



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

HPC Parallel Programming Models for Real-Time Control Systems

Eduardo Quiñones

{eduardo.quinones@bsc.es}

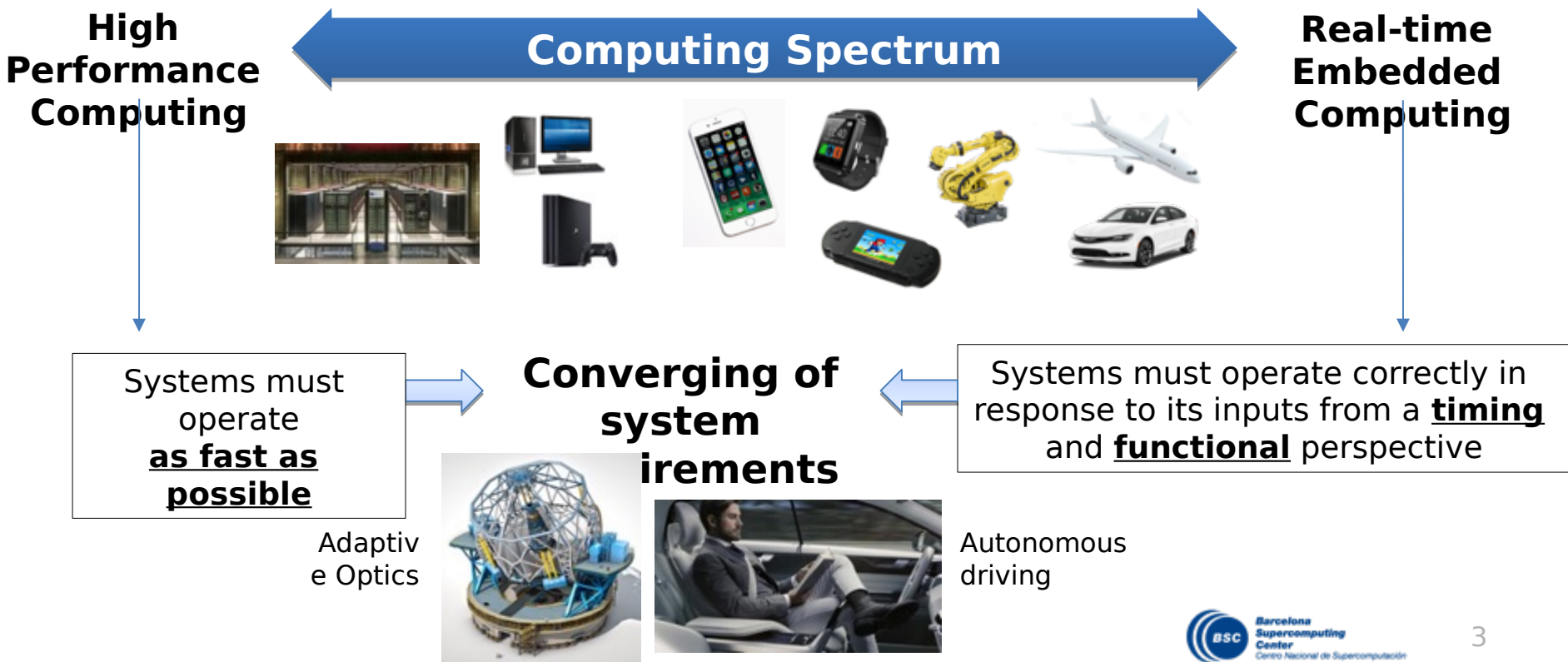
5th workshop on Real-time Control for Adaptive
Optics

Paris, 25 October 2018

Agenda

- The importance of parallel programming models
- The OpenMP parallel programming model
 - Timing guarantees and functional correctness in OpenMP
 - The real-time control (RTC) loop

Towards the Convergence of High-Performance and Real-Time Computing Domains



Towards the Convergence of High-Performance and Real-Time Computing Domains

- **Parallel computing** is key to cope with performance requirements

HPC Domain (~300W)



NVIDIA TitanV
(5120 CUDA cores)



Intel Xeon Phi KNL
(72-core fabric)

Embedded Domain (~10-15W)



