www.bsc.es

**Barcelona
Supercomputing
Center**
*Centro Nacional de Supercomputación*

# On the Instrumentation of OpenMP and OmpSS Tasking Constructs

H.Servat, X.Teruel, G. Llort, A.Duran, J.Giménez, X.Martorell, E. Ayguadé and J. Labarta
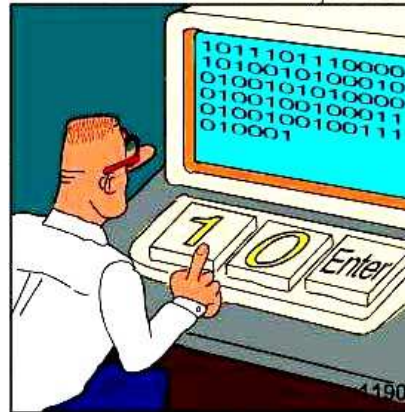
# Supercomputers today

## )) Top500 (June 2012)

| Rank | Site | Computer/Year Vendor | Cores | Rank | Site | Computer/Year Vendor | Cores |
|---|---|---|---|---|---|---|---|
| 1 | DOE/NNSA/LLNL United States | **Sequoia** - BlueGene/Q, Power BQC 16C 1.60 GHz, Custom / 2011 IBM | 1572864 | 6 | DOE/SC/Oak Ridge National Laboratory United States | **Jaguar** - Cray XK6, Opteron 6274 16C 2.200GHz, Cray Gemini interconnect, NVIDIA 2090 / 2009 Cray Inc. | 298592 |
| 2 | RIKEN Advanced Institute for Computational Science (AICS) Japan | K computer, SPARC64 VIIIfx 2.0GHz, Tofu interconnect / 2011 Fujitsu | 705024 | 7 | CINECA Italy | **Fermi** - BlueGene/Q, Power BQC 16C 1.60GHz, Custom / 2012 IBM | 163840 |
| 3 | DOE/SC/Argonne National Laboratory United States | **Mira** - BlueGene/Q, Power BQC 16C 1.60GHz, Custom / 2012 IBM | 736432 | 8 | Forschungszentrum Juelich (FZJ) Germany | **JuQUEEN** - BlueGene/Q, Power BQC 16C 1.60GHz, Custom / 2012 IBM | 131072 |
| 4 | Leibniz Rechenzentrum Germany | **SuperMUC** - iDataPlex DX360M4, Xeon E5-2680 8C 2.70GHz, Infiniband FDR / 2012 IBM | 147456 | 9 | CEA/TGCC-GENCI France | **Curie thin nodes** - Bullx B510, Xeon E5-2680 8C 2.700GHz, Infiniband QDR / 2012 Bull | 77184 |
| 5 | National Supercomputing Center in Tianjin China | **Tianhe-1A** - NUDT YH MPP, Xeon X5670 6C 2.93 GHz, NVIDIA 2050 / 2010 NUDT | 136368 | 10 | National Supercomputing Centre in Shenzhen (NSCS) China | **Nebulae** - Dawning TC3600 Blade System, Xeon X5650 6C 2.66GHz, Infiniband QDR, NVIDIA 2050 / 2010 Dawning | 120640 |

## )) Parallel programming is getting (even) harder

- More parallelism needed because of large core count
- Heterogenity brings an additional complexity layer

