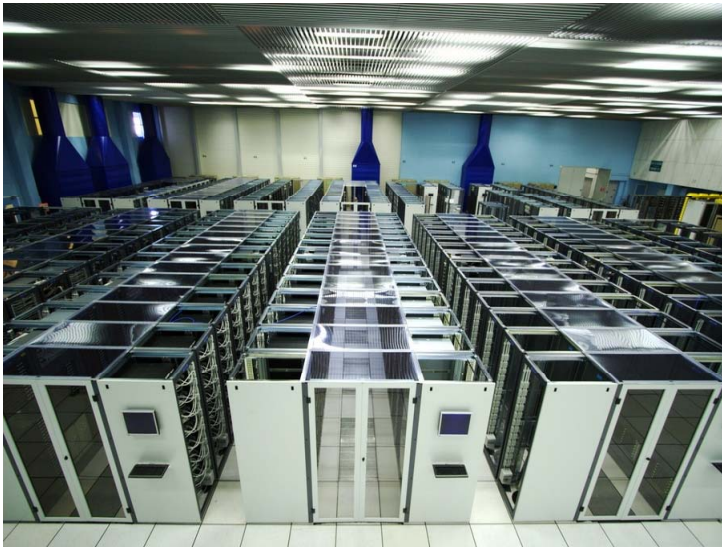
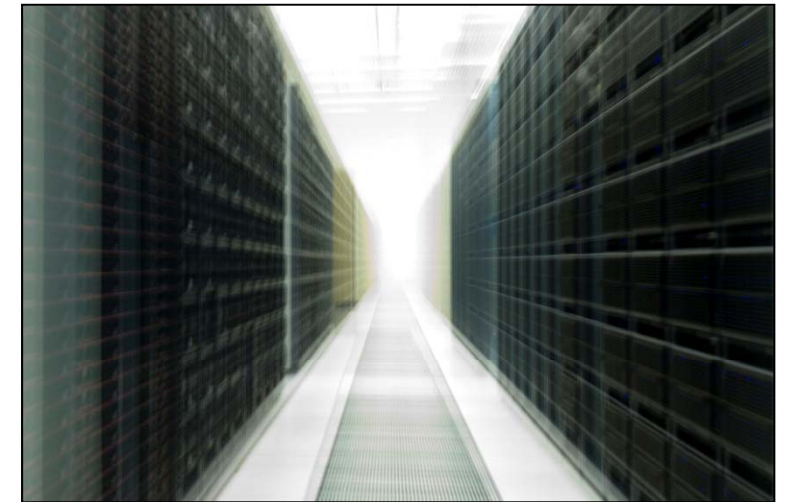


# Harnessing Future Hardware



Sverre Jarp  
CERN  
openlab  
IT Dept.

