



PPAM

2015

Krakow, Poland  
September 6-9, 2015

# Virtualizing **CUDA** enabled GPGPUs on **ARM** clusters

R. Montella<sup>1</sup>, G. Giunta<sup>1</sup>, G. Laccetti<sup>2</sup>, M. Lapegna<sup>2</sup>

C. Ferraro<sup>1</sup>, C. Palmieri<sup>1</sup>, V. Pelliccia<sup>1</sup>

<sup>1</sup>Department of Science and Technologies  
University of Napoli Parthenope

<http://hpsc.uniparthenope.it>

<sup>2</sup>Department of Mathematics and Applications  
University of Napoli Federico II

<http://dma.unina.it>



{ [raffaele.montella](mailto:raffaele.montella@uniparthenope.it), [giulio.giunta](mailto:giulio.giunta@uniparthenope.it),  
[carmine.ferraro](mailto:carmine.ferraro@uniparthenope.it), [carlo.palmieri](mailto:carlo.palmieri@uniparthenope.it), [valentina.pelliccia](mailto:valentina.pelliccia@uniparthenope.it) }  
[@uniparthenope.it](mailto:@uniparthenope.it)

{ [giuliano.laccetti](mailto:giuliano.laccetti@unina.it),  
[marco.lapegna](mailto:marco.lapegna@unina.it) }  
[@unina.it](mailto:@unina.it)



# Generic Virtualization Service

*(since March 2010)*

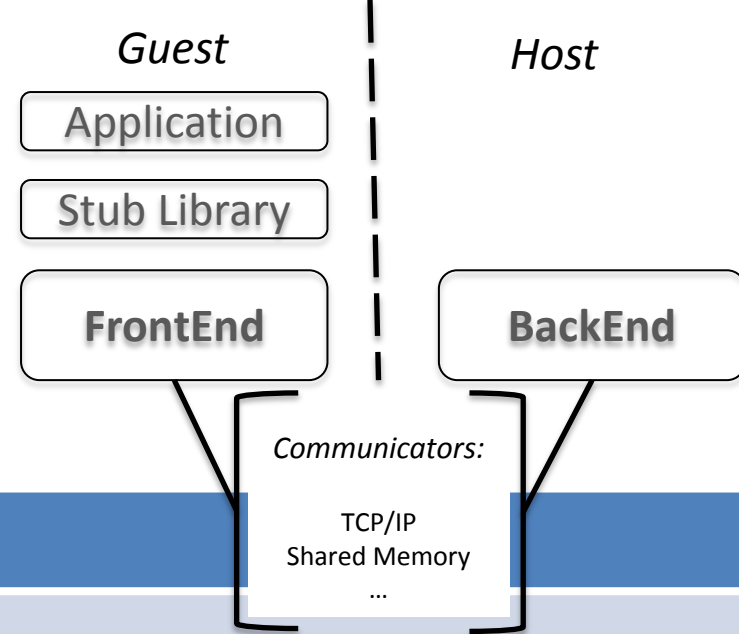
- Framework for split-driver based abstraction components
- Plug-in architecture
- **Independent from:**
  - Hypervisor (or no-hypervisor)
  - Communication
  - Target of virtualization
  - **Architecture!**
- **High performance:**
  - Enabling transparent virtualization
  - With overall performances better or not too far from un-virtualized resources

<http://hpsc.uniparthenope.it/>

**G**virtus

# The Communicator

- Provides a high performance communication between virtual machines and their hosts.
- The choice of the hypervisor deeply affects the efficiency of the communication.

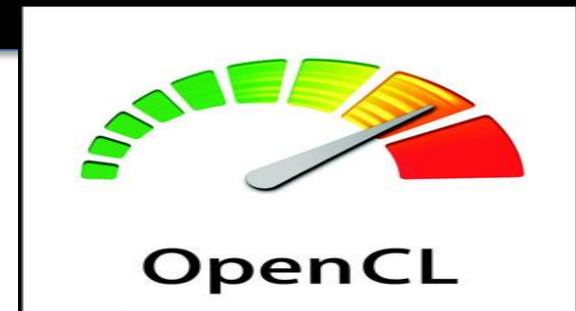


Hypervisor	FE/BE comm	Notes
No hypervisor	Unix Sockets	Used for testing purposes
Generic	TCP/IP	<ul style="list-style-type: none"> <li>• Communication testing purposes</li> <li>• <b>Remote / Distributed virtualized resources (i.e. GPUs)</b></li> <li>• <b>High Performance Internet of Things</b></li> </ul>
Xen	XenLoop	<ul style="list-style-type: none"> <li>• runs directly on the top of the hardware through a custom Linux kernel</li> <li>• provides a communication library between guest and host machines</li> <li>• implements low latency and wide bandwidth TCP/IP and UDP connections</li> <li>• app transparent and offers an automatic discovery of the supported VMs</li> </ul>
VMware	Virtual Machine Communication Interface (VMCI)	<ul style="list-style-type: none"> <li>• commercial hypervisor running at the application level</li> <li>• provides a datagram API to exchange small messages</li> <li>• a shared memory API to share data</li> <li>• an access control API to control which resources a virtual machine can access</li> <li>• and a discovery service for publishing and retrieving resources</li> </ul>

