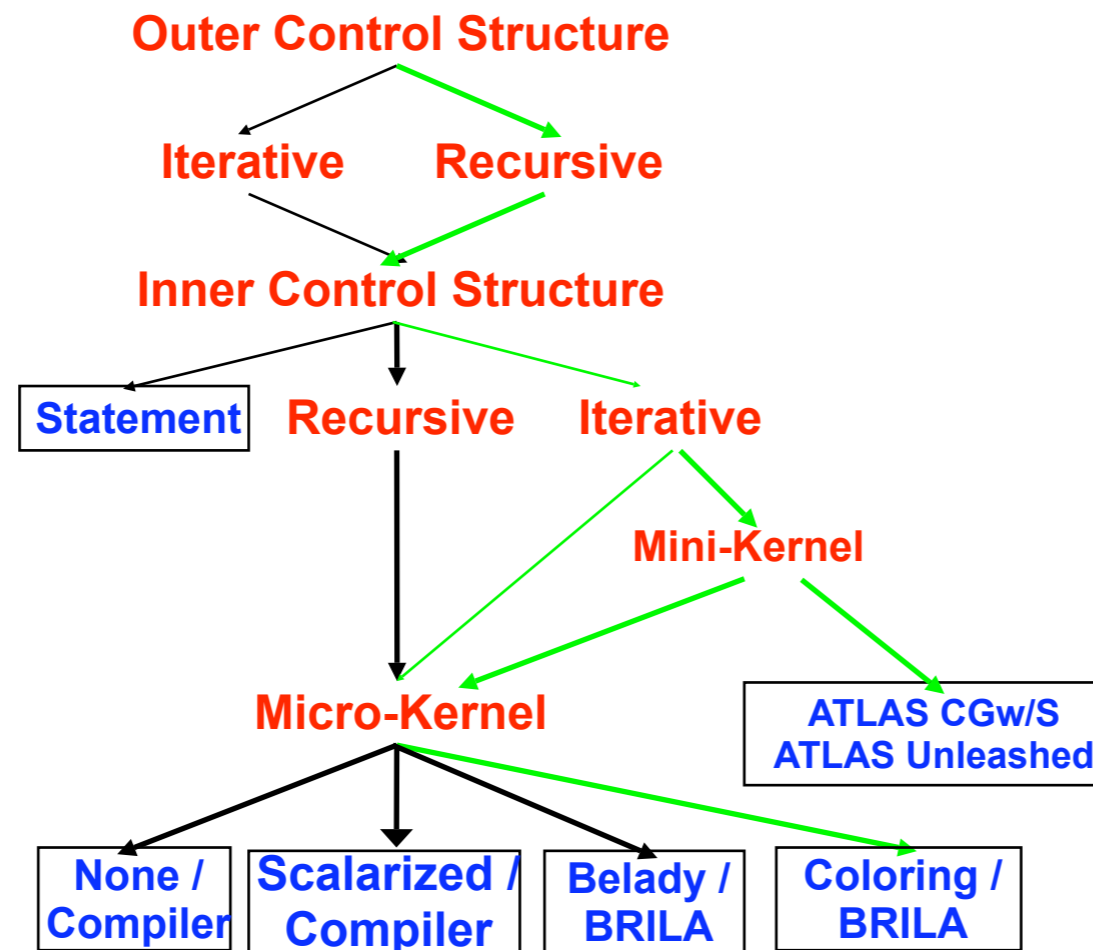
- ▣ Divide all dimensions in half
- ▣ Bilardi, *et al.*: Use grey-code ordering



No. of misses, with tall-cache assumption:

$$Q(n) = \left\{ \begin{array}{ll} 8 \cdot Q\left(\frac{n}{2}\right) & \text{if } n > \sqrt{\frac{M}{3}} \\ 3n^2 & \text{otherwise} \end{array} \right\} \leq \Theta\left(\frac{n^3}{L\sqrt{M}}\right)$$

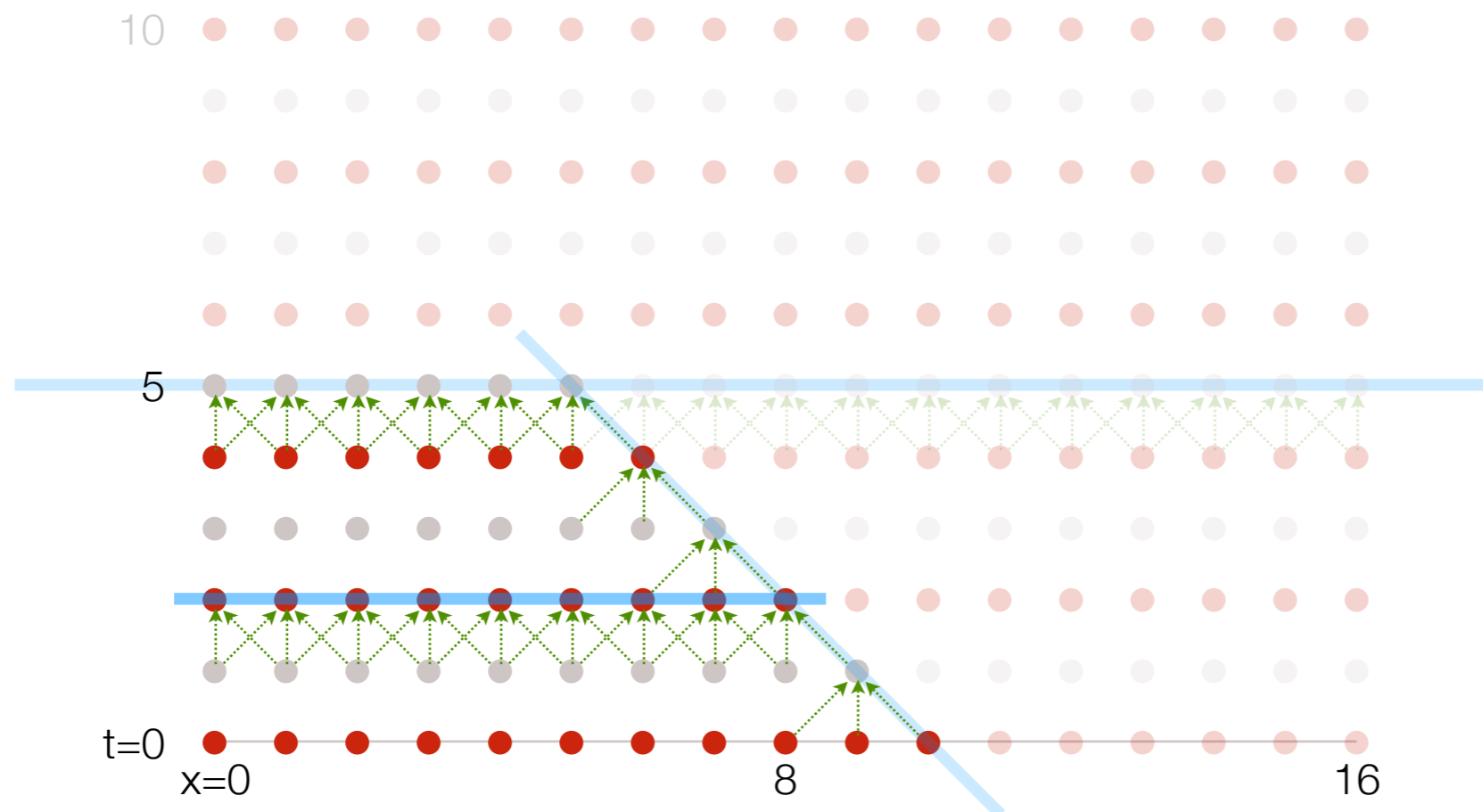
Performance-engineering challenges



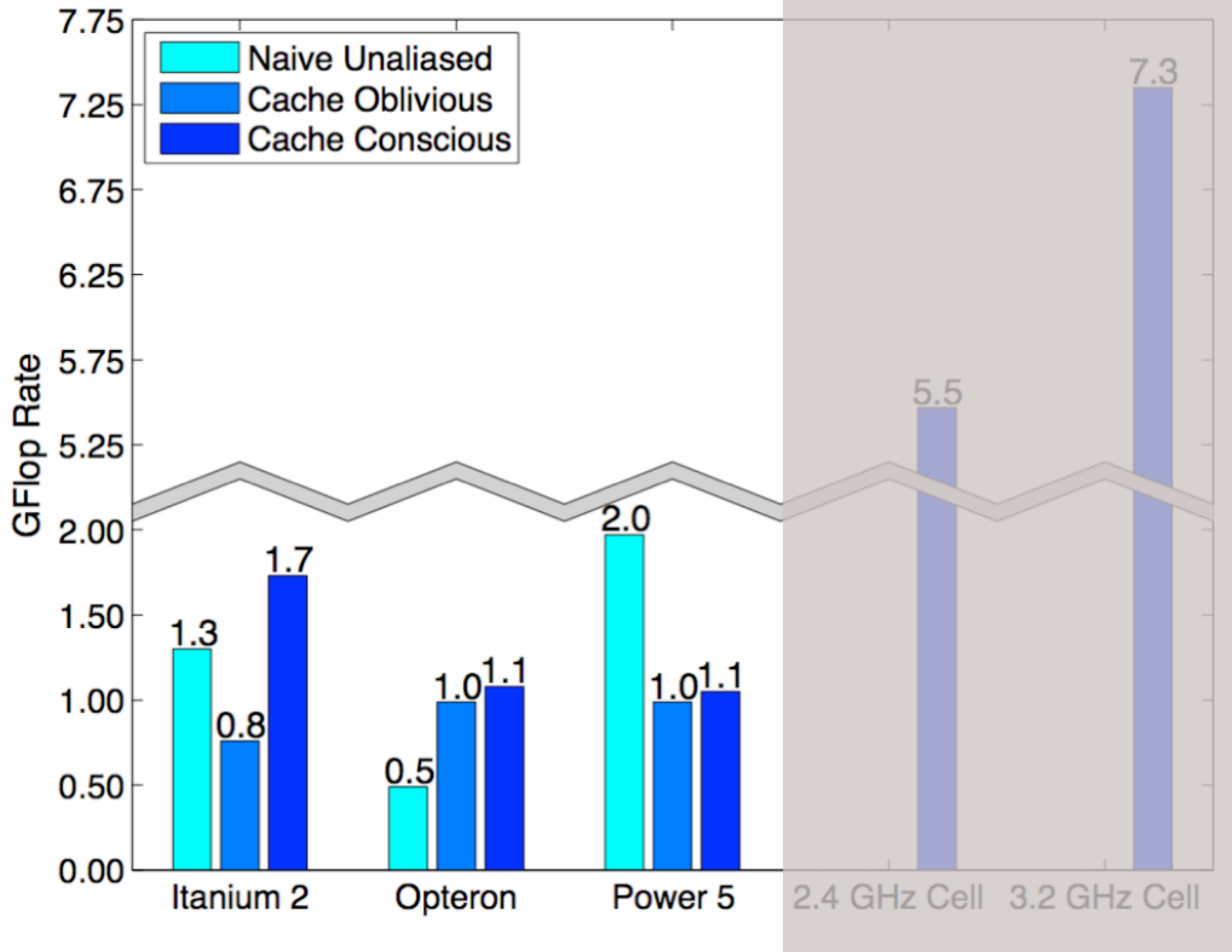
Cache-oblivious stencil computation

Theorem [Frigo & Strumpen (ICS 2005)]:
 $d = \text{dimension} \Rightarrow$

$$Q(n, t; d) = O\left(\frac{n^d \cdot t}{M^{\frac{1}{d}}}\right)$$



Cache-conscious algorithm



Source: Datta, et al. (2007)



Survey of autotuning



Early idea seedlings

- **Polyalgorithms:** John R. Rice

- (1969) “A polyalgorithm for the automatic solution of nonlinear equations”
- (1976) “The algorithm selection problem”

- **Profiling** and feedback-directed compilation

- (1971) D. Knuth: “An empirical study of FORTRAN programs”
- (1982) S. Graham, P. Kessler, M. McKusick: gprof
- (1991) P. Chang, S. Mahlke, W-m. W. Hwu: “Using profile information to assist classic code optimizations”

- Code generation from **high-level representations**

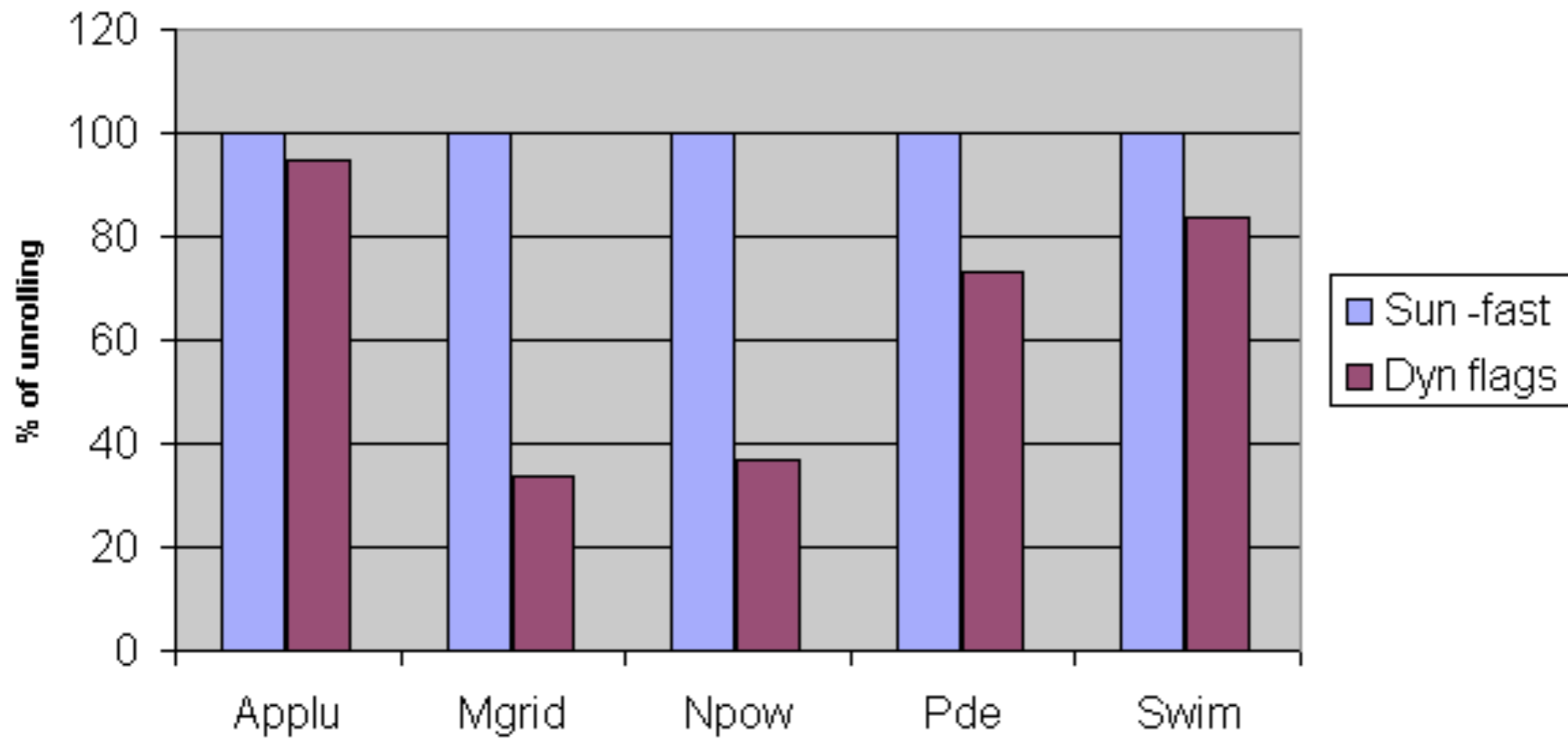
- (1989) J. Johnson, R.W. Johnson, D. Rodriguez, R. Tolimieri: “A methodology for designing, modifying, and implementing Fourier Transform algorithms on various architectures.”
- (1992) M. Covell, C. Myers, A. Oppenheim: “Computer-aided algorithm design and arrangement” (1992)



Why doesn't the compiler do the dirty work?