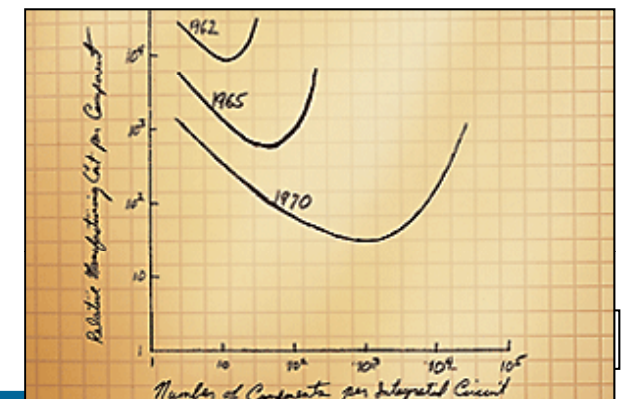
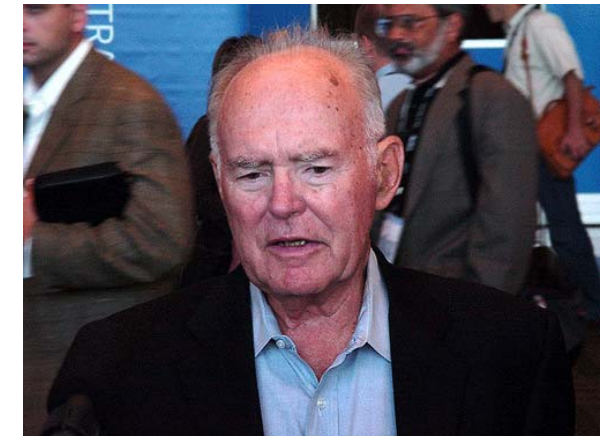
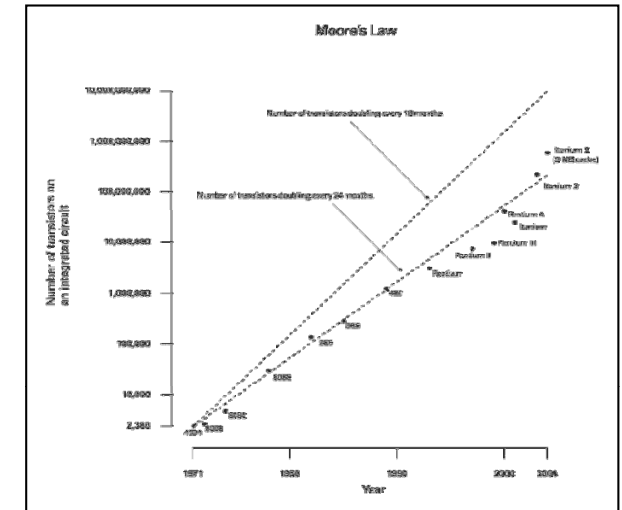


CERN

Tracking Workshop – GSI, Darmstadt, 9 – 11 June 2010

# The driving force: Moore's law

- We continue to double the number of transistors every other year(\*)
  - Latest consequence
    - Single core
      - Multicore
      - Manycore



# Real consequence of Moore's law

- We are being “**snowed under**” by transistors:
  - More (and more complex) execution units
    - Hundreds of new instructions
  - Longer SIMD/SSE vectors
  - More hardware threading
  - More and more cores
- In order to profit we need to “think parallel”
  - Data parallelism
  - Task parallelism

# “Intel platform 2015” (and beyond)

- **Today’s silicon processes:**

- 45 nm
- 32 nm

- **On the roadmap:**

- 22 nm (2011/12)
- 16 nm (2013/14)

