

How to Write Code that Will Survive the Many-Core Revolution

Write Once, Deploy Many(-Cores)



CAPS worldwide ecosystem



Open**MP**

hmpp

Competence

OpenACC.

DIRECTIVES FOR ACCELERATORS



- OpenMP ARB Accelerator program subcommittee
- OpenStandard Initiative
- HMPP Competences Centers in Europe and Asia

December 2011

www.caps-entreprise.com

Foreword



- "How to write code that will survive the many-core revolution?" is being setup as a collective initiative of HMPP Competence Centers
 - Gather competences from universities across the world (Asia, Europe, USA) and R&D projects: H4H, Autotune, TeraFlux
 - Visit <u>http://competencecenter.hmpp.org</u> for more information





Trend from Top500: Cores per Socket





Trend from Top500: Accelerators



December 2011

IDC @SC11: Study Results



- Many HPC codes aren't seeing a speed up with new hardware systems (due to many-core, lower bandwidth, lower memory/core, etc.)
 - Many applications will need a major redesign
 - o Multi-core will cause many issues to "hit-the-wall"
 - So GPUs can offer a speed-up

• BUT

- GPU still need to be more easy to program
- $\circ~$ Future portability is a key concern







SC11: Consensus about directives ?

- **OpenACC** initiative (2011)
 - CAPS, Cray, NVIDIA, PGI
 - A first common syntax for accelerator regions
 - Visit <u>http://www.openacc-standard.com</u> for more information

- OpenHMPP initiative (2010)
 - Directive Open Standard for many-core programming
 - Complient with OpenACC syntax
 - Topics not covered by OpenACC: *data-flow extension, tracing* interface, auto-tuning APIs...
 - Visit http://www.openhmpp.org for more information



DIRECTIVES FOR ACCELERATORS

OpenACC.









Many parallelism forms are needed to deal with

- Increasing number of processing units
- Some form of vector computing (AVX or SSE instructions)

What to consider by writing your code



- Computing power comes from parallelism
 - Hardware (frequency increase) to software (parallel codes) shift
 - Driven by energy consumption
 - Heterogeneity is the source of efficiency
 - Few large fast OO cores combined with many smaller cores (e.g. APUs)
- Fast moving hardware targets environment

 e.g. fast GPU improvements (RT and HW), new massively parallel CPU
 Write codes that will last many architecture generations
- Keeping a unique version of the codes, preferably monolanguage, is a necessity
 - Reduce maintenance cost
 - Directive-based approaches suitable
 - Preserve code assets



Software Main Driving Forces



- ALF Amdahl's Law is Forever
 - A high percentage of the execution time has to be parallel
 - Many algorithms/methods/techniques will have to be reviewed to scale



• Data locality is expected to be the main issue

Moving data will always suffer latency

One MPI Process per Core Approach



• Benefit

• No need to change existing code

Issues

- Extra latency compared to shared memory use
 - MPI implies some copying required by its semantics (even if efficient MPI implementations tend to reduce them)
- Excessive memory utilization
 - Partitioning for separate address spaces requires replication of parts of the data.
 - When using domain decomposition, the sub-grid size may be so small that most points are replicated (i.e. ghost zone)
 - Memory replication implies more stress on the memory bandwidth which finally produces a weak scaling
- Cache trashing between MPI processes
- Heterogeneity management
 - How are MPI processes linked to accelerator resources?
 - How to deal with different core speed?
 - The balance may be hard to find between the optimal MPI process that makes good use of the CPU core and the use of the accelerators that may be more efficiently used with a different ratio



Thread Based Parallelism Approach



- Benefits
 - Some codes are already ported to OpenMP/threads APIs
 - o Known APIs
- Issues
 - Data locality and affinity management
 - Data locality and load balancing (main target of thread APIs) are in general two antagonistic objectives
 - Usual thread APIs make it difficult / not direct to express data locality affinity
 - Reaching a tradeoff between vector parallelism (e.g. using the AVX instruction set), thread parallelism and MPI parallelism
 - Vector parallelism is expected to impact more and more on the performance
 - Current threads APIs have not been designed to simplify the implementation of such tradeoff
 - Threads granularity has to be tuned depending on core characteristics (e.g. SMT, heterogeneity)
 - Thread code writing style does not make it easy to tune



