

CAPS

## Agenda

- GPGPU & Hybrid Parallel computing
- HMPP Concepts and Overview
- Starting with HMPP
- Addressing hardware accelerators with HMPP
- HMPP Toolchain
- HMPP Runtime
- Managing Data with HMPP
- Grouping Codelets
- Sharing Data with HMPP
- HMPP Features & Roadmap

MAISON DE LA SIMULATION 2012

[www.caps-entreprise.com](http://www.caps-entreprise.com) 2

The slide features a background graphic of a grid of grey squares that curves and recedes into the distance, creating a 3D effect. The CAPS logo is in the top right corner, and the text 'Agenda' is in the top left. A large, semi-transparent watermark 'MAISON DE LA SIMULATION 2012' is oriented diagonally across the center. At the bottom, the website 'www.caps-entreprise.com' and the number '2' are visible.

CAPS

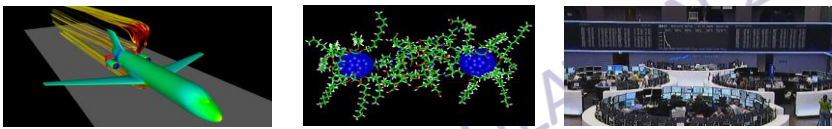
## GPGPU & Parallel Hybrid Computing

www.caps-entreprise.com 3

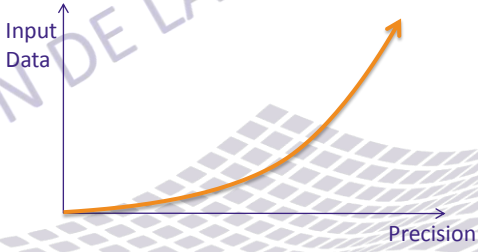
CAPS

## Industry and Business Facts

- Modeling & Simulation are pervasive



- Precision is the key to success



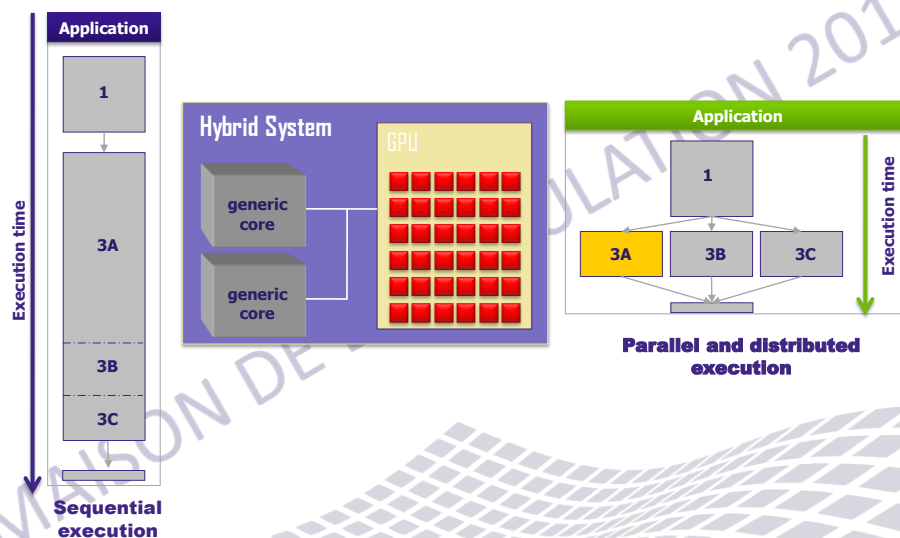
www.caps-entreprise.com 4

## Why Hybrid Computing?



- More precision implies more data to process
  - But still in the same time, so more speed in computations!
- Current technologies reached a limit (in terms of frequency)
- One solution is to use parallelism (increase # of cores)
  - We do more things in parallel instead of doing it quicker
- The efficiency scale evolves according economic issues
  - Performance / Power consumption
  - Performance / Development time
  - ...
- Mainstream applications will rely on these multicore / manycore architectures
- Various heterogeneous hardware
  - General purpose cores
  - Application specific cores (e.g. GPUs)

