

**Latinamerican School for
Computational Materials
Science.**



Optimization techniques

Stefano Cozzini

CNR-INFM DEMOCRITOS, Trieste



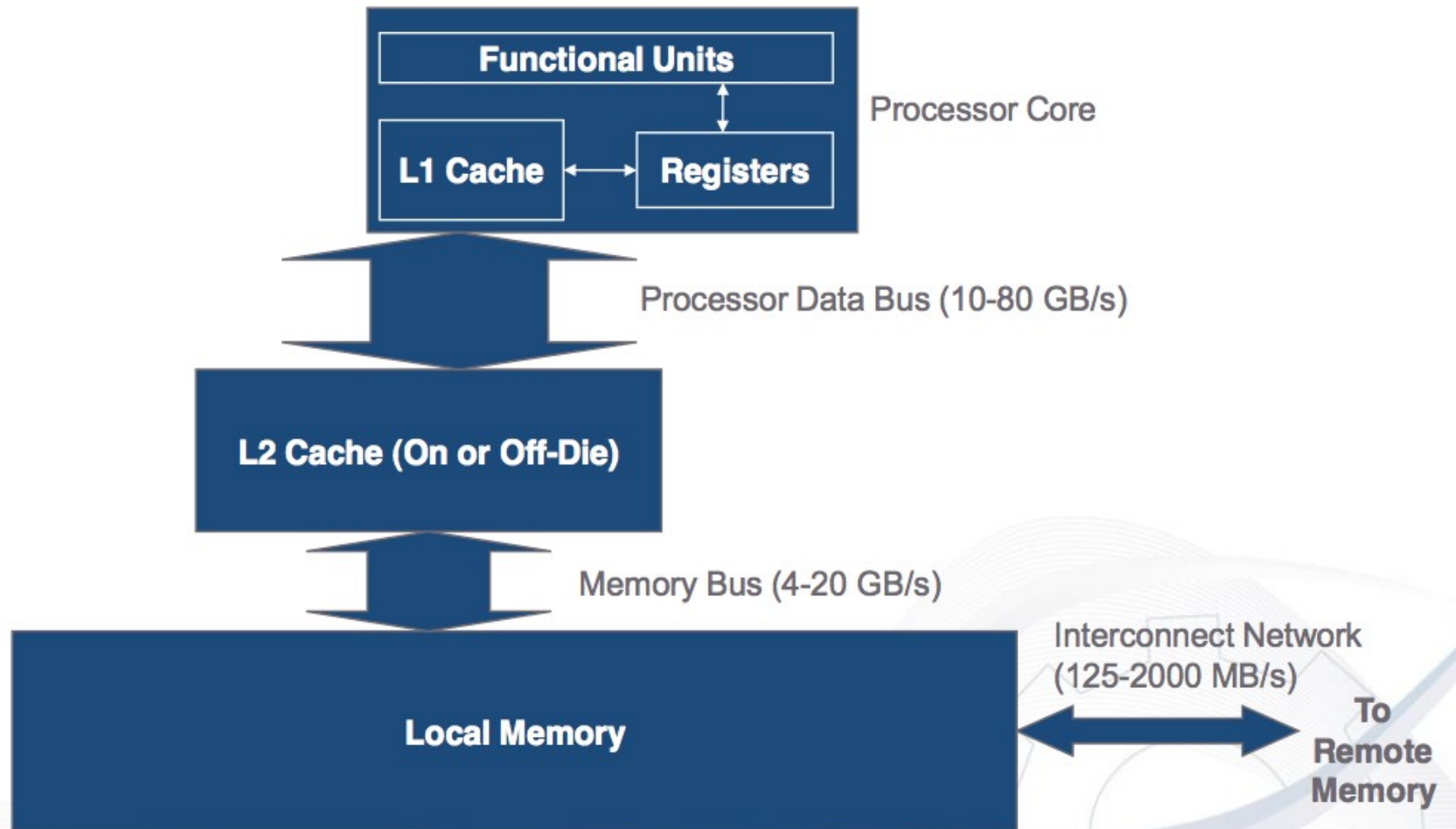
Agenda:

- Introduction
- Performance and evaluating process (Profiling and timing your code)
- Optimization techniques
- General Performance techniques
 - Use of Libraries: see next lecture..