# An application: Virtualizing gpGPUs

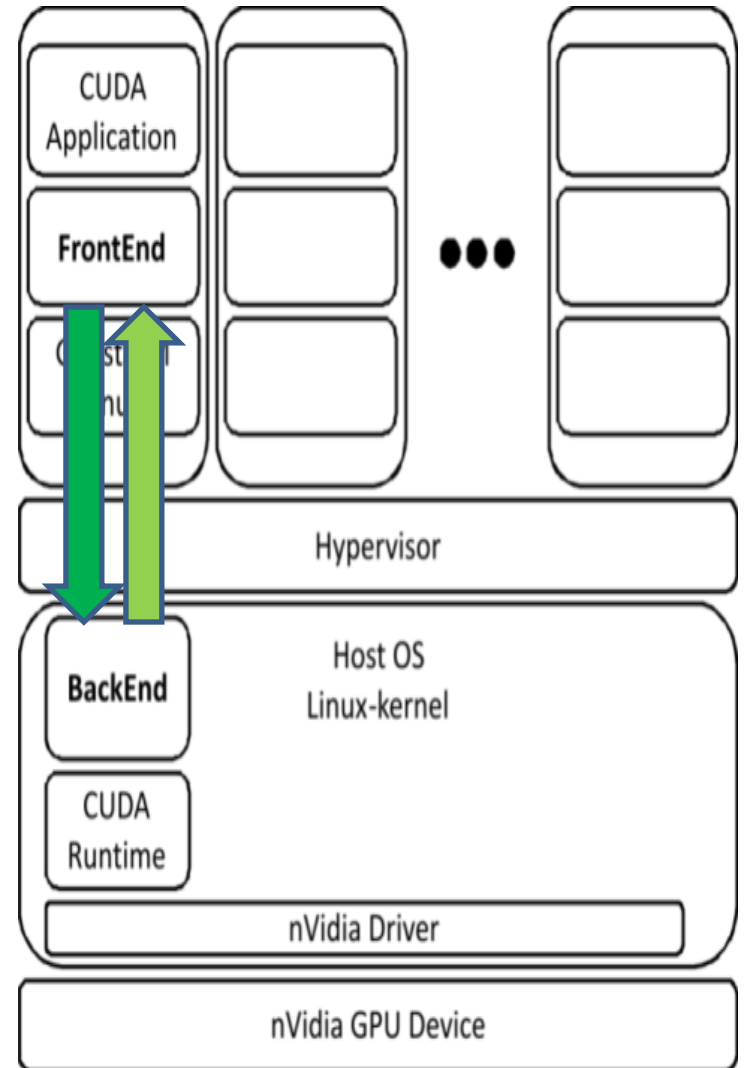
- GPUs

- Hypervisor independent
- Communicator independent
- GPU independent
- Programming model independent



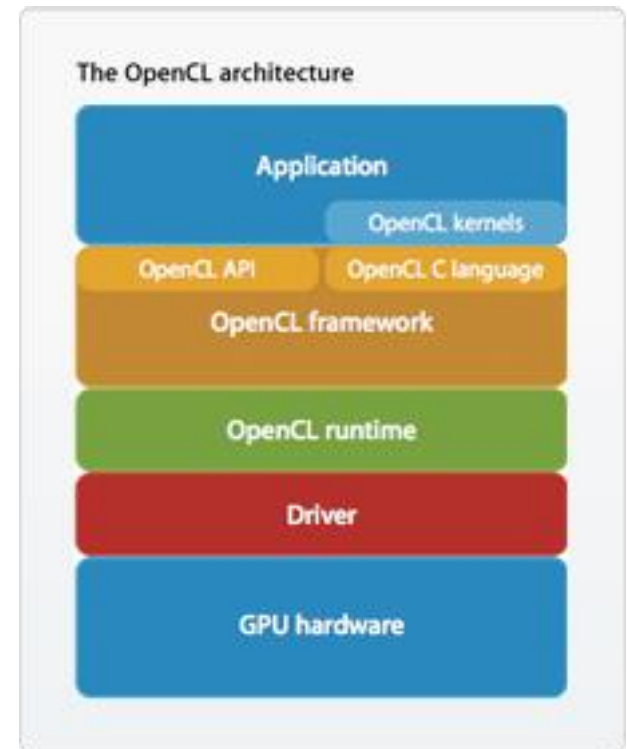
# GVirtuS – pre OpenCL

- Full support to
  - CUDA 3.x drivers
  - CUDA 3.x runtime
- Partially supporting (more work is needed)
  - OpenGL integration
- **Limitations:**  
No support after CUDA 3.x because NVIDIA issues



