Roberto De Pietri

Esperienze openACC e openMP a Parma

Problems porting application to many-core and GPU systems

R De Pietri, R. Alfieri, M. Borelli, A.Feo, P. Leoni

Summary - purpose of the research

- To check if it is easy to create a simulation program that can be easily ported to different architecture
- How to write a program that run in many-core, GPU, standard CPU and communicating through the MPIlibrary
- * How well it performs
- * Which ones are the best strategy to tune it ...
- * MINIMAL MODIFICATION TO STANDARD CODE

Problem type we are looking for

- * Time evolution of Partial Differential Equations on a cartesian grid
- * The evolution use differential operation with a fixed stencil size
- * The variables that correspond to the "physical space" need to be partitioned to be executed on different physical computing nodes that comunicate which each other



The conditions we would like to match

* 3d grids

- variable number of ghost size
- different number of floating point
 variables associated to each grid
 point
- different number of floating
 points operation needed for the
 update of the variables associated
 to each grid-point.
- When a single physical node do not have all the memory needed to store the configuration.

- A grid size of 1000x1000x1000 with 10
 variables for each grid point and 3 time levels requires 300GBytes of memory
- * IF the update of a grid variable for each point requires 50 Flops for a time step, the total **number** of floating point operations that would be required is 1 TFlops Usually we need to do, at least, 10000-20000 time steps....
- clock frequency is not going up and indeed the # of floating point operations for functional units will not increase.
- * We need to partition the computation !

The "game of life"





- * Just a simple update rule
- * + a local computation that can get easily vectorized by compiler Nc -> sum = a[i]*b[i]
- Total computations (26GFlop/5.28 TFlop(Nc=1000))



Our main update routine

```
// Compute Internals
                                    #pragma acc parallel present(grid[nrows+2][ncols+2],\
                                           next grid[nrows+2][ncols+2], sum, A[0:ncomp], B[0:ncomp]) \
 Total computations
                                            async(2) num gans(100) vector length(16)
                                      #pragma acc loop gang independent
                                      #pragma omp for private(i,j,k,neighbors,sum)
  # update=10
                                      for (i=rmin int; i<=rmax int; i++) {</pre>
                                        #pragma acc loop worker independent
                                        for (j=cmin int; j<=cmax int; j++) {</pre>
  max grid 17000x17000
                                          #pragma ivdep
                                          #pragma vector aligned
                                          #pragma acc loop vector independent reduction(+: sum) \
* vector computation =
                                          for (k=0; k < ncomp; k++)
                                                                    sum += A[k] + B[k];
  ncomp*2 Flops
                                          neighbors = grid[i+1][j+1] + grid[i+1][j] + grid[i+1][j-1]
                                           + grid[i][j+1] + grid[i][j-1] + grid[i-1][j+1]+grid[i-1][j]
                                           + grid[i-1][j-1];
* rule = 10 Flops
                                          if ( ( neighbors > 3.0 ) || ( neighbors < 2.0 ) )
                                            next grid[i][j] = 0.0;
                                          else if ( neighbors == 3.0 )
  ncomp=0 : 26GFlop
                                            next grid[i][j] = 1.0;
                                          else
                                            next grid[i][j] = grid[i][j];
  ncomp=1000 : 5.28TFlop
                                       // end Compute Internale
                                    }
```

The Nvidia KEPLER K20- GPU

SMX

* GPU are much different from CPU and need a lot of effort to write a program that run efficiently on this architecture (CUDA/OpenCL?)