Introduction

- Discuss how to measure performances of your cluster/system
- Discuss performance tuning techniques common to most modern architecture (mainly 32/64 bit commodity processor)
- Using optimization techniques users have control over
 - Code modification
 - Compiler options
- Optimization is a dirty work (and dangerous one for your code...)
- Compiler is your best friend..

Memory hierarchy



Locality of Reference

- **Most programs have a high degree of *locality* in their accesses**
- Memory hierarchy tries to exploit locality
- **Temporal locality:**
Recently referenced items (instr or data) are likely to be *referenced again* in the near future:
 - iterative loops, subroutines, local variables
 - working set* concept
- **Spatial locality:**
programs access data which is *near to each other*:
 - operations on tables/arrays
 - cache line size* is determined by spatial locality
- **Sequential locality:**
processor executes instructions in *program order*:
 - branches/in-sequence ratio is typically 1 to 5

Performance Evaluation process

- **Monitoring System:**
 - Use monitoring tools to better understand your machine's limits and usage
 - is the system limit well suited to run my application ?
 - Observe both overall system performance and single-program execution characteristics. Monitoring your own code
 - Is the system doing well ? Is my program running in a pathological situation ?
- **Monitoring your own code:**
 - Timing the code:
 - timing a whole program (time command `:/usr/bin/time`)
 - timing a portion (all portions) of the program
 - Profiling the program

Useful Monitoring Commands (Linux)

- **Uptime(1)** returns information about system usage and user load
- **ps(1)** lets you see a “snapshot” of the process table
- **top** process table dynamic display
- **free** memory usage
- **vmstat** memory usage monitor

```
top - 08:48:53 up 6 days, 18:35, 5 users, load average: 0.98, 0.55, 0.22
Tasks: 91 total, 2 running, 89 sleeping, 0 stopped, 0 zombie
Cpu(s): 99.3% us, 0.7% sy, 0.0% ni, 0.0% id, 0.0% wa, 0.0% hi, 0.0% si
Mem: 507492k total, 440336k used, 67156k free, 57928k buffers
Swap: 2048248k total, 108456k used, 1939792k free, 99908k cached
```

PID	USER	PR	NI	VIRT	RES	SHR	S	%CPU	%MEM	TIME+	COMMAND
18806	cozzini	25	0	2164	464	396	R	90.8	0.1	3:52.85	a.out
9874	moses	15	0	211m	94m	18m	S	7.6	19.0	54:25.92	firefox-bin
9578	root	15	0	194m	38m	4836	S	1.0	7.7	16:00.59	X
18807	cozzini	16	0	2320	956	756	R	0.3	0.2	0:00.53	top

Monitoring your own code (time)

NAME

time - time a simple command or give resource usage

SYNOPSIS

```
time [options] command [arguments...]
```

DESCRIPTION

The `time` command runs the specified program `command` with the given arguments. When `command` finishes, `time` writes a message to standard output giving timing statistics about this program ..

```
----->time ./a.out  
[program output]
```

```
real    0m1.361s  
user    0m0.770s  
sys     0m0.590s
```

user time: Cpu-time dedicated to your program
sys time: time used by your program to execute system calls
real time: total time aka walltime

User/System/Walltime

- *Real time (or wall clock time)* is the total elapsed time from start to end of a timed task
- *CPU user time* is the time spent executing in user space
 - Does not include time spent in system (OS calls) and time spent executing other processes
- *CPU system time* is the time spent executing system calls (kernel code)
 - System calls for I/O, devices, synchronization and locking, threading, memory allocation
 - Typically does not include process waiting time for non-ready devices such as disks
- CPU user time + CPU system time < real time
 - CPU percentage spent on process = $100\% * (\text{user} + \text{system}) / \text{real}$

a top disaster: swapping..

- virtual or swap memory:

This memory, is actually space on the hard drive. The operating system reserves a space on the hard drive for “swap space”.