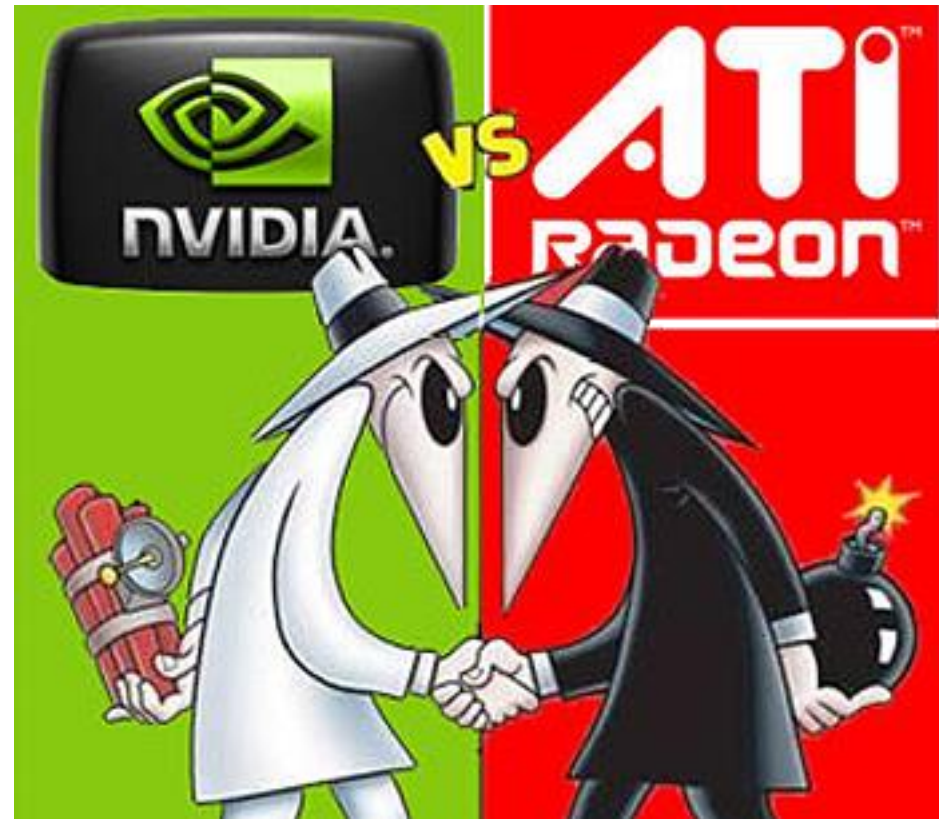
# GVirtuS - OpenCL

- Currently ongoing!
- Benchmarks available:
  - Matrix Multiply
  - Vector Product
  - Histogram
- Why OpenCL instead of CUDA:
  - Open Platform
  - Hardware independent ( CPUs, GPUs )
  - Widely supported



# But the world evolved differently...

- OpenCL didn't behaved as expected...
- ...people uses CUDA...
- ...we have to be back to the CUDA plugin.





# How it works

0x00f9ab15

0x00f9ab16

0x00f9ab17

0x00f9ab18

0x0023FF72

0x0023FF73

0x0023FF74

0x0023FF75

```
char *p=  
my_malloc(S);
```

```
char *my_malloc(  
long SIZE) {  
  
Function *f=  
new Function(  
"my_malloc");  
  
f->pushLong(SIZE);  
  
f->execute();  
}  
  
return f->pullPtr();
```

Application

GVirtuS Front-End

GVirtuS  
Communicator

```
char *my_malloc(  
Function *f) {  
  
long size=f-  
>pullLong();  
char *p=malloc(size);  
}  
  
f->pushPtr(p);  
}
```

GVirtuS Back-End

Drivers

p=0x0023FF74

- Simplified example: allocate memory on the host



# How it works

0x0023FF73

0x0023FF74

0x0023FF75

0x0023FF76

0x0023FF77

p=0x0023FF74

&result=0x0023FF94

```
success=
my_to_host(
p,"string to
copy"
);
```

```
int *my_to_host(
char *p, char *src) {

Function *f=
new Function(
"my_to_host");

f->pushPtr(p);
f->pushString(src);
f->execute();
}

return f->pullInt();
```

```
char *my_to_host(
Function *f) {

char *p=f->pullPtr();
char *src=pullString();

int
result=to_host(p,src);

f->pushInt()
}
```

Application

GVirtuS Front-End

GVirtuS  
Communicator

GVirtuS Back-End

Drivers

success=0x00

result=0x00

- Simplified example: accessing host memory

~~memcpy(p, "string to copy", 15)~~

# GVirtuS and CUDA!

Consumer – The CUDA app

```
#include <stdio.h>
#include <cuda.h>
int main(void) {
    int n;
    cudaGetDeviceCount(&n);
    printf("Number of CUDA GPU(s): %d\n", n);
    return 0;
}
```

**GVirtuS Frontend**

```
cudaError_t cudaGetDeviceCount(int *count) {
    Frontend *f = Frontend::GetFrontend();
    f->AddHostPointerForArguments(count);
    f->Execute("cudaGetDeviceCount");
    if(f->Success())
        *count =
            *(f->GetOutputHostPointer<int>());
    return f->GetExitCode();
}
```

Producer – The CUDA device host

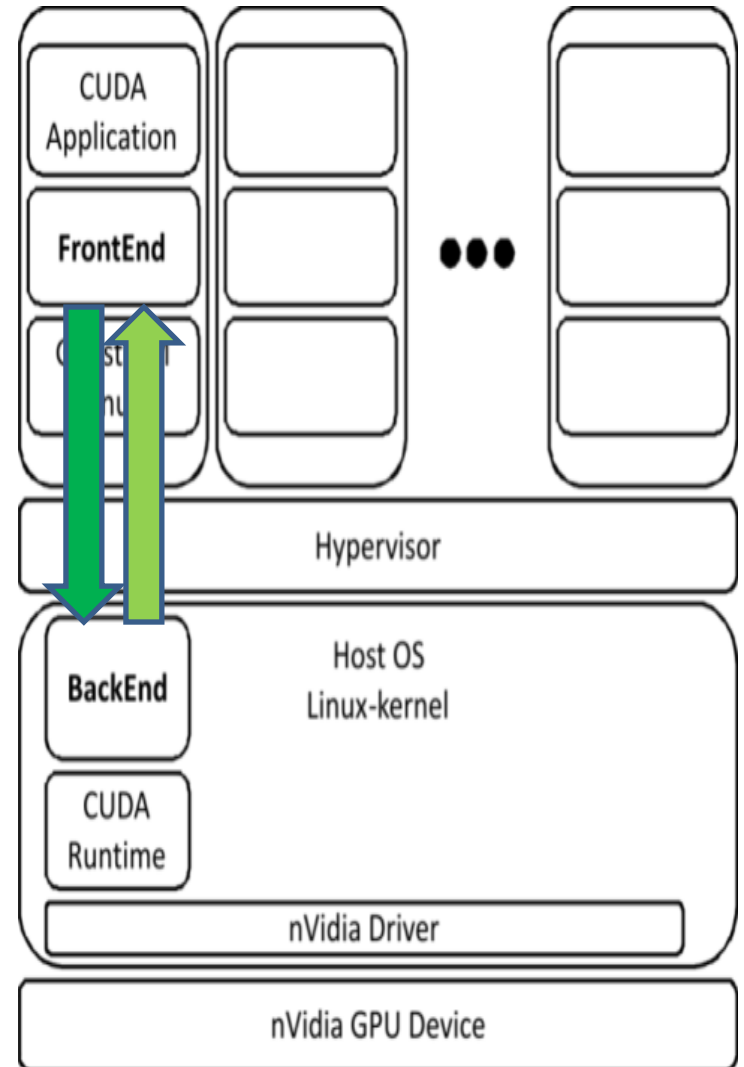
**GVirtuS Backend**

**Process Handler**

```
Result *handleGetDeviceCount(
    CudaRtHandler * pThis,
    Buffer *input_buffer) {
    int *count = input_buffer->Assign<int>();
    cudaError_t exit_code;
    exit_code = cudaGetDeviceCount(count);
    Buffer *out = new Buffer();
    out->Add(count);
    return new Result(exit_code, out);
}
```

