



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

Tutorial OmpSs: Development methodology and infrastructure

PUMPS 2013 tutorial

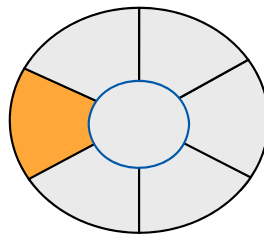
Hybrid and Heterogeneous Parallel Programming with MPI/OmpSs for Exascale Systems

Rosa M Badia, Xavier Martorell

Tutorial OmpSs

« Agenda	10:00 – 11:00	Tasking in OpenMP 3.0 and 4.0	60 min
	11:00 – 11:15	Coffee break	15 min
	11:15 – 12:15	Introduction to OmpSs programming model <ul style="list-style-type: none">• Introduction to StarSs• OmpSs syntax• Simple examples• Development methodology and infrastructure	60 min
	12:15– 12:45	Practical: heat equation example and divide-and-conquer (part I)	30 min
	12:45 – 14:00	Lunch	75 min
	14:00 – 15:00	Practical: heat equation example and divide-and-conquer (part I)	90 min
	15:00 – 15:30	Programming using a hybrid MPI/OmpSs approach	15 min
	15:30 – 17:00	Practical: heat equation example and matrix-multiply	105 min

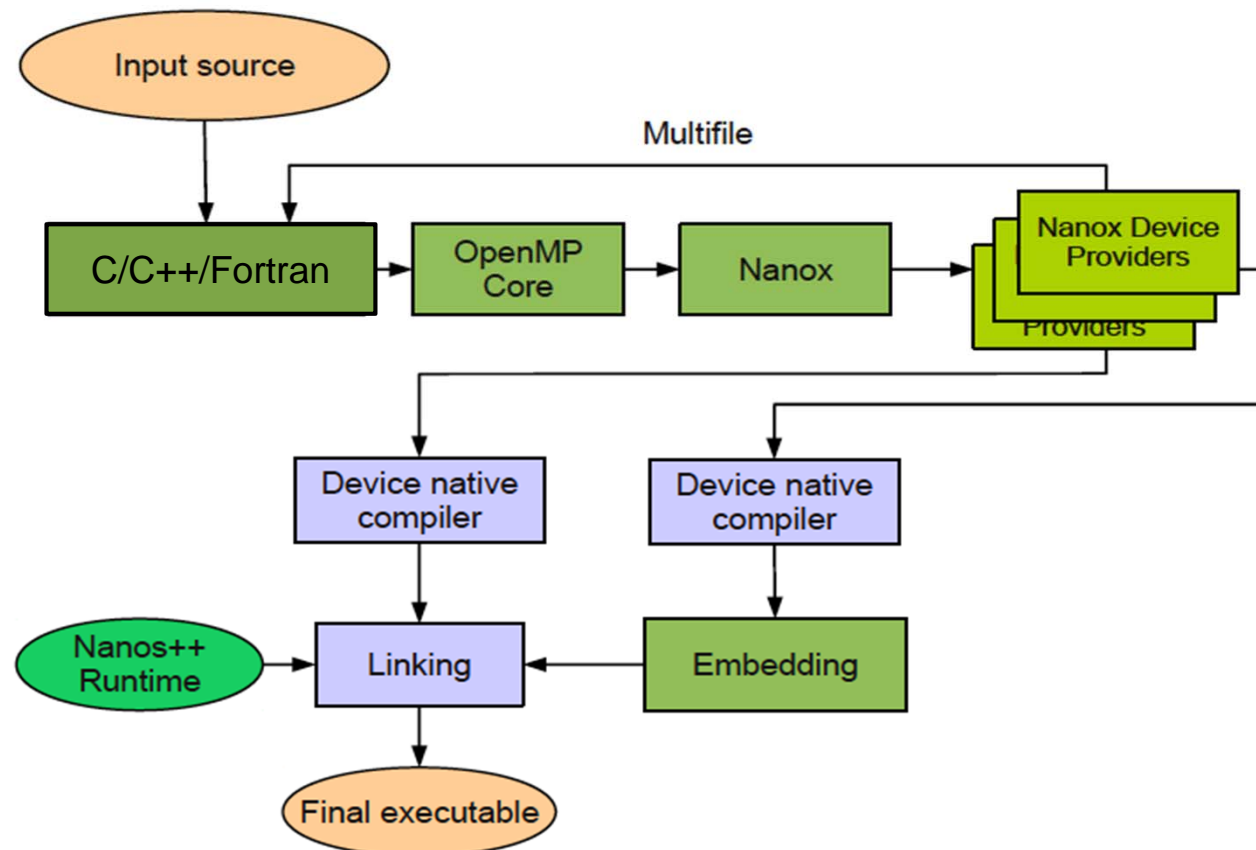
OmpSs compiler and runtime



Mercurium Compiler

- Recognizes constructs and transforms them to calls to the runtime
- Manages code restructuring for different target devices