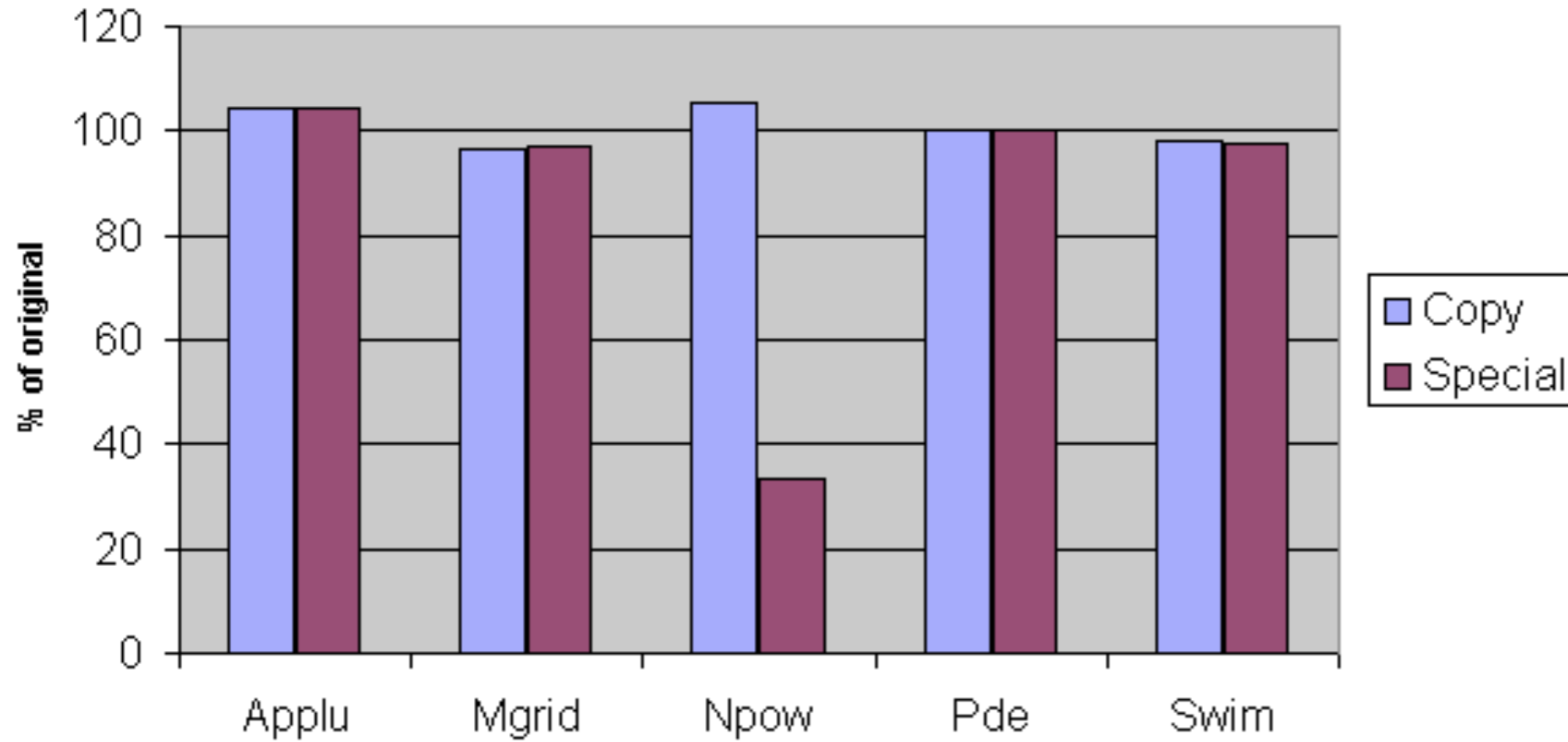
- Why doesn't the compiler do all of this?
 - Analysis
 - Over-specified dependencies
 - Correctness requirements
 - Limited access to relevant run-time information
 - Architecture: Realistic hardware models?
 - Engineering: Hard to modify a production compiler

Flag Selection



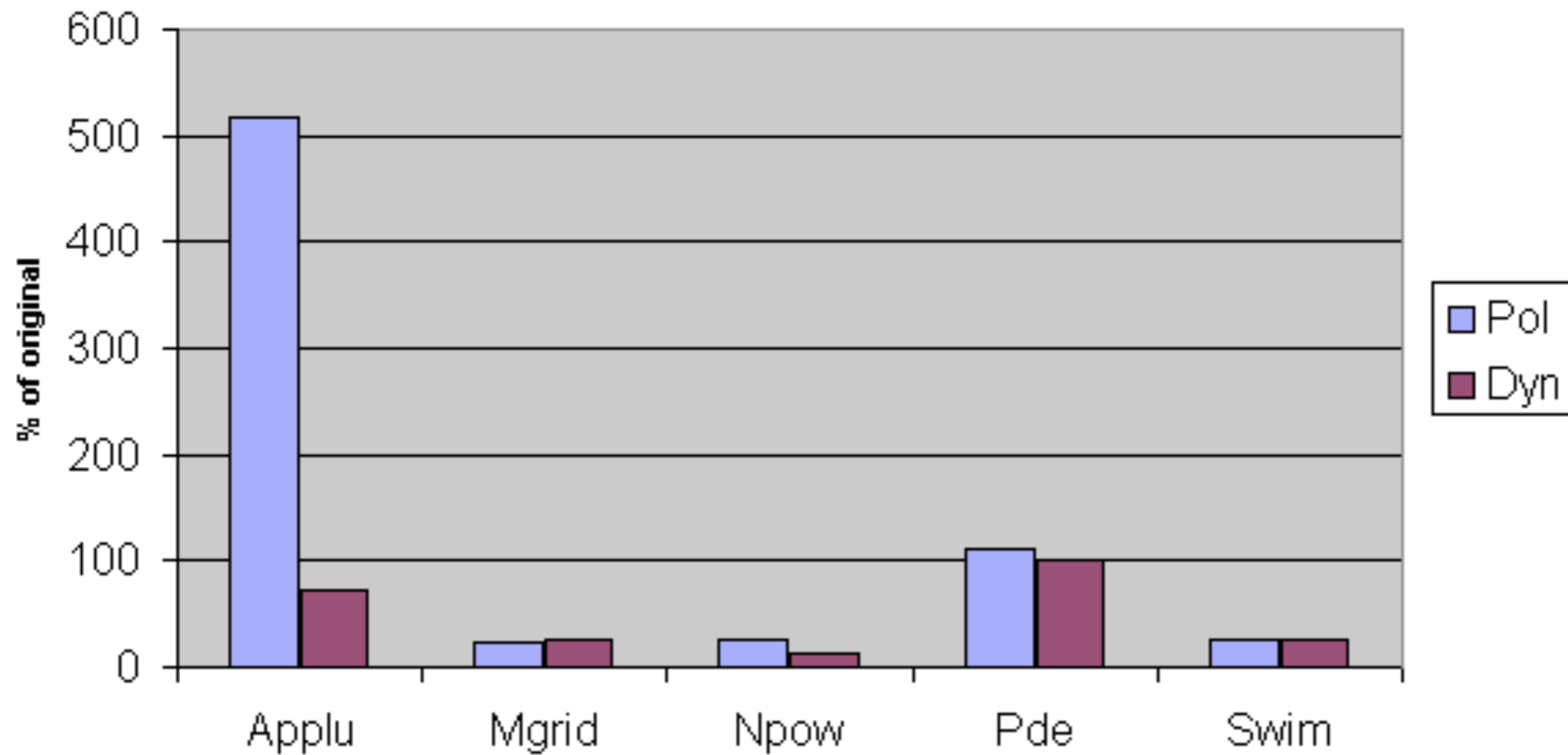
Source: Voss, ADAPT compiler project: <http://www.eecg.toronto.edu/~voss/AdaptPage/results.html>

Specialization



Source: Voss, ADAPT compiler project: <http://www.eecg.toronto.edu/~voss/AdaptPage/results.html>

Parallel On/Off



Source: Voss, ADAPT compiler project: <http://www.eecg.toronto.edu/~voss/AdaptPage/results.html>



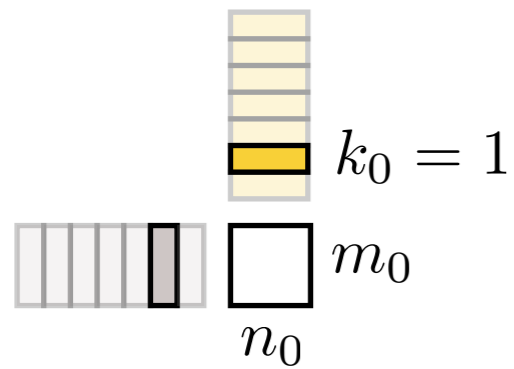
Automatic performance tuning, or “autotuning”

- Two-phase methodology for producing automatically tuned code
 - Given: Computational kernel or program; inputs; machine
 - Identify and generate a parameterized space of candidate implementations
 - Select the fastest one using empirical modeling and automated experiments
- “Autotuner” = System that implements this
 - Usually domain-specific (exception: “autotuning/iterative compilers”)
 - Leverage back-end compiler for performance and portability

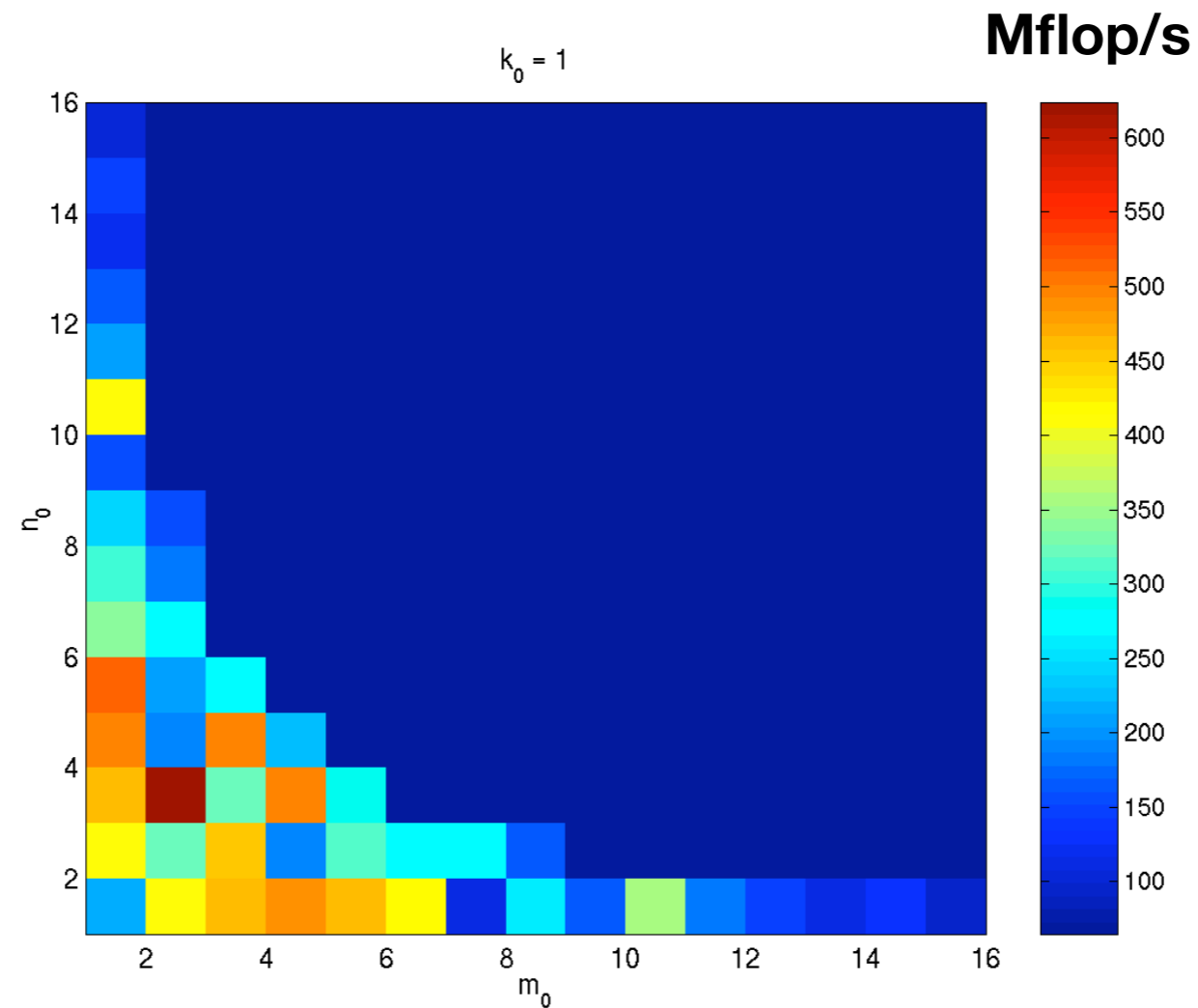
How an autotuner differs from a compiler (roughly)

	Compiler	Autotuner
Input	General-purpose source code	Specification
Code generation time	User responsive	Long, but amortized
Implementation selection	Static analysis; some run-time profiling/feedback	Automated empirical models and experiments

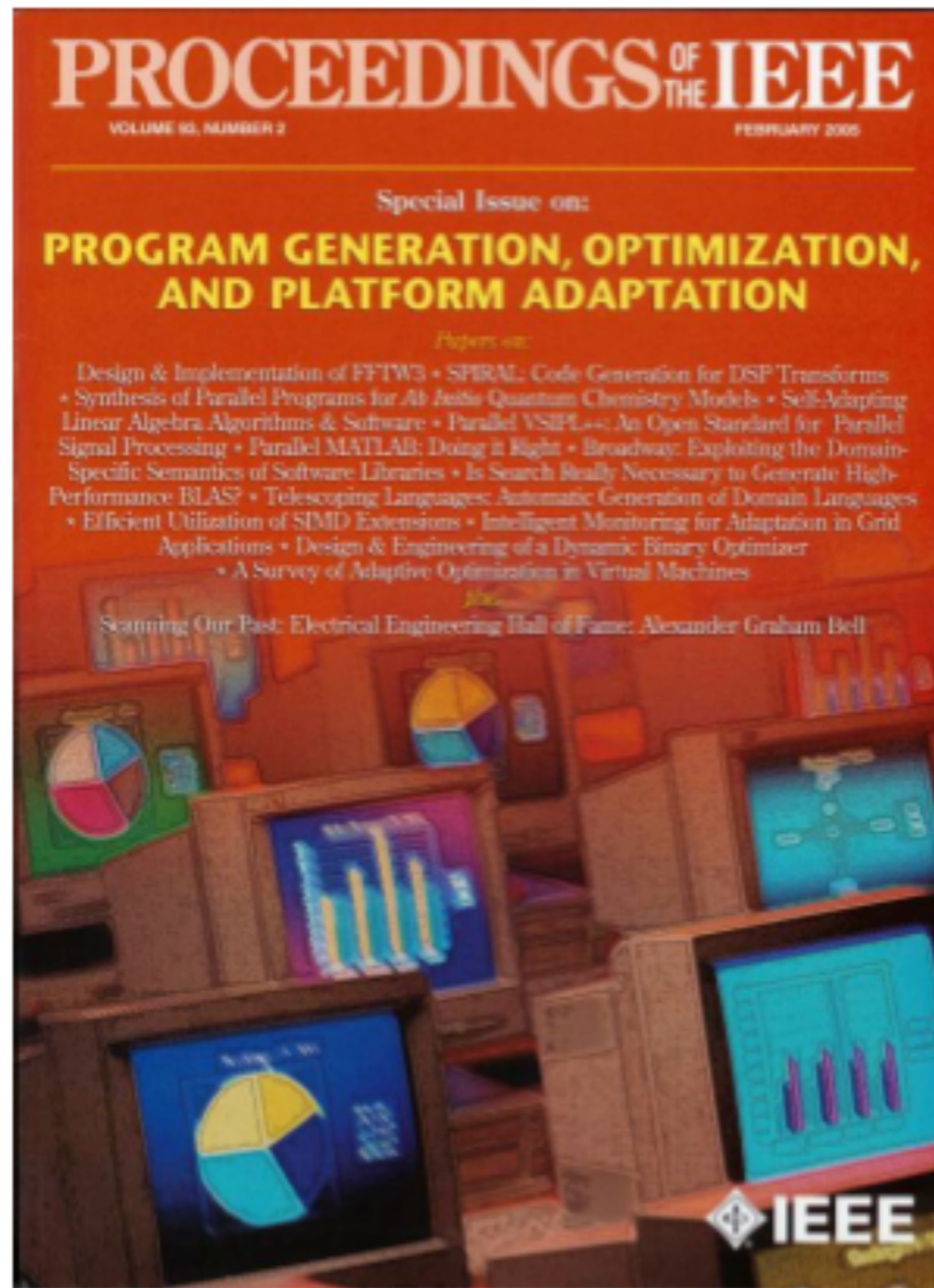
Example: What a search space looks like



- Platform: Sun Ultra Ili
- 16 double regs
- 667 Mflop/s peak
- Unrolled, pipelined inner-kernel
- Sun cc v5.0 compiler



Source: PHiPAC Project at UC Berkeley (1997)



Proceedings of the IEEE special issue, Feb. 2005





Dense linear algebra



PHiPAC (1997)

- Portable High-Performance ANSI C [Bilmes, Asanovic, Chin, Demmel (1997)]
 - Coding guidelines: C as high-level assembly language
 - Code generator for multi-level cache- and register-blocked matrix multiply
 - Exhaustive search over all parameters
 - Began as class project which beat the vendor BLAS

PHiPAC coding guideline example: Removing false dependencies

- Use local variables to remove false dependencies

```
a[i] = b[i] + c;  
a[i+1] = b[i+1] * d;
```


**False read-after-write hazard
between a[i] and b[i+1]**



```
float f1 = b[i];  
float f2 = b[i+1];
```

```
a[i] = f1 + c;  
a[i+1] = f2 * d;
```

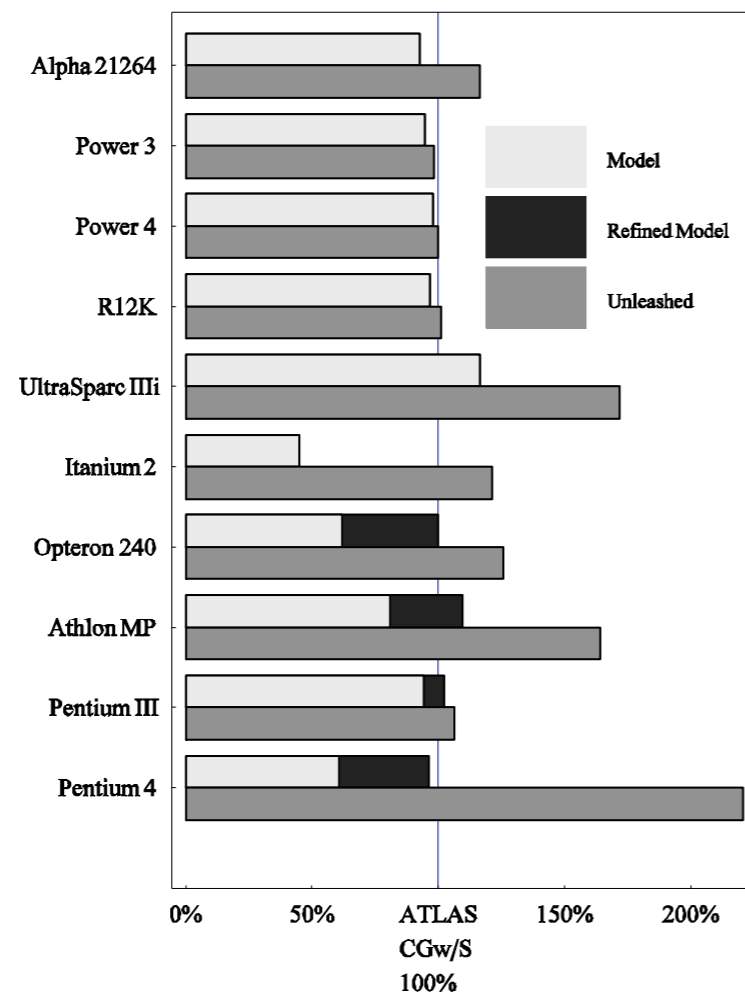
In C99, may declare a & b unaliased
("restrict" keyword)



ATLAS (1998)

- “Automatically Tuned Linear Algebra Software” — [R.C. Whaley and J. Dongarra (1998)]
 - Overcame PHiPAC shortcomings on x86 platforms
 - Copy optimization, prefetch, alternative schedulings
 - Extended to full BLAS, some LAPACK support (e.g., LU)
- Code generator (written in C, output C w/ inline-assembly) with search
 - Copy optimization prunes much of PHiPAC’s search space
 - “Simple” line searches
 - See: iterative floating-point kernel optimizer (iFKO) work

