# Parallel and Heterogeneous Computing

**Trond Hagen, SINTEF ICT**
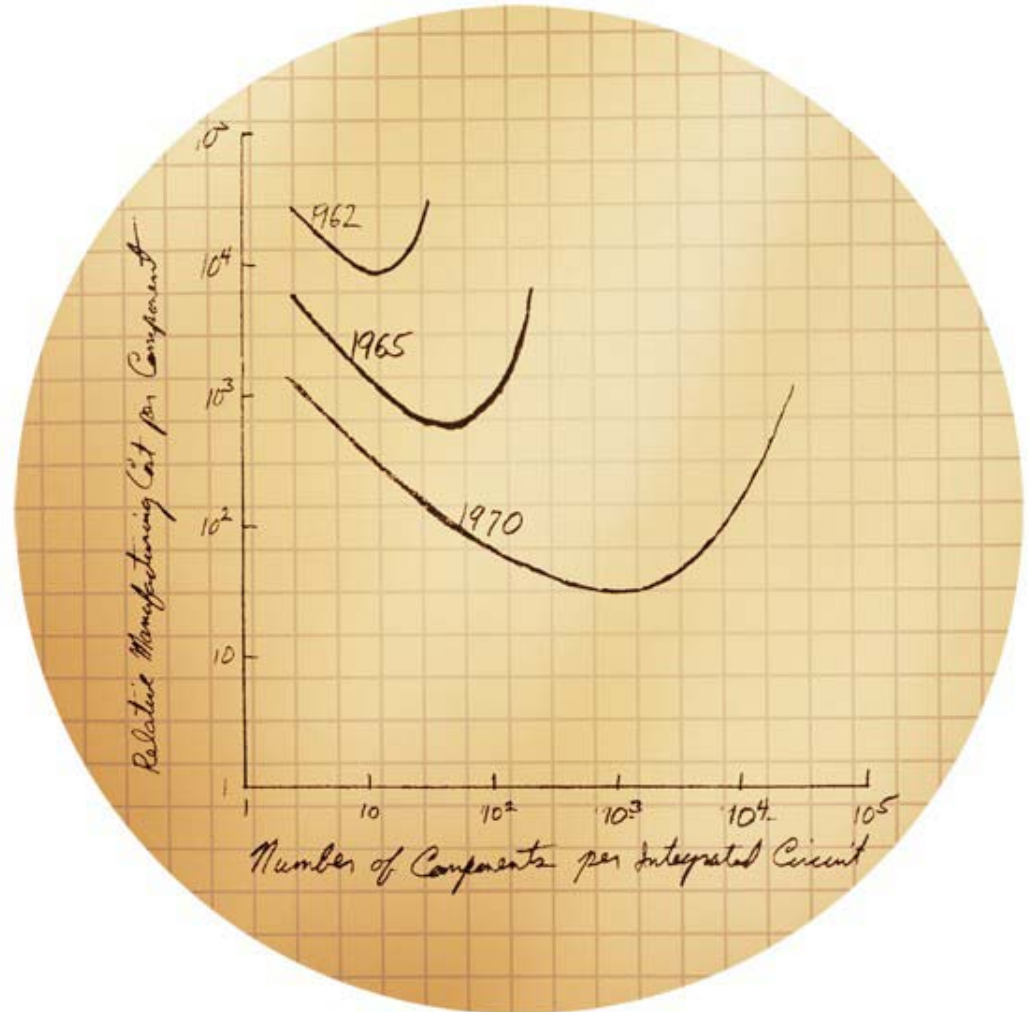
**Collab Workshop, April 12, 2010**

# Outline

- Hardware Evolution

- Parallel Computing

- Graphics Processing Units

- Algorithm Design

- GPU-based Cloud Computing

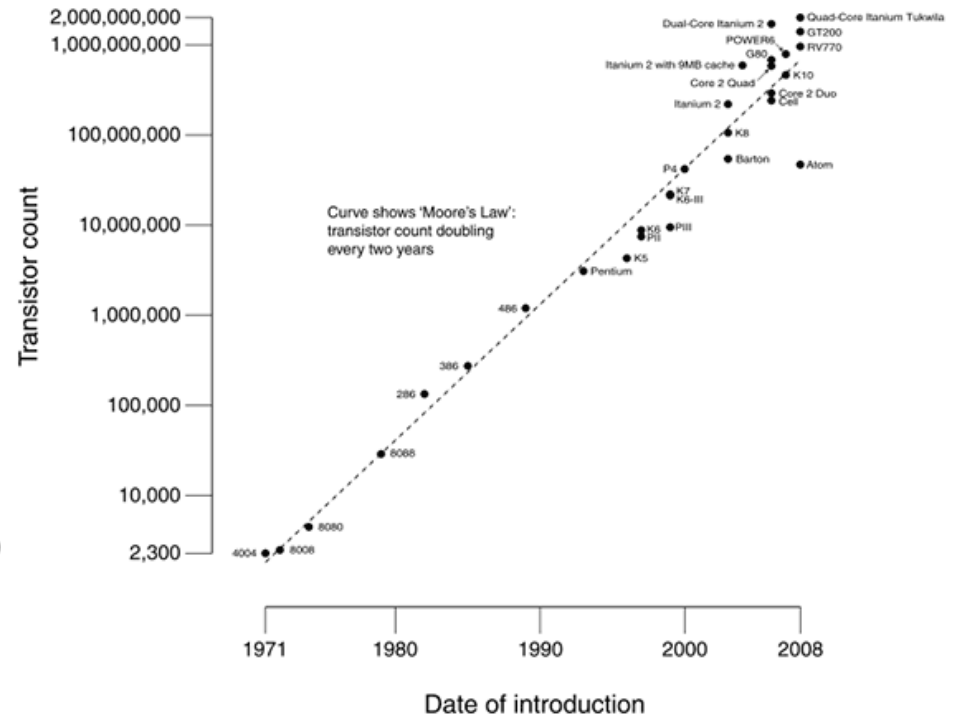# Hardware Evolution

# Moore's Law

- *"The number of transistors on an integrated circuit for minimum component cost doubles every 24 months"*
  *– Gordon Moore.*
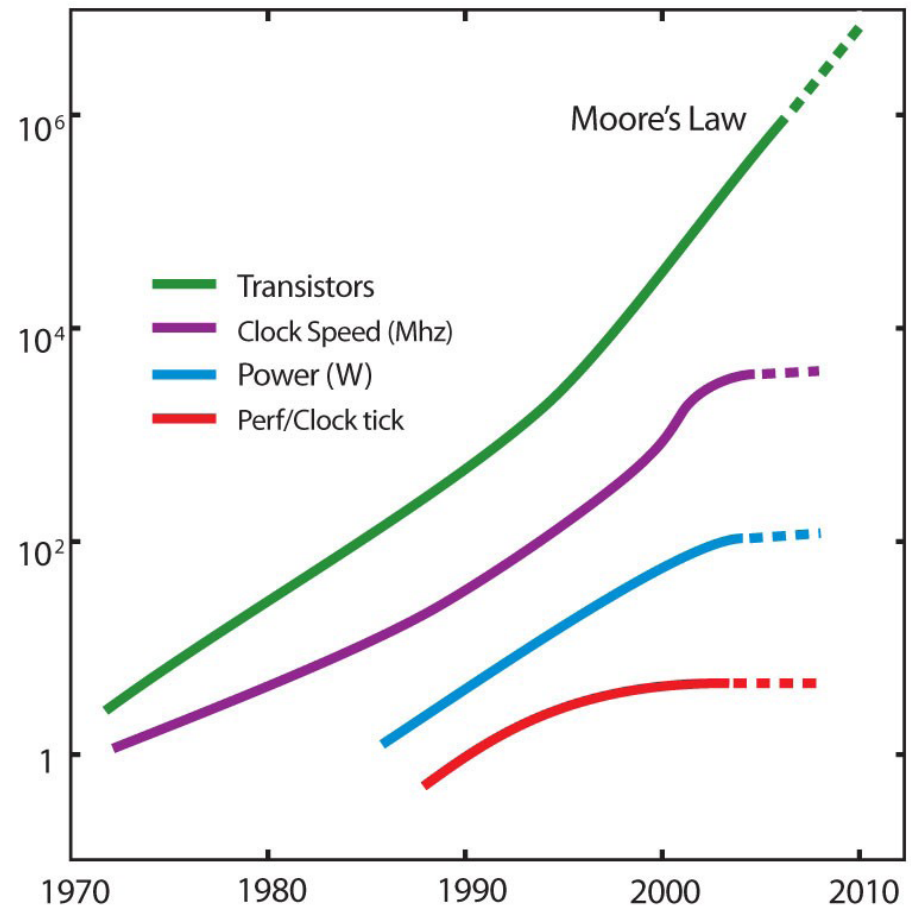
# Brief history of computing: 1970-2004