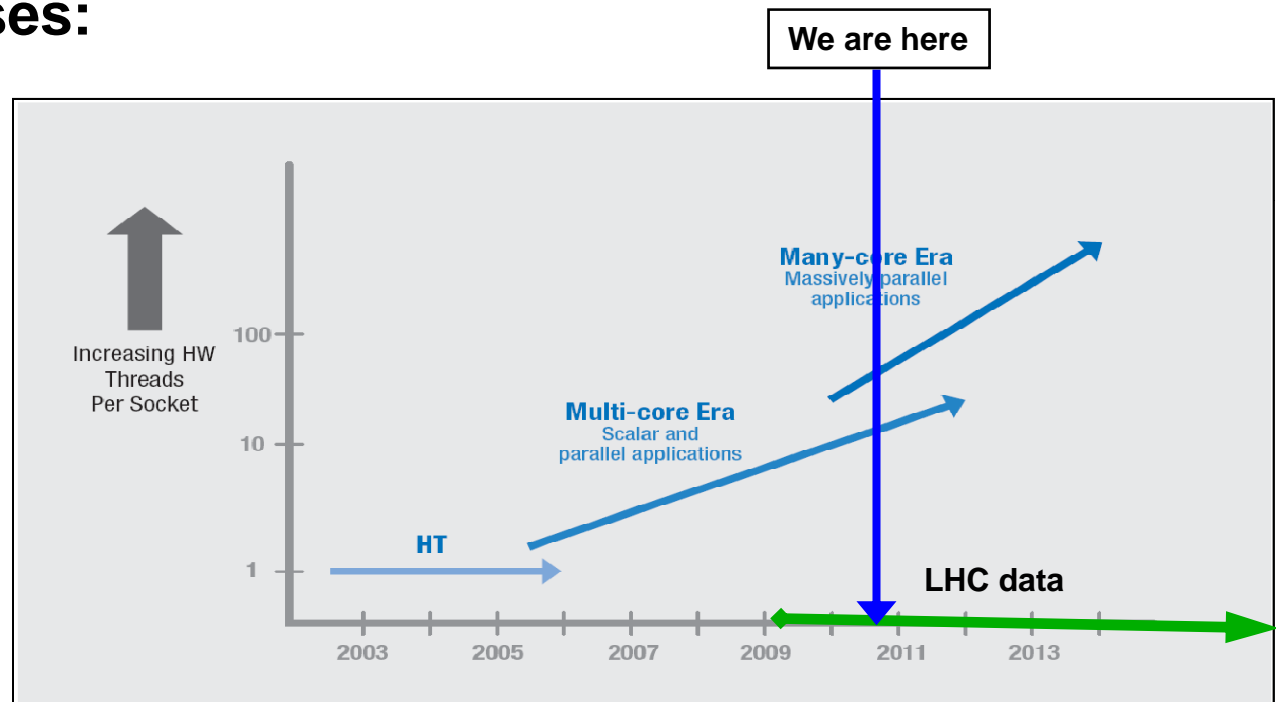
- **In research:**

- 11 nm (2015/16)
- 8 nm (2017/18)

– Source: Bill Camp/Intel HPC

- **Each generation will push the core count:**

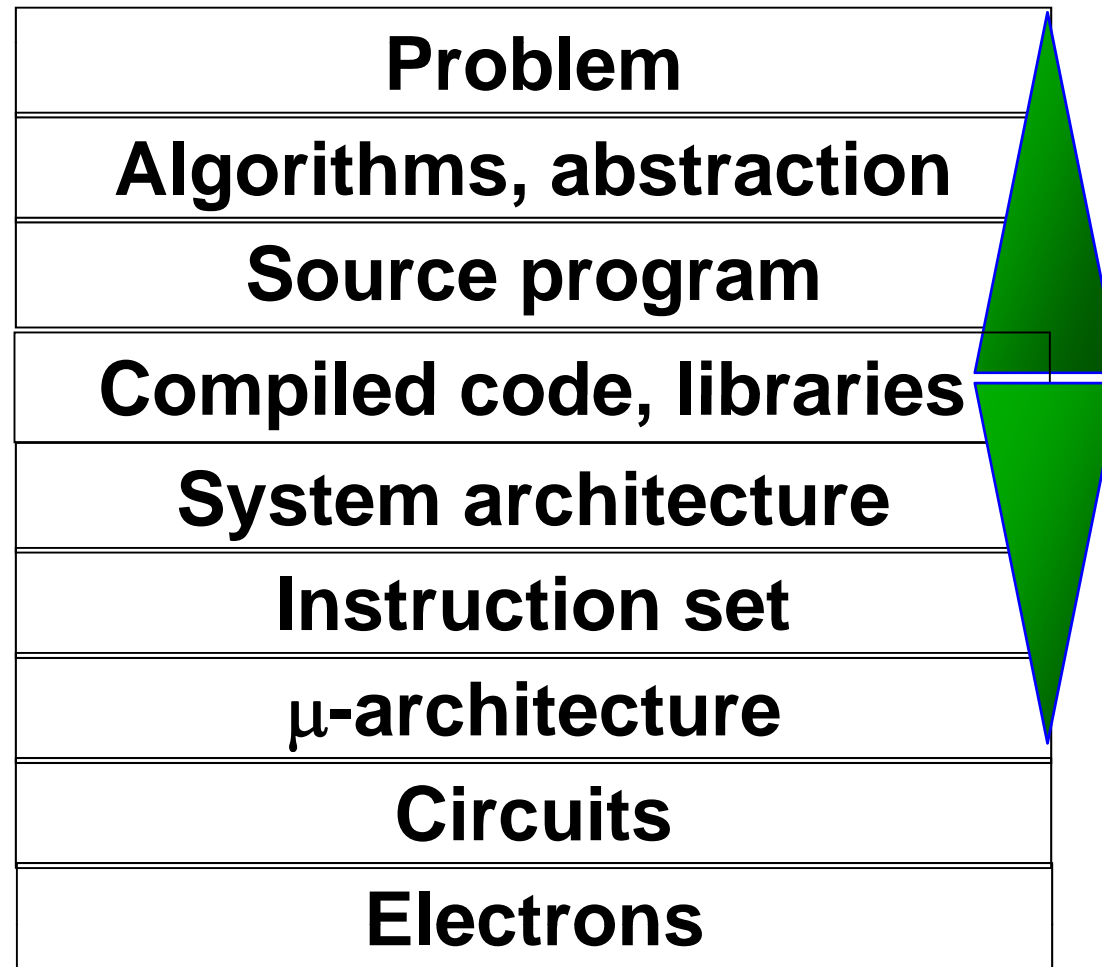
- **We are entering the many-core era (whether we like it or not) !**