Search vs. modeling



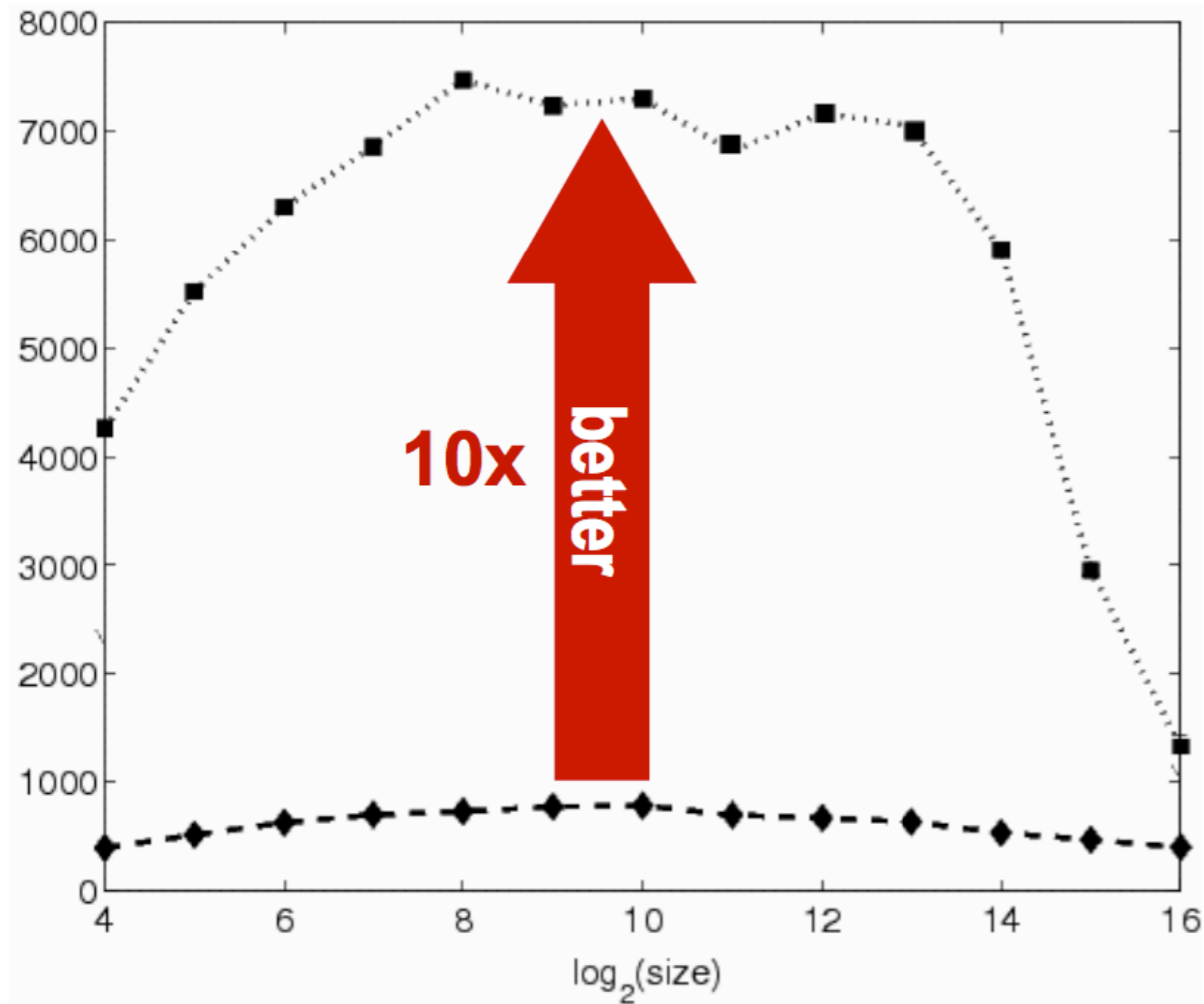
- Yotov, *et al.* “Is search really necessary to generate high-performance BLAS?”
- “Think globally, search locally”
 - Small gaps \Rightarrow local search
 - Large gaps \Rightarrow refine model
- “Unleashed” \Rightarrow hand-optimized plug-in kernels



Signal processing

Motivation for performance tuning

pseudo
Mflop/s




**best available
implementation
(FFTW, Intel IPP, Spiral)**

**roughly the same
operations count**

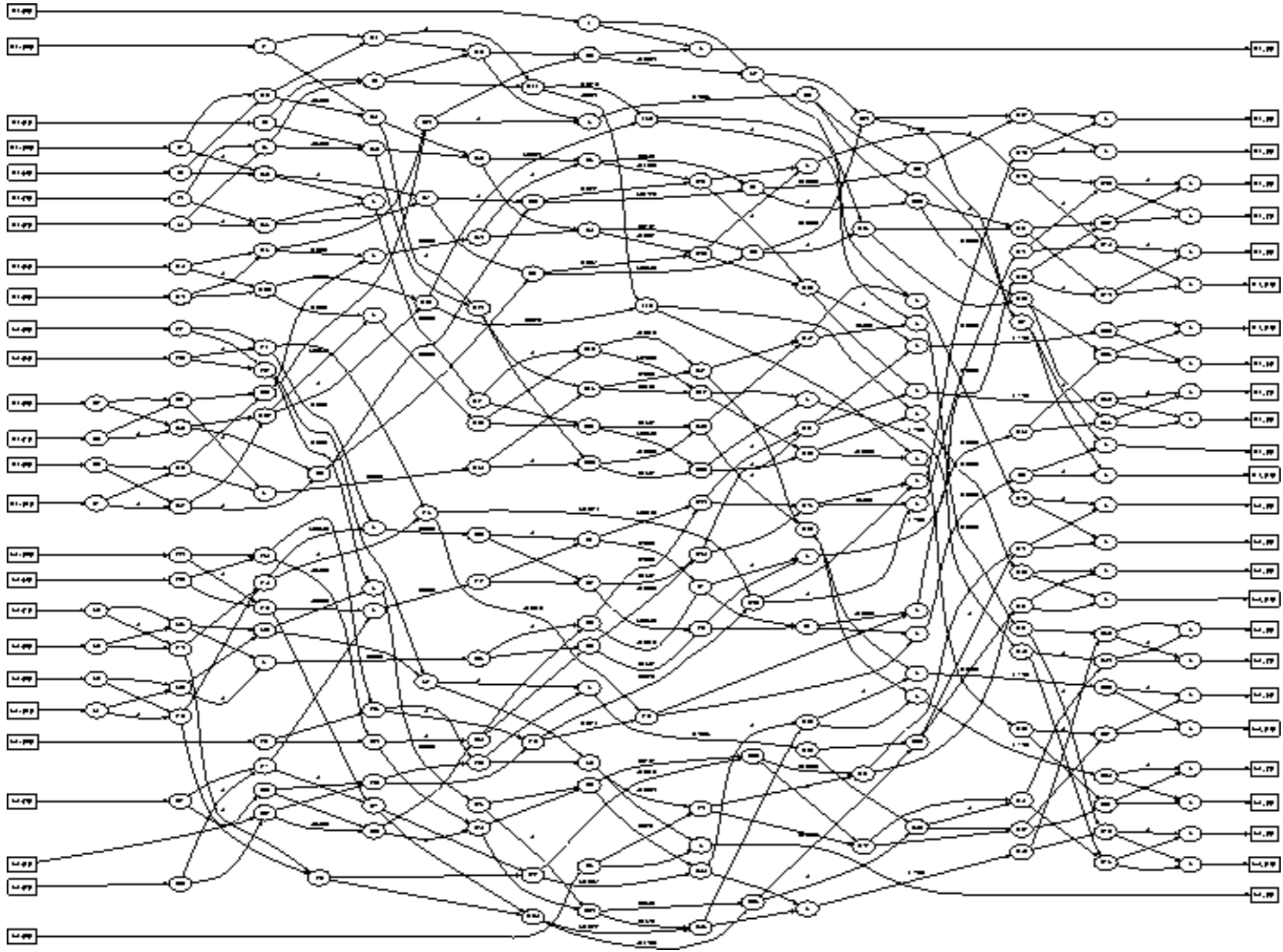
**reasonable
implementation
(Numerical recipes,
GNU scientific library)**

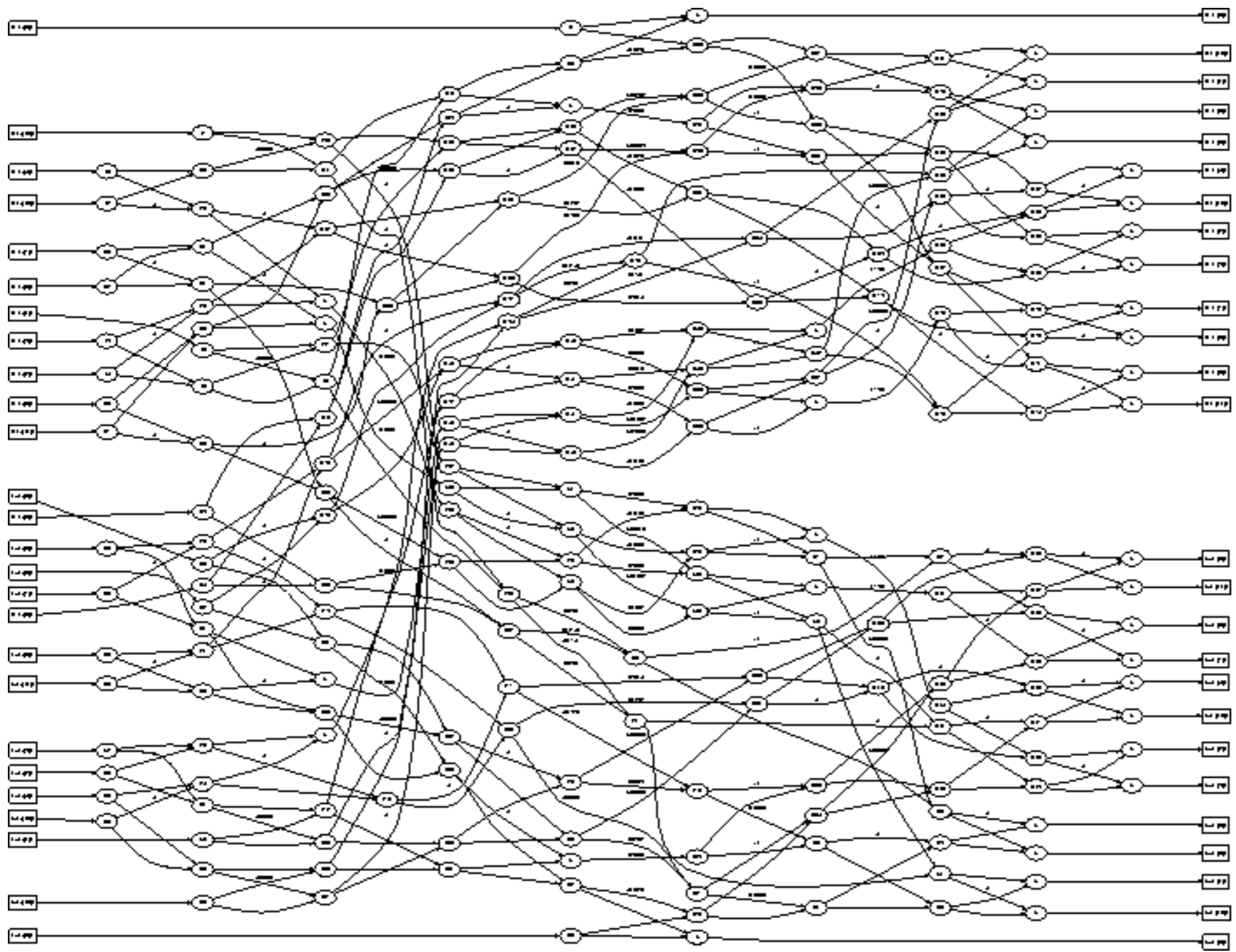
Source: J. Johnson (2007), CScADS autotuning workshop



FFTW (1997)

- “Fastest Fourier Transform in the West” [M. Frigo, S. Johnson (1997)]
- “**Codelet**” generator (in OCaml)
 - Explicit represent a small fixed-size transform by its computation DAG
 - Optimize DAG: Algebraic transformations, constant folding, “DAG transposition”
 - Schedule DAG cache-obliviously and output as C source code
- **Planner**: At run-time, determine which codelets to apply
- **Executor**: Perform FFT of a particular size using plan
- Efficient “plug-in” assembly kernels





Cooley-Tukey FFT algorithm

$$y[k] \leftarrow \text{DFT}_N(x, k) \equiv \sum_{j=0}^{N-1} x[j] \cdot \omega_N^{-kj} \quad x, y \in \mathbb{C}^N$$

$$\omega_N \equiv e^{2\pi\sqrt{-1}/N}$$

$$N \equiv N_1 \cdot N_2$$



Cooley-Tukey FFT algorithm

$$y[k] \leftarrow \text{DFT}_N(x, k) \equiv \sum_{j=0}^{N-1} x[j] \cdot \omega_N^{-kj} \quad x, y \in \mathbb{C}^N$$

$$\omega_N \equiv e^{2\pi\sqrt{-1}/N}$$

$$N \equiv N_1 \cdot N_2$$

↓

$$0 \leq k_1 < N_1 \quad \text{and} \quad 0 \leq k_2 < N_2$$

$$y[k_1 + k_2 \cdot N_1] \leftarrow \sum_{n_2=0}^{N_2-1} \left[\left(\sum_{n_1=0}^{N_1-1} x[n_1 \cdot N_2 + n_2] \cdot \omega_{N_1}^{-k_1 n_1} \right) \cdot \omega_N^{-k_1 n_2} \right] \cdot \omega_{N_2}^{-k_2 n_2}$$

N_1 -point DFT

Twiddle

N_2 -point DFT

Cooley-Tukey FFT algorithm: Encoding in the codelet generator

$$y[k] \leftarrow \text{DFT}_N(x, k) \equiv \sum_{j=0}^{N-1} x[j] \cdot \omega_N^{-kj} \quad x, y \in \mathbb{C}^N$$

$$y[k_1 + k_2 \cdot N_1] \leftarrow \sum_{n_2=0}^{N_2-1} \left[\underbrace{\left(\sum_{n_1=0}^{N_1-1} x[n_1 \cdot N_2 + n_2] \cdot \omega_{N_1}^{-k_1 n_1} \right)}_{N_1\text{-point DFT}} \cdot \underbrace{\omega_N^{-k_1 n_2}}_{\text{Twiddle}} \right] \cdot \omega_{N_2}^{-k_2 n_2}$$

N_2 -point DFT

let dftgen(N, x) \equiv fun $k \rightarrow \dots$ # $\text{DFT}_N(x, k)$

let cooley_tukey(N_1, N_2, x) \equiv

let $\hat{x} \equiv$ fun $n_2, n_1 \rightarrow x(n_2 + n_1 \cdot N_2)$ in

let $\mathbf{G}_1 \equiv$ fun $n_2 \rightarrow \text{dftgen}(N_1, \hat{x}(n_2, --))$ in

let $\mathbf{W} \equiv$ fun $k_1, n_2 \rightarrow \mathbf{G}_1(n_2, k_1) \cdot \omega_N^{-k_1 n_2}$ in

let $\mathbf{G}_2 \equiv$ fun $k_1 \rightarrow \text{dftgen}(N_2, \mathbf{W}(k_1, --))$

in

fun $k \rightarrow \mathbf{G}_2(k \bmod N_1, k \text{ div } N_1)$

**(Functional
pseudo-code)**

Planner phase

Published in *Proc. IEEE*, vol. 93, no. 2, pp. 216–231 (2005).

```
fftw_plan plan;  
fftw_complex in[n], out[n];
```

```
/* plan a 1d forward DFT: */  
plan = fftw_plan_dft_1d(n, in, out,  
                        FFTW_FORWARD, FFTW_PATIENT);
```

Initialize in[] with some data...

```
fftw_execute(plan); // compute DFT
```

Write some new data to in[] ...

```
fftw_execute(plan); // reuse plan
```

**Assembles plan
using dynamic
programming**

Fig. 8. Example of FFTW's use. The user must first create a plan, which can be then used for many transforms of the same size.

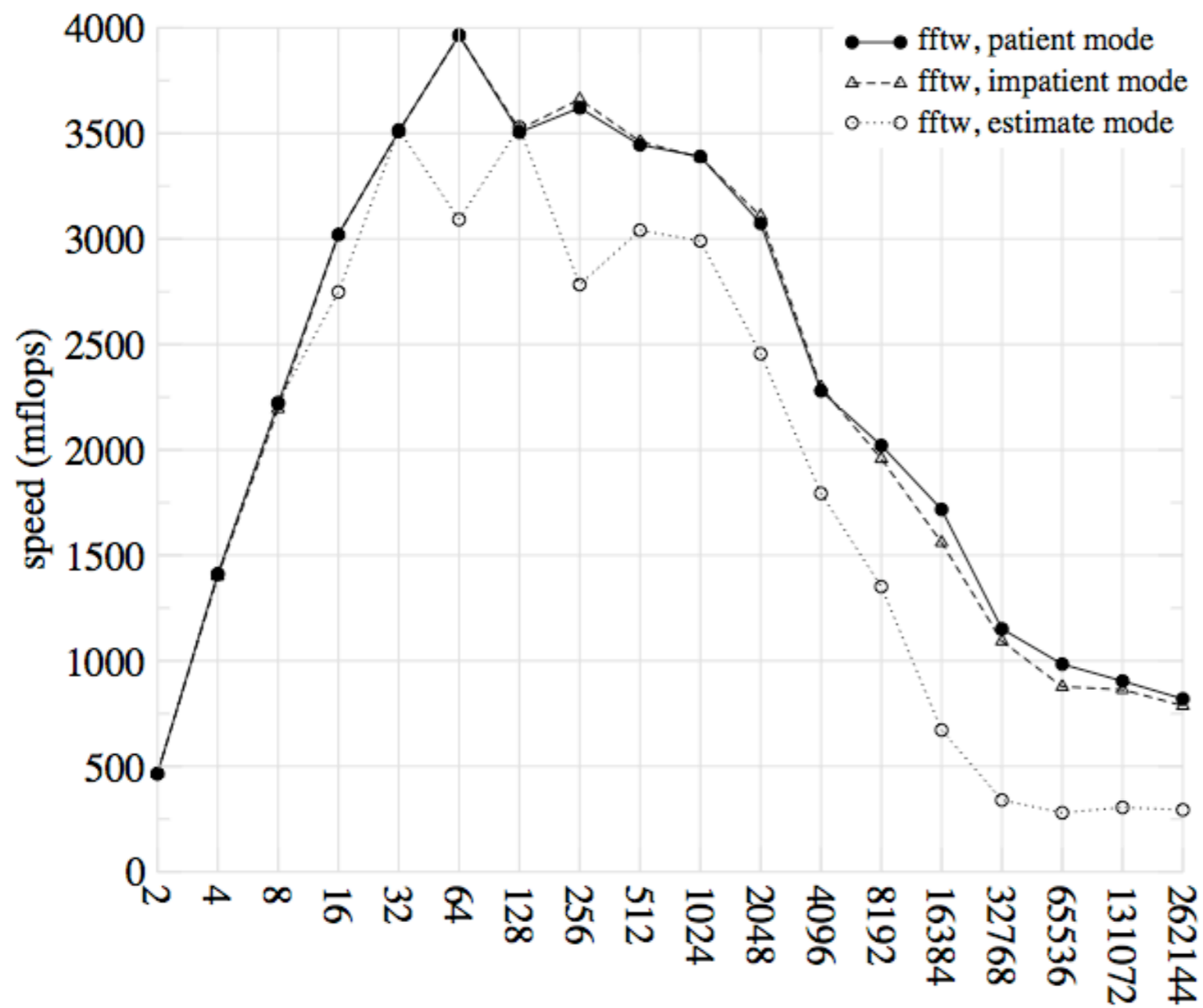


Fig. 9. Effect of planner tradeoffs: comparison of patient, impatient, and estimate modes in FFTW for double-precision 1d complex DFTs, power-of-two sizes, on a 2 GHz PowerPC 970 (G5). Compiler and flags as in Fig. 4.