- 1971 (4004)
  2300 transistors
  **1 x 0.000740 GHz**
- 2004 (Pentium 4 Prescott)
  125 000 000 transistors
  **1 x 4 GHz**
- 2008 (Core i7 Quad)
  731 000 000 transisitors
  **4 x 3.33 GHz**
- 2010+ (Itanium Tukwila/Poulson)
  2 000 000 000+ transistors
  **8 x 2 GHz?**

CPU Transistor Counts 1971-2008 & Moore's Law

Curve shows 'Moore's Law':
transistor count doubling
every two years

Transistor count

Date of introduction

SINTEF

# What happened in 2004?

- Pentium 4 original target frequency: **10 GHz**

- Increasing the frequency has several implications (three walls):
  - Memory
  - Instruction Level Parallelism
  - Power

SINTEF

# Memory Wall

- Memory speeds have not increased as fast as core frequencies
  - A processor can wait through hundreds of clock cycles if it has to get data or instructions from main memory

- Larger caches combined with instruction level parallelism can reduce the memory-wait time

# Instruction Level Parallelism Wall

- Difficulty to find enough parallelism in the instructions stream of a single process to keep higher performance processor cores busy.

- ILP:
  - Instruction pipelining (execution of multiple instructions can be partially overlapped)
  - Superscalar (parallel instruction pipelines)
  - Branch prediction (predict outcome of decisions)
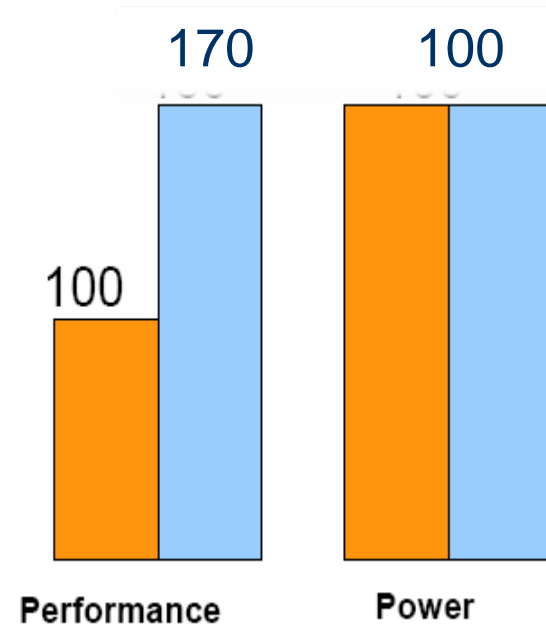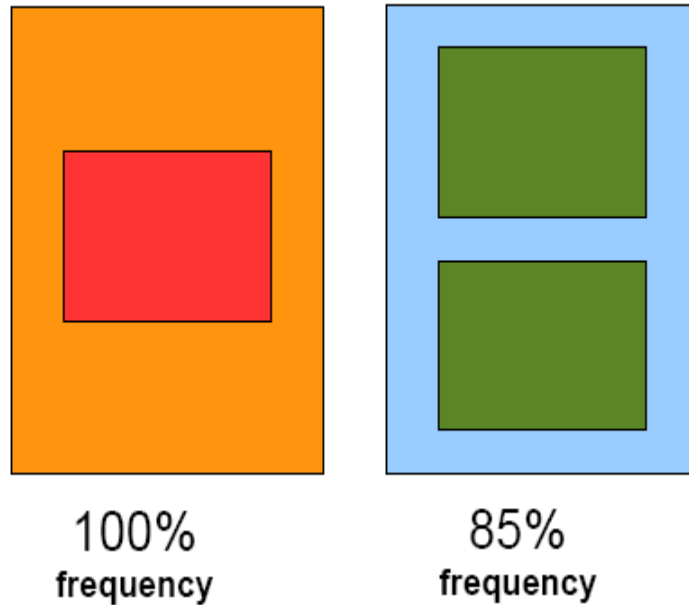  - Out-of-order execution

SINTEF

# Power Consumption Wall

- When the frequency increases, the power consumption increases disproportionately

- Power density relative to **cube of frequency.**
    - Frequency limited to around 4 GHz

$$P = C \cdot \rho \cdot f \cdot V_{dd}^2 \approx V_{dd}^3$$

$P$ is the power density in watts per unit area, $C$ is the total capacitance, $\rho$ is the transistor density, $f$ is the processor frequency, and $V_{dd}$ is the supply voltage.

# How can multi-core help?



- Two cores running at 85% frequency vs. one core at 100% frequency:
  - 100% power consumption
  - 170% performance

# Why not 100-core Pentium IV?

- Deep pipelines,
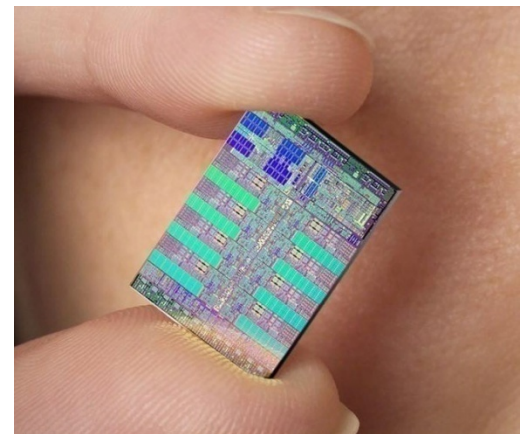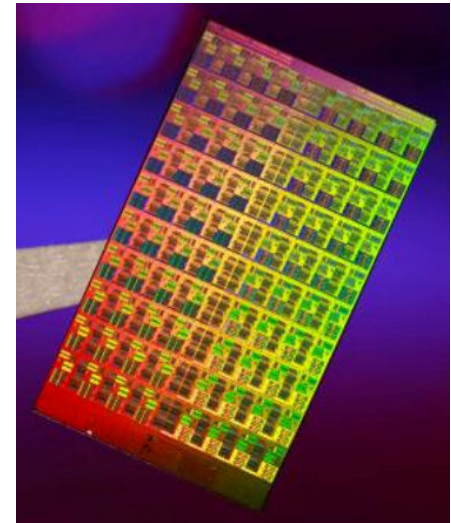- heavy ILP use, and
- huge caches drain a **lot of power.**

*About 50% of the cost of a supercomputer is the electricity bill!*



- CPUs have low CPI (cycles per instruction), but high "WPI" (Watts per Instruction)
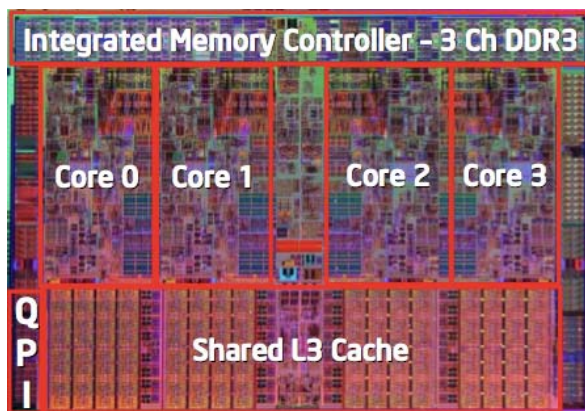- Accelerator cores have shallow pipelines, low or no ILP, and small or no caches!