- time to access virtual memory **VERY** large:
- this time is done by the system not by your program !
- sometimes the system assumes a killer to kill your program.. (oom killer)

```
top - 08:57:02 up 6 days, 19:35, 7 users, load average: 2.77, 0.73, 0.25
Tasks: 86 total, 2 running, 84 sleeping, 0 stopped, 0 zombie
Cpu(s): 0.3% us, 4.8% sy, 0.0% ni, 0.0% id, 94.2% wa, 0.6% hi, 0.0% si
Mem: 507492k total, 506572k used, 920k free, 196k buffers
Swap: 2048248k total, 941984k used, 1106264k free, 4740k cached
```

PID	USER	PR	NI	VIRT	RES	SHR	S	%CPU	%MEM	TIME+	COMMAND
11656	cozzini	18	0	2172m	408m	260	D	4.3	82.4	0:03.75	a.out
33	root	15	0	0	0	0	D	0.7	0.0	0:00.54	kswapd0
3195	root	15	0	20696	1432	1140	D	0.3	0.3	0:06.81	clock-applet

top disaster example (1)

```
[cozzini@stroligo optimization]$ /usr/bin/time ./a.out
provide an integer (suggested range 100-250)
larger values can be very memory and time-consuming
300
inizialisation time= 11.787208
10.86user 0.98system 0:14.22elapsed 83%CPU (0avgtext+0avgdata 0maxresident)k
0inputs+0outputs (5major+106090minor)pagefaults 0swaps
```

```
[cozzini@stroligo optimization]$ /usr/bin/time ./a.out
provide an integer (suggested range 100-250)
larger values can be very memory and time-consuming
320
Command terminated by signal 2
0.18user 1.81system 0:29.27elapsed 6%CPU (0avgtext+0avgdata 0maxresident)k
0inputs+0outputs (5846major+170788minor)pagefaults 0swaps
```

top disaster example (2)

```
[cozzini@stroligo optimization]$ /usr/bin/time ./a.out <300 &
[cozzini@stroligo optimization]$ free
```

	total	used	free	shared	buffers	cached
Mem:	507492	484916	22576	0	1156	10172
-/+ buffers/cache:		473588	33904			
Swap:	2048248	78108	1970140			

```
[cozzini@stroligo optimization]$ /usr/bin/time ./a.out <320 &
[cozzini@stroligo optimization]$ free
```

	total	used	free	shared	buffers	cached
Mem:	507492	506412	1080	0	252	3936
-/+ buffers/cache:		502224	5268			
Swap:	2048248	546348	1501900			

Timing portion of the code

- Record the time before portion A
- execute portion A
- record the time after portion A
- print/save the difference in time for subsequent analysis
- C function to compute time:
 - `clock`
- Fortran90 function to compute time:
 - `cpu_time routine (f95)`

```
clock_t c0, c1;  
c0 = clock();  
section to code..  
c1= clock();  
cputime = (c1 - c0)/(CLOCKS_PER_SEC );
```

```
call cpu_time(t0)  
  
section to code..  
call cpu_time(t1)  
cputime = (t1 - t0)
```

Well written codes have their own timing report..

```
!!Specific TIMING for section: MD INTEGRATION !!!!!!!!!!!
```

```
!Serial subroutines :
```

```
c
```

!section	times	avg-time	max(PE)	min(PE)
!vscale	2	0.0600	0.0600(0)	0.0600(0)
!scanpairs_prot	100	72.5500	72.5500(0)	72.5500(0)
!vertest_prot	100	2.2600	2.2600(0)	2.2600(0)
!link_list	7	70.2900	70.2900(0)	70.2900(0)
!spme_prot	100	727.8700	727.8700(0)	727.8700(0)
!fill_charge_gri	100	214.5000	214.5000(0)	214.5000(0)
!fft_back	100	79.2700	79.2700(0)	79.2700(0)
!scalar_sum	100	43.2400	43.2400(0)	43.2400(0)
!fft_forw	100	78.8400	78.8400(0)	78.8400(0)
!grad_sum	100	303.6600	303.6600(0)	303.6600(0)
!ewcorr_prot	100	24.4000	24.4000(0)	24.4000(0)
!ewald3_prot	5870100	15.4300	15.4300(0)	15.4300(0)
!pair_force_prot	100	0.0000	0.0000(0)	0.0000(0)
!srfew2_prot	5855979	817.5300	817.5300(0)	817.5300(0)
!dihfrfc_prot	100	3.0800	3.0800(0)	3.0800(0)