## You said Parallel Hybrid Application?



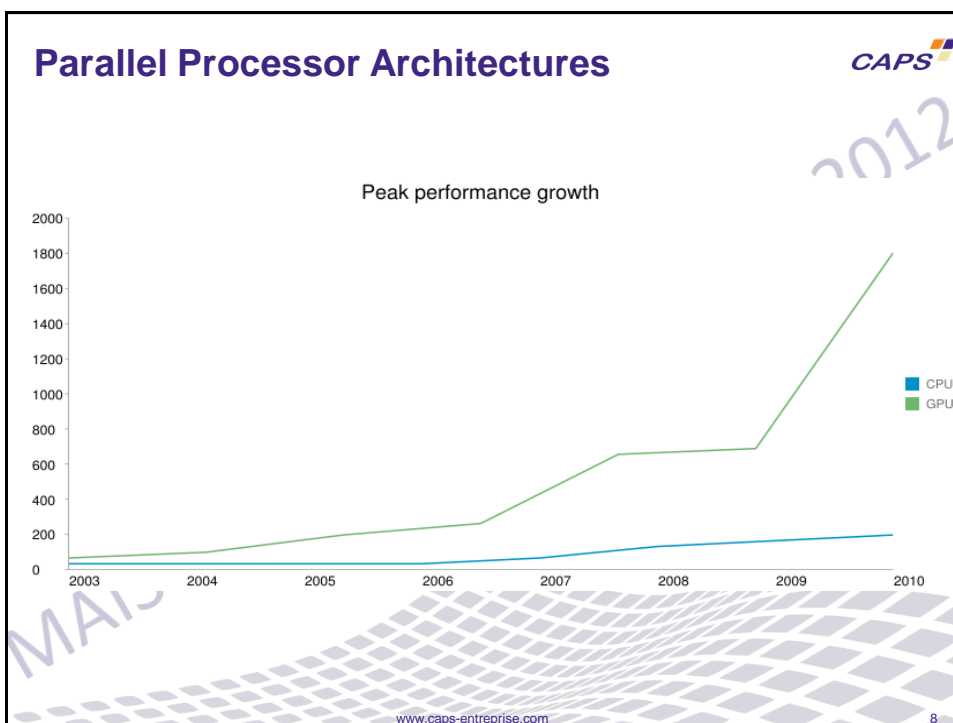
## Multiple Parallelism Levels

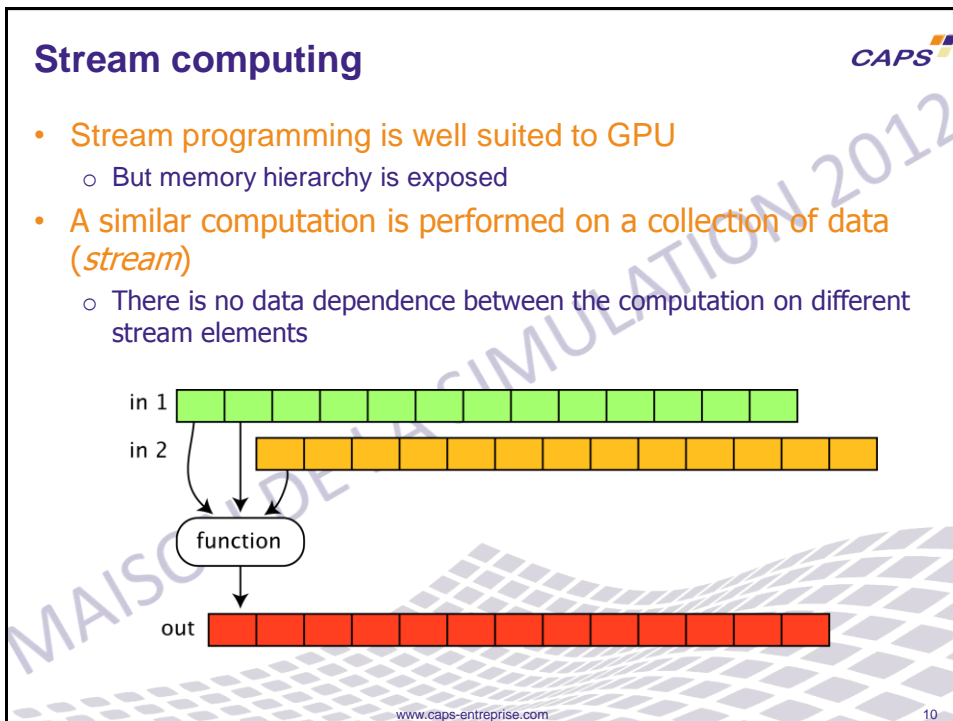
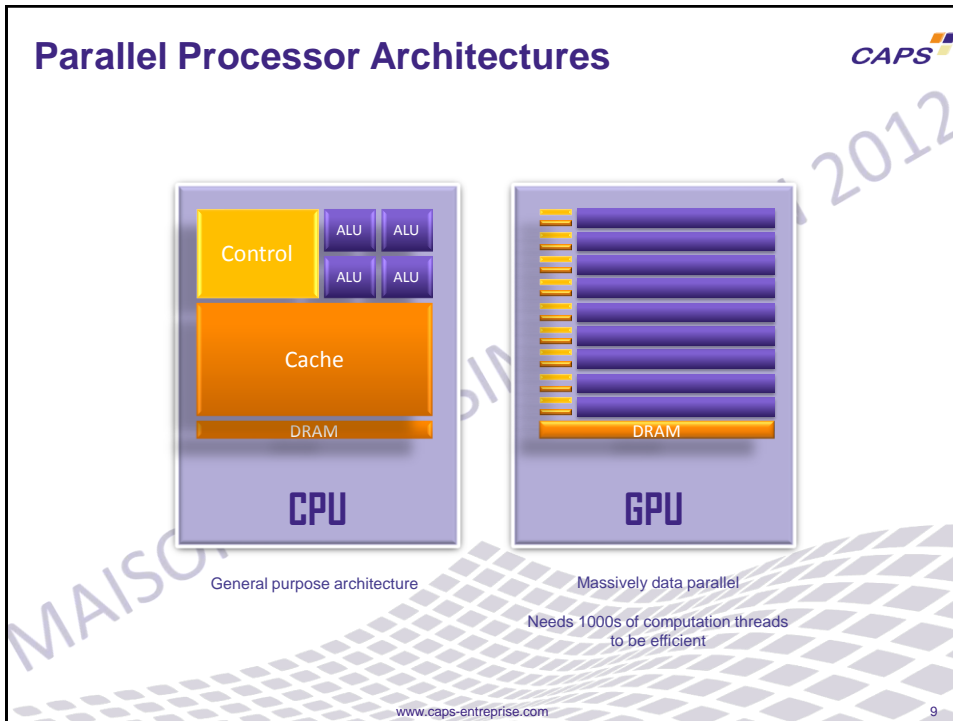
CAPS