# GVirtuS – **post** OpenCL

- On going support to
  - CUDA 6.x drivers
  - CUDA 6.x runtime
- Partially supporting (more work is needed)
  - OpenGL integration
- **Limitations:**  
For now many, but challenging!



# Scenarios and prototypal applications

- Development workstation/1
  - i7-940 2\*NVIDIA Tesla C1060
- Development workstation/2
  - i7-940 2\*NVIDIA Titan X
- The AWS EC2 GPU machine
  - g2.2xlarge 1\*NVIDIA K250
- ARM based computing node/1
  - Udo0 Cortex A9 32bit - nogpu
- ARM based computing node/2
  - NVIDIA Jetson Cortex A15 64bit – gpu/nogpu



# Development Workstation

asyncAPI	WORKS
cdpSimplePrint	WORKS
cdpSimpleQuicksort	requires GPU devices with compute SM 3.5 or higher
clock	WORKS
cplIntegration	WORKS
cppOverload	requires GPU devices with compute SM 2.0 or higher
cudaOpenMP	ERRORS
inlinePTX	WORKS
matrixMul	WORKS
matrixMulCUBLAS	Comparing CUBLAS Matrix Multiply with CPU results: FAIL
matrixMulDrv	WORKS
simpleAssert	requires a GPU with compute capability 2.0 or later
simpleAtomicIntrinsics	WORKS
simpleCallback	ERRORS
simpleCubemapTexture	requires SM 2.0 or higher
simpleIPC	ERRORS
simpleLayeredTexture	simpleLayeredTexture requires SM >= 2.0
simpleMPI	WORKS

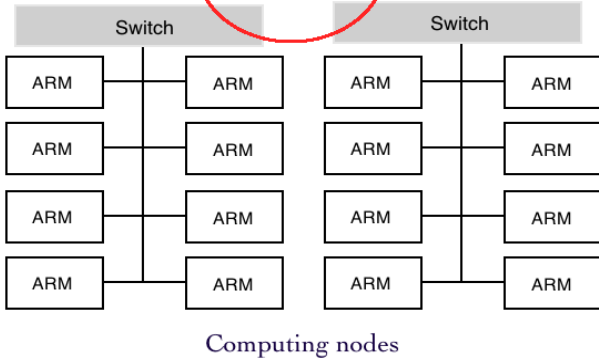
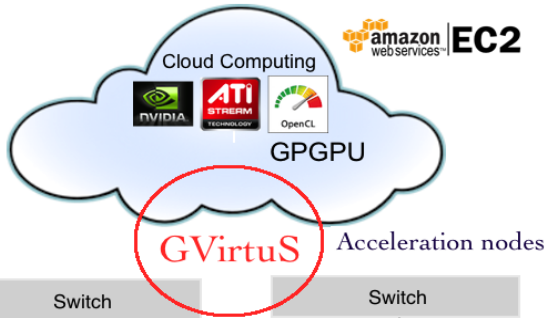
simpleMultiCopy	WORKS
simpleMultiGPU	WORKS
simpleOccupancy	cudaOccupancyMaxActiveBlocksPerMultiprocessor
simpleP2P	Two or more GPUs with SM 2.0 or higher capability are required
simplePitchLinearTexture	*** Error: cudaMemcpy2D() not yet implemented!
simplePrintf	WORKS
simpleSeparateCompilation	requires a GPU with compute capability 2.0
simpleStreams	WORKS
simpleSurfaceWrite	__cudaRegisterSurface
simpleTemplates	WORKS
simpleTexture	simpleTexture completed, returned ERROR!
simpleTextureDrv	WORKS
simpleVoteIntrinsics	WORKS
simpleZeroCopy	Device 0 does not support mapping CPU host memory!
template	WORKS
template_runtime	WORKS
UnifiedMemoryStreams	requires Compute Capability of SM 3.0
vectorAdd	WORKS
vectorAddDrv	WORKS

- Current GVirtuS CUDA 6.5 results running 0\_Simple SDK examples.

Test	With GVirtuS	Without GVirtuS
Matrix Multiplication	0.139 sec	0.092 sec
Vector Addition	0.063 sec	0.059 sec
Sorting Networks	8.787 sec	8.539 sec

- Backend and Frontend on x86\_64 (Tcplp)

# AWS EC2 K250 GPU



Test	With GVirtuS	Without GVirtuS
Matrix Multiplication	38.236 sec	0.098 sec
Vector Addition	3.298 sec	0.057 sec
Sorting Networks	2min24.648 sec	8.482 sec

Backend: EC2 K250  
Frontend: X86\_64

Test	With GVirtuS	Without GVirtuS
Matrix Multiplication	50.607 sec	N/A (Can't run)
Vector Addition	3.368 sec	N/A (Can't run)
Sorting Networks	4min17.547 sec	N/A (Can't run)

Backend: EC2 K250  
Frontend: ARM Cortex A9

- Elastic GPU sharing.