Analysis Techniques

- there are three generally available techniques for analyzing code performance:
 - Compiler reports and listings
 - Profiling
 - Hardware performance counters

Compiler Reports and Listings

- Compilers on most modern high performance computers are capable of doing a wide range of optimizations,
- By default, compilers generally do not describe in much detail what kinds of optimizations they were able to perform on a given piece of code.
- However, many compilers will optionally generate *optimization reports* and/or *listing files*.
 - Optimization reports are typically sent to **stderr** at compile time and contain messages describing what optimizations could or could not be applied at various points in the source code.
 - Listing files usually consist of a listing of the source code with messages about optimizations interspersed through the listing.

Reporting and Listing Compiler Options

GNU compilers

None

PGI compilers

-Minfo=option[,option,...]

Prints information to `stderr` on `option`; `option` can be one or more of **time**, **loop**, **inline**, **sym**, or **all**

-Mneginfo=option[,option]

Prints information to `stderr` on why optimizations of type `option` were not performed; `option` can be **concur** or **loop**

-Mlist

Generates a listing file

Intel compilers

-opt_report

Generates an optimization report on `stderr`

-opt_report_file filename

Generates an optimization report to **filename**

Profiling

- Profiling is an approach to performance analysis in which the amount of time spent in sections of code is measured (using either a sampling technique or on entry/exit of a code block) and presented as a histogram.
- This allows a developer to identify the routines which are taking the most execution time, as these are typically the best candidates for optimization.
- Profiling can be done at varying levels of granularity:
 - Subroutine
 - Basic block
 - Source code line
- Profiling usually requires special compilation.
 - The specially compiled executable will generate a file containing execution profile data as it runs.
 - This data file can be analyzed after the code is run
 - a profiling analysis program should be employed

Profiling Compiler Options

FORTTRAN INTEL:

- prof-dir <d> specify directory for profiling output files (*.dyn and *.dpi)
- prof-file <f> specify file name for profiling summary file
- prof-gen instrument program for profiling
- prof-use enable use of profiling information during optimization
- qp compile and link for function profiling with UNIX gprof tool
- p same as -qp

FORTTRAN PGF90

-Mprof[=dwarf|func|hwcts|lines|mpich1|mpich2|time]

- dwarf Generate additional code for profiling
- func Add limited DWARF info for third party profilers
- hwcts Function-level profiling
- lines PAPI-based profiling using hardware counters, 64-bit only
- mpich1/2 Line-level profiling
- time Use profiled MPI communication library; implies -Mmpich1/2
- Sample-based instruction-level profiling

GNU:

- p Generate extra code to write profile information suitable for the analysis program gprof.
- pg Generate extra code to write profile information suitable for the analysis program gprof.

Hardware Performance Counters

- Most modern microprocessors have one or more event counters which can be used to count low level processor events such as floating point operations, cache line misses, and total instructions executed.
- The output from these counters can be used to infer how well a program is utilizing the processor.
- In many cases, there are utilities for accessing these hardware counters, through either a library or a command line timing interface.

How to optimize...

- Iterative optimization
 - 1. Check for correct answers (program must be correct!)
 - 2. Profile to find the hotspots, e.g. most time-consuming routines
 - 3. Optimize these routines using compiler options, compiler directives (pragmas), and source code modifications
 - Repeat 1-3
- Optimizing the hotspots of a program improves overall performance
- Programs with “flat profiles” (flat timing histogram)
 - Programs with lots of routines that each take a small amount of time are difficult to optimize

Optimization Techniques

- There are basically three different categories:
 - **Improve memory performance (The most important)**
 - Better memory access pattern
 - Optimal usage of cache lines (improve spatial locality)
 - Re-usage of cached data (improve temporal locality)
 - **Improve CPU performance**
 - Create more opportunities to go superscalar (high level)
 - Better instruction scheduling (low level)
 - **Use already highly optimized libraries/subroutines**

Where to optimize ?