# Think in serial

# Execute in parallel
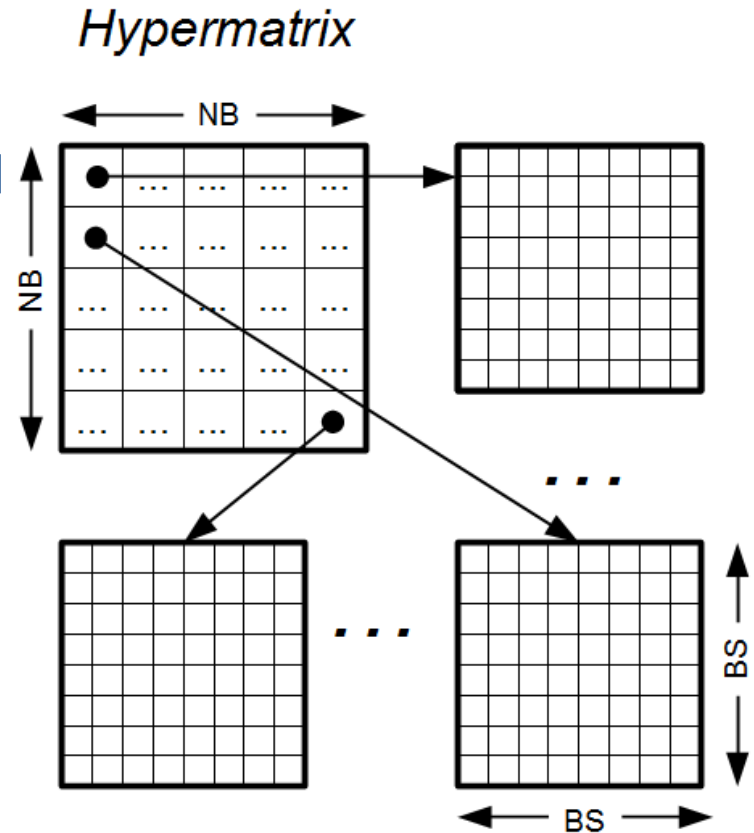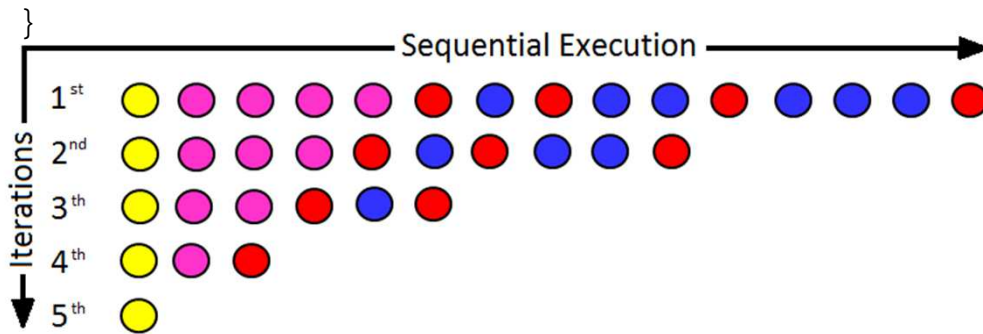
```
void Cholesky( float *A, int NB) {
  int i, j, k;
  for (k=0; k<NB; k++) {
    spotrf (A[k*NB+k]);
    for (i=k+1; i<NB; i++)
      strsm (A[k*NB+k], A[k*NB+i]);
    for (i=k+1; i<NB; i++) {
      for (j=k+1; j<i; j++)
        sgemm( A[k*NB+i], A[k*NB+j], A[j*NB+i]);
      ssyrk (A[k*NB+i], A[i*NB+i]);
    }
  }
}
```

OUT: A[0]

IN: A[0]   OUT: A[1] … A[4]

Sequential Execution →

Iterations

1st 2nd 3th 4th 5th

**Hypermatrix**

NB

NB

BS

BS

**‖ Iters = 5, Critical path = 35**

Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

# Cholesky (OpenMP 2.5)

```
void Cholesky( float *A ) {
  int NB, i, j, k;
  for (k=0; k<NB; k++) {
    ● spotrf (A[k*NB+k]);
#pragma omp parallel for
    for (i=k+1; i<NB; i++)
      ● strsm (A[k*NB+k], A[k*NB+i]);
    for (i=k+1; i<NB; i++) {
#pragma omp parallel for
      for (j=k+1; j<i; j++)
        ● sgemm( A[k*NB+i], A[k*NB+j], A[j*NB+i]);
      ● ssyrk (A[k*NB+i], A[i*NB+i]);
    }
  }
}
```
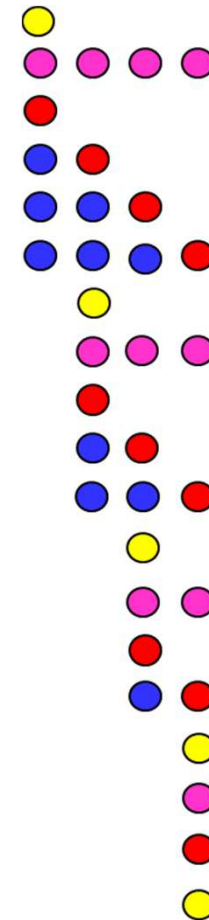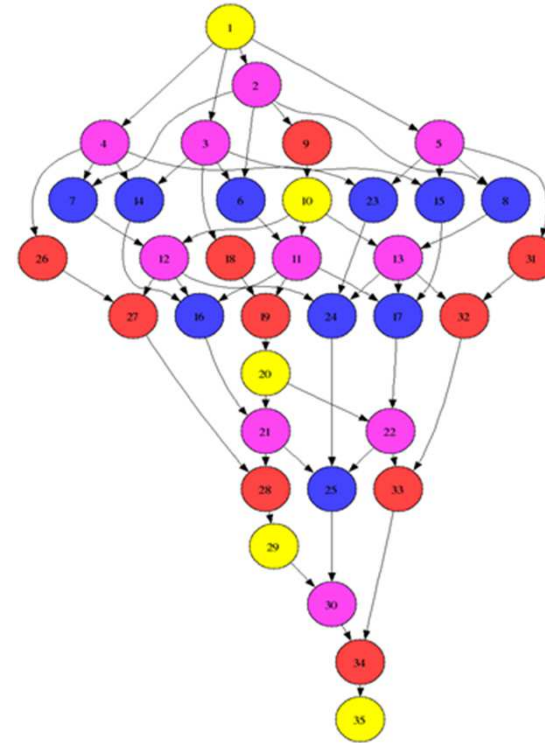
《 Iters = 5, Critical path = 25

```
void Cholesky( float *A ) {
  int NB, i, j, k;

  for (k=0; k<NB; k++) {
    ● spotrf (A[k*NB+k]);
#pragma omp parallel for
    for (i=k+1; i<NB; i++)
      ● strsm (A[k*NB+k], A[k*NB+i]);
    for (i=k+1; i<NB; i++) {
      for (j=k+1; j<i; j++)
#pragma omp task
        ● sgemm( A[k*NB+i], A[k*NB+j], A[j*NB+i]);
#pragma omp task
      ● ssyrk (A[k*NB+i], A[i*NB+i]);
#pragma omp taskwait
    }
  }
}
```