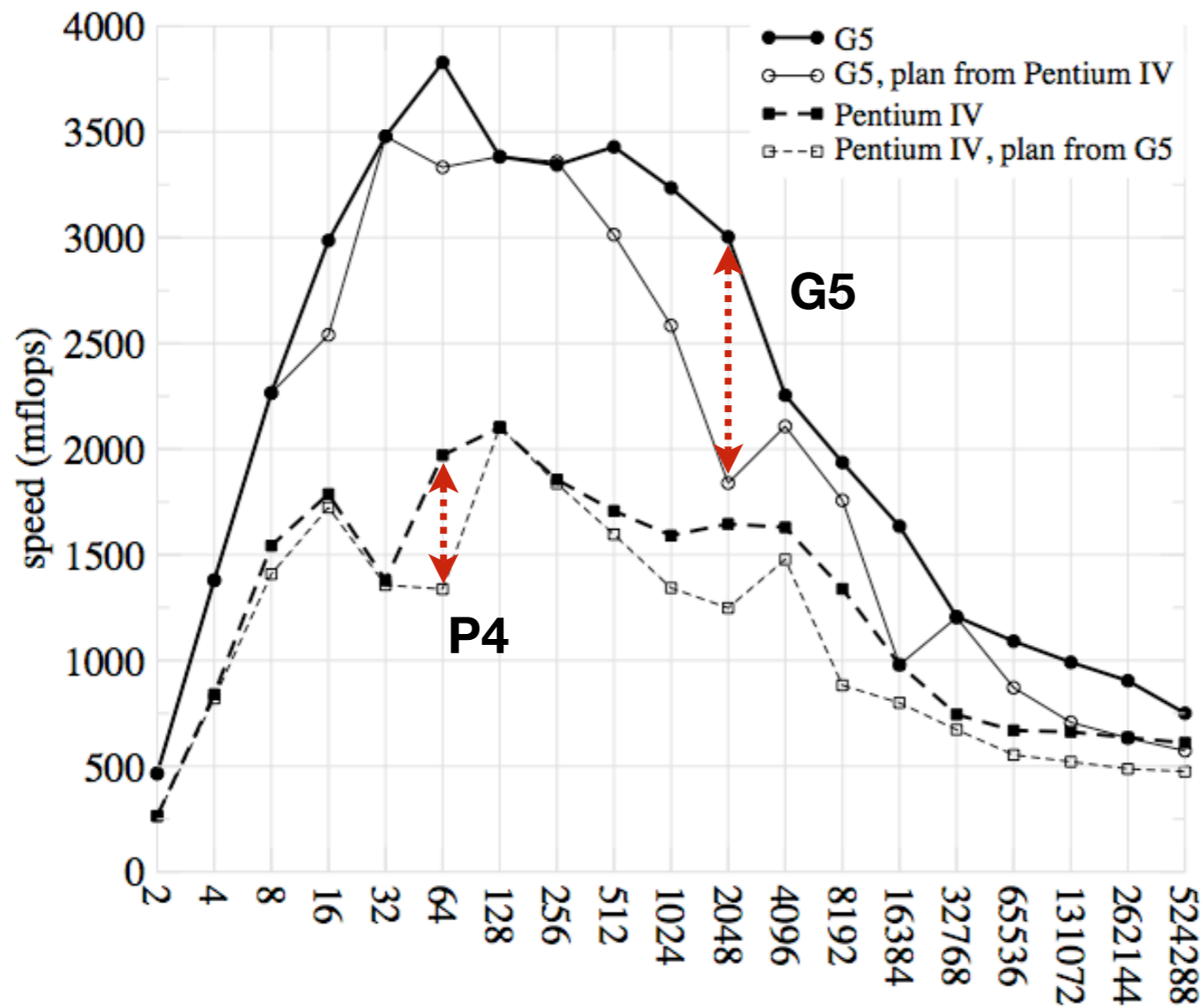
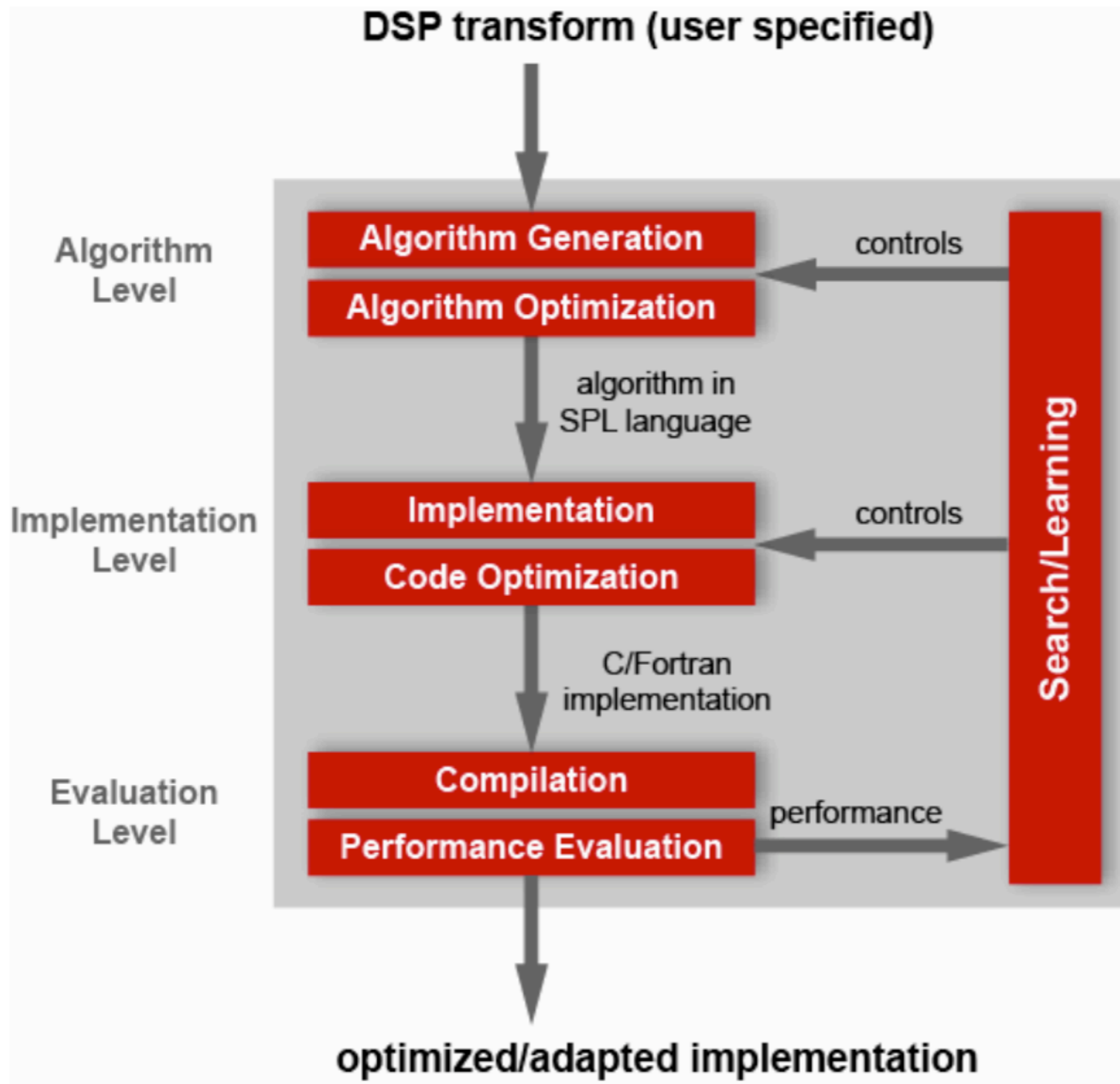


Fig. 10. Effects of tuning FFTW on one machine and running it on another. The graph shows the performance of one-dimensional DFTs on two machines: a 2 GHz PowerPC 970 (G5), and a 2.8 GHz Pentium IV. For each machine, we report both the speed of FFTW tuned to that machine and the speed tuned to the *other* machine.

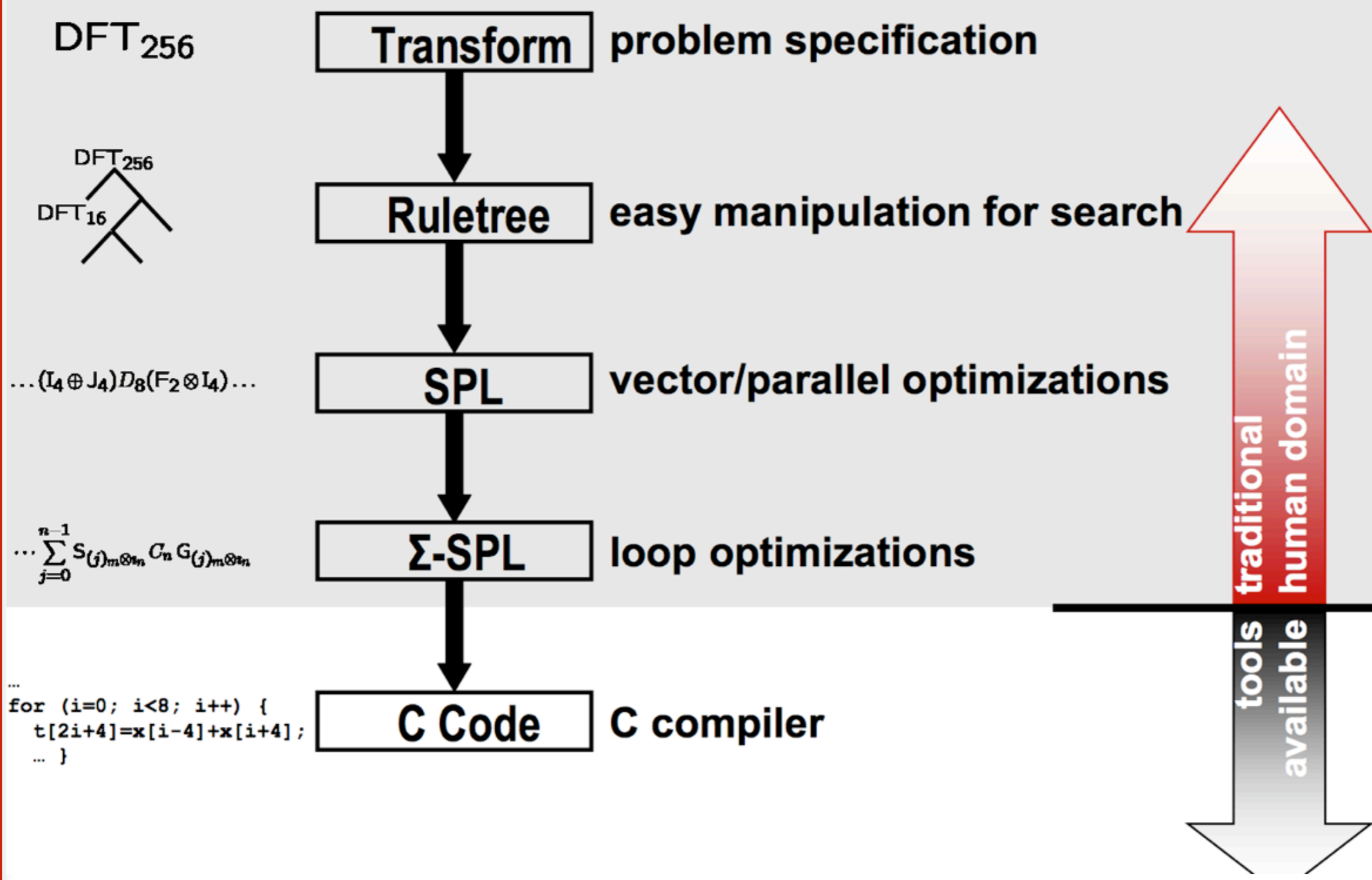


SPIRAL (1998)

- Code generator
 - Represent linear transformations as formulas
 - Symbolic algebra + rewrite engine transforms formulas
- Search using variety of techniques (more later)



Source: J. Johnson (2007), CScADS autotuning workshop



Source: J. Johnson (2007), CScADS autotuning workshop

High-level representations and rewrite rules

$$\mathbf{DFT}_N \equiv \left[\omega_N^{kl} \right]_{0 \leq k, l < N}$$

$$\mathbf{DCT-2}_N \equiv \left[\cos \frac{(2l+1)k\pi}{2N} \right]_{0 \leq k, l < N}$$

⋮

$n = k \cdot m$:

$$\implies \mathbf{DFT}_n \rightarrow (\mathbf{DFT}_k \otimes I_m) T_m^n (I_k \otimes \mathbf{DFT}_m) L_k^n$$

$n = k \cdot m, \gcd(k, m) = 1$:

$$\implies \mathbf{DFT}_n \rightarrow P_n (\mathbf{DFT}_k \otimes \mathbf{DFT}_m) Q_n$$

p is prime :

$$\implies \mathbf{DFT}_p \rightarrow R_p^T (I_1 \oplus \mathbf{DFT}_{p-1}) D_p (I_1 \oplus \mathbf{DFT}_{p-1}) R_p$$

⋮

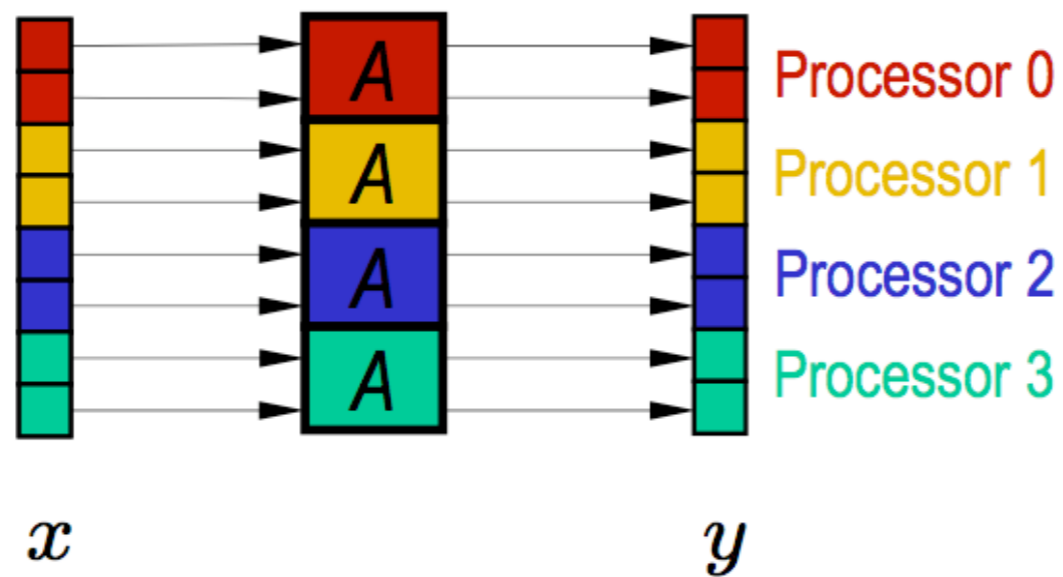
$$\mathbf{DFT}_2 \rightarrow \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$$



High-level representations expose parallelism

$$(I_4 \otimes A) \cdot \begin{bmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \end{bmatrix} = \begin{bmatrix} A & & & \\ & A & & \\ & & A & \\ & & & A \end{bmatrix} \cdot \begin{bmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \end{bmatrix} = \begin{bmatrix} AX_1 \\ AX_2 \\ AX_3 \\ AX_4 \end{bmatrix}$$

A applied 4 times independently



High-level representations expose parallelism

$$\begin{aligned} \left(\begin{bmatrix} a & b \\ c & d \end{bmatrix} \otimes I_2 \right) \cdot \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} &= \begin{bmatrix} a \cdot I_2 & b \cdot I_2 \\ c \cdot I_2 & d \cdot I_2 \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} \\ &= \begin{bmatrix} a \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + b \begin{bmatrix} x_3 \\ x_4 \end{bmatrix} \\ c \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + d \begin{bmatrix} x_3 \\ x_4 \end{bmatrix} \end{bmatrix} \end{aligned}$$