An Approach for Portable Many-Core Codes CAPS



1 - Express Parallelism, not Implementation CAPS

- Rely on code generation for implementation details
 - Usually not easy to go from a low level API to another low level one
 - Tuning has to be possible from the high level
 - But avoid relying on compiler advanced techniques for parallelism discovery, ...
 - You may have to change the algorithm!
- An example with HMPP





HMPP Hybrid Compiler



16

Rich set of directives for performance



- Manycore programming directives in legacy code
 - Declare and generate GPU versions of computations (codelets)
 - Optimize data movement
 - Distribute computations over CPU cores & GPUs
- Tuning of GPU kernels
 - Advanced code optimizations
 - Control mapping of computations
 - Fully exploit GPU stream architecture



Performance does matter



En average, HMPP reaches Cuda performance +/- 10%



2 x Intel(R) Xeon(R) X5560 @ 2.80GHz (8 cores) - MKL NVidia Tesla C2050, ECC activated – HMPP, CUBLAS, MAGMA





Performances (max) NO TRANSFERTS

rgie atomique • energies a	lemeilves		T in s				
Scalar	18457.52	3030.9	4668.67	7531.1	1054.4	15.68	1933.82
OMP=8	5040.1	379.09	1479.91	1302.64	1248.35	6.76	348.32
НМРР	820.34	56.88	388.23	156.61	81.58	2.04	29.7
	1267.66	75.01	458.95	135.61	66.99	66.51	531.3
	ALL	NOISE	DIFFUS	KERSBS	SHIFT	BOUND	FLUX
OMP / SEQ	3.66	8.00	3.15	5.78	0.84	2.32	5.55
HMPP / SEQ	22.50	53.29	12.03	48.09	12.92	7.69	65.11
CUDA / SEQ	14.56	40.41	10.17	55.53	15.74	0.24	3.64
Diffus = F KERSBS	dup FT FW + di = KER + S	iffrac + FFTE BS	3W		Geom	128 x 128 No I/O 1 MPI	3 x 256
GCdV		CEA,	DAM, DIF, F-91	297 Arpajon			2

Guillaume Colin de Verdière, Onera XtremCFD Workshop, 7th of October, 2011

2 – Do not hide parallelism



- Do not hide parallelism by awkward coding
 - Data structure aliasing, ...
 - Deep routine calling sequences
 - Separate concerns (functionality coding versus performance coding)
- Data parallelism when possible
 - Simple form of parallelism, easy to manage
 - o Favor data locality
 - o But sometimes too static
- Kernels level
 - Expose massive parallelism
 - Ensure that data affinity can be controlled
 - Make sure it is easy to tune the ratio vector / thread parallelism



Data Structure Management



- Data locality
 - Makes it easy to move from one address space to another one
 - Makes it easy to keep coherent
- Do not waste memory
 - Memory per core ratio not improving
- Choose simple data structures
 - Enable vector/SIMD computing
 - Use library friendly data structures
 - May come in multiple forms, e.g. sparse matrix representation
- For instance consider "data collections" to deal with multiple address spaces or multiple devices or parts of a device
 - o Gives a level a adaptation for dealing with heterogeneity
 - Load distribution over the different devices is simple to express



HMPP 3.0 Map Operation on Data Collection CAPS



3 - Debugging Issues



- Avoid assuming you won't have to debug!
- Keep serial semantic
 - For instance, implies keeping serial libraries in the application code
 - $\circ~$ Directives based programming makes this easy
- Ensure validation is possible even with rounding errors
 - Reductions, ...
 - Aggressive compiler optimizations
- Use defensive coding practices
 - $\circ~$ Events logging, parameterize parallelism, add synchronization points,
 - Use debuggers (e.g. Allinea DDT)