	Annual increase	Typical value in 2006
Single-chip floating-point performance	59%	4 GFLOP/s
Front-side bus bandwidth	23%	1 GWord/s = 0.25 word/flop
DRAM latency	(5.5%)	70 ns = 280 FP ops = 70 loads

Optimization Techniques for memory

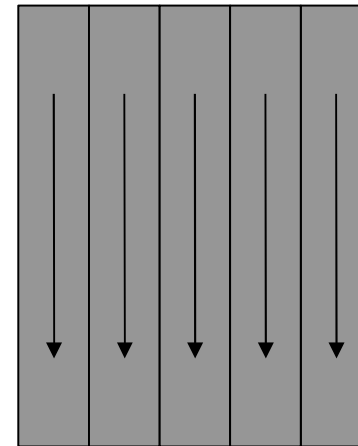
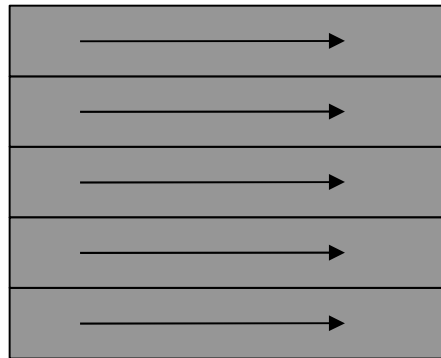
- Loop Interchanges
- Effective Reuse of Data Cache
- Loop Unrolling
- Loop Fusion/Fission
- Prefetching
- Floating Point Division

Storage in Memory

The storage order is language dependent:

Fortran stores “column-wise”

C stores “row-wise”



Accessing elements in storage order greatly enhances the performance for problem sizes that do not fit in the cache(s)
(spatial locality: **stride 1** access)

Loop Interchange

- Basic idea: In a nested loop, examine and possibly change the order of the loop
- Advantages:
 - Better memory access patterns (leading to improved cache and memory usage)
 - Elimination of data dependencies (to increase the opportunities for CPU optimization and parallelization)
- Disadvantage:
 - May make a short loop innermost (which is not good for optimal performances)

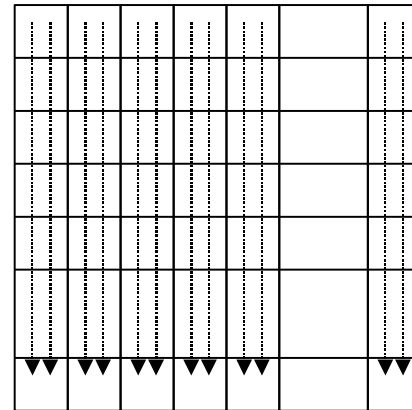
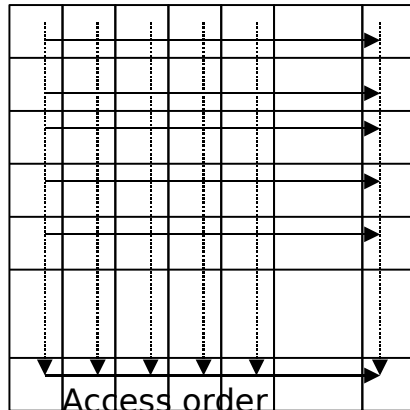
Loop Interchange - Example 1

Original

```
DO i=1,N
  DO j=1,M
    C(i,j)=A(i,j)+B(i,j)
  END DO
END O
```

Interchanged loops

```
DO j=1,M
  DO i=1,N
    C(i,j)=A(i,j)+B(i,j)
  END DO
END DO
```



→ Storage order

→ Access order

Loop Interchange in C

In C, the situation is exactly the opposite
interchange

```
for (j=0; j<M; j++)  
  for (i=0; i<N; i++)  
    C[i][j] = A[i][j] + B[i][j];
```

index reversal

```
for (i=0; i<N; i++)  
  for (j=0; j<N; j++)  
    C[i][j] = A[i][j] + B[i][j];
```

```
for (j=0; j<M; j++)  
  for (i=0; i<N; i++)  
    C[j][i] = A[j][i] + B[j][i];
```

- The performance benefit is the same in this case
- In many practical situations, loop interchange is much easier to achieve than index reversal