# Accelerator cores

- **Modern accelerators are power efficient!**
  - Intel Teraflops research chip (AKA Polaris, 80-core): 1 teraflops using 97 watt (Mattson 2008)
  - Core i7 Quad offers 100 gigaflops for 130 watt...
  - GPUs, FPGAs, and the Cell BE are examples of commodity accelerators
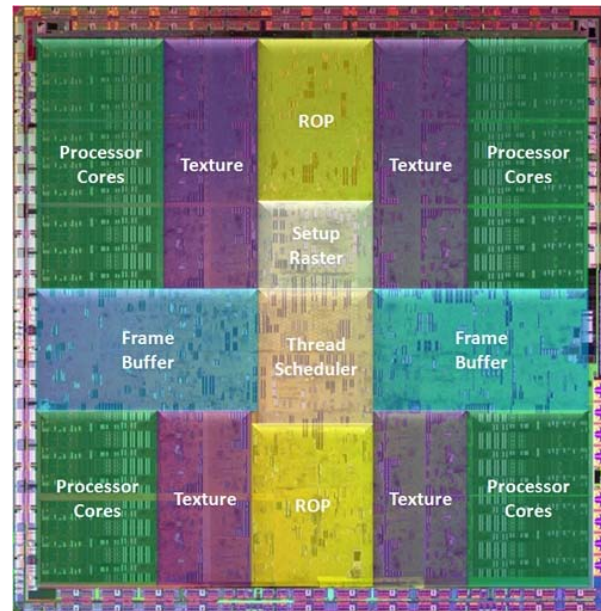
# Heterogeneous Computer

■ A *heterogeneous computer* is a tightly coupled system of processing units with distinct characteristics.

■ A modern desktop or laptop computer is an example of such a system, as most systems include both a task-parallel, multi-core CPU and one or more data-parallel processors in the form of programmable graphics processing units(GPUs).



CPU (Intel Nehalem)                    GPU (NVIDIA GeForce 280)

# Parallel Computing

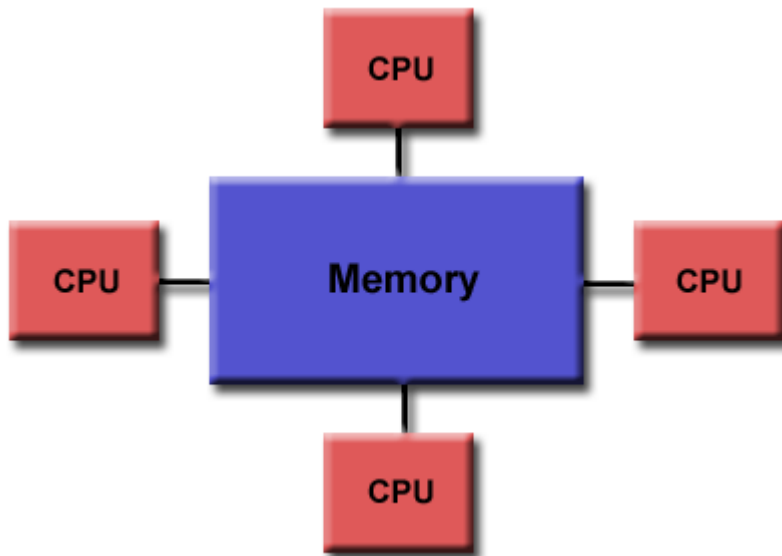# Task vs. Data Parallelization

Task parallelization:

- Each processor executes a different thread (or process) on the same or different data.
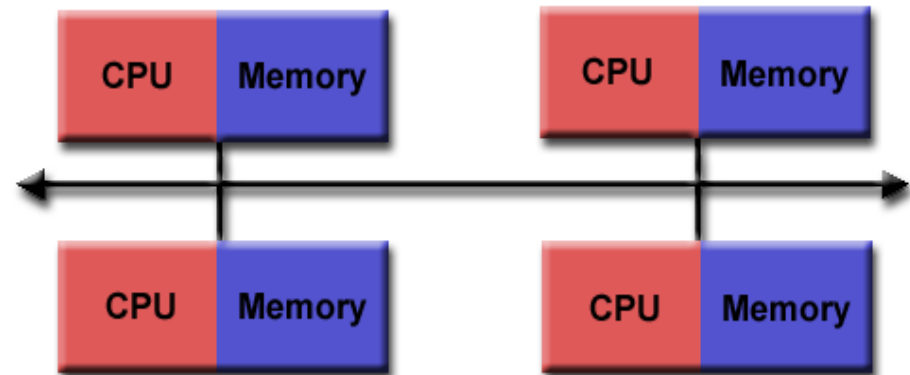
- CPUs

Data parallelization:

- Each processor performs the same task on different pieces of distributed data

- GPUs

SINTEF

# Shared Memory vs. Distributed Memory

- Multiple processors can operate independently but share the same memory resources
- OpenMP

- Distributed memory systems require a communication network to connect inter-processor memory
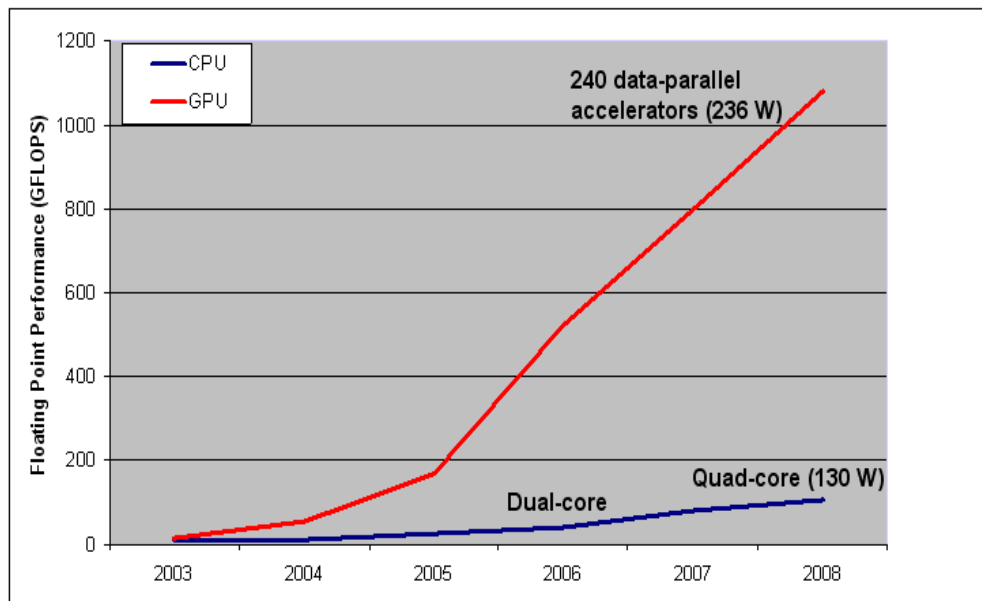- MPI

