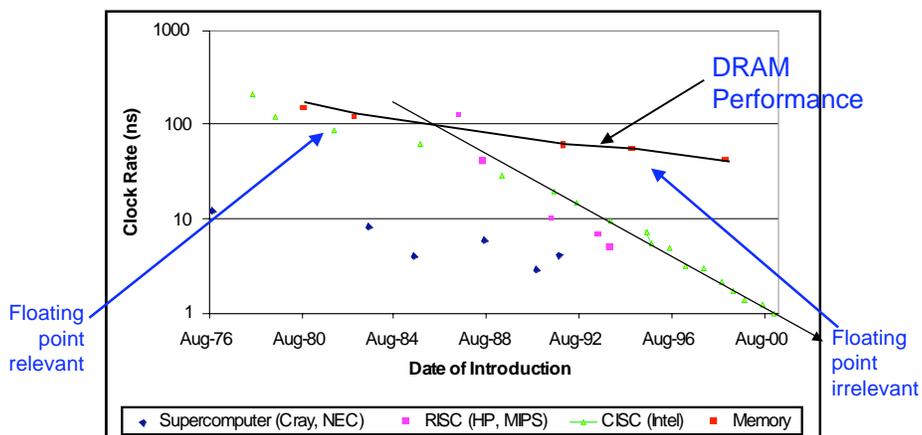


Overcoming the Barriers to Sustained Petaflop Performance

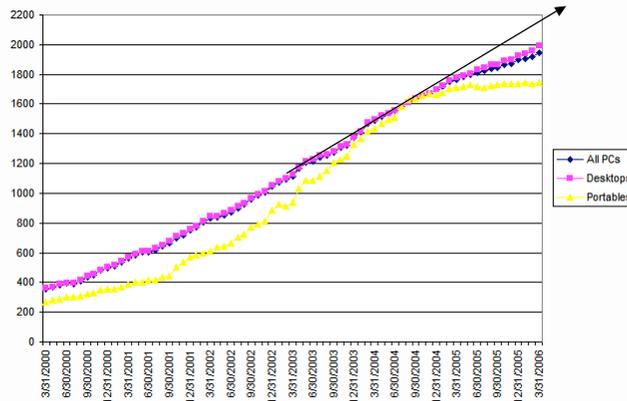
William D. Gropp
Computer Science
www.cs.uiuc.edu/homes/wgropp



Why is achieved performance on a single node so poor?



CPU performance is not doubling any more

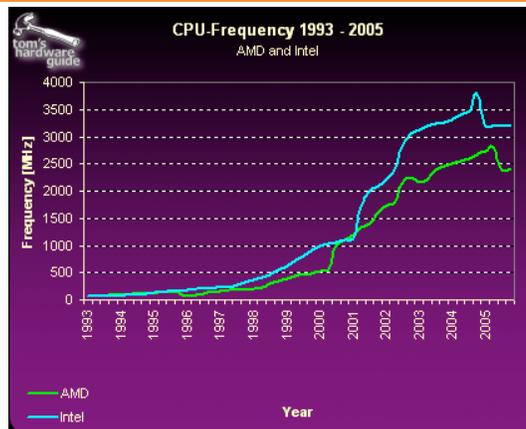


- Average CPU clock speeds (from <http://www.pcpitstop.com/research/cpu.asp>)



3

Peak CPU speeds are stable



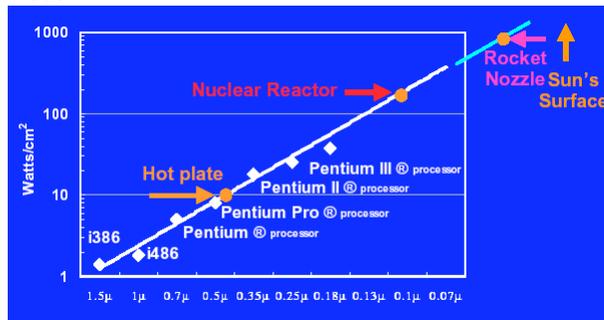
- From http://www.tomshardware.com/2005/11/21/the_mother_of_all_cpu_charts_2005/



4

Why are CPUs not getting faster?

- Power dissipation problems will force more changes
 - ◆ Current trends imply chips with energy densities greater than a nuclear reactor
 - ◆ Already a problem: Recalls of recent Mac laptops because they could overheat.
 - ◆ Will force new ways to get performance, such as extensive parallelism
 - ◆ Note highly parallel nodes already common - routers (Cisco has 188 cores) GPUs, ...