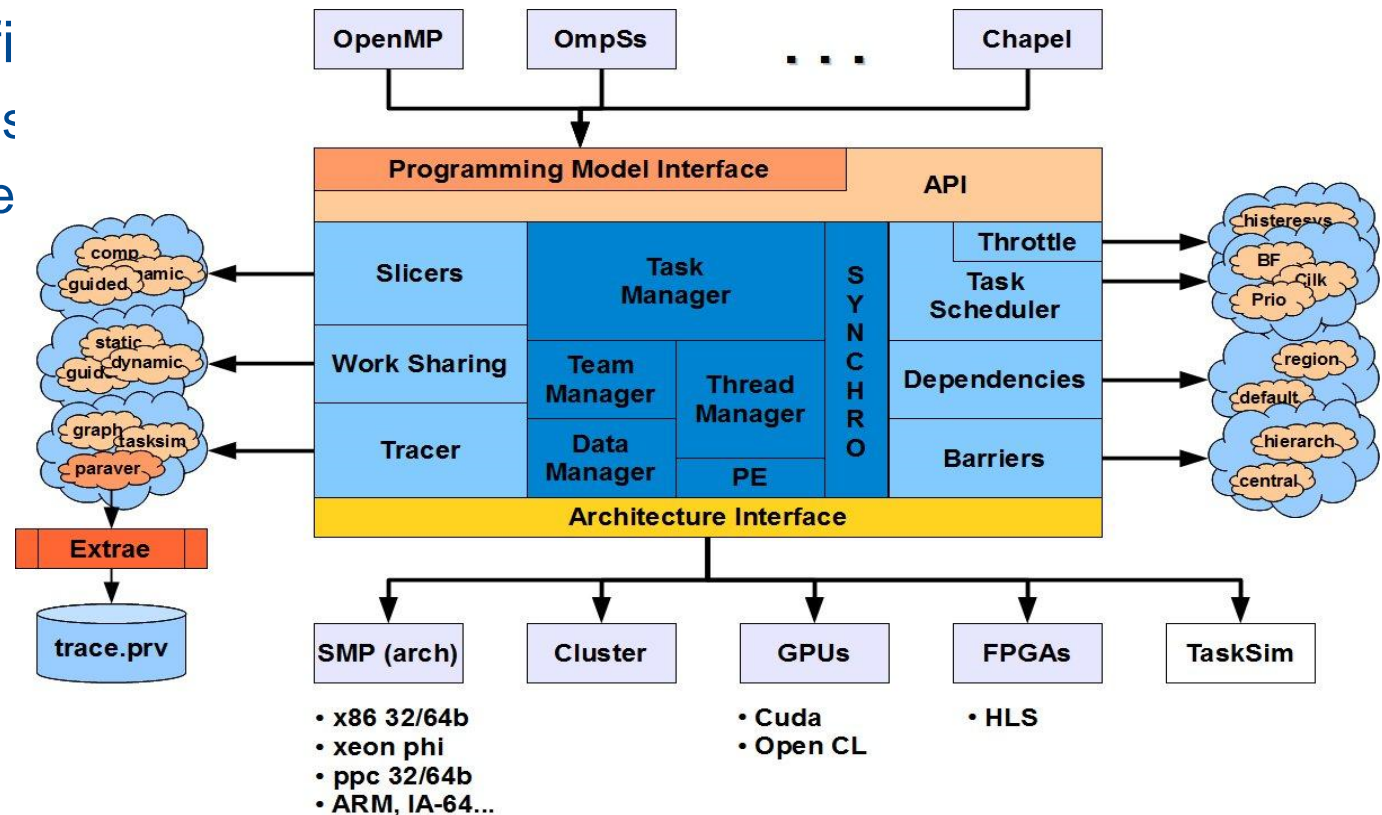
- Device-specific handlers
- May generate code in a separate file
- Invokes different back-end compilers
 - nvcc for NVIDIA



The NANOS++ Runtime

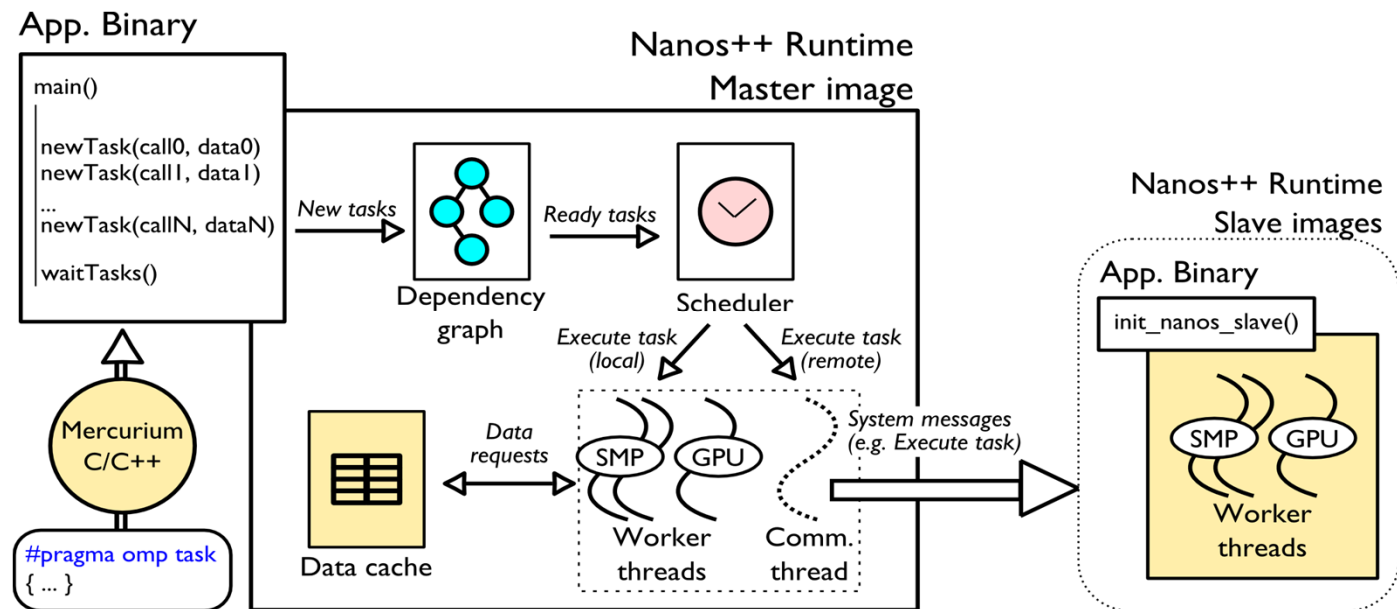
« Nanos++

- Common execution runtime (C, C++ and Fortran)
- Target specific features
- Task creation, dependency management, resilience, ...
- Task scheduling (BF, Cilk, Priority, Socket, ...)
- Data management: Unifi
 - Transparently manages
 - ... and data transfer be