Shared memory system

Distributed memory system

# Graphics Processing Units

# Graphics Processing Unit (GPU)

- Background: Game industry
- Massive parallel architecture (>500 cores per chip)
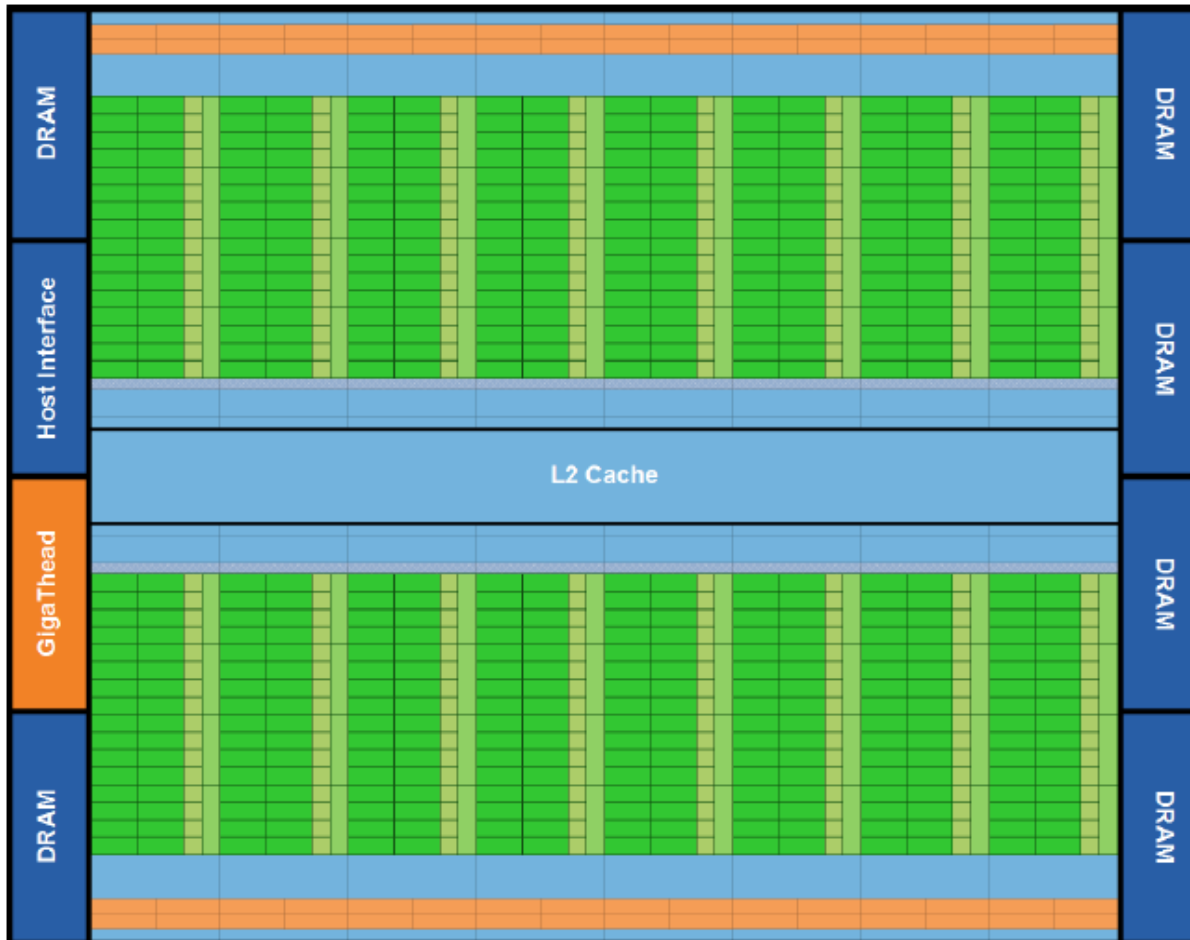- Typically 10-50 times speedup compared to CPU applications

# NVIDIA Fermi GPU

- 3 billion transistors

- 512 cores

- Up to 1 Terabyte of GPU memory

- 1.5 Teraflops performance

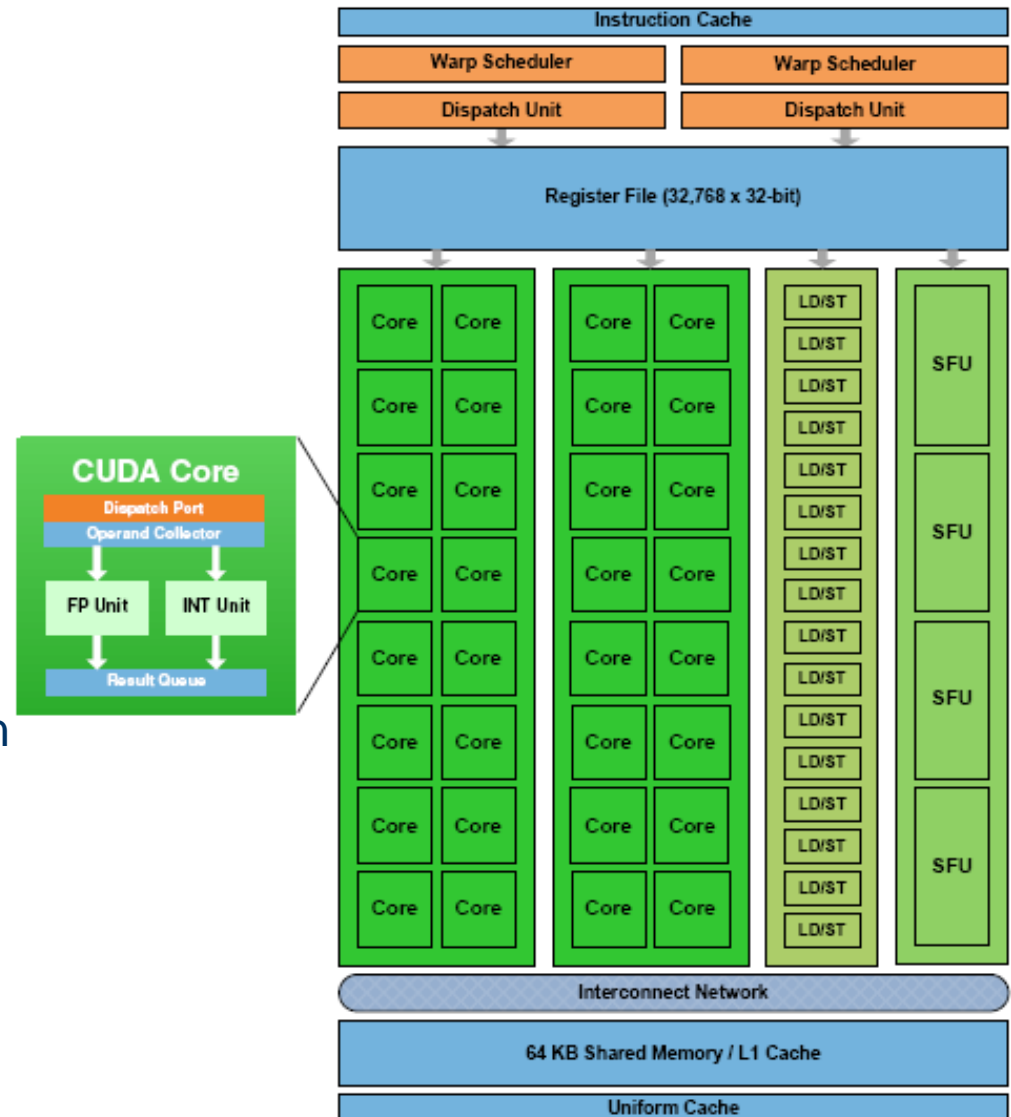- 40 nm manufacturing process

NVIDIA GTX 480

# NVIDIA Fermi Architecture



16 streaming multiprocessors are positioned around a common L2 cache.
Each SM is a vertical rectangular strip that contain an orange portion (scheduler and dispatch), a green portion (execution units), and light blue portions (register file and L1 cache).

# Fermi Architecture – Streaming Multiprocessor

- 32 cores per SM

- 64 KB of RAM for shared memory and L1 cache

- Double precision speed at 50% of peak single precision

- New IEEE 754-2008 floating point standard, surpassing even the most advanced CPUs.