Loop Interchange – Mnemonic rule

- With **row-major**, the column or "**rightmost**" index varies most quickly (C/C+)
- With **column-major**, the row of "**leftmost**" index varies most quickly. (Fortran/F90)

Loop Interchange - Example 2

```
DO i=1,300
  DO j=1,300
    DO k=1,300
      A (i,j,k) = A (i,j,k)+ B (i,j,k)* C (i,j,k)
    END DO
  END DO
END DO
```

Loop order	x335 (P4 2.4Ghz)	x330 (P3 1.4Ghz)
i j k	8.77	9.06
i k j	7.61	6.82
j i k	2	2.66
j k i	0.57	1.32
k i j	0.9	1.95
k j i	0.44	1.25

Timings are in seconds

Loop Interchange Compiler Options

GNU compilers:

None

PGI compilers:

-Mvect

Enable vectorization, including loop interchange

Intel compilers:

-O3

Enable aggressive optimization, including loop transformations

TEST IF WHAT THEY CLAIM TO DO IS WHAT THEY ACTUALLY DO

Prefetching

- *Prefetching* is the retrieval of data from memory to cache before it is needed in an upcoming calculation.
- This is an example of general optimization technique called **latency hiding**:
 - communications and calculations are overlapped and occur simultaneously.
- The actual mechanism for prefetching varies from one machine to another.

Prefetching Compiler Options

GNU:

`-fprefetch-loop-arrays`

If supported by the target machine, generate instructions to prefetch memory to improve the performance of loops that access large arrays.

PGI:

`-Mprefetch[=option:n] -Mnoprefetch`

Add (don't add) prefetch instructions for those processors that support them (Pentium 4, Opteron); `-Mprefetch` is default on Opteron; `-Mnoprefetch` is default on other processors.

INTEL:

`-O3` Enable `-O2` optimizations and in addition, enable more aggressive optimizations such as loop and memory access transformation, and **prefetching**.

TEST IF WHAT THEY CLAIM TO DO IS WHAT THEY ACTUALLY DO

Loop Unrolling Example

- Normal loop

```
do i=1,N
  a(i)=b(i)+x*c(i)
enddo
```

- Manually unrolled loop

```
do i=1,N,4
  a(i)=b(i)+x*c(i)
  a(i+1)=b(i+1)+x*c(i+1)
  a(i+2)=b(i+2)+x*c(i+2)
  a(i+3)=b(i+3)+x*c(i+3)
enddo
```

Loop Unrolling Compiler Options

GNU compilers:

`-funroll-loops`

Enable loop unrolling

`-funroll-all-loops`

Unroll all loops; not recommended

PGI compilers:

`-Munroll`

Enable loop unrolling

`-Munroll=c:N`

Unroll loops with trip counts of at least **N**

`-Munroll=n:M`

Unroll loops up to **M** times

Intel compilers:

`-unroll`

Enable loop unrolling

`-unrollM`

Unroll loops up to **M** times

TEST IF WHAT THEY CLAIM TO DO IS WHAT THEY ACTUALLY DO

Floating Point Division

- Floating point division is an expensive operation
 - Takes 22-60 CPs to complete (average about 32 CPs)
 - Usually not pipelined
 - According to the IEEE floating point standard, divisions **must** be carried out as such and not replaced with a multiplication by a reciprocal (even for division by a constant!)
- A possible optimization technique is to “relax” the IEEE requirements and replace a division with multiplication by a reciprocal. Most compilers do this automatically at higher levels of optimization.

Floating Point Division Compiler Options

GNU:

`-funsafe-math-optimizations`

Allow optimizations for floating-point arithmetic that (a) assume that arguments and results are valid and (b) may violate IEEE or ANSI standards.

PGI:

`--Kieee -Knoieee (default)`

Perform floating-point operations in strict conformance with the IEEE 754 standard. Some optimizations are disabled with `-Kieee`, and a more accurate math library is used. The default `-Knoieee` uses faster but very slightly less accurate methods.

INTEL:

`--no-prec-div (i32 and i32em)`

Enables optimizations that give slightly less precise results than full IEEE division. With some optimizations, such as `-xN` and `-xB`, the compiler may change floating-point division computations into multiplication by the reciprocal of the denominator.