5

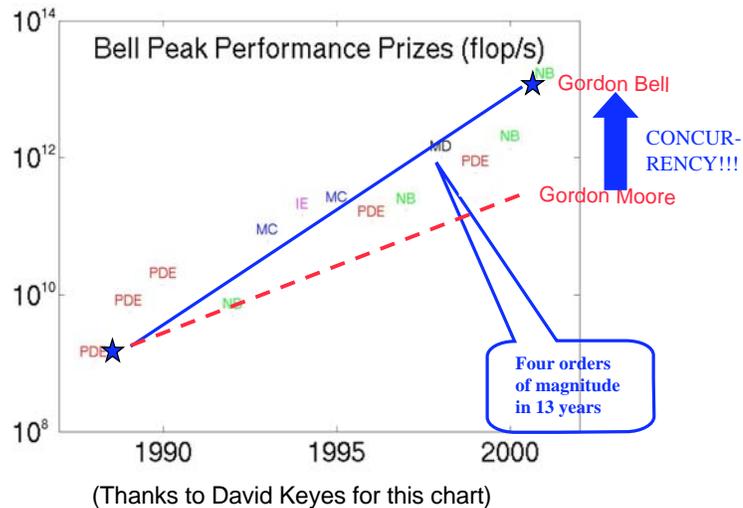
Where will we get (Sustained) Performance?

- Algorithms that are a better match for the architectures
- Parallelism at all levels
- Concurrency at all levels
- A major challenge is to realize these approaches in code
- Parallelism is not a new trend ...



6

Gordon Bell Prize outpaces Moore's Law



7

Understanding the Reasons for Low Node Performance

- Current laptops now have a peak speed (based on clock rate) of over 2 Gflops (20 Cray1s!)
- Observed (sustained) performance is often a small fraction of peak
- Why is the gap between “peak” and “sustained” performance so large?
- Lets look at a simple numerical kernel



8

Sparse Matrix-Vector Product

- Common operation for optimal (in floating-point operations) solution of linear systems
- Sample code (in C):

```
for row=1,n
  m = i[row] - i[row-1];
  sum = 0;
  for k=1,m
    sum += *a++ * x[*j++];
  y[i] = sum;
```
- Data structures are a[nnz], j[nnz], i[n], x[n], y[n]



9

Simple Performance Analysis

- Memory motion:
 - ◆ nnz (sizeof(double) + sizeof(int)) + n (2*sizeof(double) + sizeof(int))
 - ◆ Assume a perfect cache (never load same data twice; only compulsory loads)
- Computation
 - ◆ nnz multiply-add (MA)
- Roughly 12 bytes per MA
- Typical WS node can move 1-4 bytes/MA
 - ◆ Maximum performance is 8-33% of peak



10

More Performance Analysis

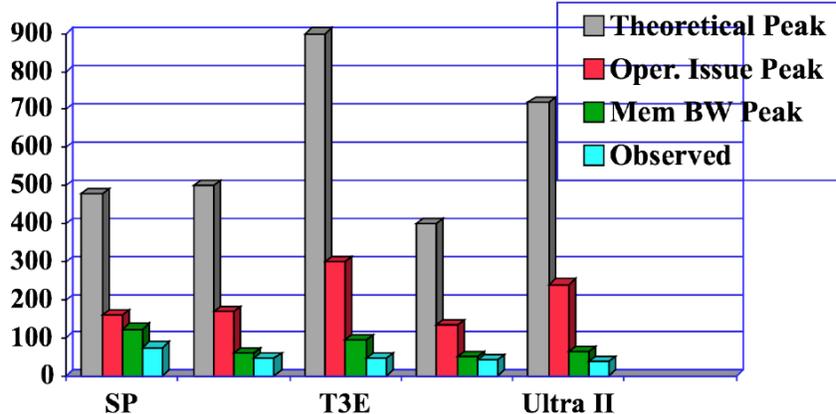
- Instruction Counts:
 - ♦ $\text{nz} (2 \times \text{load-double} + \text{load-int} + \text{mult-add}) + n (\text{load-int} + \text{store-double})$
- Roughly 4 instructions per MA
- Maximum performance is 25% of peak (33% if MA overlaps one load/store)
 - ♦ (wide instruction words can help here)
- Changing matrix data structure (e.g., exploit small block structure) allows reuse of data in register, eliminating some loads (x and j)
- Implementation improvements (tricks) cannot improve on these limits