S. Borkar et al. (Intel), "Platform 2015: Intel Platform Evolution for the Next Decade", 2005.

# Programming: A Complicated Story

- We cannot concentrate on just one layer (ignoring the others)



# Different designs: CPU or GPU

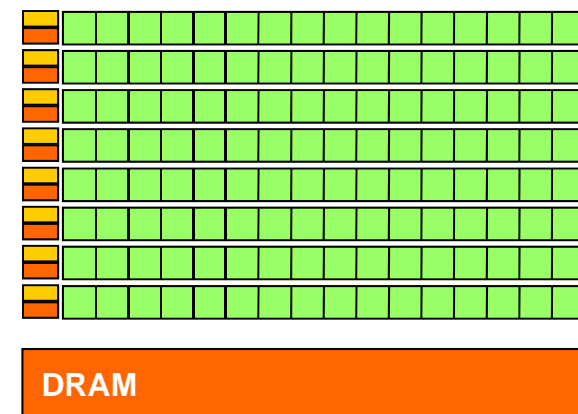
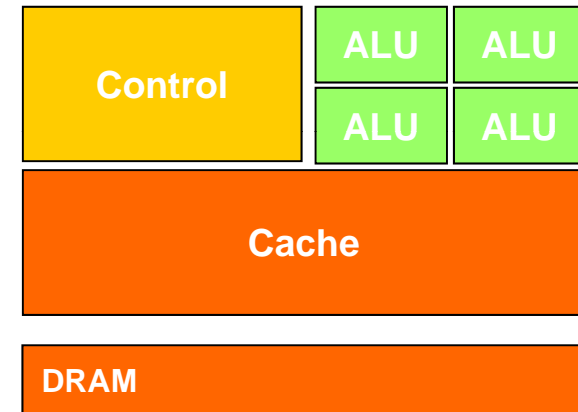
## ■ Different focus

### ■ CPU

- Optimised for low-latency access to cached data sets
- Control logic for out-of-order and speculative execution

### ■ GPU

- Optimised for data-parallel, throughput computation
- Architecture tolerant of memory latency
- More transistors dedicated to computation