Floating Point Division Example

- inverse exercise we will do in the lab..

Floating Point Division With Arrays

- Consider the following loop nest in which the array $A(i, j)$ is scaled by different factors stored in array $B(i)$:

```
do j=1,N
  do i=1,N
    A(i,j)=A(i,j)/B(i)
  enddo
enddo
```

enddo

- The compiler can do no automatic optimization to this, because $B(i)$ is not a scalar. However, you can manually do the following:
 - Create a temporary array to hold the inverses of the $B(i)$ array.
 - Replace the division in the inner loop with multiplication by the temporary array.
 - The resulting code can be unrolled and/or pipelined.



Optimization based on Microprocessor Architectures

- Pipelined Functional Units
- Superscalar Processors
 - Processors which have multiple functional units are said to be *superscalar*.
- Instruction Set Extensions
 - Newer processors have additional instructions beyond the usual floating point add and multiply instructions:
 - SSE2/SSE3/3DNow ! Etc...
 - Cat /proc/cpuinfo..

Instruction Set Extension Compiler Options

GNU:

`-mmmx/no-mmx` These switches enable or disable the use of built-in functions that allow direct access to the MMX, SSE, SSE2, SSE3 and 3Dnow extensions of the instruction set `-msse`

`-mno-sse`

`-msse2 / -mno-sse2`

`-msse3 / -mno-sse3`

`-m3dnow / -mno-3dnow`

PGI:

`--fastsse`
Chooses generally optimal flags for a processor that supports SSE instructions (Pentium 3/4, AthlonXP/MP, Opteron) and SSE2 (Pentium 4, Opteron). Use `pgf90 -fastsse -help` to see the equivalent switches.

INTEL:

`-arch SSE` Optimizes for Intel Pentium 4 processors with Streaming SIMD Extensions (SSE).

`-arch SSE2` Optimizes for Intel Pentium 4 processors with Streaming SIMD Extensions 2 (SSE2).

General techniques

- Blocking/ tiling
- !! Use of optimized libraries !!

Blocking for cache (tiling)

Blocking for cache is:

- An optimization that applies for datasets that do not entirely fit in the (2nd level) data cache
- A way to increase spatial locality of reference i.e. exploit full cache lines
- A way to increase temporal locality of reference i.e. Improves data re-usage

By way of example, let discuss the transpose of a matrix...