11

Realistic Measures of Peak Performance

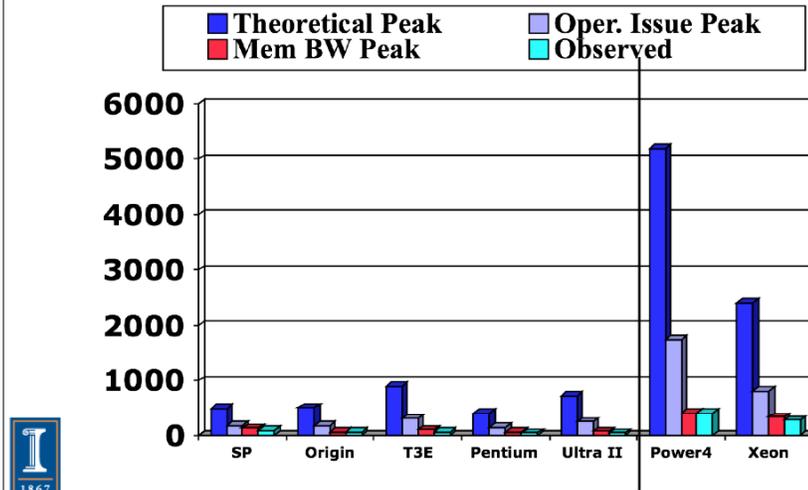
Sparse Matrix Vector Product
one vector, matrix size, $m = 90,708$, nonzero entries $\text{nz} = 5,047,120$



12

Realistic Measures of Peak Performance

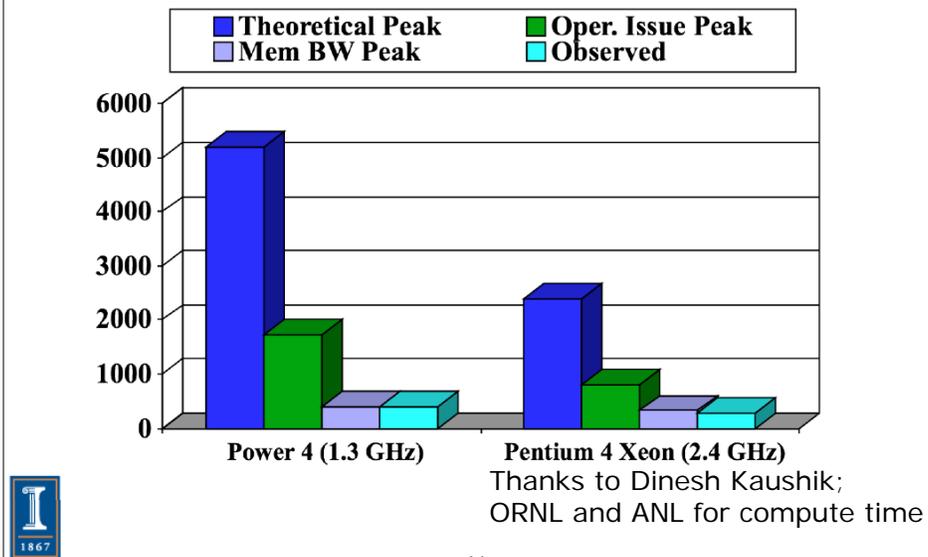
Sparse Matrix Vector Product
 one vector, matrix size, $m = 90,708$, nonzero entries $nz = 5,047,120$



13

Realistic Measures of Peak Performance

Sparse Matrix Vector Product
 One vector, matrix size, $m = 90,708$, nonzero entries $nz = 5,047,120$



14

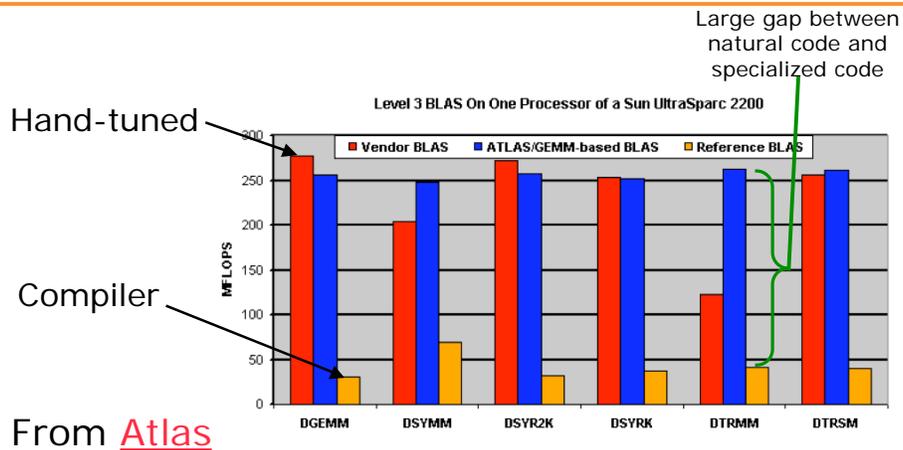
What About CPU-Bound Operations?

