A Comparison of Cache-conscious and Cache-oblivious Programs

Keshav Pingali, University of Texas, Austin

Joint work with Kamen Yotov, Goldman Sachs Tom Roeder, Cornell University John Gunnels, IBM T.J. Watson Research Center Fred Gustavson, IBM T.J.Watson Research Center

Memory Hierarchy Management

· Cache-conscious (CC) approach:

- Blocked iterative algorithms and arrays (usually)
 Code and data structures have parameters that depend on careful blocking for memory hierarchy
- Used in dense linear algebra libraries: BLAS, LAPACK
 Lots of algorithmic data reuse: O(N³) operations on O(N²) data

Cache-oblivious (CO) approach:

- Recursive algorithms and data structures (usually)
 - Not aware of memory hierarchy: approximate blocking
 - I/O optimal: Hong and Kung, Frigo and Leiserson
- Used in FFT implementations: FFTW
 - Little algorithmic data reuse: O(N(logN)) computations on O(N) data





Organization of talk

- CO and CC approaches to blocking
 control structures
 - data structures
- Non-standard view of blocking (or why CO may work well) – reduce bandwidth required from memory
- Experimental results
 - UltraSPARC IIIi
 - Itanium
 - Xeon
 - Power 5
- · Lessons and ongoing work



CO: recursive micro-kernel Internal nodes of recursion tree are recursive overhead; roughly - 100 cycles on Itanium-2 - 360 cycles on UltraSPARC IIIi \bigcirc Large overhead: for LD, roughly one internal node per leaf node Solution: Micro-kernel: code obtained by complete unrolling of recursive tree for some fixed size problem (RUxRUxRU) - Cut off recursion when sub-problem size becomes equal to micro-kernel size, and invoke micro-kernel - Overhead of internal node is recursive micro-kernel amortized over micro-kernel, rather than a single multiply-add - Choice of RU: empirical





CC algorithms: discussion

- · Iterative codes
 - Nested loops
- Implementation of blocking
- Cache blocking
 - Mini-kernel: in ATLAS, multiply NBxNB blocks
 - Choose NB so NB² + NB + 1 <= C_{L1}
- Register blocking
 - Micro-kernel: in ATLAS, multiply MUx1 block of A with 1xNU block of B into MUxNU block of C
 - Choose MU,NU so that MU + NU +MU*NU <= NR

Organization of talk

- CO and CC approaches to blocking – control structures
 - data structures
- Non-standard view of blocking – reduce bandwidth required from memory
- Experimental results
 - UltraSPARC IIIi
 - Itanium
 - Xeon
 - Power 5
- · Lessons and ongoing work

Blocking

- Microscopic view
 - Blocking reduces expected latency of memory access
- Macroscopic view
 - Memory hierarchy can be ignored if
 - memory has enough bandwidth to feed processor
 - data can be pre-fetched to hide memory latency
 - Blocking reduces bandwidth needed from memory
- Useful to consider macroscopic view in more detail



sqrt(capacity(L)/3) ≥ NB_L ≥ 4/Bandwidth(L,L+1) (levels L,L+1)



Lessons

· Reducing bandwidth requirements

- Block size does not have to be exact
- Enough for block size to lie within an interval that depends on hardware parameters
- If upper bound on NB is more than twice lower bound, divide and conquer will automatically generate a block size in this range
- → approximate blocking CO-style is OK

Reducing latency

- Accurate block sizes are better
- If block size is chosen approximately, may need to compensate with prefetching

Organization of talk

- Non-standard view of blocking

 reduce bandwidth required from memory
- CO and CC approaches to blocking
 - control structures
 - data structures
- Experimental results
 - UltraSPARC IIIi
 - Itanium
 - Xeon
 - Power 5
- · Lessons and ongoing work

UltraSPARC IIIi

• Peak performance: 2 GFlops (1 GHZ, 2 FPUs)

• Memory hierarchy:

- Registers: 32
- L1 data cache: 64KB, 4-way
- L2 data cache: 1MB, 4-way
- Compilers
 - C: SUN C 5.5







Lessons

- Bottom-line on UltraSPARC:
 - Peak: 2 GFlops
 - ATLAS: 1.75 GFlops
 - Best CO strategy: 700 MFlops
- Similar results on other machines:
 - Best CO performance on Itanium: roughly 2/3 of peak
- · Conclusion:
 - Recursive micro-kernels are not a good idea





Lessons

- Two hardware constraints on size of micro-kernels:
 I-cache limits amount of unrolling
 - Number of registers
- Iterative micro-kernel: three degrees of freedom (MU,NU,KU)
 - Choose MU and NU to optimize register usage
 - Choose KU unrolling to fit into I-cache
- Recursive micro-kernel: one degree of freedom (RU)
- But even if you choose rectangular tiles, all three degrees of freedom are tied to both hardware constraints









- Vendor BLAS gets highest performance
- Pre-fetching boosts performance by roughly 40%
- Iterative code: pre-fetching is well-understood
- Recursive code: not well-understood

Summary

- Iterative approach has been proven to work well in practice
 Vendor BLAS, ATLAS, etc.
 - But requires a lot of work to produce code and tune parameters
 - Implementing a high-performance CO code is not easy – Careful attention to micro-kernel and mini-kernel is needed
- Using fully recursive approach with highly optimized recursive micro-kernel, we never got more than 2/3 of peak.
- Issues with CO approach

•

- Recursive Micro-Kernels yield less performance than iterative ones using same scheduling techniques
- Pre-fetching is needed to compete with best code: not well-understood in the context of CO codes

Ongoing Work

- · Explain performance of all results shown
- Complete ongoing Matrix Transpose study
- · Proteus system and BRILA compiler
- I/O optimality:
 - Interesting theoretical results for simple model of computation
 - What additional aspects of hardware/program need to be modeled for it to be useful in practice?

