



PARALLEL COMPUTING LABORATORY

Auto-tuning Stencil Codes for Cache-Based Multicore Platforms

Kaushik Datta Dissertation Talk December 4, 2009



Motivation



- Multicore revolution has produced wide variety of architectures
- Compilers alone fail to fully exploit multicore resources
- Hand-tuning has become infeasible
- ✤ We need a better solution!



Contributions



- We have created an automatic stencil tuner (auto-tuner) that achieves up to 5.4x speedups over naïvely threaded stencil code
- We have developed an "Optimized Stream" benchmark for determining a system's highest attainable memory bandwidth
- We have bound stencil performance using the Roofline Model and in-cache performance



Outline



- Stencil Code Overview
- Cache-based Architectures
- Auto-tuning Description
- Stencil Auto-tuning Results

Stencil Code Overview



- For a given point, a stencil is a fixed subset of nearest neighbors
- A stencil code updates every point in a regular grid by "applying a stencil"
- Used in iterative PDE solvers like Jacobi, Multigrid, and AMR
- Also used in areas like image processing and geometric modeling
- This talk will focus on three stencil kernels:
 - 3D 7-point stencil

Electrical Engineering and Computer Sciences

- 3D 27-point stencil
- 3D Helmholtz kernel





Adaptive Mesh Refinement (AMR)





- Al is rough indicator of whether kernel is memory or compute-bound
- Counting *only* compulsory misses:



- Stencil codes usually (but not always) bandwidth-bound
 - Long unit-stride memory accesses
 - Little reuse of each grid point
 - Few flops per grid point
- Actual AI values are typically lower (due to other types of cache misses)







- Stencil Code Overview
- Cache-based Architectures
- Auto-tuning Description
- Stencil Auto-tuning Results



SPAF SPAF SPAF SPAF

179 GB/s

21.33 GB/s

SP/ SP/

МΤ ΙĘ

Crossbar

4MB Shared L2 (16 way)

(64b interleaved)

4 Coherency Hubs

2x128b controllers

667MHz FBDIMMs

MT

90 GB/s

10.66 GB/s



Sun Niagara2

GB/s

x 6.4 (

ω Ξ

SPARC IT SPAR MT SPAR T SPA

M M M M

90 GB/s

10.66 GB/s

SPARC

Crossbar

4MB Shared L2 (16 way) (64b interleaved)

4 Coherency Hubs

2x128b controllers

667MHz FBDIMMs

SPARC

μ

179 GB/s

21.33 GB/s





Sun Niagara2

IBM Blue Gene/P

10









Intel Clovertown



Intel Nehalem



Sun Niagara2

AMD Barcelona

Socket/Core/ Thread Count









AMD Barcelona

Intel Clovertown



Intel Nehalem



Total HW

Thread Count









Intel Clovertown



Intel Nehalem



Stream Copy

AMD Barcelona

Bandwidth (GB/s)









Intel Clovertown



Intel Nehalem



AMD Barcelona

Peak DP Computation Rate (GFlop/s)







- Stencil Code Overview
- Cache-based Architectures
- Auto-tuning Description
- Stencil Auto-tuning Results

General Compiler Deficiencies

- Historically, compilers have had problems with domain-specific * transformations:
 - Register allocation (explicit temps)
 - Loop unrolling
 - Software pipelining
 - Tiling

Computer Sciences

- SIMDization
- Common subexpression elimination
- Data structure transformations
 Algorithmic transformations
- Compilers typically use heuristics (not actual runs) to determine the * best code for a platform
 - Difficult to generate optimal code across many diverse multicore architectures









- Auto-tuning became popular because:
 - Domain-specific transformations could be included
 - Runs experiments instead of heuristics
 - Diversity of systems (and now increasing core counts) made performance portability vital
- Auto-tuning is:
 - Portable (to an extent)
 - Scalable
 - Productive (if tuning for multiple architectures)
 - Applicable to many metrics of merit (e.g. performance, power efficiency)
- We let the machine search the parameter space intelligently to find a (near-)optimal configuration
- Serial processor success stories: FFTW, Spiral, Atlas, OSKI, others...