- Dense Matrix-Matrix Product
 - ◆ Probably the numerical program most studied by compiler writers
 - ◆ Core of some important applications
 - ◆ More importantly, the core operation in High Performance Linpack (HPL)
 - ◆ Should give optimal performance...



15

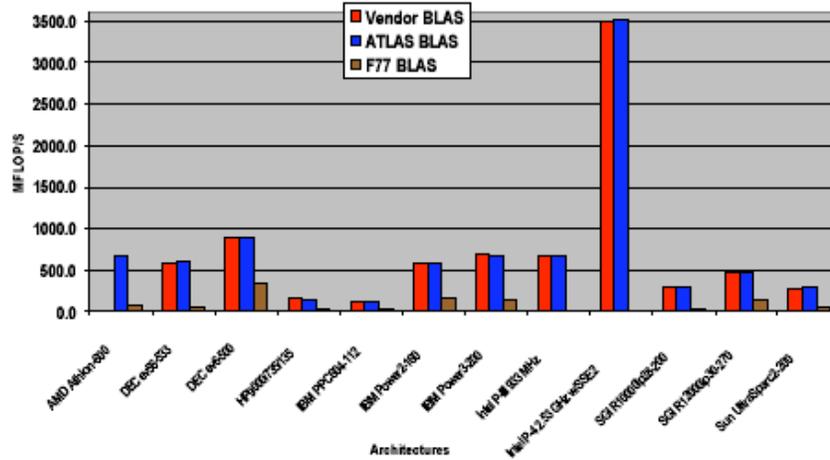
How Successful are Compilers with CPU Intensive Code?



Enormous effort required to get good performance

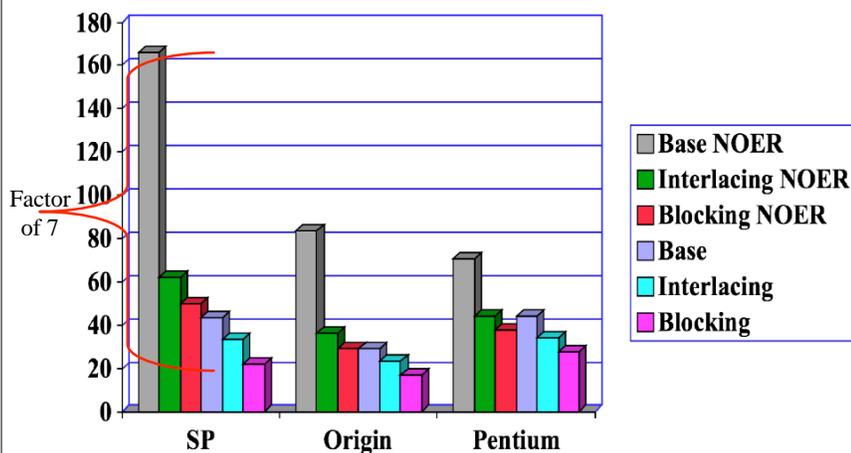
16

Very Few Compilers do well on DGEMM (n=500)



17

Effect of code transformations for uni-processor performance



18

Consequences of Memory/CPU Performance Gap

- Performance of an application may be (and often is) limited by sustained memory bandwidth (not peak memory bus speed) or latency rather than CPU clock
- “Peak” performance determined by the resource that is operating at full speed for the algorithm
 - ◆ Often memory system (e.g., see STREAM results)
 - ◆ Sometimes instruction rate/mix (including integer ops), load balance, internode communication performance, ...
- For example, sparse matrix-vector operation performance is best estimated by using STREAM performance
 - ◆ Note that STREAM performance is delivered performance to a Fortran or C program, not memory bus rate time width
 - ◆ High latency of memory and low number of outstanding loads can significantly reduce sustained memory bandwidth