Allinea DDT – Debugging Tool for CPU/GPU CAPS

Session Control Search View Help

- Debug your kernels:
 - Debug CPU and GPU concurrently
- Examine thread data
 - Display variables
- Integrated with HMPP
 - Allows HMPP directives breakpoints
 - Step into HMPP codelets





24



4 - Performance Measurements

- Parallelism is about performance!
- Track Amdahl's Law Issues
 - $\circ~$ Serial execution is a killer
 - Check scalability

December

- Use performance tools
- Add performance measurement in the code
 - Detect bottleneck a.s.a.p.
 - $\circ~$ Make it part of the validation process





HMPP Wizard



- Analyze the code to help migrate to Provide tuning advices • many-core architectures within an incremental process
- - Check coalescing of memory accesses
 - Improve parallelism ... 0



HMPP Wizard





December 2011

www.caps-entreprise.com

27

HMPP Performance Analyzer



- Dynamic analyses
- Synthetize metrics based on GPU execution profile

Home 🚯 Execution Profile 🛱 Source files 🍳 Advice 🗲 Performant	nce Analyzer 🔅 Options	CAPS		
H Filters Sesults				
<pre>143</pre>	Codelet		Tune vour	
156 DO j=1, p 157 DO i=1, m 158 prod=0 159 k=1 160 v1a = vin1(i, k) 161 v2a = vin2(k, j)	"process_matrix"	E	Codelet	
163 prod = prod + vla * v2a 164 vla = vin1(i,k) 165 v2a = vin2(k,j) 166 ENDDO 167 prod = prod + vla * v2a 168 vout(i,j) = alpha * prod + beta * vout(i,j) 169 ENDDO	Loop Nest 1 From line 146 to 152 • Kernel #1: 40.6% of the exc	ecution	for Your Hardware	
Image: Image state Image state <td>• Grid: 3D Loop gridific • GPU execution time: • Kernel name:</td> <td>ation 21850.5us</td> <td></td> <td></td>	• Grid: 3D Loop gridific • GPU execution time: • Kernel name:	ation 21850.5us		

5 - Dealing with Libraries



- Library calls can usually only be partially replaced
 - No one to one mapping between libraries (e.g.BLAS, FFTW, CuFFT, CULA,LibJacket)
 - No access to all code (i.e. avoid side effects)
 - Don't create dependencies on a specific target library as much as possible
 - Still want a unique source code
- Deal with multiple address spaces / multi-GPUs
 - Data location may not be unique (copies)
 - Usual library calls assume shared memory
 - Library efficiency depends on updated data location (long term effect)
- Libraries can be written in many different languages
 - CUDA, OpenCL, HMPP, etc.
- There is not one binding choice depending on applications/users
 - Binding needs to adapt to uses depending on how it interacts with the remainder of the code
 - Choices depend on development methodology



Library Interoperability in HMPP 3.0





Conclusion



- Software has to expose massive parallelism
 - The key to success is the algorithm!
 - The implementation has "just" to keep parallelism flexible and easy to exploit
- Directive-based approaches are currently one of the most promising track
 - Preserve code assets
 - May separate parallelism exposure from the implementation (at node level)
- Remember that even if you are very careful in your choices you may have to rewrite parts of the code
 - Code architecture is the main asset here



CAPS Many-core programming Parallelization GPGPU NVIDIA Cuda OpenHMPP Directive-based programming Code Porting Methodology **OpenACC Hybrid Many-core Programmin HPC community Petaflops** Parallel computing HPC open standard Exaflops Open CL High Performance Computing Code speedup Multi-core programming Massively parallel **HMPP** Competence Center Hardware accelerators programming DevDeck Parallel programming interface



Global Solutions for **Many-Core** Programming

http://www.caps-entreprise.com