PCI Express 3.0 Host Interface													
GigaThread Engine													
Memory Controller Memory Controller			SAX										
Memory Controller	SMX												

Instruction Cache																		
heduler		1		Warp Scheduler					Warp Scheduler					Warp Scheduler				
Dispatch			Dispatch Dispatch				Dispatch Dispatch				Dispatch Dispatch							
Register File (65.536 x 32-bit)																		
DP Unit	Core	c	ore	Core	DP Unit	LDIST	SFU	Core	Core	Core	OP Unit	Core	Core	Core	DP Unit	LDIST	SFU	
DP UNK	Core	G	ore	Core	DP Unit	LDIST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LDIST	SFU	
o DP Unit	Core	•	ore	Core	DP Unit	LDIST	SFU	Corp	Core	Core	OP Unit	Core	Core	Core	DP Unit	LDIST	SFU	
DP UNR	Core		~	Core	DP Unit	LDIST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LDIST	SFU	
				Core	DP Unit	LDIST	SFU	Core	Core	Core	OF Line	Core	Core	Core	DP Unit	LDIST	SFU	
	٦	1		Core	DP Unit	LDIST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LDIST	SFU	
	man		•	Core	DP Unit	LDIST	SFU	Core	Core	Core	OP Unit	Core	Core	Core	DP Unit	LD:\$T	SFU	
	ory co	2	•	Core	DP Unit	LDIST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LDIST	SFU	
		ntroller	•	Core	DP Unit	LDIST	SFU	Core	Core	Core	OP Unit	Core	Core	Core	DP Unit	LDIST	SFU	
			•	Core	DP Unit	LOIST	SFU	Core	Core	Core	OP Unit	Gore	Core	Core	DP Unit	LDIST	SFU	
	Ī	Morric	•	Core	DP Unit	LOVOT	sru	Core	Core	Core	OP Unit	Core	Core	Core	DP Unit	LDIST	sru	
HERE			•	Core	DP Unit	LDIST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LDIST	SFU	
	1	2	•	Core	OP Unit	LDIST	SFU	Core	Core	Core	OP Unit	Core	Core	Core	DP Unit	LDIST	SFU	
		trolle	•	Core	OP Unit	LDIST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LDIST	SFU	
			•	Core	OP Unit	LDIST	SFU	Core	Core	Core	OP Unit	Core	Core	Core	DPUnit	LDIST	SFU	
		1	•	Core	OP Unit	LDIST	SFU	Core	Core	Core	OP Unit	Core	Core	Core	DP Unit	LDIST	SFU	
	and it			Interconnect Network														
	Nuo.	ery C		64 KB Shared Memory / L1 Cache														
	Contra	Contro		48 KB Road-Only Data Cache														
	0.101			Tex		Tex	(Tex		Tex	1		Tex		Tex		
				Tex		Tex	(Tex		Tex	:		Tex		Tex		
			-	_														

OpenACC vs OpenMP

- OpenACC is a programming standard for parallel computing developed by Cray, CAPS, Nvidia and PGI. The standard is designed to simplify parallel programming of heterogeneous CPU/GPU systems.
- * Like in OpenMP, the programmer can annotate C, C++ and Fortran source code to identify the areas that should be accelerated using PRAGMA compiler directives and additional functions.
- * OpenMP (from 4.0) will be available not only on the CPU but also on GPU we need to check it out ... (gcc-4.9)

OpenMP-different compiler

Nthreads=16 eurora

ICC:
 icc -xavx -fopenmp
 -vec-report2

GNU:
 gcc -mavx -fopenmp
 -ftree-vectorize -ffast-map
 -ftree-vectorizer-verbose=1

* PGI:
pgcc -tp=sandybridge-64
-Mvect=simd:256 -mp=numao
-fast -Minfo=vec,mp



OpenMP-different compiler

- Nthreads=16 eurora
- ICC:
 icc -xavx -fopenmp
 -vec-report2

GNU: gcc -mavx -fopenmp -ftree-vectorize -ffast-ma -ftree-vectorizer-verbose

* PGI:

•

pgcc -tp=sandybridge-64 -Mvect=simd:256 -mp= -fast -Minfo=vec,mp



OpenMP-different compiler

- Nthreads=16 eurora
- ICC:
 icc -xavx -fopenmp
 -vec-report2

GNU: gcc -mavx -fopenmp -ftree-vectorize -ffast-ma -ftree-vectorizer-verbose

* PGI:

•

pgcc -tp=sandybridge-64 -Mvect=simd:256 -mp= -fast -Minfo=vec,mp



OpenACC – #kernel vs #parallel

- 2014-03-31 life_hpc2 on eurora NCOMP=1000 # pragma acc kernels 120 kernels directives 100 Within a kernels construct parallel directives compiler uses classical 60 40 automatic parallelization 20 technology to identify paralle 4000 2000 6000 8000 10000 12000 14000 16000 18000 loops. Lattice Size
- 2014-03-31 life_hpc2 on eurora NCOMP=0 20 # pragma acc parallel parallel directives This is effectively the same⁵ behavior as with OpenMP with the nowait clause on all worksharing loops.. 2000 4000 6000 10000 12000 8000 14000 18000 16000 Lattice size

http://insidehpc.com/2012/08/30/michael-wolfe-on-openacc-kernels-and-parallel-constructs/ https://www.pgroup.com/lit/articles/insider/v4n2a1.htm



https://www.pgroup.com/lit/articles/insider/v4n2a1.htm

```
// Compute Internals
                             #pragma acc kernels present(grid[nrows+2][ncols+2], \
                                         next grid[nrows+2][ncols+2], sum, A[0:ncomp], B[0:ncomp]) async(2)
                               #pragma acc loop gang(100) independent
                               for (i=rmin int; i<=rmax int; i++) { // righe</pre>
                                 #pragma acc loop workers independent
                            // Compute Internals
                            #pragma acc parallel present(grid[nrows+2][ncols+2],\
* # pragma acc kei
                                   next_grid[nrows+2][ncols+2], sum, A[0:ncomp], B[0:ncomp]) \
                                    async(2) num gans(100) vector length(16)
   Within a le
                              #pragma acc loop gang independent
                              #pragma omp for private(i,j,k,neighbors,sum)
   compiler uses cla
                              for (i=rmin int; i<=rmax int; i++) {</pre>
                                #pragma acc loop worker independent
   automatic parall
                                for (j=cmin int; j<=cmax int; j++) {</pre>
                                  #pragma ivdep
   technology to ide
                                  #pragma vector aligned
                                  #pragma acc loop vector independent reduction(+: sum) private(sum)
                                  for (k=0; k < ncomp; k++) sum += A[k] + B[k];
   loops.
                                  // LIFE
                                  neighbors = grid[i+1][j+1] + grid[i+1][j] + grid[i+1][j-1]
* # pragma acc par
                                            + grid[i][j+1] + grid[i][j-1] + grid[i-1][j+1]+grid[i-1][j]
                                            +grid[i-1][j-1];
                                  if ( (neighbors > 3.0 ) || (neighbors < 2.0 ) )
   This is effectively
                                    next grid[i][j] = 0.0;
                                  else if ( neighbors == 3.0 )
   behavior as with
                                    next grid[i][j] = 1.0;
                                  else
   the nowait claus
                                    next grid[i][j] = grid[i][j];
   sharing loops..
                               // end Compute Internals
```

http://insidehpc.com/2012/08/30/nhenaer-worre-on-openaee-kerners-and-paraner-construct

https://www.pgroup.com/lit/articles/insider/v4n2a1.htm

Overlap communication/computation



Overlap communications and computations

* The diagram below show how the communications between the GPU and the CPU and the computations got split



The comparison of the three "devices"

- 17000x17000 is the maximum size we can load on K20 GPU
- Nc=0:26GFlop
- * Nc=1000 : 5.28TFlop
- * K20 is declared of having a peak performance of 1.17 Tflops
- The single node on *Eurora* 16*2.1*8= 0.268 Tflops



Zefiro .vs. Eurora

The same test on ZEFIRO





4000

6000

8000

Lattice size

10000

12000 14000 16000 18000

2000

Test MPI – "eurora" -1-4-8 nodi

 First test on EURORA with some MPI communications (still a small number of nodes —- icc)



Conclusions.....

- OpenMP / OpenACC are very promising compiler technologies for creating portable programs that run on different "future" architectures for HPC systems
- Still difficult to get improvement with respect to OpenMP on a full node.
- * Performance are still not satisfactory !
- * We hope to get better performances going to more Physical systems where the computational load is better balanced.