- Newly designed integer ALU optimized for 64-bit and extended precision operations.
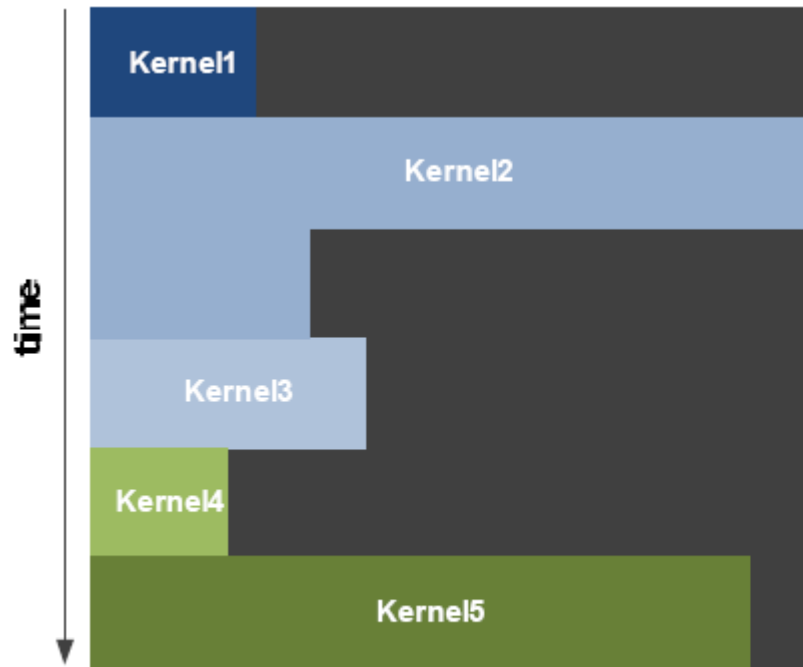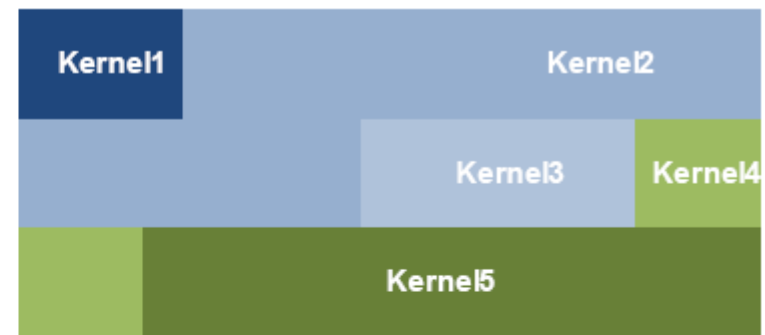
# Cached Memory Hierarchy

- First GPU to support a true cache hierarchy in combination with on-chip shared memory

- L1 cache per SM (32 cores)

- Unified L2 cache (768KB)
  - Fast, coherent data sharing across all cores.

# Hardware Thread Scheduler

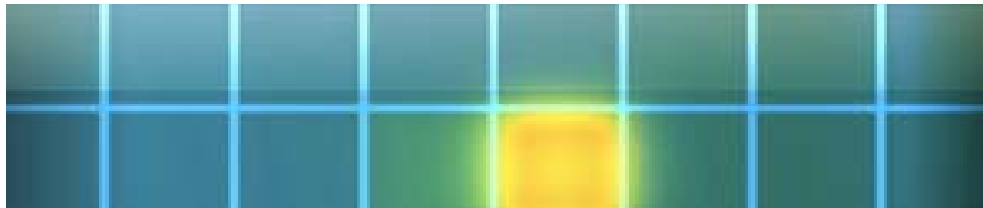- Concurrent kernel execution and faster contex switch



Serial Kernel Execution

Concurrent Kernel Execution

# Full ECC Support

- First GPU architecture to support ECC

- Detects and corrects errors before system is affected

- Protects register files, shared memories, L1 and L2 cache, and DRAM
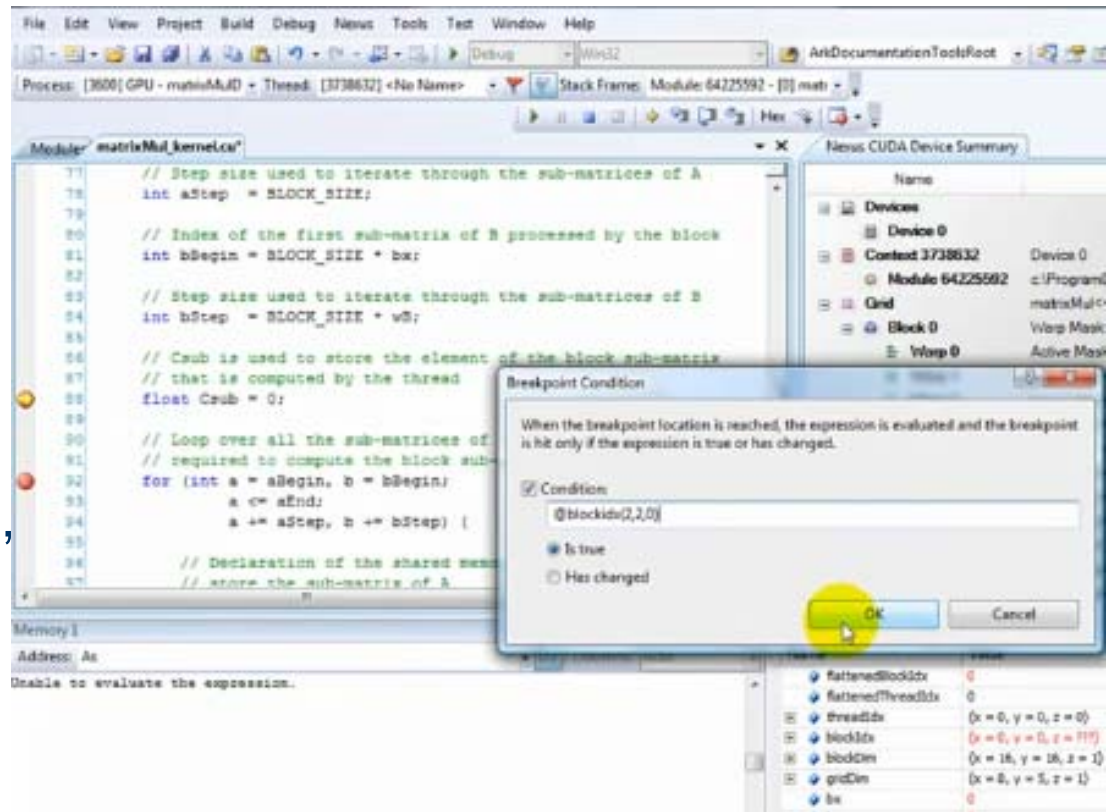
# Enhanced Software Support

- **Full C++ support**
  - Virtual functions
  - Try / catch hardware support

- **System call support**
  - printf() etc.

SINTEF

# Programming Frameworks

- Compute Unified Device Architecture (CUDA)
  - Developed by NVIDIA
  - Only available from NVIDIA GeForce 8 series and above

- Open Computing Language (OpenCL)
  - Specified by the Khronos group
  - Cross platform

- DirectCompute
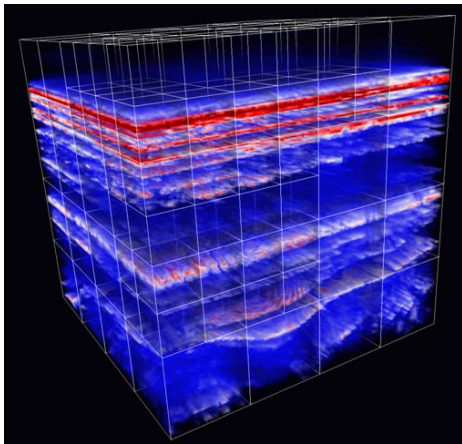  - Part of Microsoft DirectX

SINTEF

# NVIDIA Parallel Nsight

- Debug, profile, and analyze GPU code

- Integrated with Microsoft Visual Studio 2008

- Will support CUDA C, OpenCL, DirectCompute, Direct3D, and OpenGL.