Runtime structure behaviour: task handling

- Task generation
- Data dependence analysis
- Task scheduling



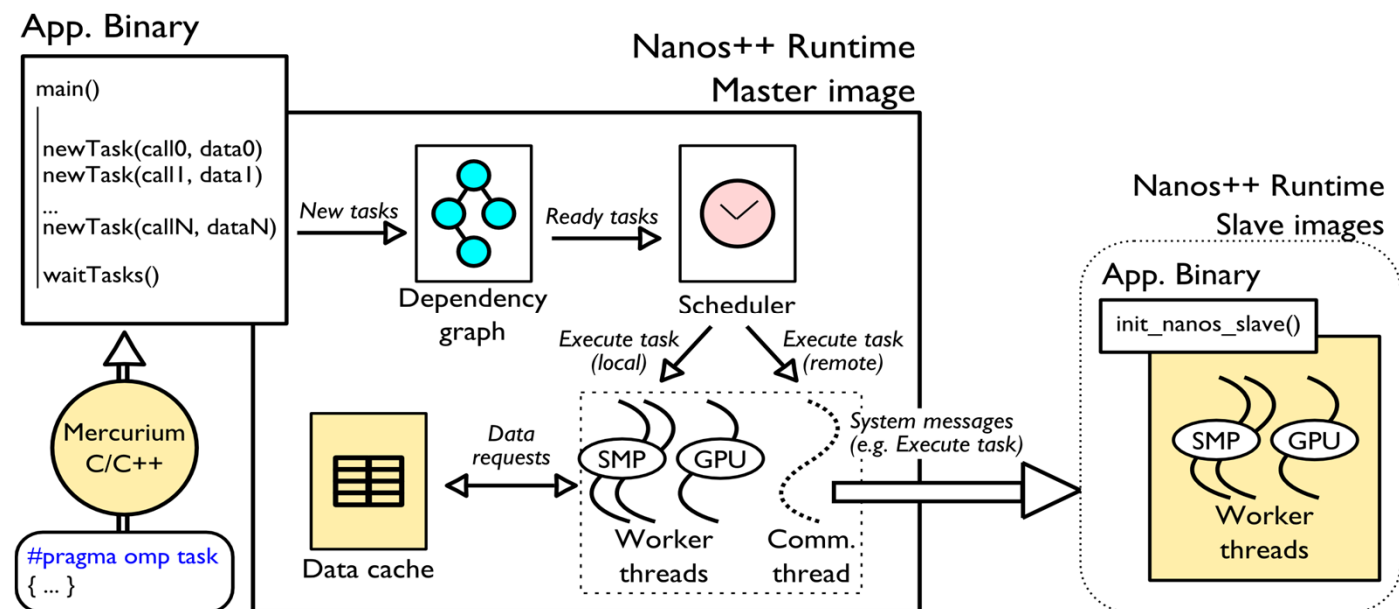
Runtime structure behaviour: coherence support

☞ Different address spaces managed with:

- A hierarchical directory
- A software cache per each:
 - Cluster node
 - GPU

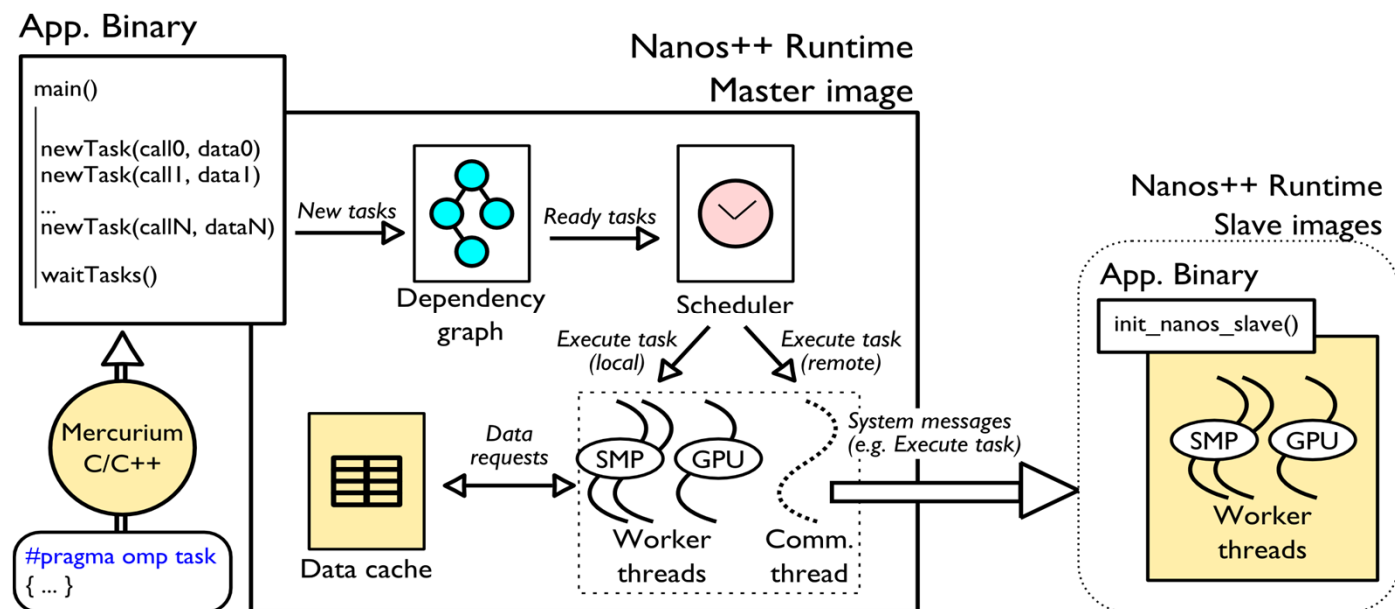
☞ Data transfers between different memory spaces only when needed