- Stencil Code Overview
- Cache-based Architectures
- Auto-tuning Description
 - Identify motif-specific optimizations
 - Generate code variants based on these optimizations
 - Traverse parameter space for best configuration
- Stencil Auto-tuning Results

Problem Decomposition (across an SMP)









Core Blocking

Electrical Engineering and Computer Sciences

• Allows for domain decomposition and cache blocking

Thread Blocking

• Exploit caches shared among threads within a core

Register Blocking

Loop unrolling in any of the three dimensions
Makes DLP/ILP explicit

- This decomposition is universal across all examined architectures
- Decomposition does not change data structure
- Need to choose best block sizes for each hierarchy level



Data Allocation





NUMA-Aware Allocation

• Ensures that the data is colocated on same socket as the threads processing it



Array Padding

- Alters the data placement so as to minimize conflict misses
- Tunable parameter



Bandwidth Optimizations





Software Prefetching

• Helps mask memory latency by adjusting look-ahead distance

• Can also tune number of software prefetch requests



- Eliminates cache line fills on a write miss
- Reduces memory traffic by 50% on write misses!
- Only available on x86 machines

In-Core Optimizations







- Single instruction processes multiple data items
- Non-portable code

Electrical Engineering and Computer Sciences

> c = a+b; d = a+b; e = c+d; c = a+b; d = a+b; c = a+b; d = a+b; d = a+b; d = a+b; d = a+b; c = a+b; d = a+b; d = a+b; e = c+d; c = a+b; c = a+b; d = a+b; e = c+d; c = a+b; c

- Reduces flops by removing redundant expressions
- icc and gcc often fail to do this







- Stencil Code Overview
- Cache-based Architectures
- Auto-tuning Description
 - Identify motif-specific optimizations
 - Generate code variants based on these optimizations
 - Traverse parameter space for best configuration
- Stencil Auto-tuning Results



- Hand-tuned code only performs well on a single platform
- Perl code generator can produce many different code variants for performance portability
- Intelligent code generator can take pseudo-code and specified set of transformations to produce code variants
 - Type of domain-specific compiler







- Stencil Code Overview
- Cache-based Architectures
- Auto-tuning Description
 - Identify motif-specific optimizations
 - Generate code variants based on these optimizations
 - Traverse parameter space for best configuration
- Stencil Auto-tuning Results

Electrical Engineering and Computer Sciences

- We introduced 9 different optimizations, each of which has its own set of parameters
- Exhaustive search is *impossible*
- To make problem tractable, we:
 - Used expert knowledge to order the optimizations
 - Applied them consecutively
- Every platform had its own set of best parameters



BERKELEY PAR LAB







- Stencil Code Overview
- Cache-based Architectures
- Auto-tuning Description
- Stencil Auto-tuning Results
 - 3D 7-Point Stencil (Memory-Intensive Kernel)
 - 3D 27-Point Stencil (Compute-Intensive Kernel)
 - 3D Helmholtz Kernel

Computer Sciences **BERKELEY PAR LAB** The 3D 7-point stencil performs: * 8 flops per point 16 or 24 Bytes of memory traffic per point Al is either 0.33 or 0.5 (w/ cache bypass) * This kernel should be memory-bound on most * architectures: Helmholtz kernel Memory Computation Ideal Arithmetic Intensity Bound Bound ' 2 7-point stencil 27-point stencil

3D 7-Point Stencil Problem

Electrical Engineering and

We will perform a single out-of-place sweep of this stencil over a 256³ grid



Naïve Stencil Code



- We wish to exploit multicore resources
- First attempt at writing parallel stencil code:
 - Use pthreads
 - Parallelize in least contiguous grid dimension
 - Thread affinity for scaling: multithreading, then multicore, then multisocket





Naïve Performance (3D 7-Point Stencil)

8

16

Naive







Perf. Limit (blue=comp., red=bandwidth)







Electrical Engineering and Computer Sciences

IBM Blue Gene/P













Intel Clovertown



IBM Blue Gene/P



Intel Nehalem



Sun Niagara2



AMD Barcelona

Parallel Scaling Speedup Over Single Core Performance

How much improvement is there? (3D 7-Point Stencil)





Electrical Engineering and Computer Sciences

0.4





IBM Blue Gene/P





Sun Niagara2



Tuning Speedup Over Best Naïve Performance



Electrical Engineering and Computer Sciences











- Stencil Code Overview
- Cache-based Architectures
- Auto-tuning Description
- Stencil Auto-tuning Results
 - 3D 7-Point Stencil (Memory-Intensive Kernel)
 - 3D 27-Point Stencil (Compute-Intensive Kernel)
 - 3D Helmholtz Kernel