SIMD-vectorizable



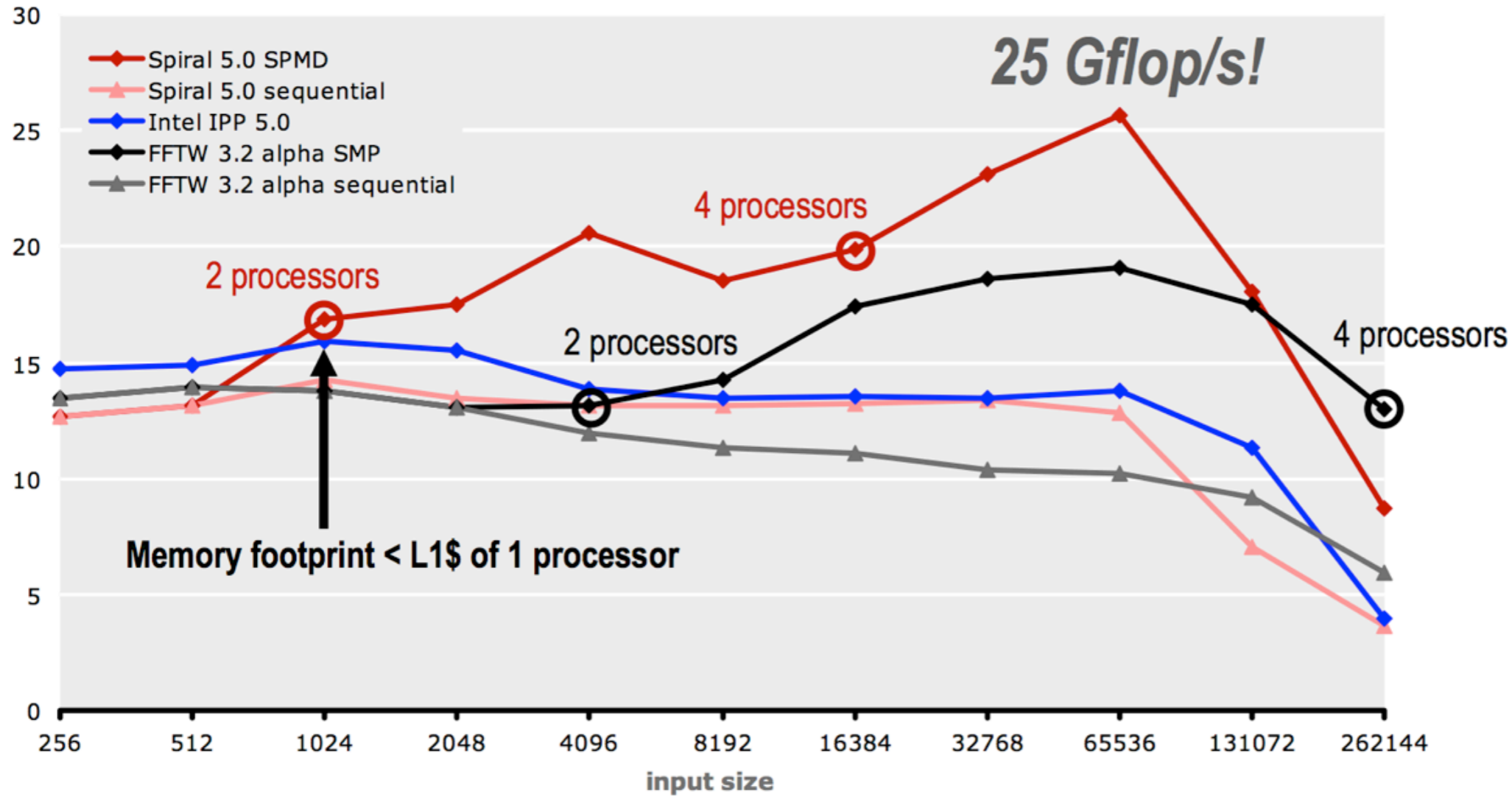


Search in SPIRAL

- Search over ruletrees, i.e., possible formula expansions
- Empirical search
 - Exhaustive
 - Random
 - Dynamic programming
 - Evolutionary search
 - Hill climbing
 - Machine learning methods

Example: SMP + vectorization results

DFT (single precision): on 3 GHz 2 x Core 2 Extreme
performance [Gflop/s]



4-way vectorized + up to 4-threaded + adapted to the memory hierarchy

Source: F. Franchetti (2007), CScADS autotuning workshop



Administrivia



Upcoming schedule changes

- Some adjustment of topics (TBD)
- Tu 3/11 — Project proposals due
- Th 3/13 — SIAM Parallel Processing (attendance encouraged)
- Tu 4/1 — No class
- Th 4/3 — Attend talk by Doug Post from DoD HPC Modernization Program



Homework 1:


Parallel conjugate gradients

- Put name on write-up!
- Grading: 100 pts max
 - Correct implementation — 50 pts
 - Evaluation — 30 pts
 - Tested on two samples matrices — 5
 - Implemented and tested on stencil — 10
 - “Explained” performance (e.g., per proc, load balance, comp. vs. comm) — 15
 - Performance model — 15 pts
 - Write-up “quality” — 5 pts



Projects

- **Proposals due Tu 3/11**
- Your goal should be to do something useful, interesting, and/or publishable!
 - Something you're already working on, suitably adapted for this course
 - Faculty-sponsored/mentored
 - Collaborations encouraged



My criteria for “approving” your project

- “Relevant to this course:” Many themes, so think (and “do”) broadly
 - Parallelism and architectures
 - Numerical algorithms
 - Programming models
 - Performance modeling/analysis



General styles of projects

- Theoretical: Prove something hard (high risk)
- Experimental:
 - Parallelize something
 - Take existing parallel program, and improve it using models & experiments
 - Evaluate algorithm, architecture, or programming model

Examples

- *Anything of interest to a faculty member/project outside CoC*
- Parallel sparse triple product ($R^*A^*R^T$, used in multigrid)
- Future FFT
- Out-of-core or I/O-intensive data analysis and algorithms
- Block iterative solvers (convergence & performance trade-offs)
- Sparse LU
- Data structures and algorithms (trees, graphs)
- Look at mixed-precision
- Discrete-event approaches to continuous systems simulation
- Automated performance analysis and modeling, tuning
- “Unconventional,” but related
 - Distributed deadlock detection for MPI
 - UPC language extensions (dynamic block sizes)
 - Exact linear algebra





Sparse linear algebra



Key distinctions in autotuning work for sparse kernels

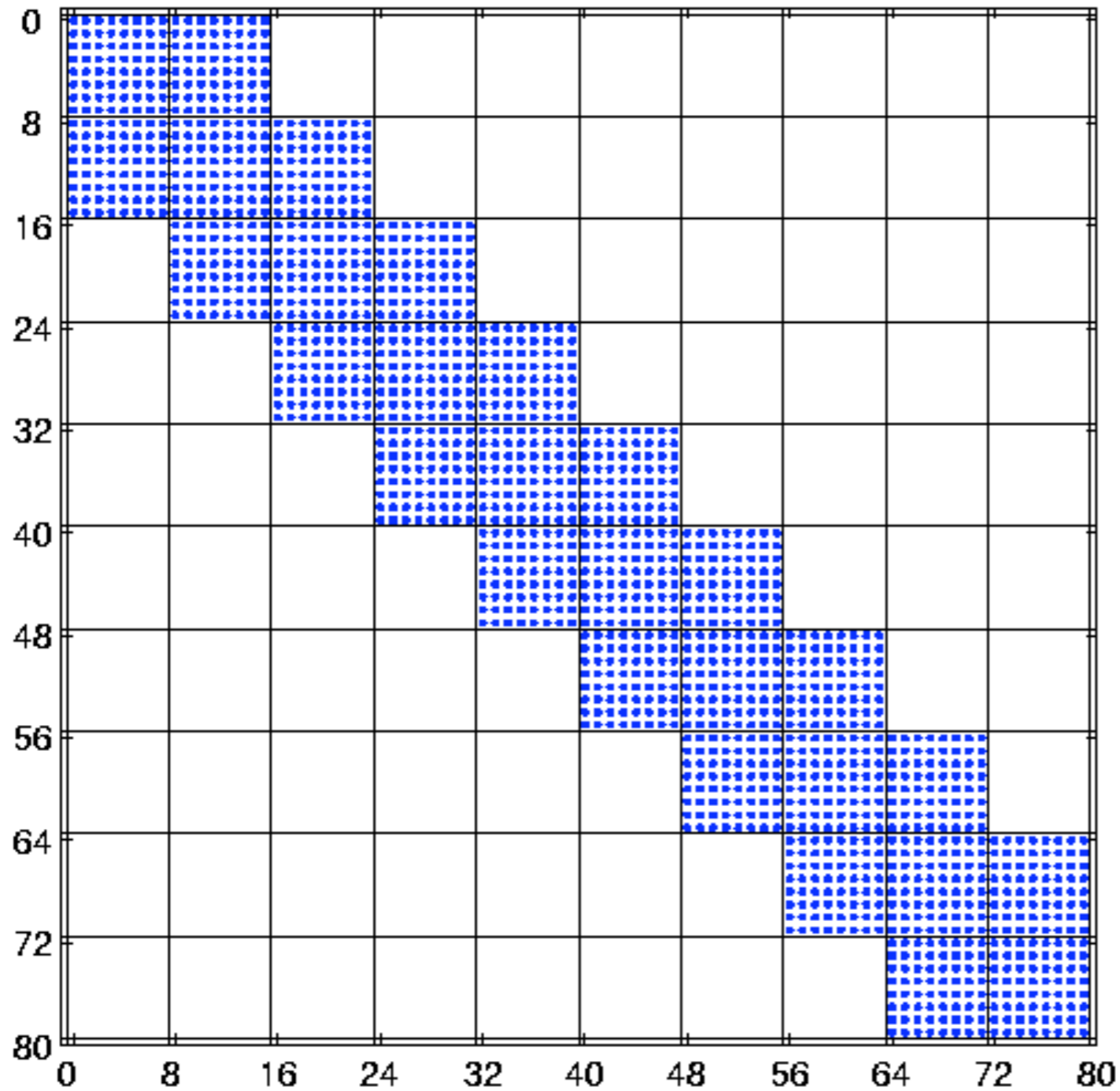
- Data structure transformations
 - Recall HW1
 - Sparse data structures require meta-data overhead
 - Sparse matrix-vector multiply (SpMV) is memory bound
 - Bandwidth limited \Rightarrow minimize data structure size
- Run-time tuning: Need lightweight techniques
- Extra flops pay off



Sparsity (1998) and OSKI (2005)

- Berkeley projects (BeBOP group: Demmel & Yelick; Im, Vuduc, *et al.*)
 - PHiPAC \Rightarrow SPARSITY \Rightarrow OSKI
 - On-going: See multicore optimizations by Williams, *et al.*, in SC 2007
- Motivation: Sparse matrix-vector multiply (SpMV) \leq 10% peak or less
 - Indirect, irregular memory access
 - Low q vs. dense case
 - Depends on machine and matrix, possibly unknown until run-time

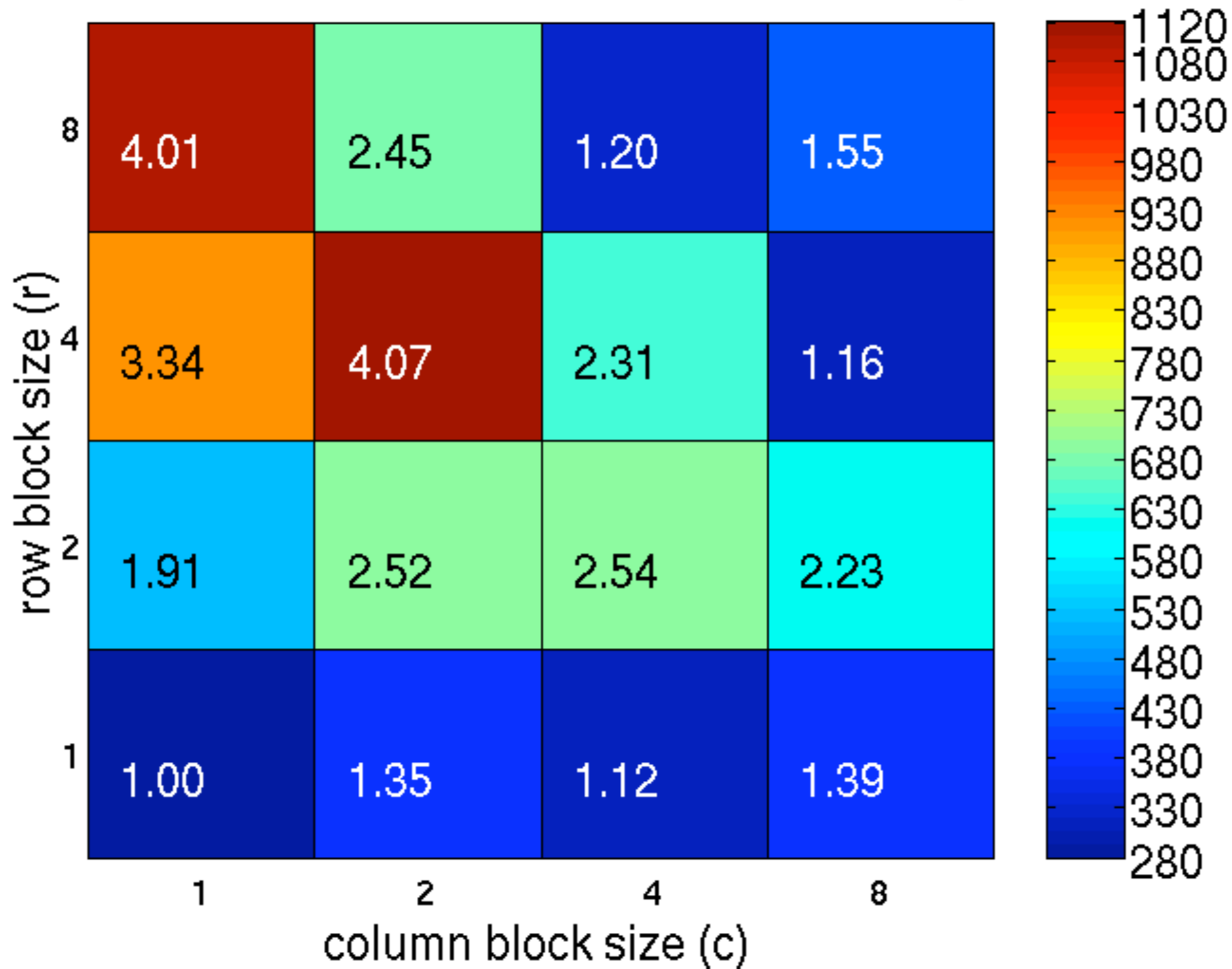
Matrix 02-raefsky3



1792 ideal nz + 0 explicit zeros = 1792 nz

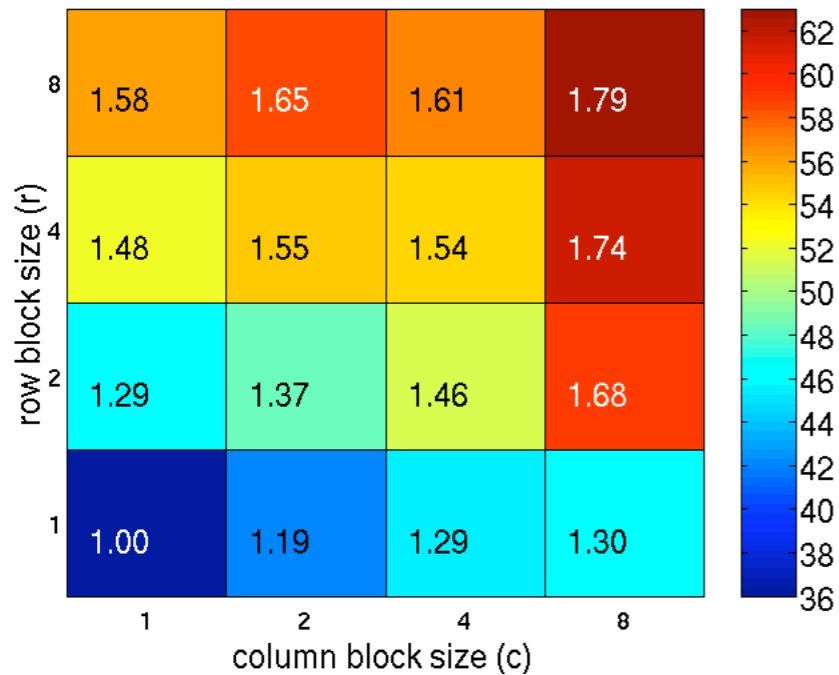


900 MHz Itanium 2, Intel C v8: ref=275 Mflop/s

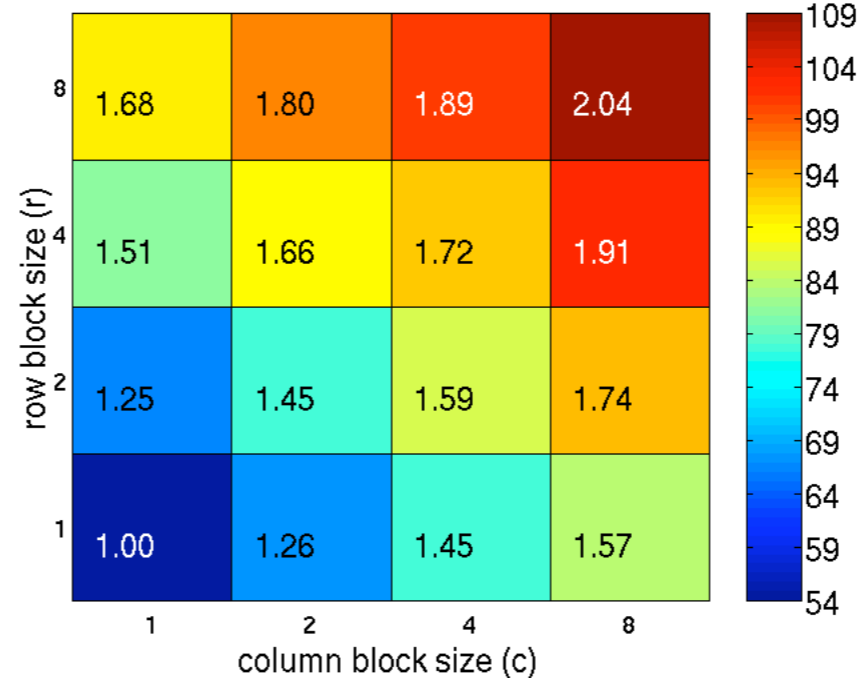




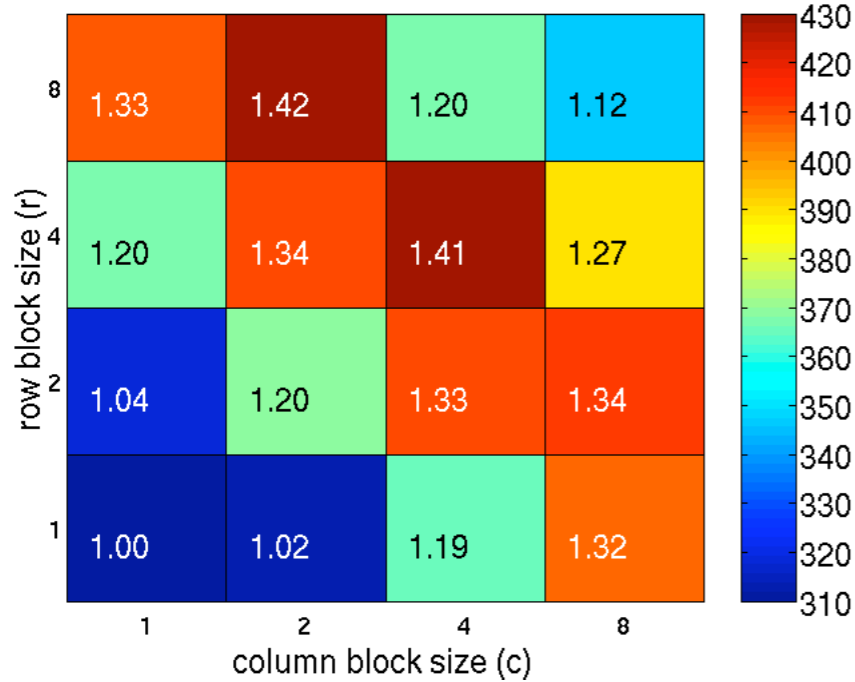
333 MHz Sun Ultra 2i, Sun C v6.0: ref=35 Mflop/s