19

Performance for Real Applications

- Dense matrix-matrix example shows that even for well-studied, compute-bound kernels, compiler-generated code achieves only a small fraction of available performance
 - ◆ “Fortran” code uses “natural” loops, i.e., what a user would write for most code
 - ◆ Others use multi-level blocking, careful instruction scheduling etc.
- Algorithms design also needs to take into account the capabilities of the system, not just the hardware
 - ◆ Example: Cache-Oblivious Algorithms (<http://supertech.lcs.mit.edu/cilk/papers/abstracts/abstract4.html>)
- Adding concurrency (whether multicore or multiple processors) just adds to the problems...



20

Possible solutions

- Single, integrated system
 - ◆ Best choice when it works
 - E.g., Matlab, R
- Current Terascale systems and many proposed petascale systems exploit hierarchy
 - ◆ Successful at many levels
 - Cluster hardware
 - OS scalability
 - ◆ We should apply this to productivity software
 - The problem is hard
 - Apply classic and very successful Computer Science strategies to address the complexity of generating fast code for a wide range of user-defined data structures.
- How can we apply hierarchies?
 - ◆ One approach is to use libraries
 - Limited by the operations envisioned by the library designer
 - ◆ Another is to enhance the users ability to express the problem in source code



25

Annotations

- Aid in the introduction of hierarchy into the software
 - ◆ Its going to happen anyway, so make a virtue of it
- Create a “market” or ecosystem in transformation tools
- Longer term issues
 - ◆ Integrate annotation language into “host” language to ensure type safety, ensure consistency (both syntactic and semantic), closer debugger integration, additional optimization opportunities through information sharing, ...



26

Observations

- Much use of mechanical transformations of code to achieve better performance
 - ◆ Compilers do not do this well
 - Too many other demands on the compiler
- Use of carefully crafted algorithms for specific operations such as allreduce, matrix-matrix multiply
 - ◆ Far more challenging than the performance transformations
- Increasing acceptance of some degree of automation in creating code
 - ◆ ATLAS, PhiPAC, TCE
 - ◆ Source-to-source transformation systems
 - E.g., ROSE, Aspect Oriented Programming (asod.net)



27

Getting Performance Out of Source Code

- Let the compiler do it
 - ◆ Too difficult (the compiler has too much to worry about)
- Improve the language so that the compiler can do it
- Build better tools to work with the language and the compiler



28

Potential challenges faced by languages

1. Time to develop the language.
 2. Divergence from mainstream compiler and language development.
 3. Mismatch with application needs.
 4. Performance.
 5. Performance portability.
 6. Concern of application developers about the success of the language.
- Understanding these provides insights into potential solutions
 - Annotations can *complement* language research by using the principle of *separation of concerns*
 - The annotation approach can be applied to *new* languages, as well



29

Key Observations

- 90/10 rule
 - ◆ current languages adequate for 90% of code
 - ◆ 10% of code causes 90% of trouble
- Memory hierarchy issues a major source of problems
 - ◆ Significant effort is put into relatively mechanical transformations of code
 - ◆ Other transformations are avoided because of their negative impact on the readability and maintainability of the code.
 - Example is loop fusion for routines that sweep over a mesh to apply different physics. Fusion, needed to reduce memory bandwidth requirements, breaks modularity of routines written by different groups.
- Coordination of distributed data structures another major source of problems
 - ◆ But the need for performance encourages a global/local separation
 - Reflected in PGAS languages
- New languages may help, but not anytime soon
 - ◆ New languages will never be the entire solution
 - ◆ Applications need help now



30

One Possible Approach

- Use annotations to augment existing languages
 - ◆ Not a new approach; used in HPF, OpenMP, others
 - ◆ Some applications already use this approach for performance portability
 - WRF weather code
- Annotations do have limitations
 - ◆ Fits best when most of the code is independent of the parts affected by the annotations
 - ◆ Limits optimizations that are available to approaches that augment the language (e.g., telescoping languages)
- But they also have many advantages...