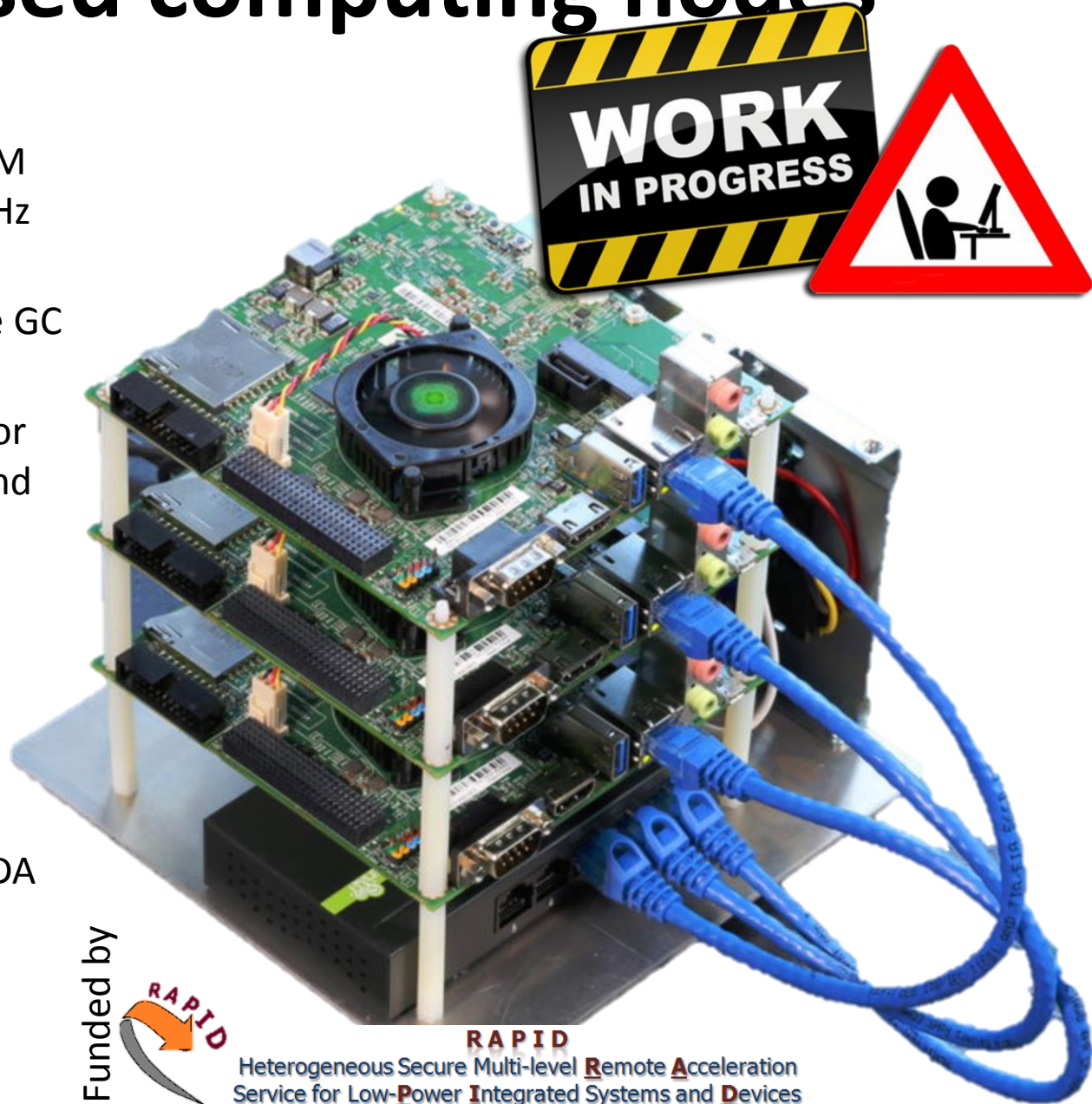
# ARM based computing nodes

- UDOO

- CPU Freescale i.MX 6 ARM Cortex-A9 Quad core 1GHz
- GPU Vivante GC 2000 + Vivante GC 355 + Vivante GC 320
- Integrated accelerators for 2D, OpenGL® ES2.0 3D and OpenVG™
- RAM DDR3 1GB

- NVIDIA Jetson

- Tegra K1 SOC
- Kepler GPU with 192 CUDA cores
- 4-Plus-1 quad-core ARM Cortex A15 CPU
- 2 GB x16 memory 64 bit