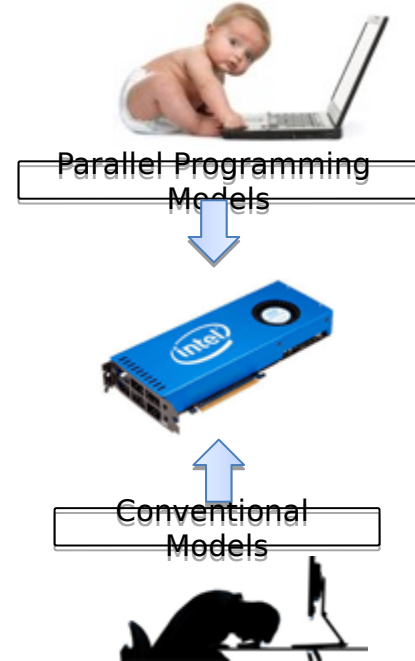
NVIDIA Tegra
(256 CUDA Cores)



Kalray MPPA
(256-core fabric)

Parallel Programming Models

- Mandatory to enhance **productivity**
 - **Programmability.** Provides an abstraction to express parallelism while hiding processor complexities
 - Defines parallel regions and synchronization mechanisms
 - **Portability/scalability.** Allows executing the same source code in different parallel platforms
 - **Performance.** Rely on run-time mechanisms to exploit the performance capabilities of parallel platforms
 - Including accelerators devices (e.g., FPGAs, GPUs, DSPs)



Parallel programming models comparison

Type	Language	✓ Strengths	✗ Weaknesses
<i>Hardware Centric</i>	Intel® TBB	<ul style="list-style-type: none"> - Highly tunable - High-level (task concept) 	<ul style="list-style-type: none"> - Portability - Mapping thread/core not part of the model
	NVIDIA® CUDA	<ul style="list-style-type: none"> - Highly tunable - Wrappers for many languages 	<ul style="list-style-type: none"> - Low level (explicit data management) - Restricted to NVIDIA GPUs
<i>Application Centric</i>	OpenCL	<ul style="list-style-type: none"> - Automatic vectorization - Executes in host and accelerator 	<ul style="list-style-type: none"> - Low level (explicit data management) - Full rewriting
	Pthreads	<ul style="list-style-type: none"> - Full execution control (thread concept) - Dynamic creation/destruction of threads 	<ul style="list-style-type: none"> - Low level (reductions, work distribution, synchronization, etc. by hand)

OpenMP

- **Mature language** constantly reviewed (last release Nov 2018, v5.0)
 - Defacto industrial standard in HPC shared memory processor architectures
- **Performance** and **efficiency**
 - Tantamount to other models (e.g. TBB, CUDA, OpenCL and MPI)
 - Support for fine-grain data- and task-parallelism
 - Features an advanced accelerator model for heterogeneous computing
- **Portability**
 - Supported by many chip and compiler vendors (Intel, IBM, ARM, NVIDIA, TI)
- **Programmability**
 - Currently available for C, C++ and Fortran (**#pragma omp**)
 - Allows incremental parallelization
 - Can be easily compiled sequentially (easing debugging)
- Very active research community

Principle behind OpenMP

- *Developers specify what the application does and not how it is done*
 - Computation is not fully controlled by the programmer but by the parallel framework

- **Complicates deriving timing analysis and functional correctness!**

OpenMP can guarantee:

- ***Time predictability.*** Reasoning about the timing behaviour of the parallel execution
- ***Correctness.*** Ensuring a correct operation in response to its inputs



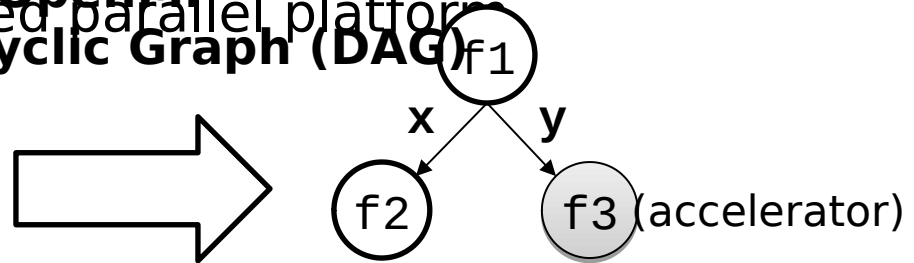
"I'm a software engineer, so I can confirm it works by magic."