- Write-through
- Write-back



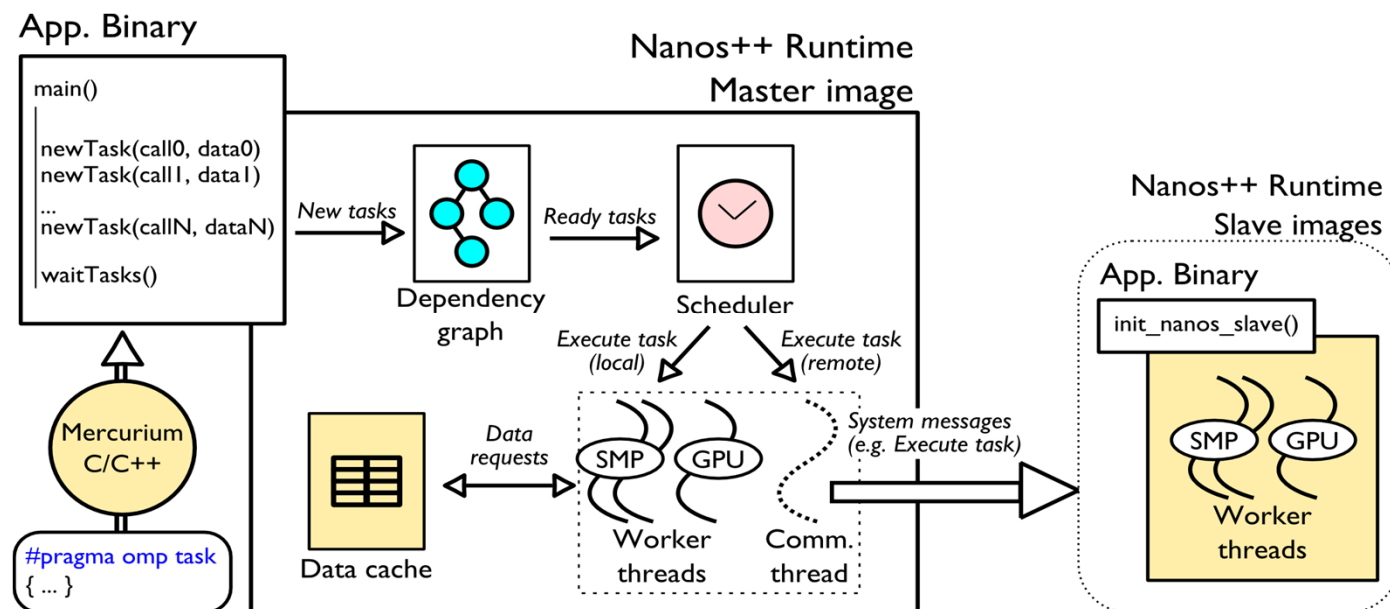
Runtime structure behaviour: GPUs

- Automatic handling of Multi-GPU execution
- Transparent data-management on GPU side (allocation, transfers, ...) and synchronization
- One manager thread in the host per GPU. Responsible for:
 - Transferring data from/to GPUs
 - Executing GPU tasks
 - Synchronization
- Overlap of computation and communication
- Data pre-fetch



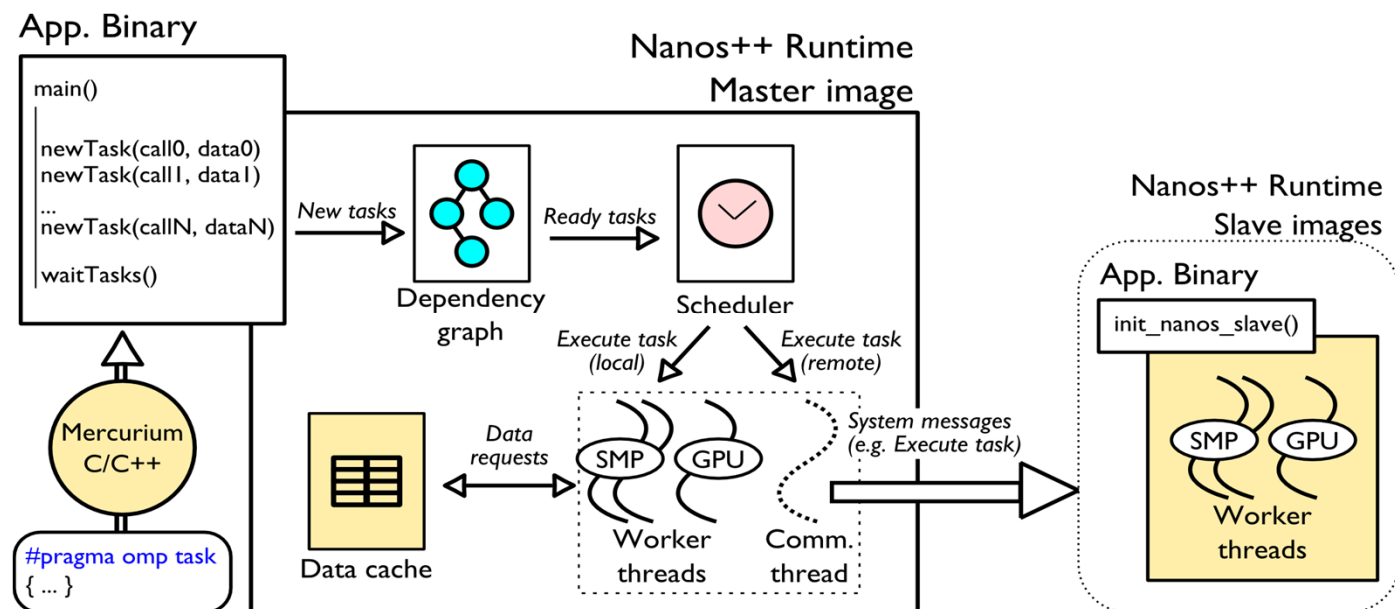
Runtime structure behaviour: clusters

- ❧ One runtime instance per node
 - One master image
 - N-1 slave images
- ❧ Low level communication through active messages
- ❧ Tasks generated by master
 - Tasks executed by worker threads in the master
 - Tasks delegated to slave nodes through the communication thread
- ❧ Remote task execution:
 - Data transfer (if necessary)
 - Overlap of computation with communication
 - Task execution
 - Local scheduler



Runtime structure behavior: clusters of GPUs

- Composes previous approaches
- Supports for heterogeneity and hierarchy:
 - Application with homogeneous tasks: SMP or GPU
 - Applications with heterogeneous tasks: SMP and GPU
 - Applications with hierarchical and heterogeneous tasks:
 - I.e., coarser grain SMP tasks
 - Internally generating GPU tasks



Compiling

« Compiling

```
mcc --ompss -c bin.c
```

« Linking

```
mcc --ompss -o bin bin.o
```

« where frontend can be:

mcc	C
mcxx	C++
mnvcc	CUDA & C
mnvcxx	CUDA & C++
mfc	Fortran

Compiling

« Compatibility flags:

- -l, -g, -L, -I, -E, -D, -W

« Other compilation flags:

-k	Keep intermediate files
--debug	Use Nanos++ debug version
--instrumentation	Use Nanos++ instrumentation version
--version	Show Mercurium version number
--verbose	Enable Mercurium verbose output
--Wp,flags	Pass flags to preprocessor (comma separated)
--Wn,flags	Pass flags to native compiler (comma separated)
--Wl,flags	Pass flags to linker (comma separated)
--help	To see many more options :-)

Executing

⌘ No LD_LIBRARY_PATH or LD_PRELOAD needed

```
./bin
```

⌘ Adjust number of threads with OMP_NUM_THREADS

```
OMP_NUM_THREADS=4 ./bin
```

Nanos++ options

- Other options can be passed to the Nanos++ runtime via `NX_ARGS`

```
NX_ARGS="options" ./bin
```

<code>--schedule=name</code>	Use name task scheduler
<code>--throttle=name</code>	Use name throttle-policy
<code>--throttle-limit=limit</code>	Limit of the throttle-policy (exact meaning depends on the policy)
<code>--instrumentation=name</code>	Use name instrumentation module
<code>--disable-yield</code>	Nanos++ won't yield threads when idle
<code>--spins=number</code>	Number of spin loops when idle
<code>--disable-binding</code>	Nanos++ won't bind threads to CPUs
<code>--binding-start=cpu</code>	First CPU where a thread will be bound
<code>--binding-stride=number</code>	Stride between bound CPUs

Nanox helper

« Nanos++ utility to

- list available modules:

```
nanox --list-modules
```

- list available options:

```
nanox --help
```

Tracing

« Compile and link with --instrument

```
mcc --ompss --instrument -c bin.c
```

```
mcc -o bin --ompss --instrument bin.o
```

« When executing specify which instrumentation module to use:

```
NX_INSTRUMENTATION=extrae ./bin
```

« Will generate trace files in executing directory

- 3 files: prv, pcf, rows
- Use paraver to analyze

Reporting problems

« Compiler problems

- <http://pm.bsc.es/projects/mcxx/newticket>

« Runtime problems

- <http://pm.bsc.es/projects/nanox/newticket>

« Support mail

- pm-tools@bsc.es

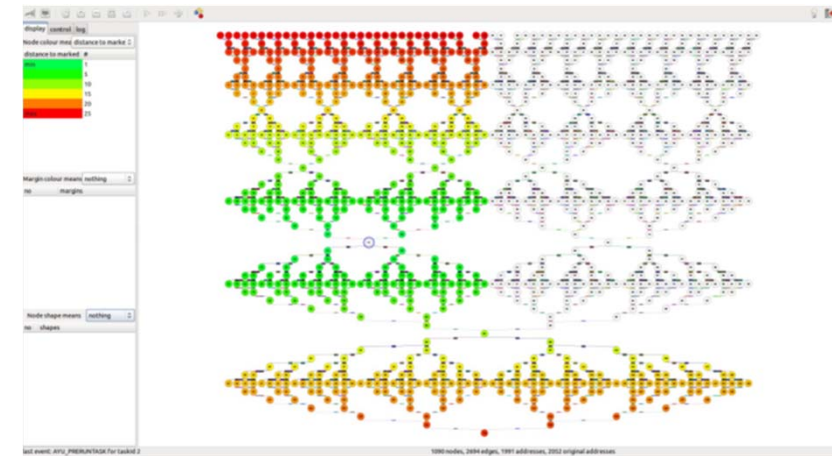
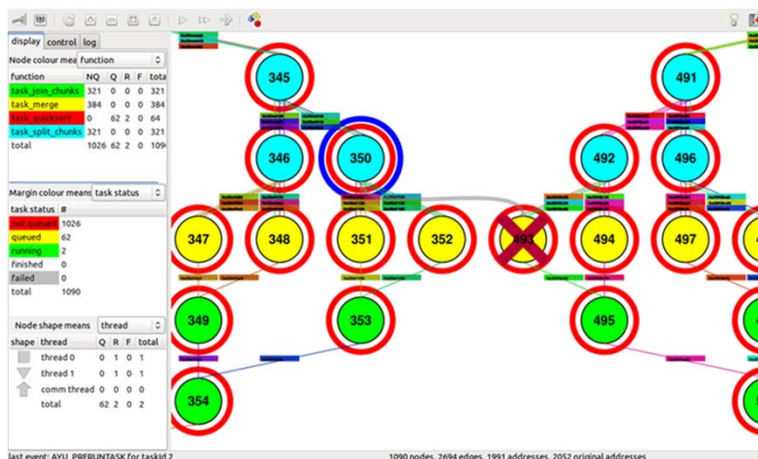
« Please include snapshot of the problem

Programming methodology

- ⌘ Correct sequential program
- ⌘ Finding tasks with Tareador
- ⌘ Debugging with Ayudame/Temanejo
- ⌘ Incremental taskification
 - Test every individual task with forced sequential in-order execution
 - → 1 thread, scheduler = FIFO, throttle=1
- ⌘ Single thread out-of-order execution
- ⌘ Increment number of threads
 - Use taskwaits to force certain levels of serialization

Debugging: AYUDAME/TEMANEJO

- « Leverage probe hooks provided by compiler and runtime
- « Task based debugging:
 - Display graph
 - Control execution environment (#threads,...)
 - Breakpoints at tasks
- « Interface to instruction level debugger (gdb)



Visualizing Paraver tracefiles

- « Set of Paraver configuration files ready for OmpSs. Organized in directories
 - **Tasks: related to application tasks**
 - Runtime, nanox-configs: related to OmpSs runtime internals
 - **Graph_and_scheduling: related to task-graph and task scheduling**
 - DataMgmt: related to data management
 - CUDA: specific to GPU