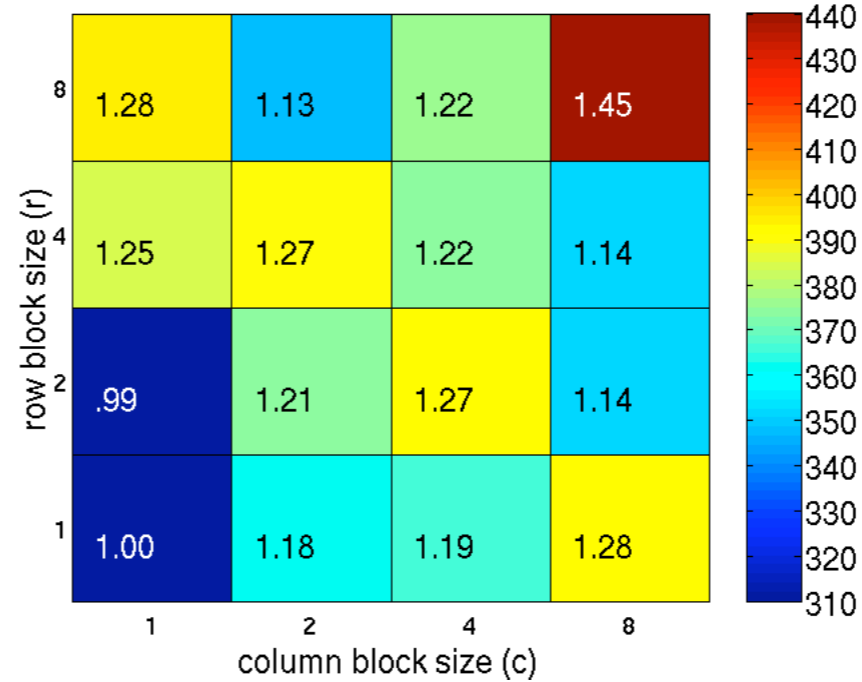
900 MHz Ultra 3, Sun CC v6: ref=54 Mflop/s



2 GHz Pentium M, Intel C v8.1: ref=308 Mflop/s

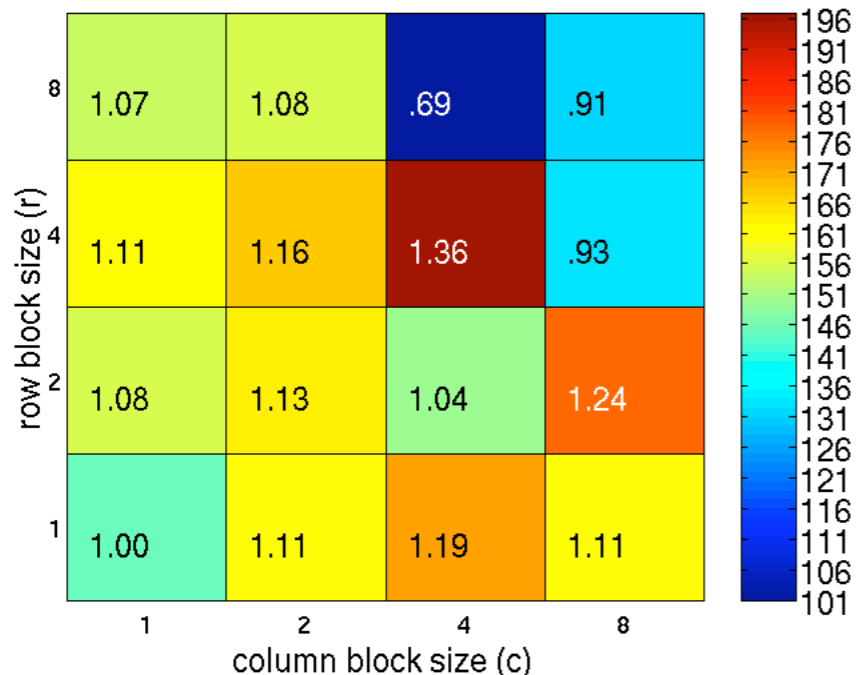


1.4 GHz Opteron, gcc 3.4.2: ref=308 Mflop/s

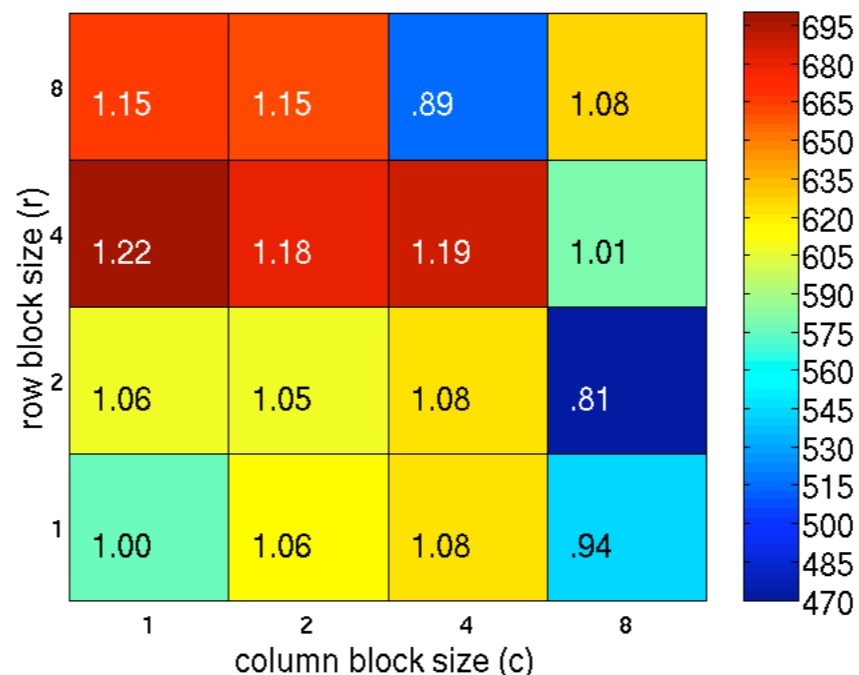




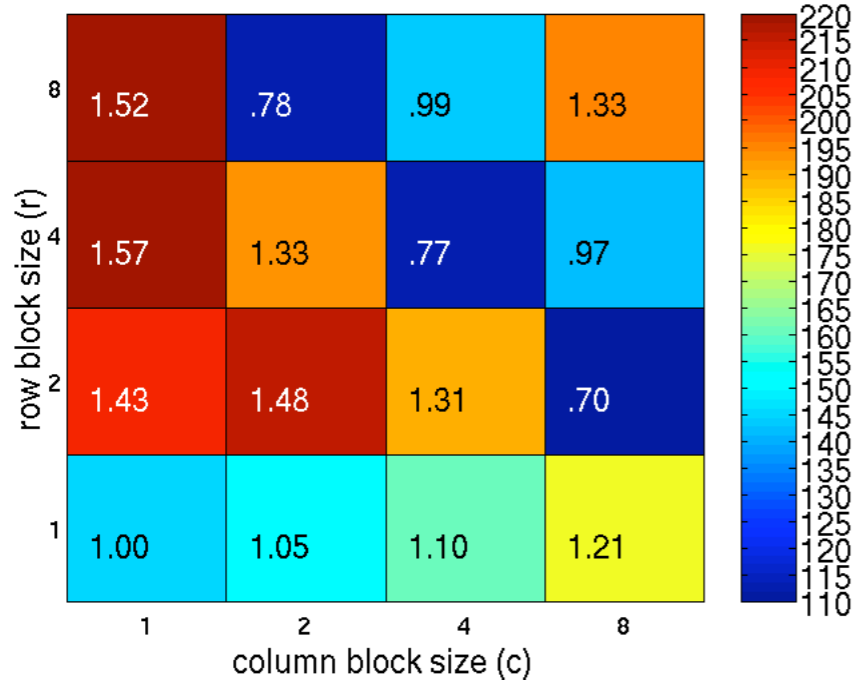
375 MHz Power3, IBM xlc v6: ref=145 Mflop/s



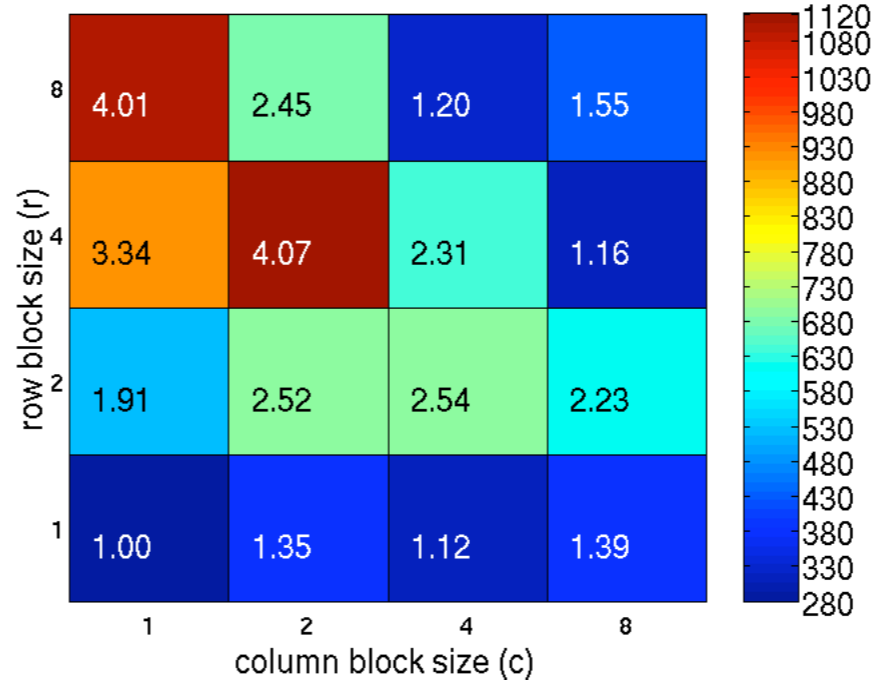
1.3 GHz Power4, IBM xlc v6: ref=577 Mflop/s



800 MHz Itanium, Intel C v7: ref=146 Mflop/s

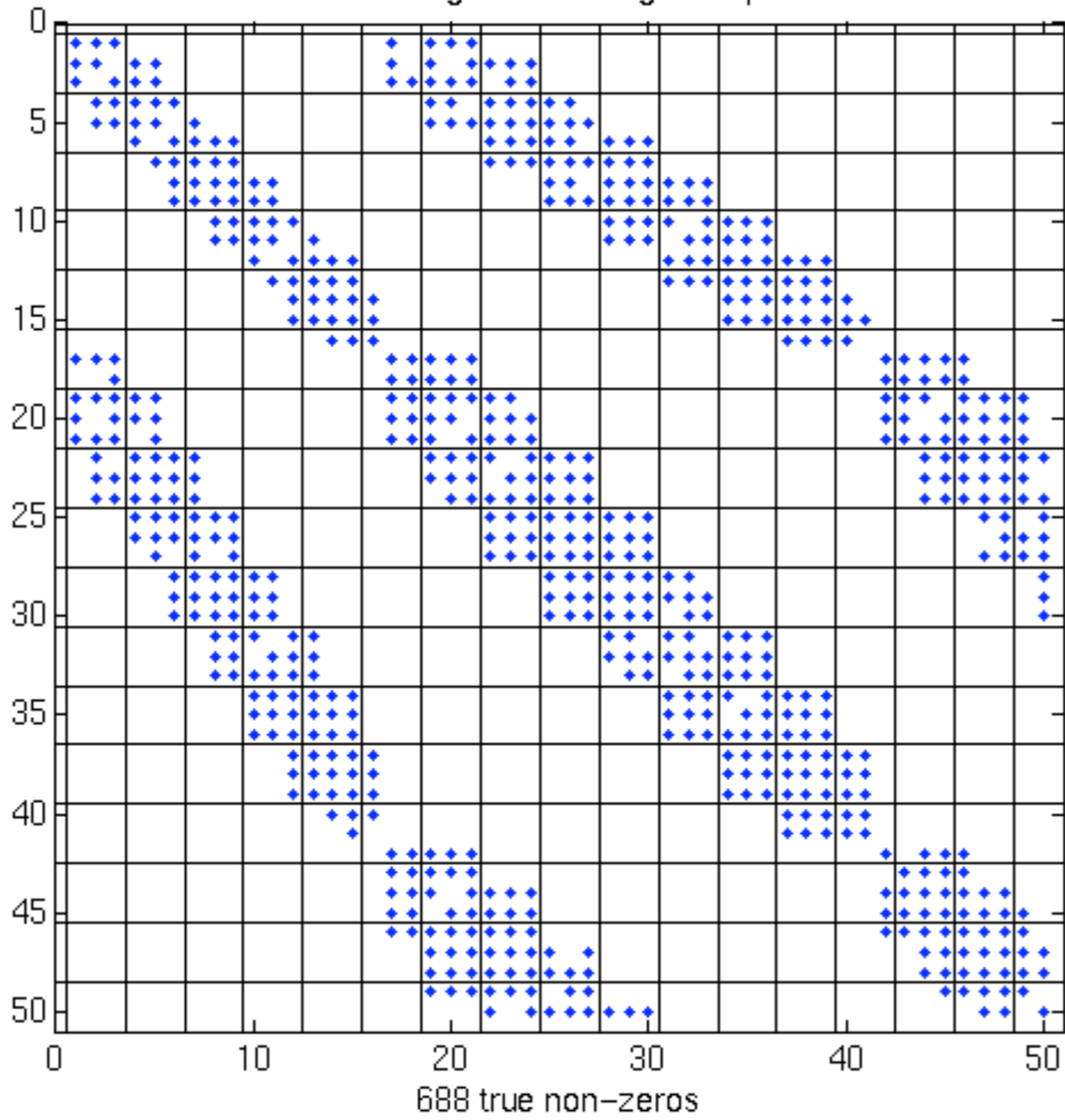


900 MHz Itanium 2, Intel C v8: ref=275 Mflop/s

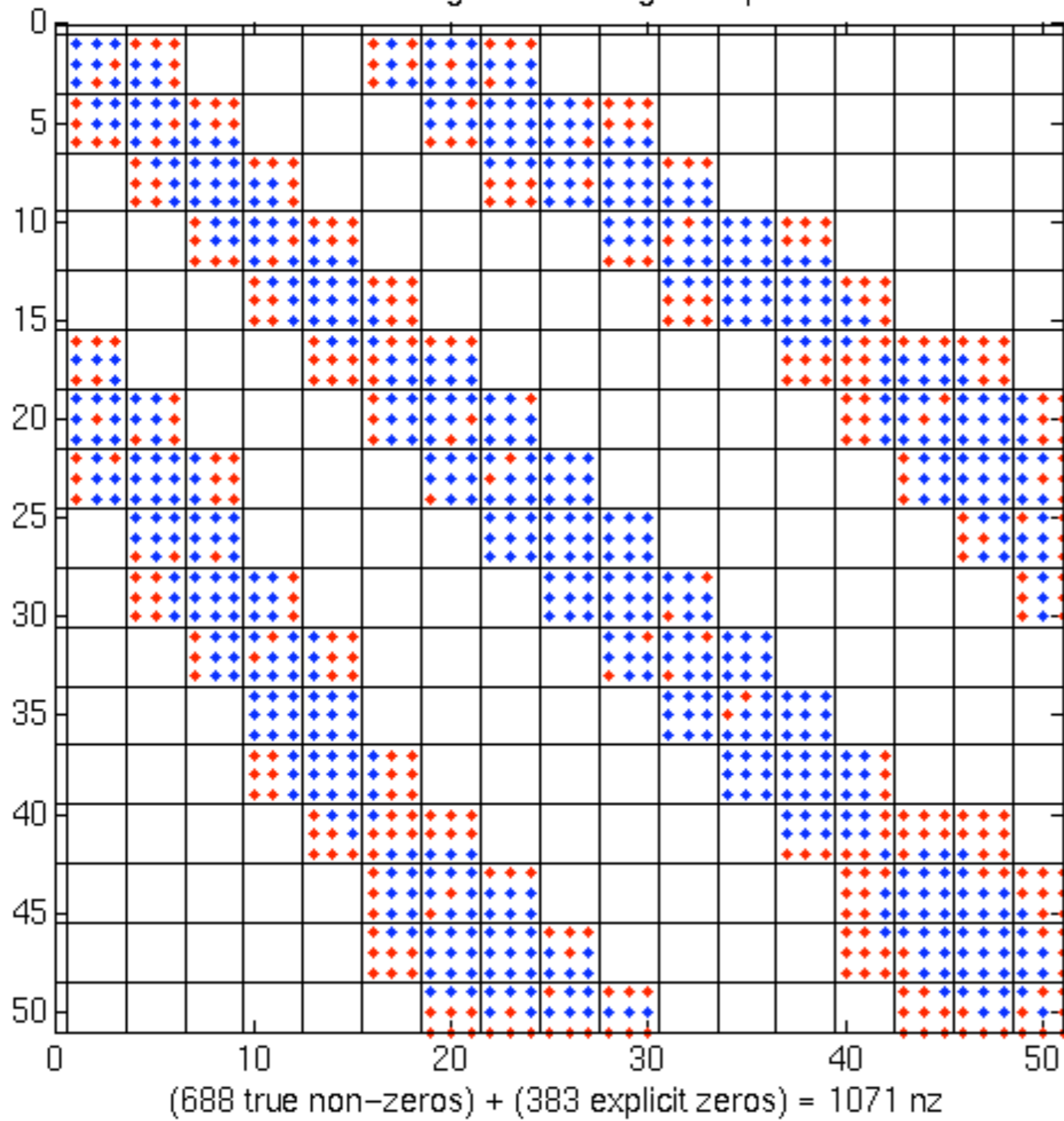




3 x 3 Register Blocking Example



3 x 3 Register Blocking Example



50% extra zeros

**1.5x faster
(2/3 time) on
Pentium III**



How OSKI tunes

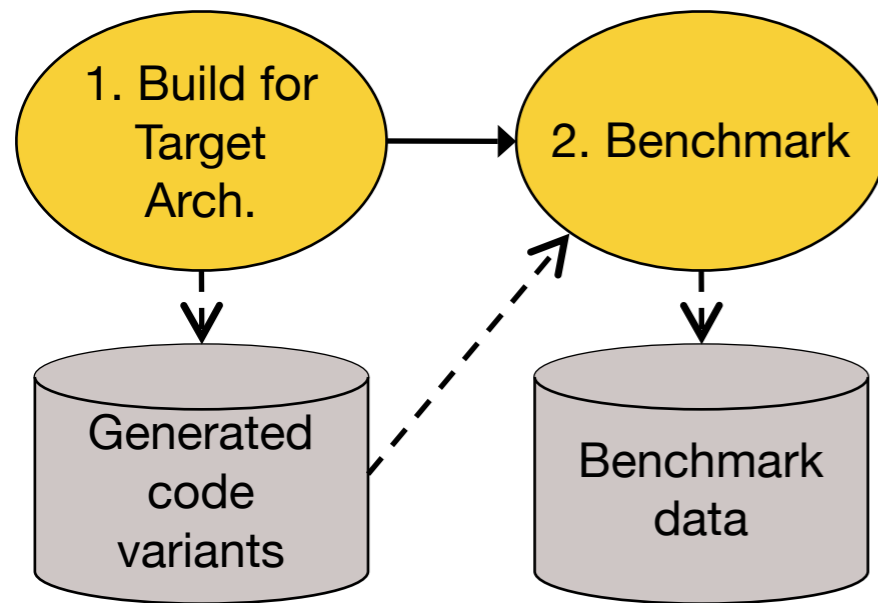
Library Install-Time (offline) ← → Application Run-Time