- Adding a new layer of specific hardware is adding a new workload to the developer
- Programming various hardware components of a node cannot be done separately

www.caps-entreprise.com

7





## Current GPUs for GPU Computing



- **NVIDIA GeForce 200, for instance GTX 285**
  - Max device memory size is 4 GB
  - Peak FP SP is 1 teraflops
  - Peak FP DP is 0.1 teraflops
  - Memory bandwidth is 159 GB/s
  - Power consumption 204 Watts
  - CUDA and OpenCL
- **AMD Firestream, for instance Radeon HD 5870**
  - Max device memory size is 2 GB
  - Peak FP SP is 2,5 teraflops
  - Peak FP DP is 0.5 teraflops
  - Memory bandwidth is 153 GB/s
  - Power consumption 228 Watts
  - OpenCL and CAL/IL

## Offloading computations



- **Host: General purpose cores**
  - Share a main memory
  - Core ISA provides fast SIMD instructions
- **Device: Streaming cores**
  - GPU, DSP, FPGA... (vector, SIMD)
  - Application specific architectures ("narrowband")
  - Can be extremely fast
- **Hundreds of GigaOps**
  - But not easy to leverage
  - Restriction to one platform is not acceptable

## Key Issues

CAPS

- Using accelerators makes you stick to a proprietary language/environment/technology
- Huge potential performance but accelerators are far from host memory
  - Data must be copied on the remote device
  - Due to narrowband links between CPU/HWA, data transfers are critical
- Therefore rethink the computation organization/algorithm

The diagram illustrates the data flow between a HOST and a DEVICE. The HOST is represented by a purple box containing 'code' and 'memory'. The DEVICE is represented by a green box containing 'code' and 'memory'. Two horizontal arrows with question marks connect the 'code' sections and the 'memory' sections of the HOST and DEVICE, indicating the critical nature of data transfers between them.

www.caps-entreprise.com 13

## Manycore Challenges

CAPS

- Programming
  - Medium
- Resources management
  - Medium
- Application deployment
  - Hard
- Portable performance
  - Extremely hard

The graph shows the progression of manycore challenges over time. The x-axis represents time (t) with markers for GPGPU 2004, 2008, and End user 2009-2010. The y-axis represents the difficulty of the challenges. Four blue bars represent the challenges: Programming, Ressource Allocation, Deployment, and Portable performance. A red arrow labeled 'Move with the hardware' points upwards and to the right, indicating that as hardware evolves, these challenges become more difficult.

www.caps-entreprise.com 14



## Why hybrid/heterogeneous may be unavailable



- Share some goals with embedded systems that lead to the same technological issues in HPC
  - Optimizing energy consumption leads to specialized architectures
- No common programming API
  - APIs always make some underlying architecture assumptions
  - No low level programming API common to all devices
  - An API mainly addresses a specific hardware component, as a consequence we need many

www.caps-entreprise.com

15

**Accelerator Programming Model**      Parallelization      CAPS

Directive-based programming      GPGPU      Manycore programming

Hybrid Manycore Programming      HPC community

Petaflops      Parallel computing      HPC open standard

Multicore programming      Exaflops      NVIDIA CUDA

Code speedup      Hardware accelerators programming

High Performance Computing

Parallel programming interface

Massively parallel

Open CL

<http://www.caps-entreprise.com>  
<http://twitter.com/CAPSentreprise>  
<http://www.openhmpp.org>