

Barcelona Supercomputing Center Centro Nacional de Supercomputación

The OmpSs Programming Model encorei

Jesus Labarta Director Computer Sciences Research Dept. BSC

Challenges on the way to Exascale

- Efficiency (..., power, ...)
- Variability
- Memory
- Faults
- Scale (...,concurrency, strong scaling,...)
- **Complexity** (...Hierarchy /Heterogeneity,...)



"To kill flies with

CANONS " Spanish saying



Node Pckg Chip

J. Labarta, et all, "BSC Vision towards Exascale"

Cluster



Core

IJHPCA vol 23, n. 4 Nov 2009



Is any of them more important than the others?

Which?

The sword to cut the "multicore" Gordian Knot



StarSs: a pragmatic approach

- Rationale
 - Runtime managed, asynchronous data-flow execution models are key
 - Need to provide a natural migration towards dataflow
 - Need to tolerate "acceptable" relaxation of pure models
 - Focus on algorithmic structure and not so much on resources
- StarSs: a family of task based programming models
 - Basic concept: write sequential on a flat single address space + directionality annotations
 - Order **IS** defined !!!
 - Dependence and data access related information (NOT specification) in a single mechanism
 - Think global, specify local
 - Power to the runtime !!!





StarSs: data-flow execution of sequential programs



Supercomputing Center

BSC



StarSs vs OpenMP

(



Barcelona Supercomputing

Center

BSC

StarSs: the potential of data access information

- Flat global address space seen by programmer
- Flexibility to dynamically traverse dataflow graph "optimizing"
 - Concurrency. Critical path
 - Memory access: data transfers performed by run time
- Opportunities for runtime to
 - Prefetch
 - Reuse
 - Eliminate antidependences (rename)
 - Replication management
 - Coherency/consistency handled by the runtime





Hybrid MPI/StarSs

- Overlap communication/computation
- Extend asynchronous data-flow execution to outer level
- Linpack example: Automatic lookahead

```
...
for (k=0; k<N; k++) {
   if (mine) {
      Factor panel(A[k]);
      send (A[k])
   } else {
      receive (A[k]);
      if (necessary) resend (A[k]);
   for (j=k+1; j<N; j++)</pre>
      update (A[k], A[j]);
#pragma css task inout(A[SIZE])
void Factor panel(float *A);
#pragma css task input(A[SIZE]) inout(B[SIZE])
void update(float *A, float *B);
```



void receive(float *A);

#pragma css task input(A[SIZE])

void resend(float *A);

V. Marjanovic, et al, "Overlapping Communication and Computation by using a Hybrid MPI/SMPSs Approach" ICS 2010



All that easy/wonderful?

- Difficulties for adoption
 - Chicken and egg issue users ↔ manufacturers
 - Availability.
 - Runtime implementations chasing new platforms
 - Development as we go
 - Fairly stable, minimal application update cost.
 - Happens to all models, by all developers (companies, research,...)
- Lack of program development support
 - Understand application dependences
 - Understand potential and best direction
- Difficulties of the models themselves
 - Simple concepts take time to be matured
 - As clean/elegant as we claim?
 - Legacy sequential code less structured than ideal





- Towards EXaflop applicaTions (EC FP7 Grant 261580)
- Demonstrate that Hybrid MPI/SMPSs addresses the Exascale challenges in a an productive and efficient way.
 - Deploy at supercomputing centers: Julich, EPCC, HLRS, BSC
 - Port Applications (HLA, SPECFEM3D, PEPC, PSC, BEST, CPMD, LS1 MarDyn) and develop algorithms.
 - Develop additional environment capabilities
 - tools (debug, performance)
 - improvements in runtime systems (load balance and GPUSs)
 - Support other users
 - Identify users of TEXT applications
 - Identify and support interested application developers
 - Contribute to Standards (OpenMP ARB, PERI-XML)