Funded by



**RAPID**  
Heterogeneous Secure Multi-level Remote Acceleration  
Service for Low-Power Integrated Systems and Devices

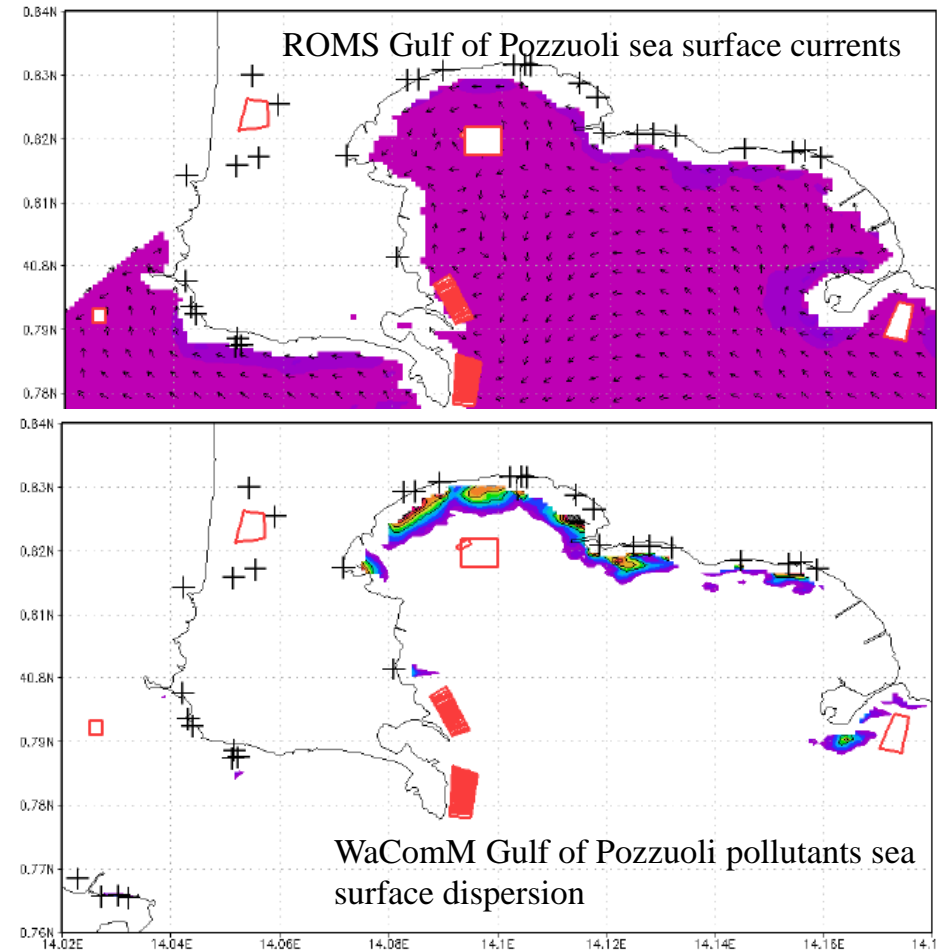


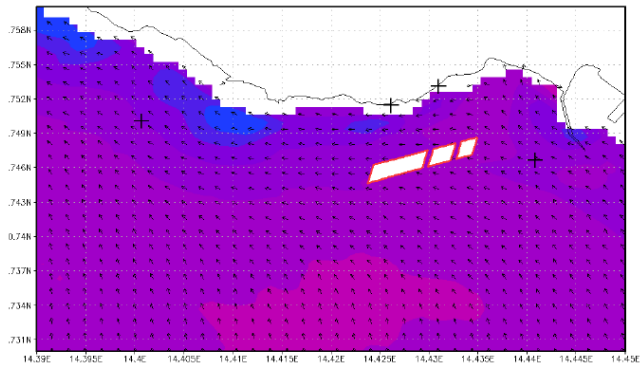
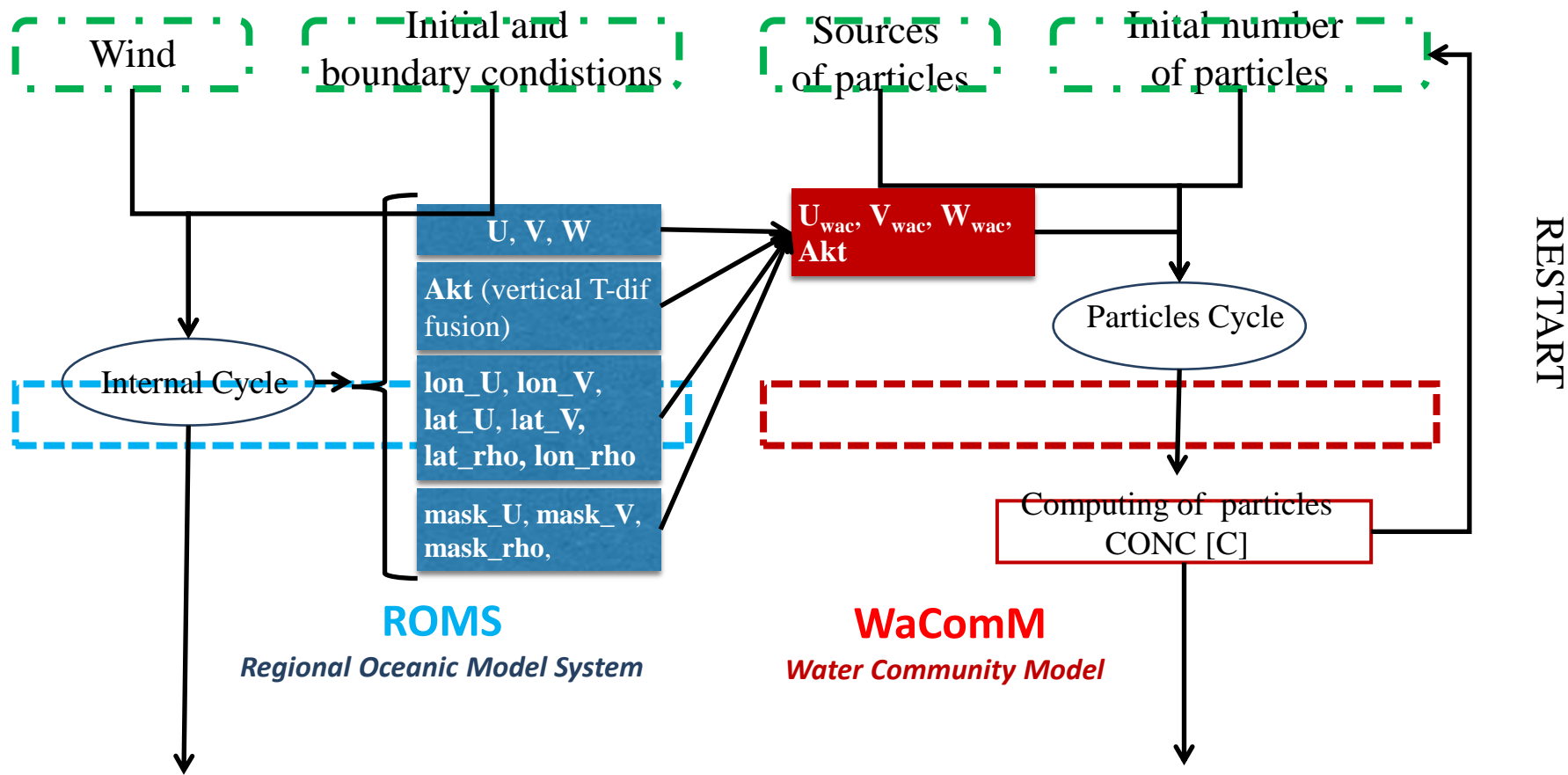
# A real life application:

## Coupling EulerianLagrangian models for offshore and coastal pollution tracing

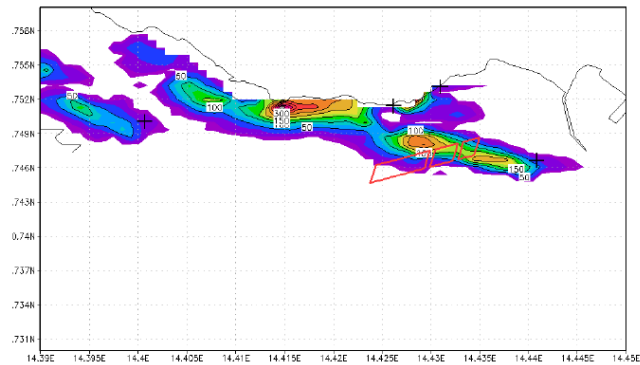
### WaComM (*Water Community Model*)

- Lagrangian 3D model
- simulate the transport and dispersion of pollutants spilled out into offshore ocean currents.
- Evolution of the *Lagrangian Assessment for Marine Pollution 3D* (LAMP3D) model.
- Restarts
- Shared memory / CUDA Hybrid parallel
- Hybrid approach (Eulerian-Lagrangian models)
- Coastal dynamics of the Bay of Naples
- Forecast the impact of pollutants spilled out from both natural and anthropic sources in high density areas of mussel culture.