# Application Showcase (SINTEF)
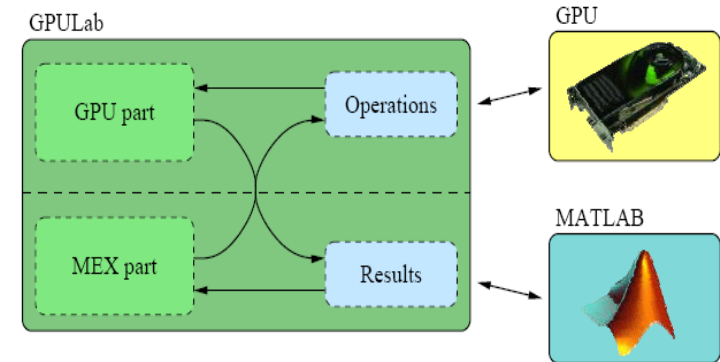
Isosurface visualization

Simulations



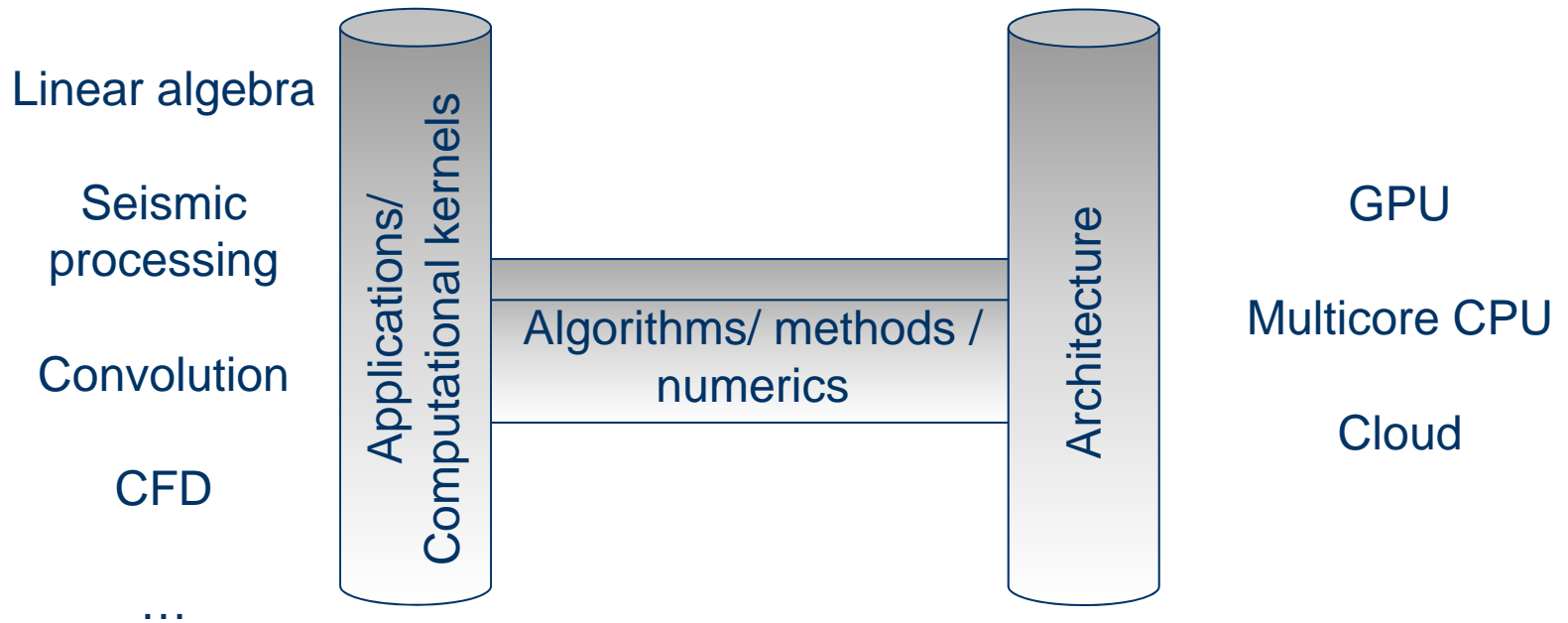Seismic/ reservoir processing



Ultrasound processing



Matlab interface for GPUs

# Algorithm Design

# Architecture and Algorithm Interaction

*The key to performance is to understand the architecture and algorithm interaction.*

Linear algebra

Seismic processing

Convolution

CFD

…

Applications/ Computational kernels

Algorithms/ methods / numerics

Architecture

GPU

Multicore CPU

Cloud

"Multicore computers shift the burden of software performance from chip designers and processor architects to software developers"
(J. Larus in *Spending Moore's Dividend*, CACM May 2009)