Name / Institution	Type of processor	Total core count
MareNostrum (BSC)	Power PC 970 MP	10240
JuGene (JSC)	32-bit PowerPC 450	294912
JuRoPA (JSC)	Dual-Quadcore Intel	26304
	Nehalem	
Laki NEC Nehalem cluster	Intel Xeon (X5560)	5600
(HLRS)	Nehalem	
Cray XT5m (HLRS)	Quad-Core AMD Opteron	896
HECToR Cray XE6 (EPCC)	12-core AMD Opteron	44544
	Magny Cours	
Hermit Cray XE6 (HLRS)	AMD Opteron 8-Core,	1344
	2GHz Magny Cours	



A PASIF



Codes being ported

- Scalapack: Cholesky factorization (UJI)
 - Example of the issues in porting legacy code
 - Demonstration that it is feasible
 - The importance of scheduling
- LBC Boltzmann Equation Solver Tool (HLRS)
 - Solver for incompressible flows based on Lattice-Boltzmann methods (LBM)
 - LBM well suited for highly complex geometries. Simplified implementation: lbc





StarSs: history/strategy/versions

Basic SMPSs

must provide directionality ∀argument Contiguous, non partially overlapped **Renaming** Several schedulers (priority, locality,...) No nesting

C/Fortran

MPI/SMPSs optims.

SMPSs regions

C, No Fortran must provide directionality ∀argument ovelaping &strided Reshaping strided accesses Priority and locality aware scheduling

Evolving research since 2005



OMPSs

C, C++, Fortran OpenMP compatibility (~) Contiguous and strided args. Separate dependences/transfers Inlined/outlined pragmas Nesting Heterogeneity: SMP/GPU/Cluster No renaming, Several schedulers: "Simple" locality aware sched,...



OmpSs

- What; Our long term infrastructure
 - "Acceptable" relaxation of basic StarSs concept
 - Reasonable merge/evolution of OpenMP
- Basic features
 - Inlined/outlined task specifications
 - Support multiple implementations for outlined tasks
 - Separation of information to compute dependences and data movement
 - Not necessary to specify directionality for an argument
 - Concurrent: Breaking inout chains (for reduction implementation)
 - Nesting
 - Heterogeneity: CUDA, OpenCL (in the pipe)
 - Strided and partially aliased arguments
 - C, C++ and Fortran





```
#pragma omp target device(cuda)
  global void cuda perlin (pixel output [], float time,
                                         int j, int rowstride)
{
     unsigned int i = blockIdx.x * blockDim.x + threadIdx.x;
     unsigned int off = blockIdx.y * blockDim.y + threadIdx.y;
       float vdx = 0.03125f;
       float vdy = 0.0125f;
       float vs = 2.0f;
      float bias = 0.35f;
      float vx = 0.0f;
      float red, green, blue;
      float xx, yy;
      float vy, vt;
                     for (j = 0; j < img height; j+=BS) {
                                                                                  // BS image rows per task
      vx = ((float) i) * vdx;
      vy = ((float) (j+off))
      vt = time * vs;
                                pixel *out = &output[j*rowstride];
      xx = vx * vs;
                                #pragma omp target device(cuda) copy deps
      yy = vy * vs;
      red = noise3(xx, vt, y
                                #pragma omp task output([rowstride*BS]out)
      green = noise3(vt, yy,
      blue = noise3(yy, xx,
      red += bias;
      green += bias:
                                      dim3 dimBlock;
      blue += bias;
      // Clamp to within [0
                                      dim3 dimGrid;
      red = (red > 1.0f) ? 1
      green = (green > 1.0f)
                                      dimBlock.x = (img width < BSx) ? img width : BSx;
      blue = (blue > 1.0f) ?
      red = (red < 0.0f) ? 0
                                      dimBlock.y = (BS < BSy) ? BS : BSy;
      green = (green < 0.0f)
      blue = (blue < 0.0f) ?
                                      dimBlock.z = 1;
      red *= 255.0f;
      green *= 255.0f;
                                      dimGrid.x = img width/dimBlock.x;
      blue *= 255.0f;
                                      dimGrid.y = BS/dimBlock.y;
      output[(off * rowstride
      output[(off * rowstride
      output[(off * rowstride
                                      dimGrid.z = 1;
      output[(off * rowstride
                                      cuda perlin <<<dimGrid, dimBlock>>> (out, time, j, rowstride);
                     #pragma omp taskwait noflush
```