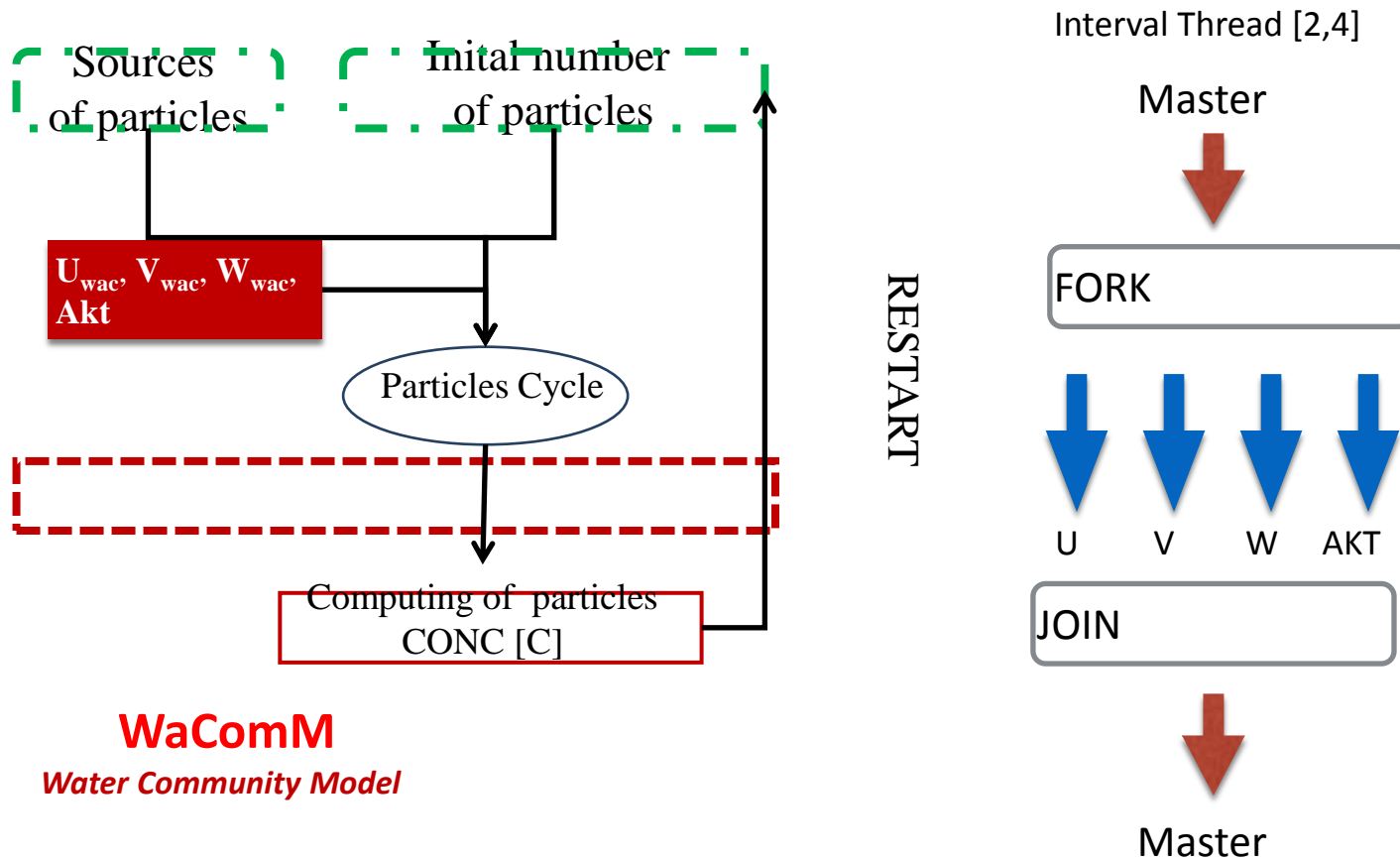


ROMS Torre del Greco (NA) sea surface currents



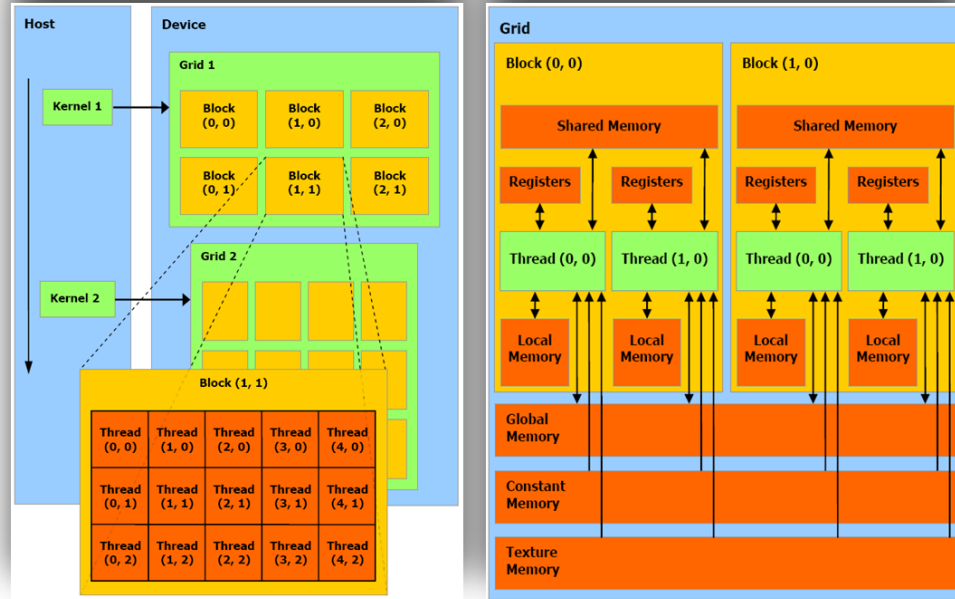
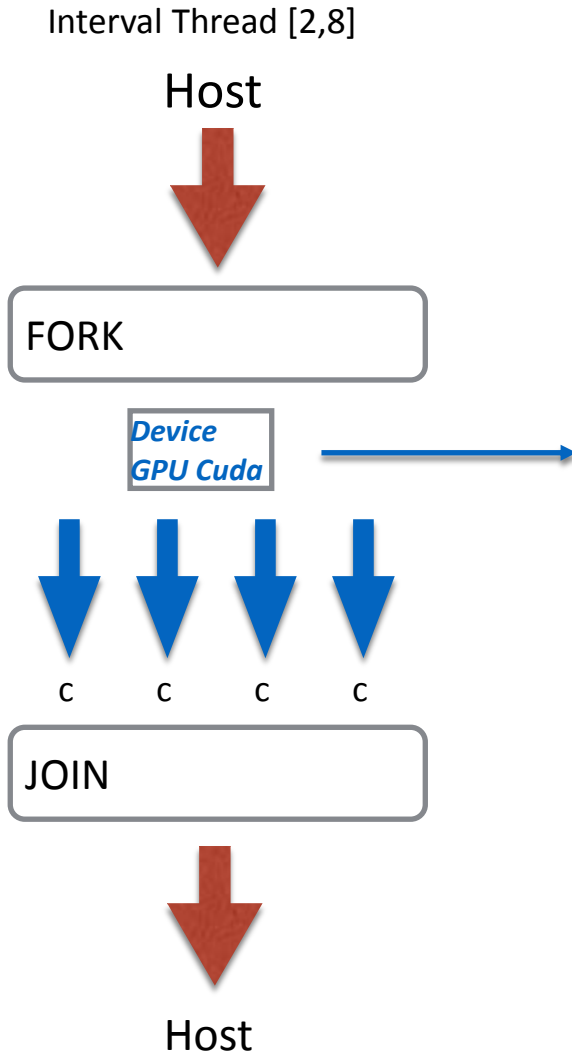
WaComM Torre del Greco (NA) pollutants sea surface dispersion

# Parallel WaComM



Parallel computations on data structures  $u$   $v$   $w$  and  $akt$ . The results obtained will be managed for the calculation of the variable  $conc$  ( $c$ ).

# Parallel Wacomm + GPU Cuda



Funded by

Modeling, forecasting and mapping of pollution by escherichia coli and salmonella in coastal areas dedicated to mussel production.

project.

IZS ME 04/12 RC/C78C120017001

# Current GVirtuS limitations

- The CUDA plugin is **not thread-safe** at the consumer/front-end side
- The management of the **pinned memory** and **unified memory** will be a serious challenge

# Conclusions

- GVirtuS works with CUDA 6.x
- Clusters of (*inexpensive*) ARM computing nodes can share one or more high-end GPUs hosted on x86\_64 machines.
- Clusters of **x86\_64** or **ARM** machines can consume GPUs available on the **cloud**
- Virtual clusters of x86\_64 machines could share **virtual GPUs**

# Future directions





# Get, test and extend

- High Performance Scientific Computing Smart-Lab:

<http://hpsc.uniparthenope.it>

- The source code public repository:

<https://bitbucket.org/montella/gvirtus-dist>

The logo for HPSC SMART LAB features the text 'HPSC' in a large, stylized font, followed by a circular icon containing a globe and the text 'SMART LAB' in a smaller font.A banner for HPSC SMART LAB with a background of binary code. The text 'HPSC SMART LAB' is in a bold, sans-serif font, and 'HIGH PERFORMANCE SCIENTIFIC COMPUTING SMART LABORATORY' is in a smaller font below it.