# Computations: Algorithms and hardware advances doubles the speedup

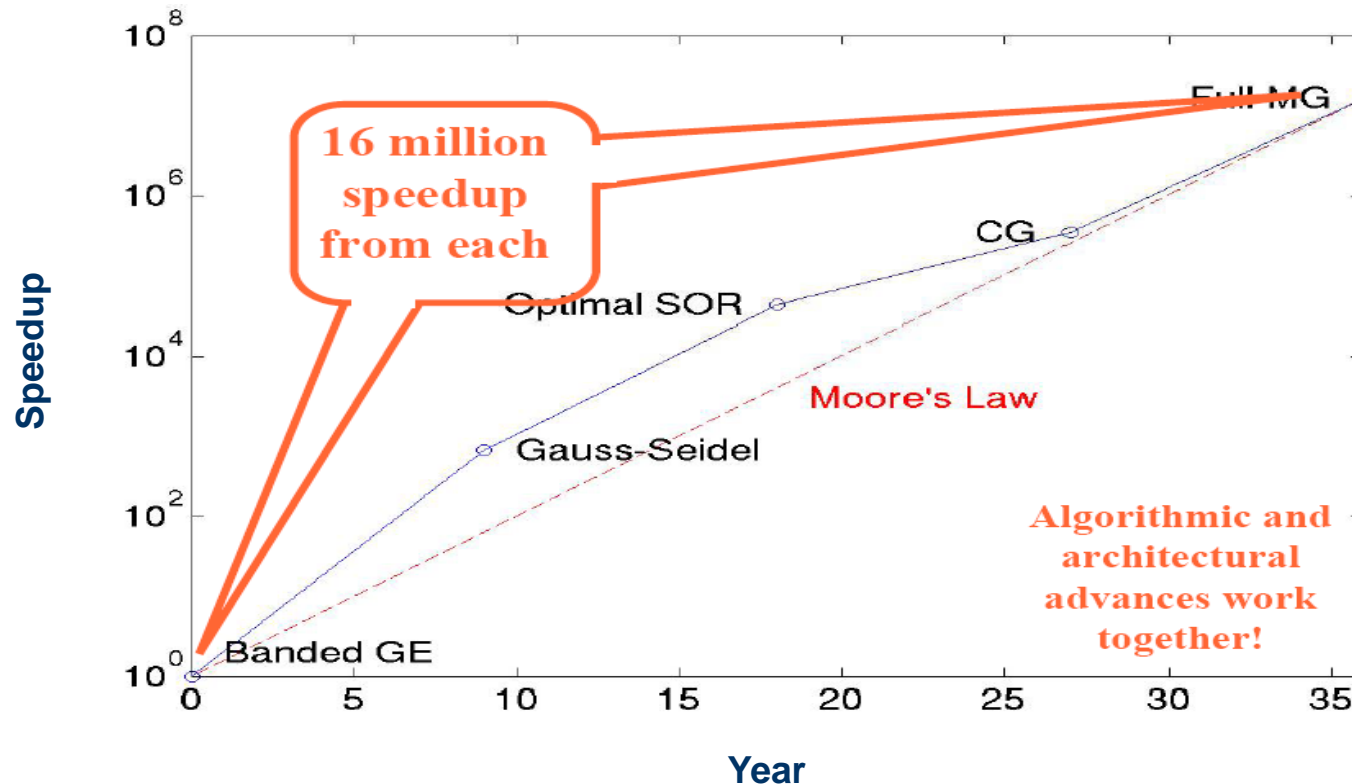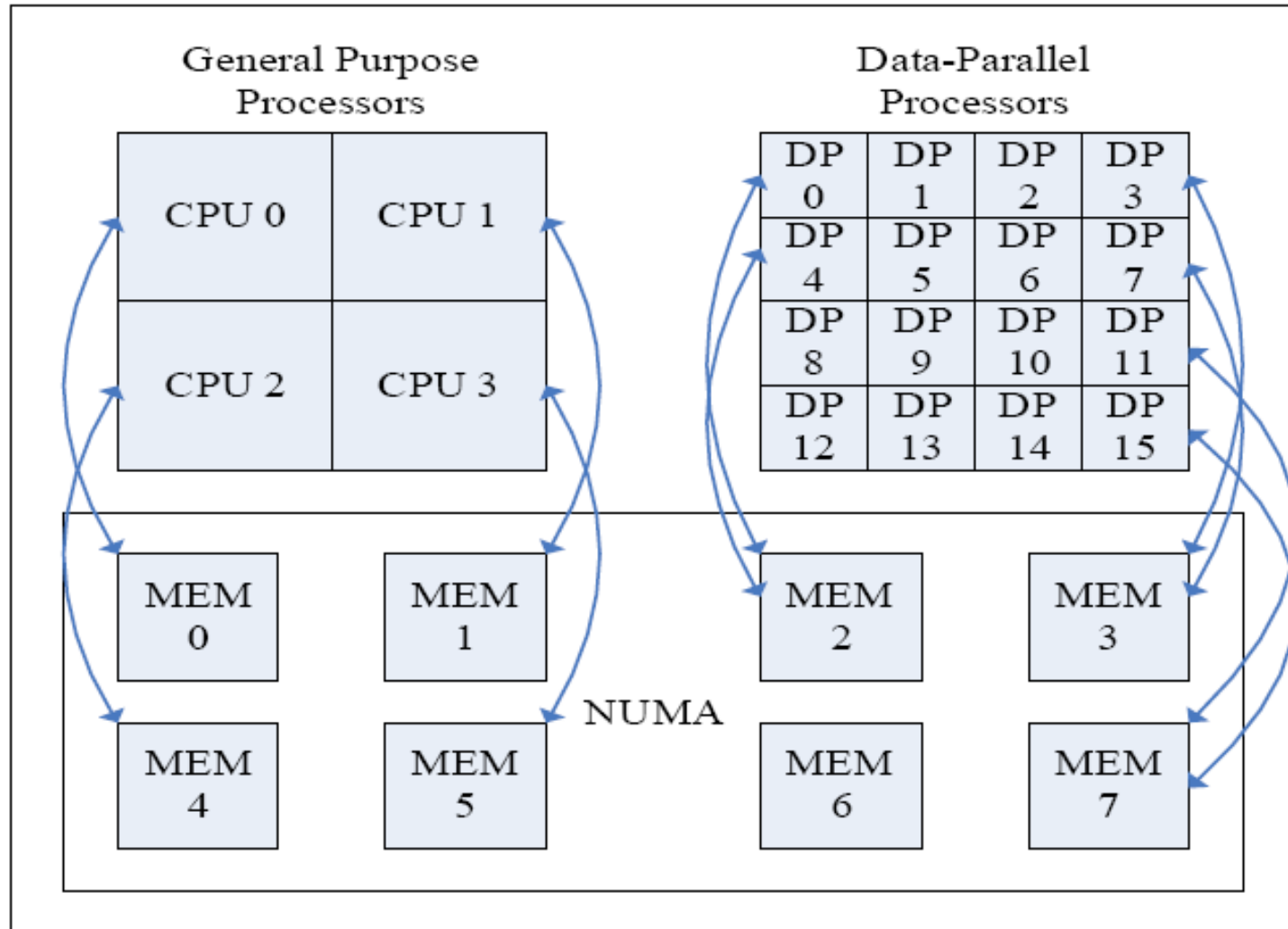**Solving linear systems of equations:**



*Image courtesy of David E. Keyes, Columbia University*

# An Idealized Heterogeneous Computer

# Consequences for Computational Algorithms

- The vast computational resources available make it crucial to design algorithms that are parallel from the ground up.

- Different sequential tasks are issued to the task-parallel processors

- Inheritably parallel tasks are executed on the data-parallel processors.

SINTEF

# NUMA memory layout of HCs

- The growth in processor performance seems to outpace the growth of memory bandwidth

- The data-flow must be organized so that processors access data that is close to it.

- Identify and minimize the dependency between data

- Use extra computations to avoid data synchronization

SINTEF

# Precision

- Double-precision numbers requires twice the bandwidth

- Double-precision numbers compute at half the speed of single-precision

- Mixed-precision algorithms will play a key role.

# What about existing code?

■ Numerical software can often gain dramatic increases in performance by refactoring just a few selected parts

```
for(…){
  for(…){
    computationalKernel();
  }
}
```
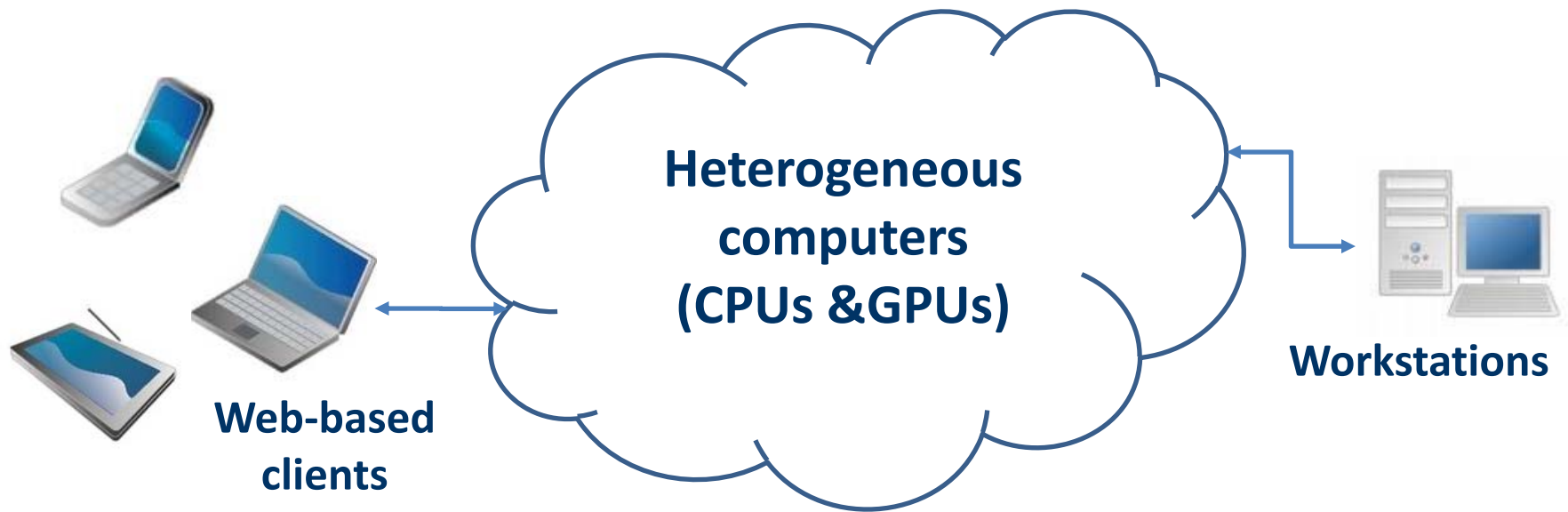
■ Such approaches will most likely reach the limits of Amdahl's law: *The theoretical maximum speedup is limited by the relative time needed for the non-parallel part of the algorithm.*
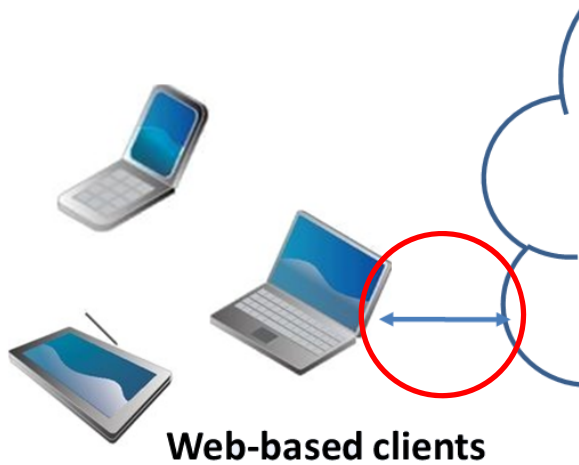
# GPU-based Cloud Computing