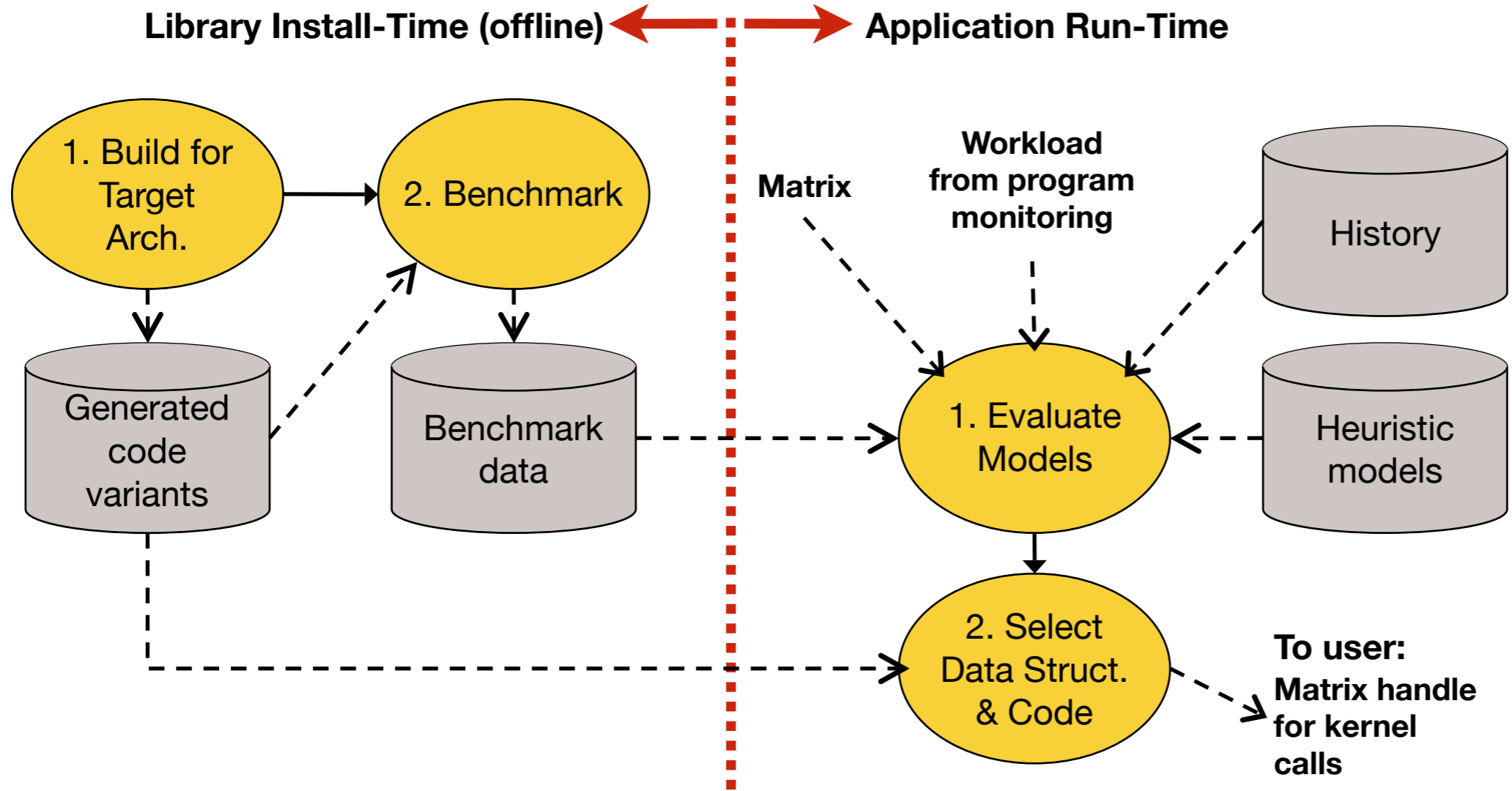
How OSKI tunes




Library Install-Time (offline) ← Application Run-Time



How OSKI tunes





Heuristic model example: Selecting a block size

- Idea: Hybrid off-line/run-time model
 - **Offline benchmark:** Measure Mflops(r, c) on dense matrix in sparse format
 - **Run-time:** Sample matrix to quickly estimate Fill(r, c)
 - Run-time **model:** Choose r, c to maximize Mflops(r, c) / Fill(r, c)
 - Accurate in practice (selects $r \times c$ with performance within 10% of best)
- Run-time **cost?**
 - Roughly 40 SpMV
 - Dominated by conversion (~80%)



Workload tuning

- Consider BiCG solver: Equal mix of A^*x and A^T*y (independent)
 - 3×1 : $A \cdot x, A^T \cdot y = 1053, 343 \text{ Mflop/s} \Rightarrow 517 \text{ Mflop/s}$
 - 3×3 : $A \cdot x, A^T \cdot y = 806, 826 \text{ Mflop/s} \Rightarrow 816 \text{ Mflop/s}$
- Higher-level operation: Fused (A^*x, A^T*y) kernel
 - 3×1 : 757 Mflop/s
 - 3×3 : 1400 Mflop/s



Tensor Contraction Engine (TCE) for quantum chemistry



Tensor Contraction Engine (TCE)

- Application domain: Quantum chemistry
 - Electronic structure calculations
 - Dominant computation expressible as a “tensor contraction”
- TCE generates a complete parallel program from a high-level spec
 - Automates time-space trade-offs
 - Output
- S. Hirata (2002), and many others
- Following presentation taken from Proc. IEEE 2005 special issue

Motivation: Simplify program development

$$\begin{aligned} \text{hbar}[a,b,i,j] = & \text{sum}[f[b,c] * t[i,j,a,c], c] - \text{sum}[f[k,c] * t[k,b] * t[i,j,a,c], k,c] + \text{sum}[f[a,c] * t[i,j,c,b], c] - \text{sum}[f[k,c] * t[k,a] * t[i,j,c,b], k,c] - \text{sum}[f[k,j] * t[i,k,a,b], k] - \text{sum}[f[k,c] * \\ & t[j,c] * t[i,k,a,b], k,c] - \text{sum}[f[k,i] * t[j,k,b,a], k] - \text{sum}[f[k,c] * t[i,c] * t[j,k,b,a], k,c] + \text{sum}[t[i,c] * t[j,d] * v[a,b,c,d], c,d] + \text{sum}[t[i,j,c,d] * v[a,b,c,d], c,d] + \text{sum}[t[j,c] * v[a,b,i,c], \\ & c] - \text{sum}[t[k,b] * v[a,k,i,j], k] + \text{sum}[t[i,c] * v[b,a,j,c], c] - \text{sum}[t[k,a] * v[b,k,j,i], k] - \text{sum}[t[k,d] * t[i,j,c,b] * v[k,a,c,d], k,c,d] - \text{sum}[t[i,c] * t[j,k,b,d] * v[k,a,c,d], k,c,d] - \text{sum}[t[j,c] \\ & * t[k,b] * v[k,a,c,i], k,c] + 2 * \text{sum}[t[j,k,b,c] * v[k,a,c,i], k,c] - \text{sum}[t[j,k,c,b] * v[k,a,c,i], k,c] - \text{sum}[t[i,c] * t[j,d] * t[k,b] * v[k,a,d,c], k,c,d] + 2 * \text{sum}[t[k,d] * t[i,j,c,b] * v[k,a,d,c], \\ & k,c,d] - \text{sum}[t[k,b] * t[i,j,c,d] * v[k,a,d,c], k,c,d] - \text{sum}[t[j,d] * t[i,k,c,b] * v[k,a,d,c], k,c,d] + 2 * \text{sum}[t[i,c] * t[j,k,b,d] * v[k,a,d,c], k,c,d] - \text{sum}[t[i,c] * t[j,k,d,b] * v[k,a,d,c], k,c,d] - \\ & \text{sum}[t[j,k,b,c] * v[k,a,i,c], k,c] - \text{sum}[t[i,c] * t[k,b] * v[k,a,j,c], k,c] - \text{sum}[t[i,k,c,b] * v[k,a,j,c], k,c] - \text{sum}[t[i,c] * t[j,d] * t[k,a] * v[k,b,c,d], k,c,d] - \text{sum}[t[k,d] * t[i,j,a,c] * v[k,b,c,d], \\ & k,c,d] - \text{sum}[t[k,a] * t[i,j,c,d] * v[k,b,c,d], k,c,d] + 2 * \text{sum}[t[j,d] * t[i,k,a,c] * v[k,b,c,d], k,c,d] - \text{sum}[t[j,d] * t[i,k,c,a] * v[k,b,c,d], k,c,d] - \text{sum}[t[i,c] * t[j,k,d,a] * v[k,b,c,d], k,c,d] \\ & - \text{sum}[t[i,c] * t[k,a] * v[k,b,c,j], k,c] + 2 * \text{sum}[t[i,k,a,c] * v[k,b,c,j], k,c] - \text{sum}[t[i,k,c,a] * v[k,b,c,j], k,c] + 2 * \text{sum}[t[k,d] * t[i,j,a,c] * v[k,b,d,c], k,c,d] - \text{sum}[t[j,d] * t[i,k,a,c] * \\ & v[k,b,d,c], k,c,d] - \text{sum}[t[j,c] * t[k,a] * v[k,b,i,c], k,c] - \text{sum}[t[j,k,c,a] * v[k,b,i,c], k,c] - \text{sum}[t[i,k,a,c] * v[k,b,j,c], k,c] + \text{sum}[t[i,c] * t[j,d] * t[k,a] * t[l,b] * v[k,l,c,d], k,l,c,d] - 2 * \\ & \text{sum}[t[k,b] * t[l,d] * t[i,j,a,c] * v[k,l,c,d], k,l,c,d] - 2 * \text{sum}[t[k,a] * t[l,d] * t[i,j,c,b] * v[k,l,c,d], k,l,c,d] + \text{sum}[t[k,a] * t[l,b] * t[i,j,c,d] * v[k,l,c,d], k,l,c,d] - 2 * \text{sum}[t[j,c] * t[l,d] * \\ & t[i,k,a,b] * v[k,l,c,d], k,l,c,d] - 2 * \text{sum}[t[j,d] * t[l,b] * t[i,k,a,c] * v[k,l,c,d], k,l,c,d] + \text{sum}[t[j,d] * t[l,b] * t[i,k,c,a] * v[k,l,c,d], k,l,c,d] - 2 * \text{sum}[t[i,c] * t[l,d] * t[j,k,b,a] * v[k,l,c,d], \\ & k,l,c,d] + \text{sum}[t[i,c] * t[l,a] * t[j,k,b,d] * v[k,l,c,d], k,l,c,d] + \text{sum}[t[i,c] * t[l,b] * t[j,k,d,a] * v[k,l,c,d], k,l,c,d] + \text{sum}[t[i,k,c,d] * t[j,l,b,a] * v[k,l,c,d], k,l,c,d] + 4 * \text{sum}[t[i,k,a,c] * \\ & t[j,l,b,d] * v[k,l,c,d], k,l,c,d] - 2 * \text{sum}[t[i,k,c,a] * t[j,l,b,d] * v[k,l,c,d], k,l,c,d] - 2 * \text{sum}[t[i,k,a,b] * t[j,l,c,d] * v[k,l,c,d], k,l,c,d] - 2 * \text{sum}[t[i,k,a,c] * t[j,l,d,b] * v[k,l,c,d], k,l,c,d] \\ & + \text{sum}[t[i,k,c,a] * t[j,l,d,b] * v[k,l,c,d], k,l,c,d] + \text{sum}[t[i,c] * t[j,d] * t[k,l,a,b] * v[k,l,c,d], k,l,c,d] + \text{sum}[t[i,j,c,d] * t[k,l,a,b] * v[k,l,c,d], k,l,c,d] - 2 * \text{sum}[t[i,j,c,b] * t[k,l,a,d] * \\ & v[k,l,c,d], k,l,c,d] - 2 * \text{sum}[t[i,j,a,c] * t[k,l,b,d] * v[k,l,c,d], k,l,c,d] + \text{sum}[t[j,c] * t[k,b] * t[l,a] * v[k,l,c,i], k,l,c] + \text{sum}[t[l,c] * t[j,k,b,a] * v[k,l,c,i], k,l,c] - 2 * \text{sum}[t[l,a] * t[j,k,b,c] \\ & * v[k,l,c,i], k,l,c] + \text{sum}[t[l,a] * t[j,k,c,b] * v[k,l,c,i], k,l,c] - 2 * \text{sum}[t[k,c] * t[j,l,b,a] * v[k,l,c,i], k,l,c] + \text{sum}[t[k,a] * t[j,l,b,c] * v[k,l,c,i], k,l,c] + \text{sum}[t[k,b] * t[j,l,c,a] * v[k,l,c,i], \\ & k,l,c] + \text{sum}[t[j,c] * t[l,k,a,b] * v[k,l,c,i], k,l,c] + \text{sum}[t[i,c] * t[k,a] * t[l,b] * v[k,l,c,j], k,l,c] + \text{sum}[t[l,c] * t[i,k,a,b] * v[k,l,c,j], k,l,c] - 2 * \text{sum}[t[l,b] * t[i,k,a,c] * v[k,l,c,j], k,l,c] \\ & + \text{sum}[t[l,b] * t[i,k,c,a] * v[k,l,c,j], k,l,c] + \text{sum}[t[i,c] * t[k,l,a,b] * v[k,l,c,j], k,l,c] + \text{sum}[t[j,c] * t[l,d] * t[i,k,a,b] * v[k,l,d,c], k,l,c,d] + \text{sum}[t[j,d] * t[l,b] * t[i,k,a,c] * v[k,l,d,c], \\ & k,l,c,d] + \text{sum}[t[j,d] * t[l,a] * t[i,k,c,b] * v[k,l,d,c], k,l,c,d] - 2 * \text{sum}[t[i,k,c,d] * t[j,l,b,a] * v[k,l,d,c], k,l,c,d] - 2 * \text{sum}[t[i,k,a,c] * t[j,l,b,d] * v[k,l,d,c], k,l,c,d] + \text{sum}[t[i,k,c,a] * \\ & t[j,l,b,d] * v[k,l,d,c], k,l,c,d] + \text{sum}[t[i,k,a,b] * t[j,l,c,d] * v[k,l,d,c], k,l,c,d] + \text{sum}[t[i,k,c,b] * t[j,l,d,a] * v[k,l,d,c], k,l,c,d] + \text{sum}[t[i,k,a,c] * t[j,l,d,b] * v[k,l,d,c], k,l,c,d] + \text{sum}[t[k,a] \\ & * t[l,b] * v[k,l,i,j], k,l] + \text{sum}[t[k,l,a,b] * v[k,l,i,j], k,l] + \text{sum}[t[k,b] * t[l,d] * t[i,j,a,c] * v[l,k,c,d], k,l,c,d] + \text{sum}[t[k,a] * t[l,d] * t[i,j,c,b] * v[l,k,c,d], k,l,c,d] + \text{sum}[t[i,c] * t[l,d] * \\ & t[j,k,b,a] * v[l,k,c,d], k,l,c,d] - 2 * \text{sum}[t[i,c] * t[l,a] * t[j,k,b,d] * v[l,k,c,d], k,l,c,d] + \text{sum}[t[i,c] * t[l,a] * t[j,k,d,b] * v[l,k,c,d], k,l,c,d] + \text{sum}[t[i,j,c,b] * t[k,l,a,d] * v[l,k,c,d], k,l,c,d] \\ & + \text{sum}[t[i,j,a,c] * t[k,l,b,d] * v[l,k,c,d], k,l,c,d] - 2 * \text{sum}[t[l,c] * t[i,k,a,b] * v[l,k,c,j], k,l,c] + \text{sum}[t[l,b] * t[i,k,a,c] * v[l,k,c,j], k,l,c] + \text{sum}[t[l,a] * t[i,k,c,b] * v[l,k,c,j], k,l,c] + \\ & v[a,b,i,j] \end{aligned}$$

Source: Baumgartner, et al. (2005)

Rewriting to reduce operation counts

Naïvely, $\approx 4 \times N^{10}$ flops

$$S_{abij} = \sum_{c,d,e,f,k,l} A_{acik} \times B_{befl} \times C_{dfjk} \times D_{cdel}$$

\Downarrow

$$S_{abij} = \sum_{c,k} \left(\sum_{d,f} \left(\sum_{e,l} B_{befl} \times D_{cdel} \right) \times C_{dfjk} \right) \times A_{acik}$$