Representation of Parallel Execution

- **Static analysis methods** to extract a complete representation of the parallel execution of OpenMP programs
 - Includes all the information for timing and functional correctness

```
#pragma omp parallel
#pragma omp single{
  int x,y;
  #pragma omp task depend(out:x,y)
  { f1(&x,&y); }
  #pragma omp task depend(in:x)
  { f2(x); }
  #pragma omp target map(to:y) depend(in:y)
  { f3(y); }
}
```

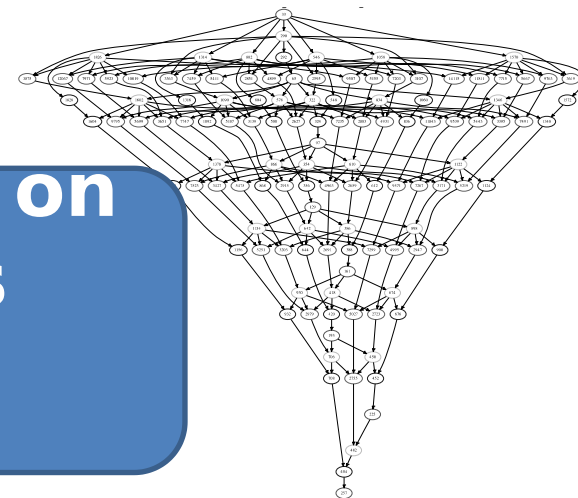
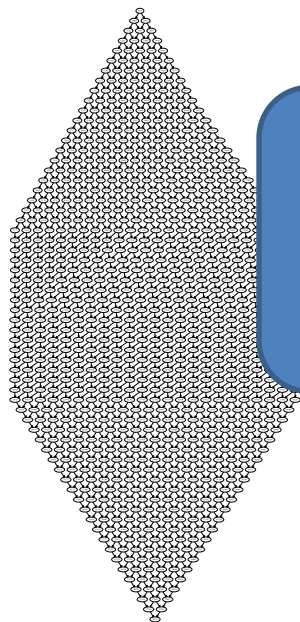
OpenMP
Direct Acyclic Graph (DAG)



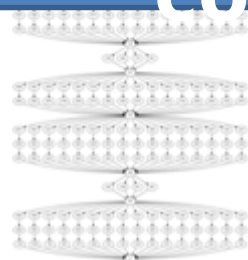
Platform Independent OpenMP-DAGs

Pedestrian detection (automotive) Infra-red sensor pre-processing (space) 3D Path Planning (avionics)

Cholesky

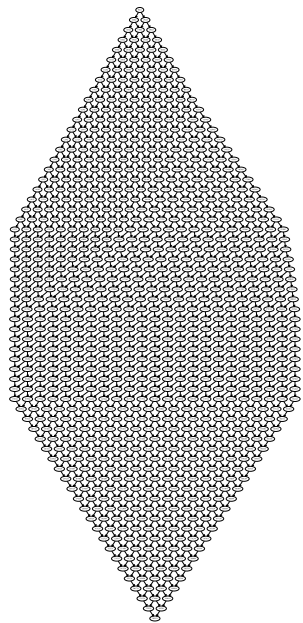


Provide guarantees on
1. Timing analysis
2. Functional correctness



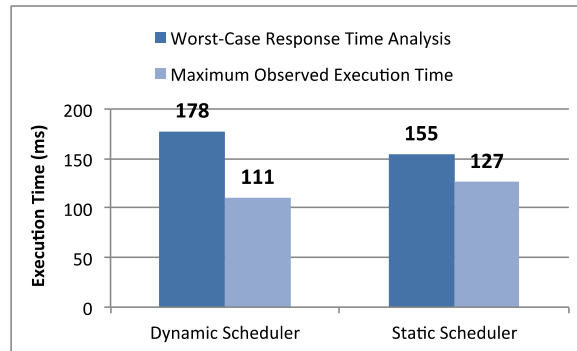
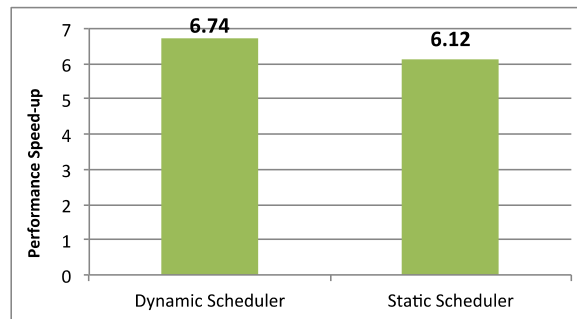
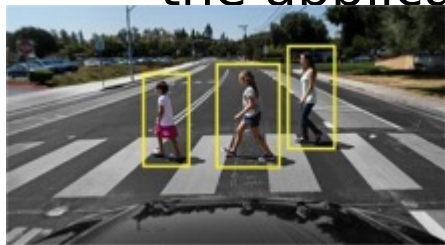
Timing Analysis

pedestrian detector
(automotive)