❰❰ Iters = 5, Critical path = 19

# Cholesky (OmpSs)

```
void Cholesky( float *A ) {
  int NB, i, j, k;

  for (k=0; k<NB; k++) {
#pragma omp task inout(A[k*NB+k])
    ● spotrf (A[k*NB+k]);
    for (i=k+1; i<NB; i++)
#pragma omp task input(A[k*NB+k]) inout(A[k*NB+i])
      ● strsm (A[k*NB+k], A[k*NB+i]);
    for (i=k+1; i<NB; i++) {
      for (j=k+1; j<i; j++)
#pragma omp task input(A[k*NB+i], A[k*NB+j]) \
                 inout(A[j*NB+i])
        ● sgemm( A[k*NB+i], A[k*NB+j], A[j*NB+i]);
#pragma omp task input (A[k*NB+i]) inout(A[i*NB+i])
        ● ssyrk (A[k*NB+i], A[i*NB+i]);
    }
  }
}
```



**《 Iters = 5, Critical path = 13** (33% shorter than OpenMP)

# OpenMP vs. OmpSs

## OpenMP – You say how to parallelize

- Based in compiler directives
- Worksharing constructs (<= 2.5): do/for loops, sections
- Tasking constructs (>= 3.0): task, taskwait, taskyield
    - While, linked lists, recursivity

## OmpSs – You say how data is used ➜ System does the rest

- Extends OpenMP
- Focus on tasks
- Heterogeneity
    - SMP, GPU
- Data dependences

```
#pragma omp task out(x)        // 1
x = 5;
#pragma omp task in(x)         // 2
printf("%d\n" , x ) ;
#pragma omp task in(x) out(y) // 3
y = x + 1;
#pragma omp task in(y)         // 4
printf ("%d\n" , y ) ;
```

One image worths a thousand words.
Do not speculate about your code perfomance.
**Look at it.**

- **((** Since 1991
- **((** Based on traces
- **((** Open-source
- **((** Core Tools
  - ▪ Extrae    – Trace generation
  - ▪ Paraver  – Trace analyzer
  - ▪ Dimemas– Message passing simulator
- **((** Detail and intelligence

# Extrae – Trace generation

**((** **Parallel programming model runtime**

 - MPI, OpenMP, Pthreads, CUDA…

**((** **Counters**

 - CPU         – PAPI and PMAPI
 - Network  – Myrinet (GM and MX)
 - OS          – Memory allocation, resource usage

**((** **Links to source**

 - Callstack at MPI calls
 - User functions selected (default none)

**((** **Periodic samples**

 - PAPI counters + callstack

**((** **User events**

# Synergy

- **What?**
  - OmpSs + CEPBA Tools

- **Why?**
  - Improve the analysis adding data only known by the parallel runtime
    - Ready task, blocked task, spins, sleeps, yields, etc.
  - Facilitate the development – analysis – optimization cycle

- **How?**
  - Implement services to communicate both suites

# Code transformations

## « OmpSs in several flavours

- **Performance**
  - Production runs
- **Instrumentation**
  - Logs information at compile and run stages

```
{
  code 1
  #pragma omp task
  {
    code 2
  }
  code 3
}
```
**User source code**

```
code2OL
{
  code 2
}

{
  code 1
  t = nanos_create_task (code2OL);
  nanos_submit_task (t);
  code 3
}
```
**Mercurium translation (Performance)**