The 3D 27-point stencil performs: * 30 flops per point 16 or 24 Bytes of memory traffic per point Al is either 1.25 or 1.88 (w/ cache bypass) ** CSE can reduce the flops/point • This kernel should be compute-bound on most * architectures: Helmholtz kernel **Memory** Computation Ideal Arithmetic Intensity Bound Bound 2 7-point stencil **27-point stencil**

3D 27-Point Stencil Problem

Electrical Engineering and Computer Sciences

We will perform a single out-of-place sweep of this stencil over a 256³ grid

BERKELEY PAR LAB



Naïve Performance (3D 27-Point Stencil)







Perf. Limit (blue=comp., red=bandwidth)

Naive

Auto-tuned Performance (3D 27-Point Stencil)





Electrical Engineering and Computer Sciences







0.00



2 Cores

IBM Blue Gene/P

4







Parallel Scaling Speedup Over Single Core Performance

39

How much improvement is there? (3D 27-Point Stencil)





IBM Blue Gene/P

Electrical Engineering and

Computer Sciences

1.4 1.2 1.0 0.8 0.6 0.4 0.2 0.0 1 2 4 8 Fully Populated Cores Intel Nehalem



Sun Niagara2



Tuning Speedup Over Best Naïve Performance



How well can we do? (3D 27-Point Stencil)











- Stencil Code Overview
- Cache-based Architectures
- Auto-tuning Description
- Stencil Auto-tuning Results
 - 3D 7-Point Stencil (Memory-Intensive Kernel)
 - 3D 27-Point Stencil (Compute-Intensive Kernel)
 - 3D Helmholtz Kernel



- The 3D Helmholtz kernel is very different from the previous kernels:
 - Gauss-Seidel Red-Black ordering
 - 25 flops per stencil

Electrical Engineering and Computer Sciences

- 7 arrays (6 are read only, 1 is read and write)
- Many small subproblems- no longer one large problem
- Ideal AI is about 0.20
- This kernel should be memory-bound on most architectures:







- 1-2 threads per problem is optimal in cases where load balancing is not an issue
- If this trend continues, load balancing will be an even larger issue in the manycore era



- This is performance of 16³ subproblems in a 0.5 GB memory footprint
- Performance gets worse with more threads per subproblem



Conclusions



- Compilers alone achieves poor performance
 - Typically achieve a low fraction of peak performance
 - Exhibit little parallel scaling
- Autotuning is essential to achieving good performance
 - 1.9x-5.4x speedups across diverse architectures
 - Automatic tuning is *necessary* for scalability
 - With few exceptions, the same code was used
- Ultimately, we are limited by the hardware
 - We can only do as well as Stream or in-core performance
 - The memory wall will continue to push stencil codes to be bandwidthbound
- When dealing with many small subproblems, fewer threads per subproblem performs best
 - However, load balancing becomes a major issue
 - This is an even larger problem for the manycore era



Future Work



Better Productivity:

- Current Perl scripts are primitive
- Need to develop an auto-tuning framework that has semantic knowledge of the stencil code (S. Kamil)
- Better Performance:
 - We currently do no data structure changes other than array padding
 - May be beneficial to store the grids in a recursive format using spacefilling curves for better locality (S. Williams?)

Better Search:

- Our current search method does require expert knowledge to order the optimizations appropriately
- Machine learning offers the opportunity for tuning with little domain knowledge and many more parameters (A. Ganapathi)





- Kathy and Jim for sure-handed guidance and knowledge during all these years
- Sam Williams for always being available to discuss research (and being an unofficial thesis reader)
- Rajesh Nishtala for being a great friend and officemate
- Jon Wilkening for being my outside thesis reader
- The Bebop group, including Shoaib Kamil, Karl Fuerlinger, and Mark Hoemmen
- The scientists at LBL, including Lenny Oliker, John Shalf, Jonathan Carter, Terry Ligocki, and Brain Van Straalen
- The members of the Parlab and Radlab, including Dave Patterson and Archana Ganapathi
- Many others that I don't have space to mention here...



Supplemental Slides







- Lawrence Berkeley Laboratory (LBL) is using stencil auto-tuning as a building block of its Green Flash supercomputer (Google: Green Flash LBL)
- Dr. Franz-Josef Pfreundt (head of IT at Fraunhofer-ITWM) used stencil tuning to improve the performance of oil exploration code