1. Dynamic and static allocation strategies of processor resources
2. Upper-bound the response time of the application

$$R_k^{ub} \leftarrow \frac{len(C_k)}{m} \cdot \frac{1}{cm(C_k)} \cdot cm(C_k) + I_k^{lp} + I_k^{lp}$$



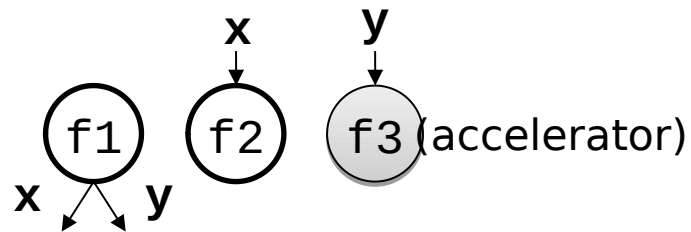
Platform dependent! (Intel 8-core processor)

Functional Correctness

- Race conditions

- Incorrect data scoping definition and usage of synchronization mechanisms

```
#pragma omp parallel {  
#pragma omp single{  
  int x,y;  
  #pragma omp task depend(out:x,y)  
  { f1(&x,&y); }  
  #pragma omp task depend(in:x)  
  { f2(x); }  
  #pragma omp target map(to:y) depend(in:y)  
  { f3(y); }  
}}
```

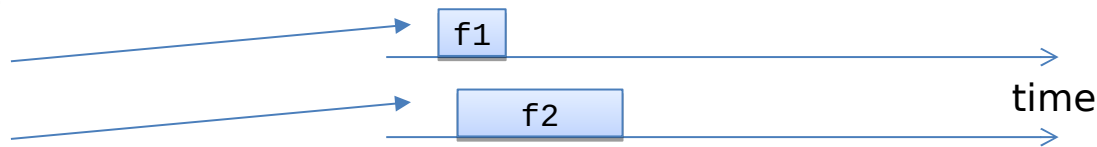


**f2 and f3 may use
incorrect values of x
and y!**

Real-time Control Loop (RTC)

- HPC Parallel programming models are time-agnostic

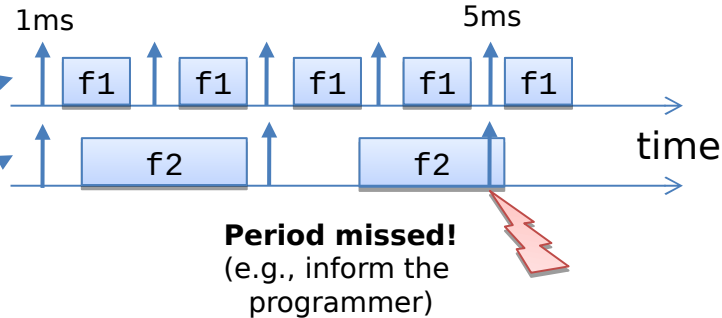
```
#pragma omp parallel {  
#pragma omp single {  
  int x,y;  
  #pragma omp task  
  { f1(&x); }  
  #pragma omp task  
  { f2(&y); }  
}}
```



Real-time Control Loop (RTC)

- Parallel programming models are time-agnostic

```
#pragma omp parallel {  
#pragma omp single{  
  int x,y;  
  while(1) {  
    if (period == 1ms){  
      #pragma omp task  
      { f1(&x); }  
    }  
    if period == 2ms)  
      #pragma omp task  
      { f2(&y); }  
  }  
}}
```



```
#pragma omp parallel {  
#pragma omp single{  
  int x,y;  
  #pragma omp task period(1ms)  
  { f1(&x); }  
  #pragma omp task period(2ms)  
  { f2(&y); }  
}}
```

Conclusions

1. Parallel programming models are fundamental for productivity in terms of programmability, portability and performance
 - They do not provide timing guarantees and functional correctness
2. OpenMP is a very convenient model for converging productivity with timing guarantees and functional correctness
 - Extracting a representation of the parallel execution (OpenMP-DAG)
 - Introducing a RTC loop within the parallel framework
3. **We are working with the language standardization committee to enhance OpenMP**



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

HPC Parallel Programming Models for Real-Time Control Systems

Eduardo Quiñones

{eduardo.quinones@bsc.es}

5th workshop on Real-time Control for Adaptive
Optics

Paris, 25 October 2018