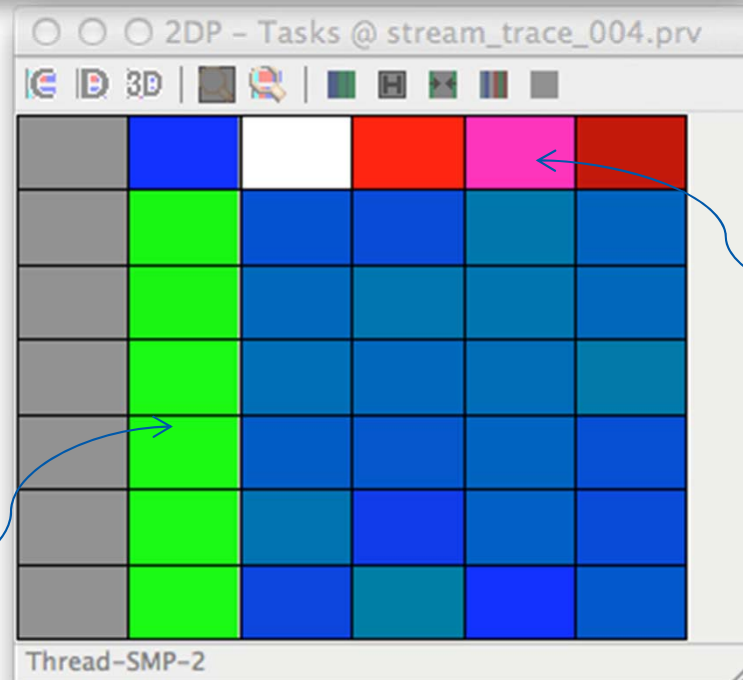
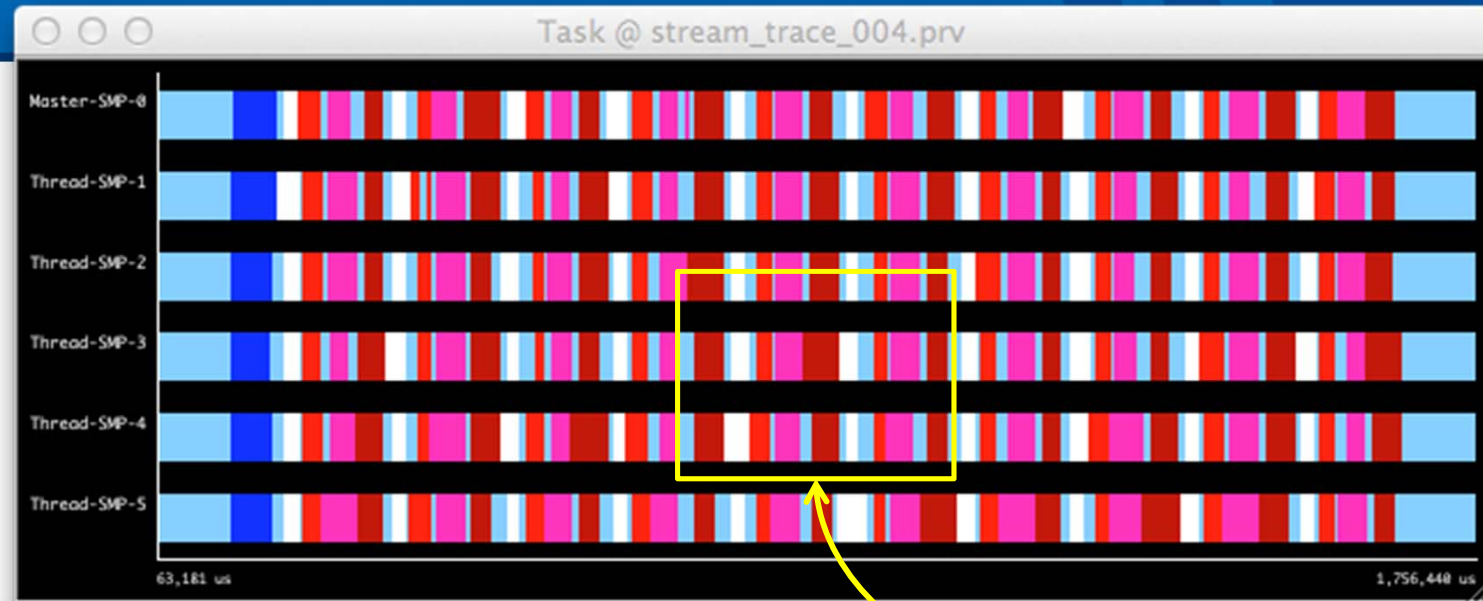
Tasks' profile

- « 2dp_tasks.cfg
- « Tasks' profile

control window:
timeline where each
color represent the
task been executed
by each thread