# Cloud Computing Characteristics

- Customers do not own the infrastructure

- Consumers pay for only what resources they use

- Dynamic hardware scalability

- Device and location independent
  - Only need an internet connection
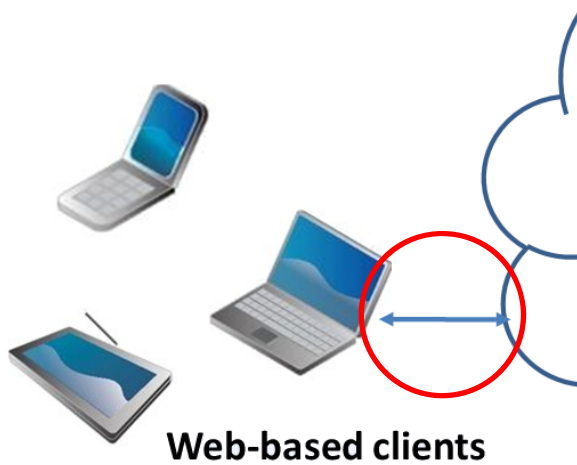
# Cloud of Heterogeneous Computers

**Heterogeneous computers (CPUs &GPUs)**

**Web-based clients**

**Workstations**

# Web-Based Clients



**Web-based clients**
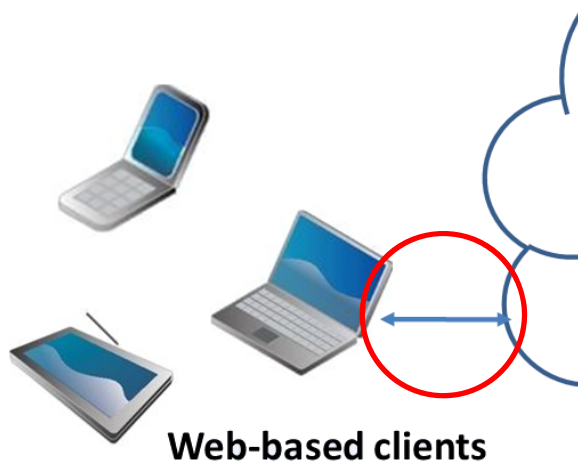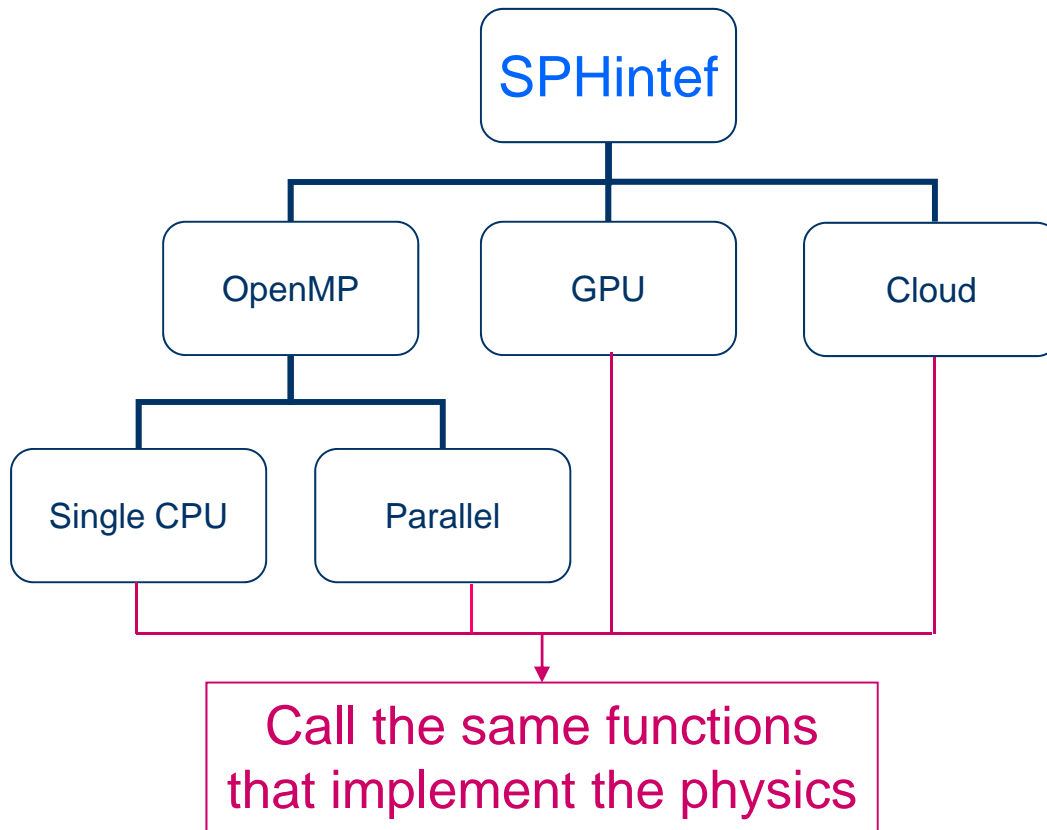
# Web-Based Clients



**Web-based clients**

- Processing in the cloud
- Streaming of results to the clients
- User input from the client to the cloud.
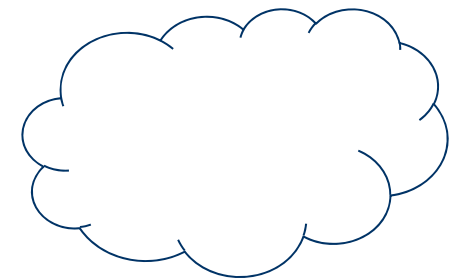
# Web-Based Clients



**Web-based clients**

- Processing in the cloud
- Streaming of results to the clients
- User input from the client to the cloud.

# Example: Code Framework – SPHintef



SPHintef

- OpenMP
- GPU
- Cloud

OpenMP:
- Single CPU
- Parallel

**Call the same functions that implement the physics**

GPU – Graphics Processing Unit

Cloud computing

Visualization of results:

# Research Challenges: Auto-Tuning in the Cloud

- Future clouds will be based on heterogeneous architectures.

- New auto-tuning algorithms must be developed
  - Traditional auto-tuning techniques do not map to new heterogeneous architectures.

- Very few generic libraries developed
  - BLAS, FFT

- Multipe levels of parallelism:
  - Distributed memory, shared memory
  - MPI, OpenMP, threads, CUDA, etc.

SINTEF

# Thanks!

# Questions?