# Seven dimensions of performance

- **First three dimensions:**

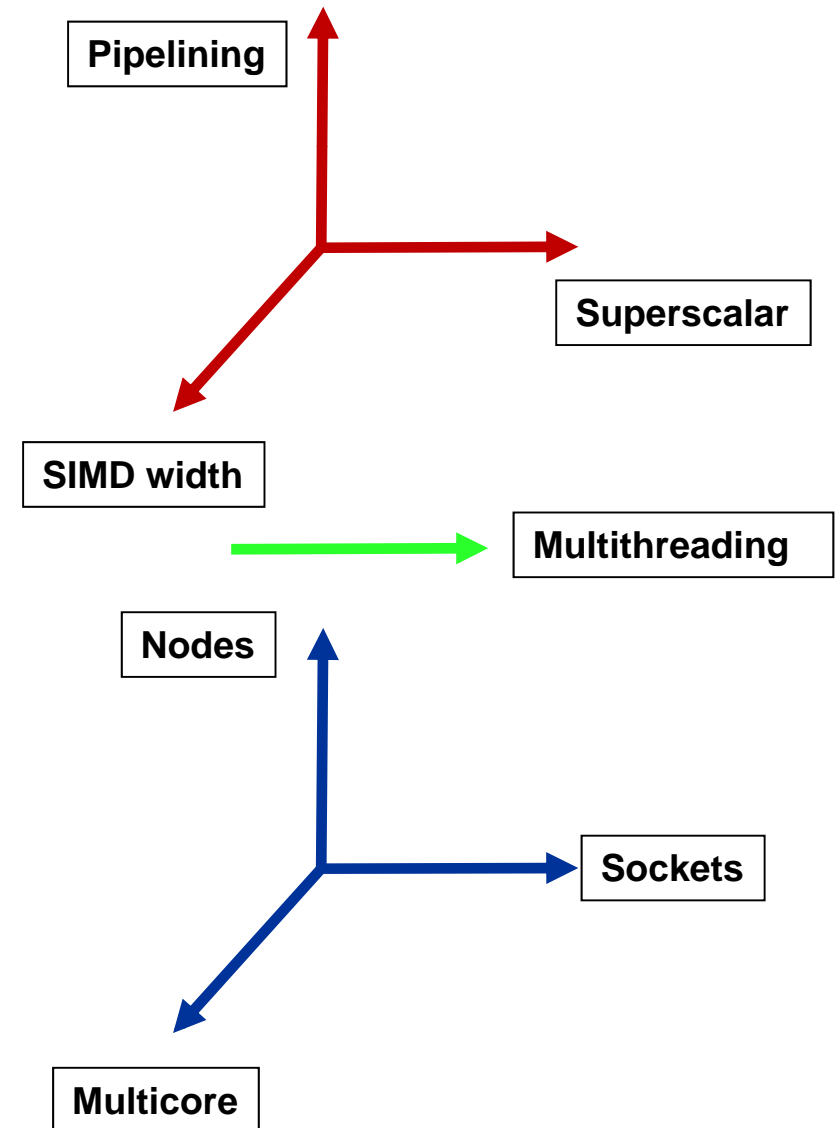
- Superscalar
- Pipelining
- Computational width/SIMD