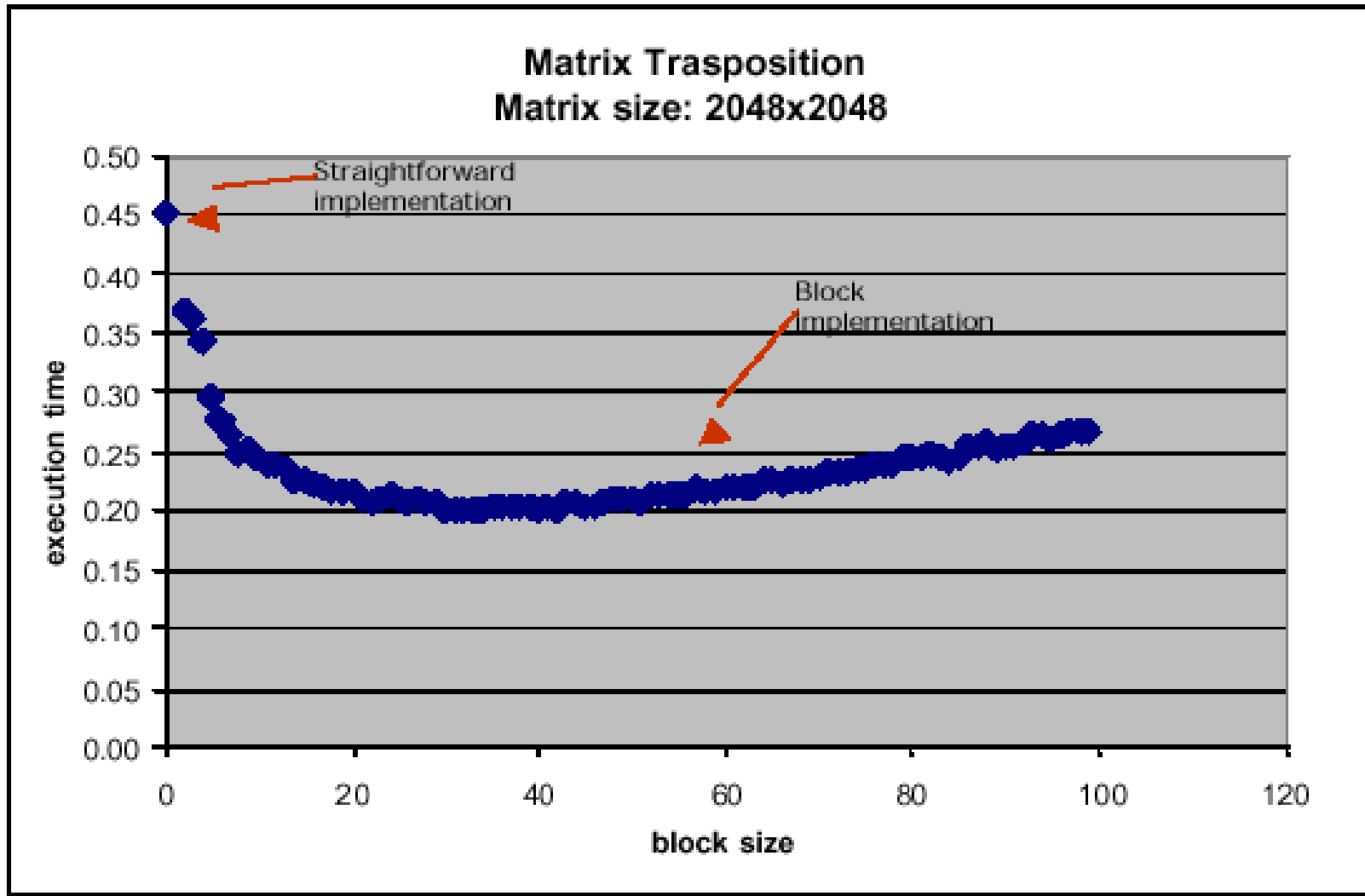
```
do i=1,n
  do j=1,n
    a(i,j)=b(j,i)
  end do
end do
```

Block algorithm for transposing a matrix:

- block data size= bsize
 - $mb=n/bsize$
 - $nb=n/bsize$
- Code is a little bit more complicated if
 - $MOD(n,bsize)$ is not zero
 - $MOD(m,bsize)$ is not zero

```
do ib = 1, nb
  ioff = (ib-1) * bsiz
  do jb = 1, mb
    joff = (jb-1) * bsiz
    do j = 1, bsiz
      do i = 1, bsiz
        buf(i,j) = x(i+ioff, j+joff)
      enddo
    enddo
    do j = 1, bsiz
      do i = 1, j-1
        bswp = buf(i,j)
        buf(i,j) = buf(j,i)
        buf(j,i) = bswp
      enddo
    enddo
    do i=1,bsiz
      do j=1,bsiz
        y(j+joff, i+ioff) = buf(j,i)
      enddo
    enddo
  enddo
enddo
```

Results... (Carlo Cavazzoni data)



Optimizing Matrix Multiply for Caches

- Several techniques for making this faster on modern processors
 - heavily studied
- Some optimizations done automatically by compiler, but can do much better
- In general, you should use optimized libraries (often supplied by vendor) for this and other very common linear algebra operations
 - BLAS = Basic Linear Algebra Subroutines
- Other algorithms you may want are not going to be supplied by vendor, so need to know these techniques

Summary

- Performance programming on uniprocessors requires
 - understanding of memory system
 - levels, costs, sizes
 - understanding of fine-grained parallelism in processor to produce good instruction mix
 - understanding your program
- Compilers are good at instruction level optimization and loop transformation
- User is responsible to present code in most natural way for compiler optimizations..
- The techniques work for any architecture, but choosing details depends on the architecture
- Blocking (tiling) is a basic approach that can be applied to many matrix algorithms