**Mercurium translation (Instrumented)**

```
nanos_create_task
{
  // Task creation
}

nanos_submit_task
{
  // Queue task into "run" queue
}
```
**Nanos++ RTL (Performance)**

**Nanos++ RTL (Instrumented)**

# Thread identification

**(( Extrae needs to know which thread is emitting the event**

- Leveraged to Nanos++ via callback

# Task suspension, resumption & migration



**((  Each task keeps an instrumentation context**

- Includes information to be backed up and restored at schedule points
  - Routine being executed, task identifier, task state…

**((  When the task migrates, also does the context**

- Enables support for untied tasks

# Tasks relations



## Draw lines relating
- Task migrations
- Data dependences
- Others: taskwait dependences, data transfers…

## Useful to study
- Execution/Critical path
- Scheduling policies

# Visualization

**◖ Thread-view (standard)**

- Determine thread usage
- Shows task migrations



Recursive multisort

Cilk_sort          Cilk_merge

**◖ Task-view (new)**

- Displays task liveness, logical order and dependences (critical path)



Cilk_sort          Cilk_merge

**◖ Complementary views**

# BOTS Benchmarks

## Floorplan



## N Queens



## Sparse LU



## Strassen

# N Queens overhead study



| | .00-.01 | .01-.02 | .02-.03 | .03-.04 | .04-.05 | .05-.06 | .06-.07 | .07-.08 | .08-.09 |
|---|---|---|---|---|---|---|---|---|---|
| Cut-off 2 | 99.99% | 0.01% | | | | | | | |
| Cut-off 3 | 99.99% | 0.01% | | | | | | | |
| Cut-off 4 | 99.99% | 0.01% | | | | | | | |
| Cut-off 5 | 44.02% | 8.51% | 8.55% | 10.76% | 11.52% | 6.94% | 4.30% | 2.75% | 1.69% |

**《 TLB misses rocket**

- Cut-off ↑ ➔ tasks ↑↑ ➔ traced events ↑↑ ➔ memory usage ↑↑

```
for ( ; ; ) {

  #pragma omp task output(*qleft,*qright)
  {
    gatherConservativeVars();
    constoprim();
    equation_of_state();
    if (H.iorder != 1) slope();
    trace();
    qleftright();
  }

  #pragma omp task input(*qleft,*qright)\
                   output(*qgdnv)
  {
    riemann()
  }

  #pragma omp task input(*qgdnv)
  {
    cmpflx();
    updateConservativeVars();
  }

}
#pragma omp taskwait
```
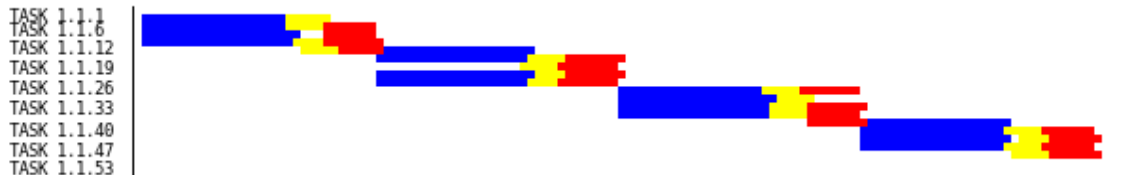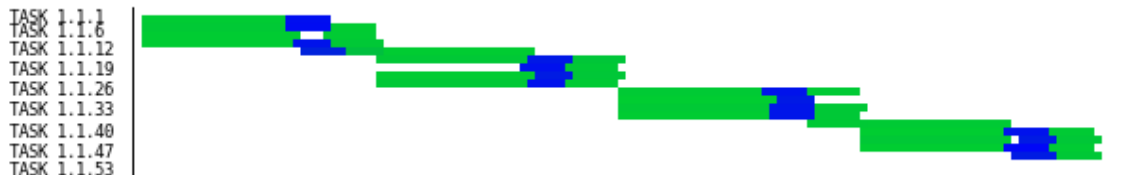


Nanos user function @ hydro.4tasks.4pes.256.task.prv

# Conclusions

**((** **Integrated environment for productive programming and analysis**

- OpenMP and OmpSs codes
- Support for tasking constructs (either tied or untied)
- Extra information: Application performance + Runtime internals

**((** **No need for "magic tricks" to instrument**

- The runtime itself produces the data automatically

**((** **New displays**

- Thread-view + Task-view
- Track task relations (migrations, dependences, etc.)
- Correlate with hardware counters

**((** **Tool for the application and the programming model developer**

# Future Work

**《** Critical path studies

**《** Scalability

- Reduce the information gathered from the runtime

**《** Support for distributed memory systems and accelerators

www.bsc.es

**Barcelona**
**Supercomputing**
**Center**
*Centro Nacional de Supercomputación*

Thank you!