Barcelona Supercomputing

Centro Nacional de Supercomputació

Center

BSC

One source → many configurations of clusters with CUDA











🔳 Task @ n	body.2n_2GPU.prv	
THREAD 1.1.1		
THREAD 1.1.2		
THREAD 1.1.3		
THREAD 1.1.4		
THREAD 1.2.1		
THREAD 1.2.2		
THREAD 1.2.3		
THREAD 1.2.4		
	9.320.196 us	35.493.973 118



🔳 Task @ nl	body.4n_2GPU.prv	
THREAD 1.1.1 THREAD 1.1.2 THREAD 1.1.2 THREAD 1.1.4 THREAD 1.2.1 THREAD 1.2.1 THREAD 1.2.2 THREAD 1.2.4 THREAD 1.2.4 THREAD 1.2.4 THREAD 1.3.1 THREAD 1.4.3 THREAD 1.4.3 THREAD 1.4.3		
	18.273.887 us	44.447.664 us

J. Bueno et al, "Productive Programming of GPU Clusters with OmpSs", IPDPS2012



Jesus Labarta. OmpSs @ EPoPPEA, January 2012

StarSs NOT only «scientific computing»

- Plagiarism detection
 - Histograms, sorting, …
- Trace browsing
 - Paraver
- Clustering algorithms
 - G-means
- Image processing
 - Tracking
- Embedded and consumer





- Discrete/atomic task
 - Run to completion task. Start and end only interaction points. No dependencies in/out to/from inside a task
 - Interactions half way through a task?
- Late dependence binding
 - Dependences are computed at task instantiation time.
 - Do we need mechanisms for later dependence computation?
- OmpSs relaxation of functional model
 - No need to specify directionality for all arguments, Commutative clause,...
 - Flexibility risk tradeoff?



Limitations?

- Limitation in data access patterns
 - Contiguous/Strided regions
 - Need/can afford further structures? Irregularly scattered, pointer traversal, nested,
 ...
- Granularity: flexibility vs. cost
 - Parallelism and lookahead more important than overhead

J.M. Perez et al, "Handling task dependencies under strided and aliased references" ICS 2010

• When: determined at instantiation time: may be too early if too much lookahead

How much of a limitation, alternatives, worthwhile? needed usage feedback



StarSs: Enabler for exascale

- Can exploit very unstructured parallelism
 - Not just loop/data parallelism
 - Easy to change structure
- Supports large amounts of lookahead
 - Not stalling for dependence satisfaction
- Allow for locality optimizations to tolerate latency
 - Overlap data transfers, prefetch
 - Reuse
- Nicely hybridizes into MPI/StarSs
 - Propagates to large scale the node level dataflow characteristics
 - Overlap communication and computation
 - A chance against Amdahl's law

- Homogenized view at heterogeneity
 - Any # and combination of CPUs, GPUs
 - Support autotuning
- Malleability: Decouple program from resources
 - Allowing dynamic resource allocation and load balance
 - Tolerate noise

Data-flow; Asynchrony

Potential is there; Can blame runtime

Compatible with proprietary low level technologies



• A change in mentality



- Deeply rooted (in or genes), but need to overcome our fears.
 - May require some effort, but it is possible and there is a lot to gain.
 - Understanding and confidence through tools will be key
 - Need education from very early levels (shape instead of reshape minds)
- Adaptability/Flexibility is key to survive in rapidly changing environments