Assuming associativity and distributivity, $\approx 6 \times N^6$ flops,
but also **requires temporary storage.**

Source: Baumgartner, et al. (2005)

Operation and storage minimization *via* loop fusion

$$T_{bcdf}^{(1)} = \sum_{e,l} B_{befl} \times D_{cdel}$$

$$T_{bcjk}^{(2)} = \sum_{d,f} T_{bcdf}^{(1)} \times C_{dfjk}$$

$$S_{abij} = \sum_{c,k} T_{bcjk}^{(2)} \times A_{acik}$$

$T1 = T2 = S = 0$

for b, c, d, e, f, l **do**

$T1[b, c, d, f] += B[b, e, f, l] \cdot D[c, d, e, l]$

for b, c, d, f, j, k **do**

$T2[b, c, j, k] += T1[b, c, d, f] \cdot C[d, f, j, k]$

for a, b, c, i, j, k **do**

$S[a, b, i, j] += T2[b, c, j, k] \cdot A[a, c, i, k]$



Operation and storage minimization *via* loop fusion

$$T_{bcdf}^{(1)} = \sum_{e,l} B_{befl} \times D_{cdel}$$

$$T_{bcjk}^{(2)} = \sum_{d,f} T_{bcdf}^{(1)} \times C_{dfjk}$$

$$S_{abij} = \sum_{c,k} T_{bcjk}^{(2)} \times A_{acik}$$

$T1 = T2 = S = 0$

for b, c, d, e, f, l **do**

$T1[b, c, d, f] += B[b, e, f, l] \cdot D[c, d, e, l]$

for b, c, d, f, j, k **do**

$T2[b, c, j, k] += T1[b, c, d, f] \cdot C[d, f, j, k]$

for a, b, c, i, j, k **do**

$S[a, b, i, j] += T2[b, c, j, k] \cdot A[a, c, i, k]$

$S = 0$

for b, c **do**

$T1f \leftarrow 0, T2f \leftarrow 0$

for d, f **do**

for e, l **do**

$T1f += B[b, e, f, l] \cdot D[c, d, e, l]$

for j, k **do**

$T2f[j, k] += T1f \cdot C[d, f, j, k]$

for a, i, j, k **do**

$S[a, b, i, j] += T2f[j, k] \cdot A[a, c, i, k]$

Time-space trade-offs

Max index of
 $a-f$: $O(1000)$
 $i-k$: $O(100)$

for a, e, c, f do

for i, j do

$$X_{aecf} += T_{ijae} \cdot T_{ijcf} \quad \leftarrow \text{“Contraction” of } T \text{ over } i, j$$

for c, e, b, k do

$$T_{cebk}^{(1)} \leftarrow f_1(c, e, b, k)$$

for a, f, b, k do

$$T_{afb}^{(2)} \leftarrow f_2(a, f, b, k)$$

Integrals, $O(1000)$ flops

for c, e, a, f do

for b, k do

$$Y_{ceaf} += T_{cebk}^{(1)} \cdot T_{afb}^{(2)} \quad \leftarrow \text{“Contraction” over } T^{(1)} \text{ and } T^{(2)}$$

for c, e, a, f do

$$E += X_{aecf} \cdot Y_{ceaf}$$

Time-space trade-offs

for a, e, c, f do ←

for i, j do

$$X_{aecf} += T_{ijae} \cdot T_{ijcf}$$

for c, e, b, k do

$$T_{cebk}^{(1)} \leftarrow f_1(c, e, b, k)$$

for a, f, b, k do

$$T_{afb k}^{(2)} \leftarrow f_2(a, f, b, k)$$

for c, e, a, f do ←

for b, k do

$$Y_{ceaf} += T_{cebk}^{(1)} \cdot T_{afb k}^{(2)}$$

for c, e, a, f do ←

$$E += X_{aecf} \cdot Y_{ceaf}$$

Max index of
 $a-f$: $O(1000)$
 $i-k$: $O(100)$

Same indices
⇒ Loop fusion candidates

Time-space trade-offs

for a, e, c, f do	\longrightarrow	for a, e, c, f do	
for i, j do		for i, j do	
$X_{aecf} += T_{ijae} \cdot T_{ijcf}$		$X_{aecf} += T_{ijae} \cdot T_{ijcf}$	
for c, e, b, k do		for a, c, e, f, b, k do	
$T_{cebk}^{(1)} \leftarrow f_1(c, e, b, k)$		$T_{cebk}^{(1)} \leftarrow f_1(c, e, b, k)$	Add extra flops
for a, f, b, k do		for a, e, c, f, b, k do	
$T_{afbk}^{(2)} \leftarrow f_2(a, f, b, k)$		$T_{afbk}^{(2)} \leftarrow f_2(a, f, b, k)$	
for c, e, a, f do	\longrightarrow	for c, e, a, f do	
for b, k do		for b, k do	
$Y_{ceaf} += T_{cebk}^{(1)} \cdot T_{afbk}^{(2)}$		$Y_{ceaf} += T_{cebk}^{(1)} \cdot T_{afbk}^{(2)}$	
for c, e, a, f do	\longrightarrow	for c, e, a, f do	
$E += X_{aecf} \cdot Y_{ceaf}$		$E += X_{aecf} \cdot Y_{ceaf}$	



Time-space trade-offs

for a, e, c, f do \longrightarrow

for i, j do

$$X_{aecf} += T_{ijae} \cdot T_{ijcf}$$

for c, e, b, k do

$$T_{cebk}^{(1)} \leftarrow f_1(c, e, b, k)$$

for a, f, b, k do

$$T_{afbk}^{(2)} \leftarrow f_2(a, f, b, k)$$

for c, e, a, f do \longrightarrow

for b, k do

$$Y_{ceaef} += T_{cebk}^{(1)} \cdot T_{afbk}^{(2)}$$

for c, e, a, f do \longrightarrow

$$E += X_{aecf} \cdot Y_{ceaef}$$

for a, e, c, f do \Leftarrow Fused

for i, j do

$$x += T_{ijae} \cdot T_{ijcf}$$

for b, k do

$$T_{cebk}^{(1)} \leftarrow f_1(c, e, b, k)$$

$$T_{afbk}^{(2)} \leftarrow f_2(a, f, b, k)$$

$$y += T_{cebk}^{(1)} \cdot T_{afbk}^{(2)}$$

$$E += x \cdot y$$



Tiled & partially fused

for a, e, c, f **do** \longrightarrow

for i, j **do**

$$X_{aecf} += T_{ijae} \cdot T_{ijcf}$$

for c, e, b, k **do**

$$T_{cebk}^{(1)} \leftarrow f_1(c, e, b, k)$$

for a, f, b, k **do**

$$T_{afbk}^{(2)} \leftarrow f_2(a, f, b, k)$$

for c, e, a, f **do** \longrightarrow

for b, k **do**

$$Y_{ceaf} += T_{cebk}^{(1)} \cdot T_{afbk}^{(2)}$$

for c, e, a, f **do** \longrightarrow

$$E += X_{aecf} \cdot Y_{ceaf}$$

for a^B, e^B, c^B, f^B **do**

for a, e, c, f **do**

for i, j **do**

$$\hat{X}_{aecf} += T_{ijae} \cdot T_{ijcf}$$

for b, k **do**

for c, e **do**

$$\hat{T}_{ce}^{(1)} \leftarrow f_1(c, e, b, k)$$

for a, f **do**

$$\hat{T}_{af}^{(2)} \leftarrow f_2(a, f, b, k)$$

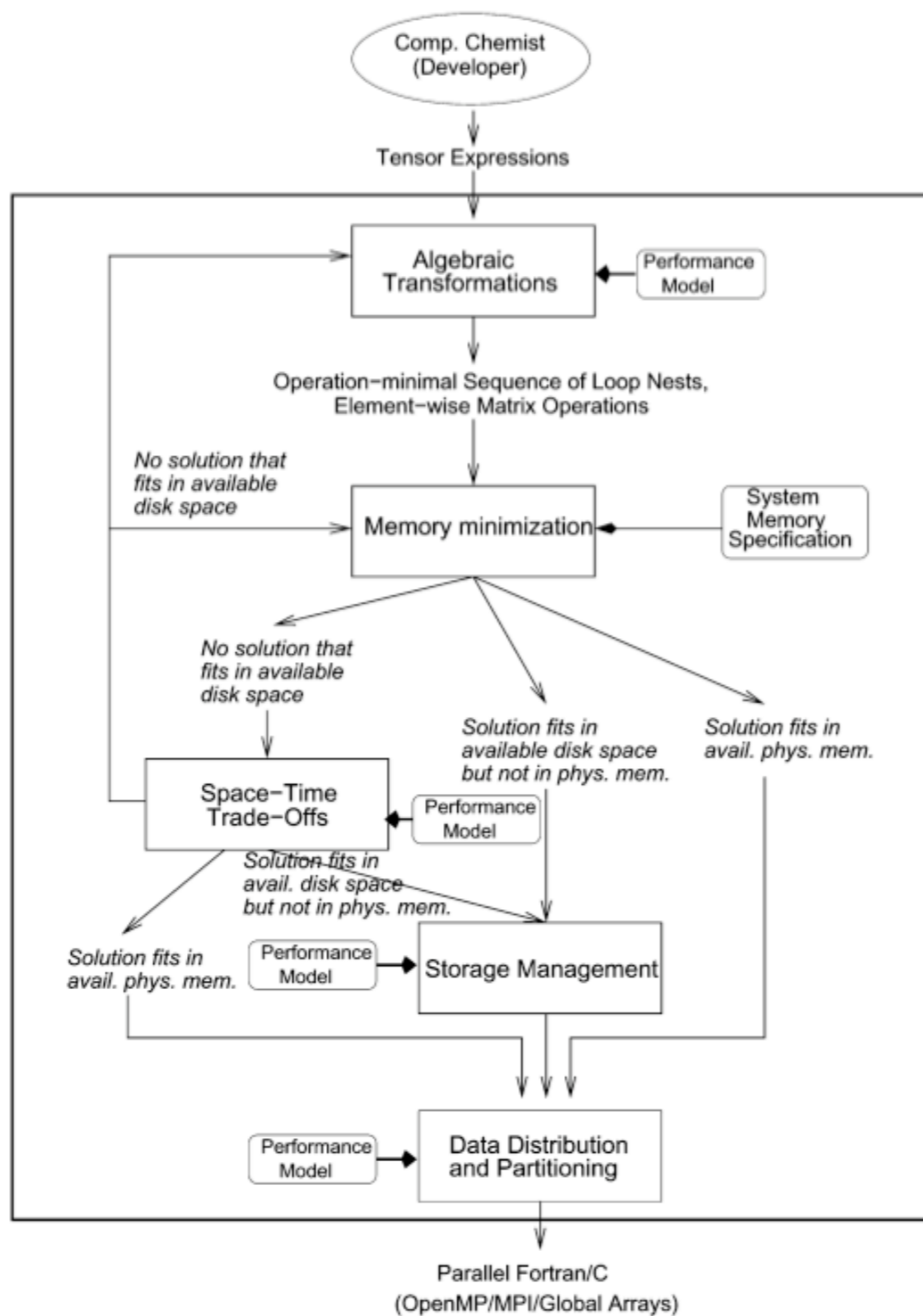
for c, e, a, f **do**

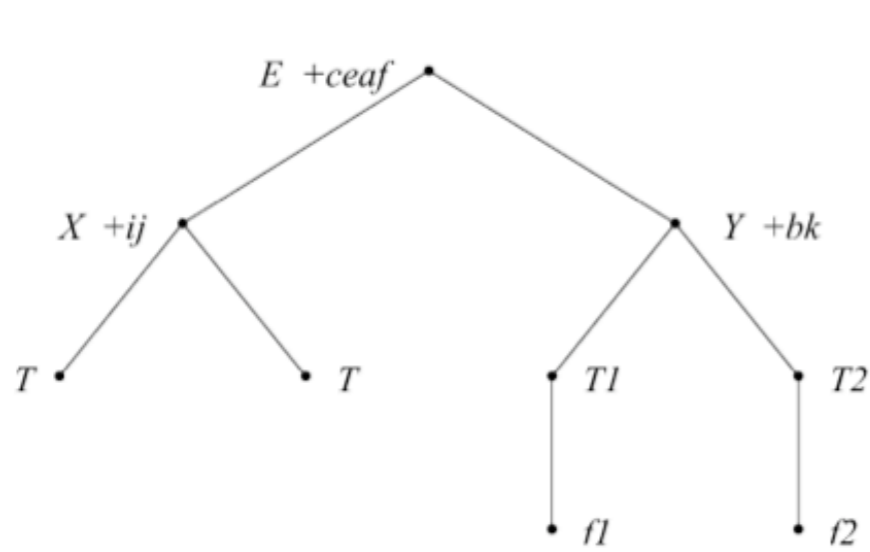
$$\hat{Y}_{ceaf} += \hat{T}_{ce}^{(1)} \cdot \hat{T}_{af}^{(2)}$$

for c, e, a, f **do**

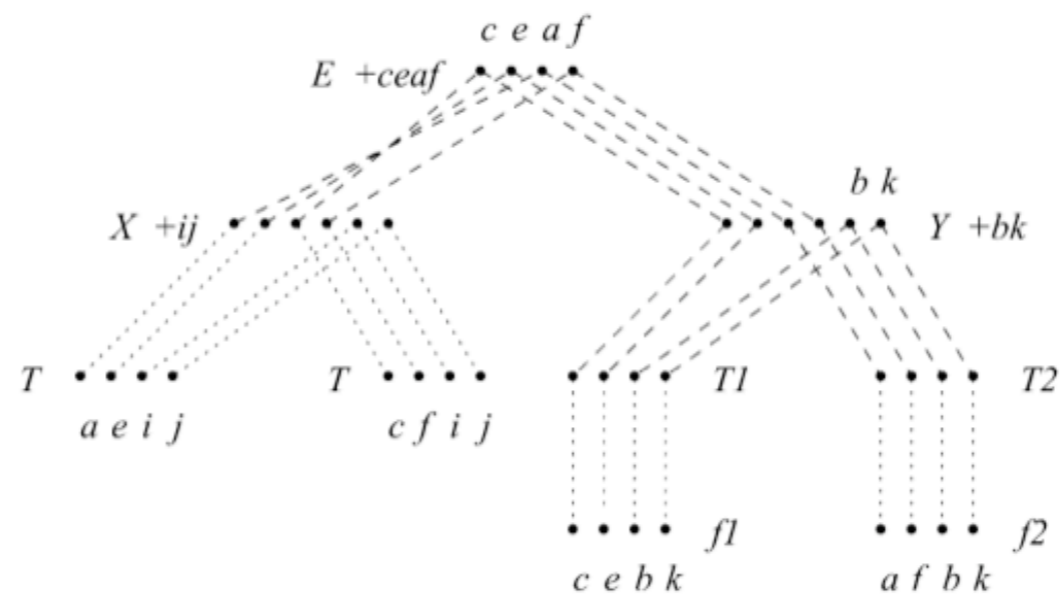
$$E += \hat{X}_{aecf} \cdot \hat{Y}_{ceaf}$$



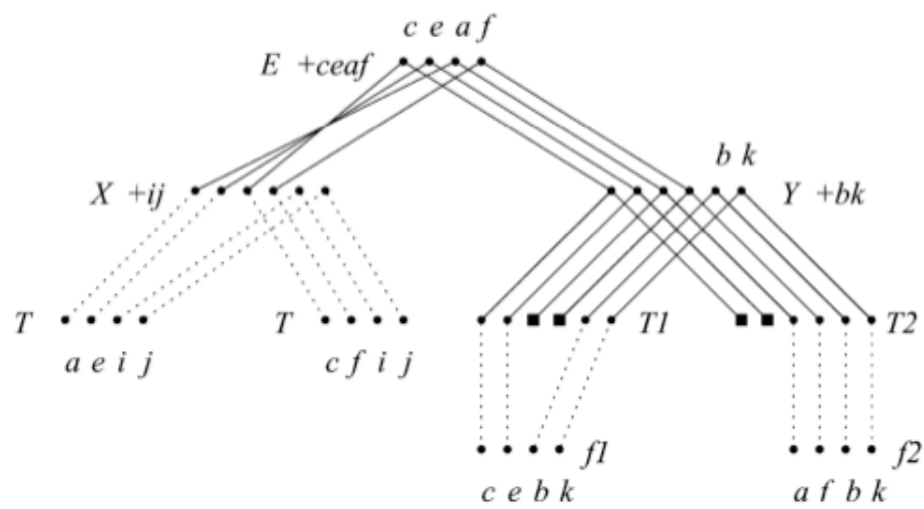




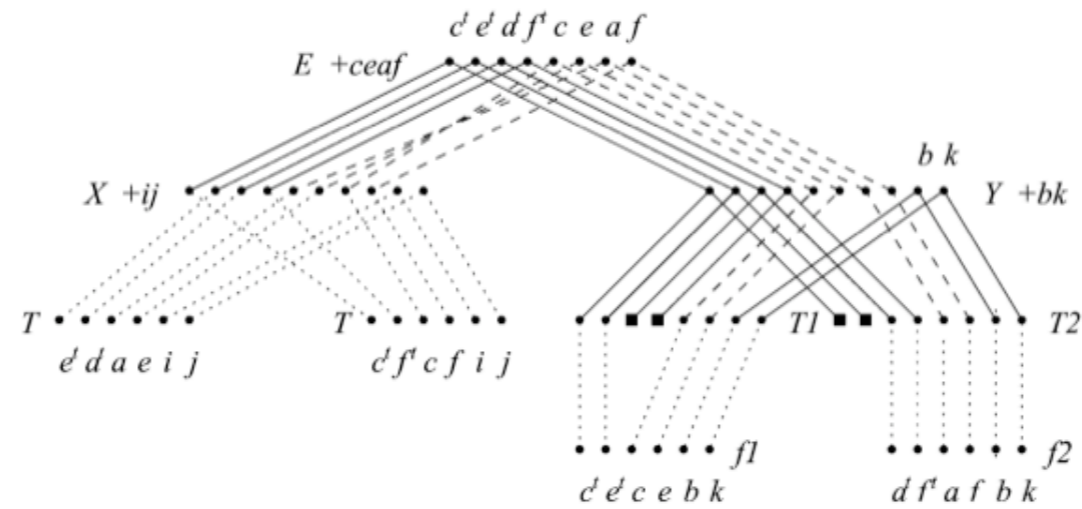
(a) Expression tree.



(b) Fusion graph.



(a) Fully fused computation from Fig. 4.



(b) Partially fused computation from Fig. 5.



Next time:
Empirical compilers and tools



“In conclusion...”



Backup slides