31

Annotations example: STREAM triad.c for BG/L

```
void triad(double *a, double *b, double *c, int n)
{
  int i;
  double ss = 1.2;
  /* --Align::var:a,b,c;; */
  for (i=0; i<n; i++)
    a[i] = b[i] + ss*c[i];
  /* --end Align */
}
```

```
void triad(double *a, double *b, double *c, int n)
{
  #pragma disjoint (*c,*a,*b)
  int i;
  double ss = 1.2;
  /* --Align::var:a,b,c;; */
  if ( ((int)(a) | (int)(b) | (int)(c) & 0xf == 0) {
    __alignx(16,a);
    __alignx(16,b);
    __alignx(16,c);
    for (i=0;i<n;i++) {
      a[i] = b[i] + ss*c[i];
    }
  }
  else {
    for (i=0;i<n;i++) {
      a[i]=b[i] + ss*c[i];
    }
  }
  /* --end Align */
}
```



32

Simple annotation example: STREAM triad.c on BG/L

Size	No Annotations (MB/s)	Annotations (MB/s)	
10	1920.00	2424.24	
100	3037.97	6299.21	
1000	3341.22	8275.86	2.5X
10000	1290.81	3717.88	
50000	1291.52	3725.48	
100000	1291.77	3727.21	2.9X
500000	1291.81	1830.89	
1000000	1282.12	1442.17	
2000000	1282.92	1415.52	
5000000	1290.81	1446.48	1.12X



33

Advantages of annotations

- These parallel the challenges for languages
 1. Speeds development and deployment by using source-to-source transformations.
 - ◆ Higher-quality systems can preserve the readability of the source code, avoiding one of the classic drawbacks of preprocessor and source-to-source systems.
 2. Leverages mainstream language developments by building on top of those languages, not replacing them.
 3. Provides a simpler method to match application needs by allowing experts to develop abstractions tuned to the needs of a class (or even a single important) application.
 - ◆ Also enables the evaluation of new features and data structures



34

Advantages of annotations (con't)

4. Provides an effective approach for addressing performance issues by permitting (but not requiring) access by the programmer to low-level details.
 - ◆ Abstractions that allow domain or algorithm-specific approaches to performance can be used because they can be tuned to smaller user communities than is possible in a full language.
5. Improves performance portability by abstracting platform-specific low-level optimization code.
 - ◆ Including heterogeneous processor system (e.g., FPGAs, mixed scalar/vector)
6. Preserves application investment in current languages.
 - ◆ Allows use of existing development tools (debuggers) and allows maintenance and development of code independent of the tools the process the annotations.



35

Is This Ugly?

- You bet!
 - ◆ But it starts the process of considering the code generation process as consisting of a hierarchy of solutions
 - ◆ Separates the integration of the tools as seen by the user from the integration as seen by "the code"
- It can evolve toward a cleaner approach, with well-defined interfaces between hierarchies, and with a compilation-based approach to provide better syntax and semantic analysis
- But only if we accept the need for a hierarchical, compositional approach.
- This complements rather than replaces advances in languages, such as global view approaches



36

What About Language and Runtime Efforts for Multicore?

- Some, like the Intel thread building blocks, address only one aspect of the programming problem
 - ◆ These are essentially a kind of annotation, exploiting the power of the C++ compiler
 - ◆ A likely approach-“peel off” the easy parallel models, leaving the hard stuff for the heroic programmers
- Others, like the HPCS languages (Chapel, X10, Fortress), add some explicit support for particular structures, but extensibility with performance is unproven.
- Transactional memory can be thought of as a generalization of the “load-reservation/store-conditional” to larger (potentially huge) memory regions. Simple cases are very attractive in transactional memory, but it is unclear how well this will work in practice.
 - ◆ Transactional memory *hardware* may help
 - ◆ Much like shared memory *hardware* is a great way to support higher level programming models
- Microsoft is exploring many approaches and has the advantage of control over their compilers and a pressing need to “solve” the productivity problem.



37

Conclusions

- It's the memory hierarchy!
- A pure, compiler based approach is not credible until
 - ◆ $\frac{\min(\text{performance of compiler on MM})}{\max(\text{performance of hand-tuned MM})} > 0.9$
 - ◆ The “condition number” of that ratio is small (less than 2)
 - ◆ Your favorite performance challenge
- Compilation is hard!
- At the node, the memory hierarchy limits performance
 - ◆ Architectural changes can help (e.g., prefetch, more pending loads/stores) but will always need algorithmic and programming help
 - ◆ Algorithms must adapt to the realities of modern architectures
- Between nodes, complexity of managing distributed data structures limits productivity, ability to adopt new algorithms
 - ◆ Domain (or better, data-structure) specific nano-languages, used as part of a hierarchical language approach, can help



38