light blue: not executing
tasks

gradient color,
indicates given estadístic:
i.e., number of tasks instances



different colours
represent different
task type

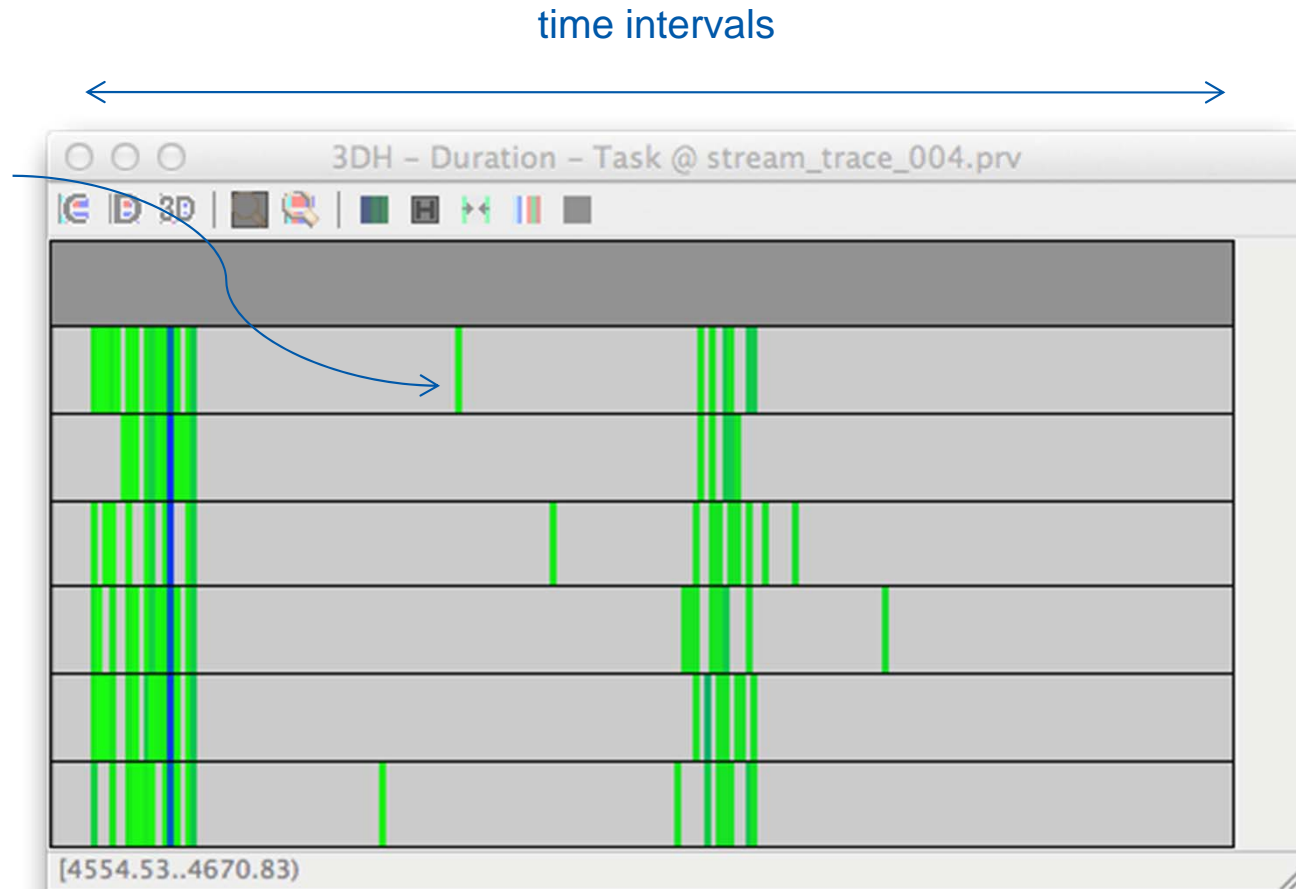
Tasks duration histogram

3dh_duration_task.cfg

gradient color,
indicates given estadístic:
i.e., number of tasks instances

threads

time intervals



Tasks duration histogram

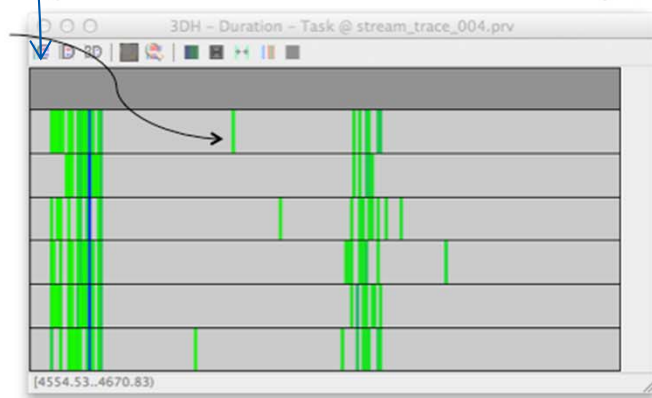
3dh_duration_task.cfg

control window:
task duration



gradient color,
indicates given estadistic:
i.e., number of tasks instances

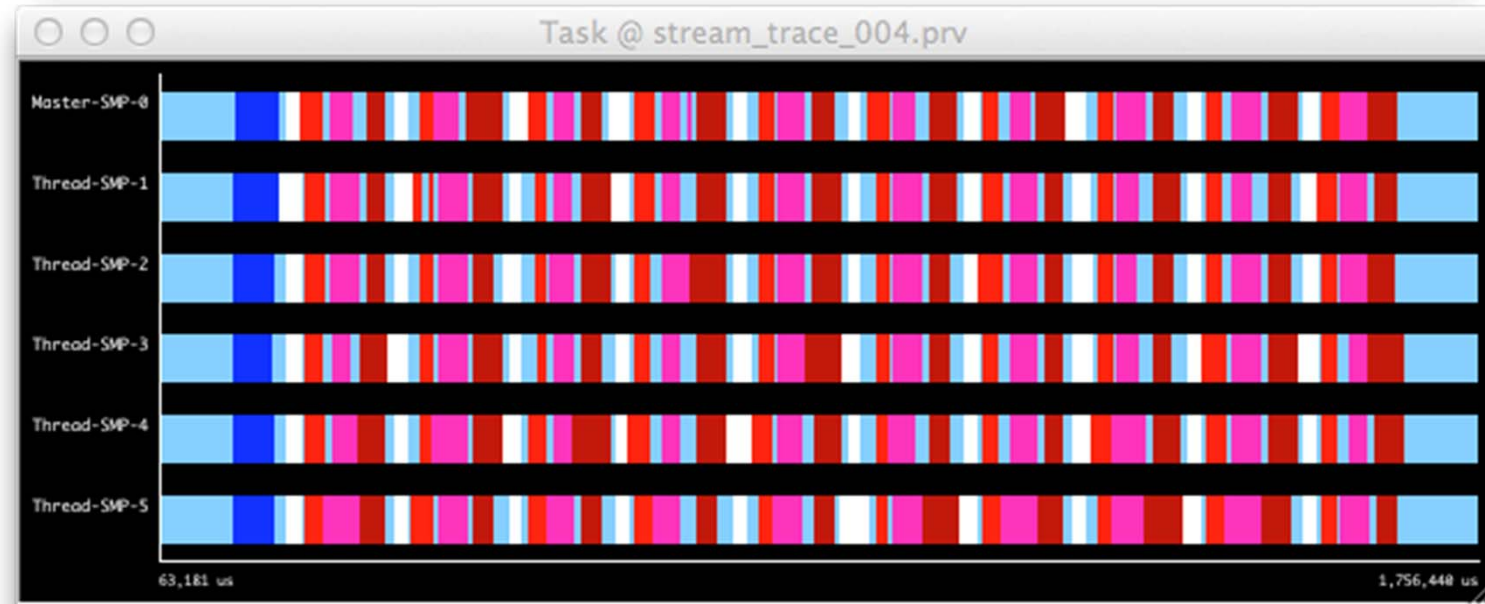
threads



Tasks duration histogram

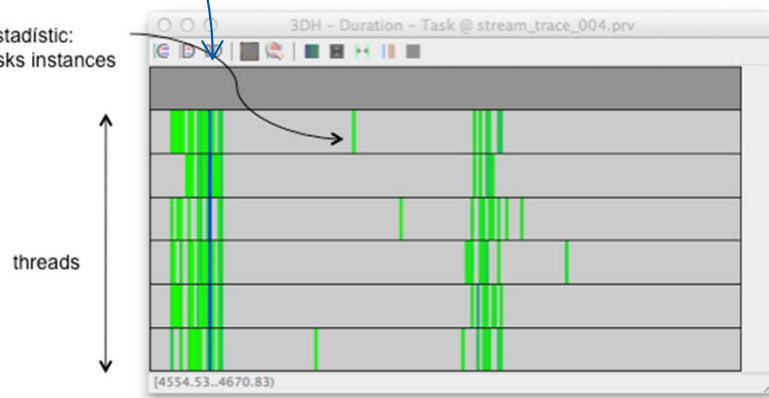
3dh_duration_task.cfg

3D window:
task type



time intervals

gradient color,
indicates given estadistic:
i.e., number of tasks instances



Tasks duration histogram

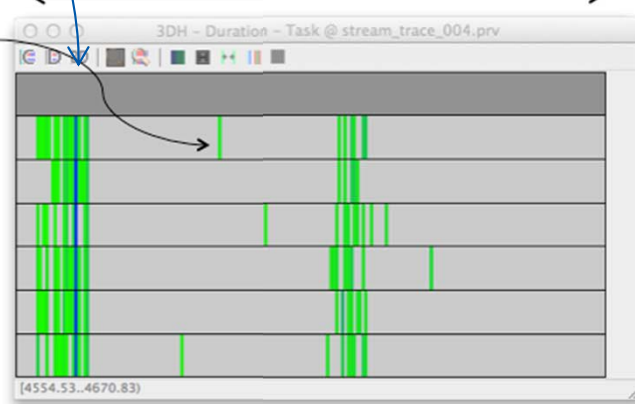
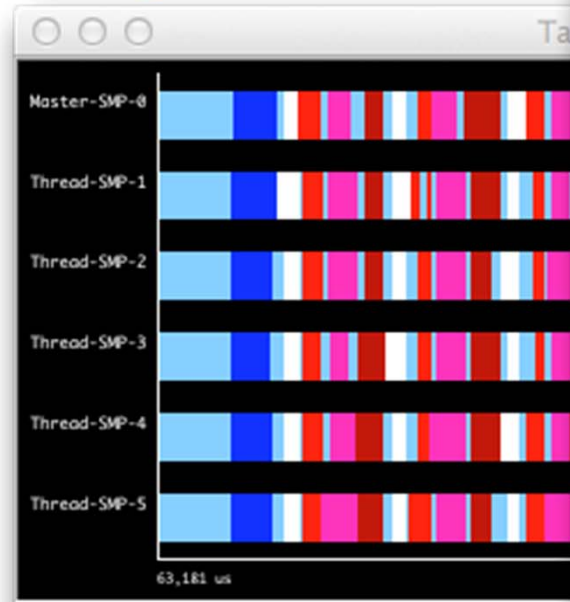
3dh_duration_task.cfg

3D window:
task type

gradient color,
indicates given estadistic:
i.e., number of tasks instances

threads

time intervals



Window browser

/Users/rosab/Documents/shared/trazas/judge/...

- Task
- 2DP - Tasks
- Task
- 2DP - Tasks
- Time btw. uf (task) events
- Task
- 3DH - Duration - Task**
- Task
- Time btw. uf (task) events
- Task
- 3DH - Duration - Task

Files & Window Properties

Statistics

Type	Semantic
Statistic	Time
Minimum Gradient	833.828
Maximum Gradient	112342.349

Data

3D

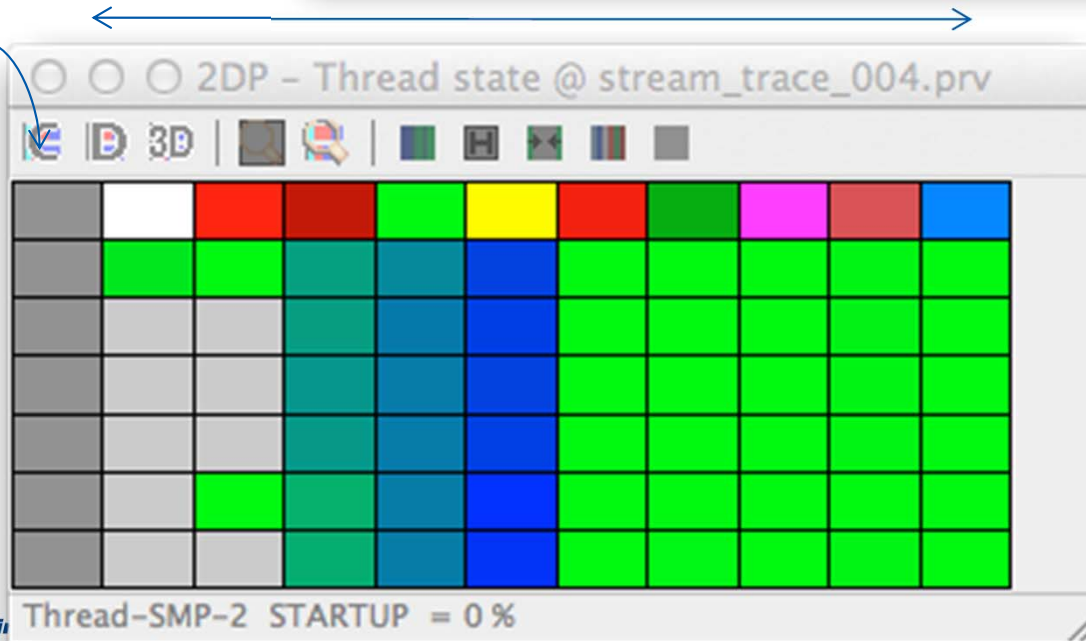
3rd Window	Task
Minimum	1
Maximum	5
Delta	1
Plane	Task 'triad_task'

chooser:
task type

Threads state profile

2dp_threads_state.cfg

control window:
timeline where each
color represent the
runtime state of each
thread





**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

Thank you!

For further information please contact
rosa.m.badia@bsc.es