- **Next dimension is a “pseudo” dimension:**

- Hardware multithreading

- **Last three dimensions:**

- Multiple cores
- Multiple sockets
- Multiple compute nodes



# In the days of the Pentium

## ■ First three dimensions:

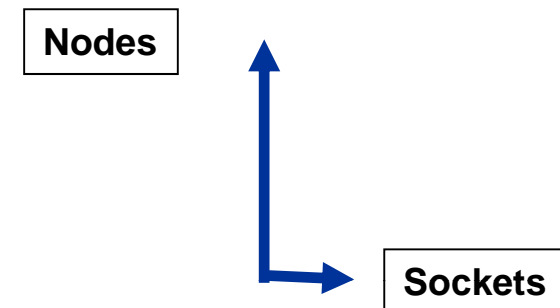
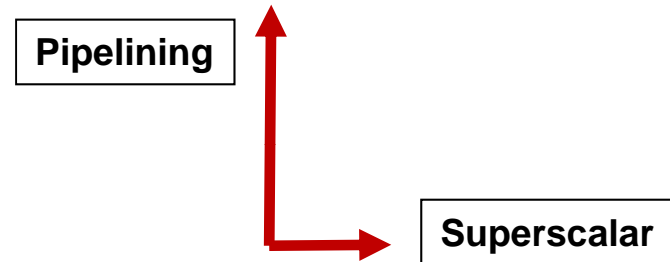
- Superscalar (**only two ports**)
- Pipelining (OK)
- **No** vectors

## ■ Next dimension is a “pseudo” dimension:

- **No** hardware multithreading

## ■ Last three dimensions:

- **No** cores
- **Hardly any** dual socket systems
- Multiple compute nodes (OK)





# Seven multiplicative dimensions:

## ■ First three dimensions:

- Superscalar
- Pipelining
- Computational width/SIMD

**Data parallelism  
inside the core  
(Vectors/Matrices)**

## ■ Next dimension is a “pseudo” dimension:

- Hardware multithreading

**Task parallelism  
across CPUs  
(Events/Tracks)**

## ■ Last three dimensions:

- Multiple cores
- Multiple sockets
- Multiple compute nodes

**Process/Task  
parallelism**

# The move to many-core systems

- **Examples of “CPU slots”: Sockets \* Cores \* HW-threads**

- Basically what you observe in “cat /proc/cpuinfo”

- **Conservative:**

- Dual-socket AMD six-core (Istanbul):  $2 * 6 * 1 = 12$

- Dual-socket Intel six-core Westmere:  $2 * 6 * 2 = 24$

- Quad-socket Intel Dunnington server:  $4 * 6 * 1 = 24$

- **Aggressive:**

- Quad-socket AMD Magny-Cours (12 core)  $4 * 12 * 1 = 48$

- Octo-socket Nehalem-EX “octo-core”:  $8 * 8 * 2 = 128$

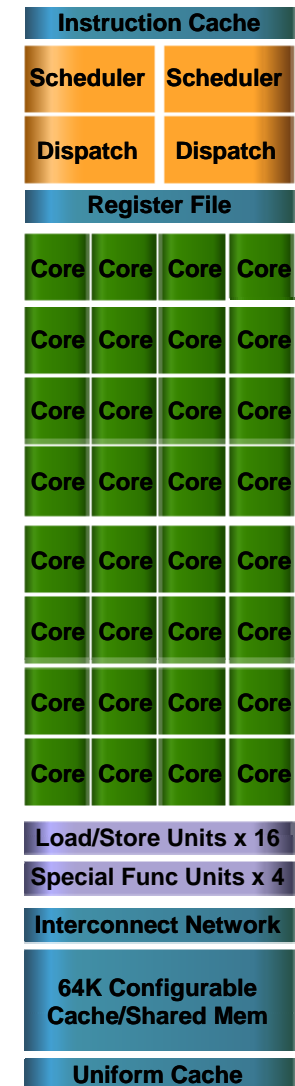
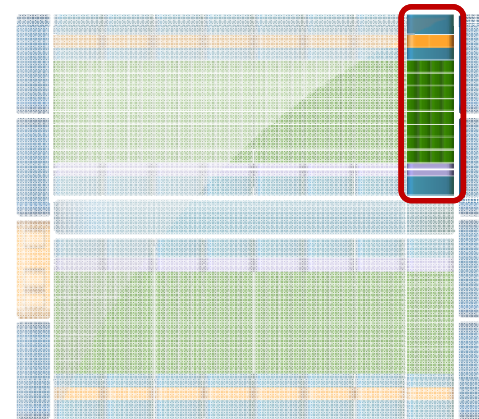
- Quad-socket Sun Niagara (T3) processors w/16 cores and 8 threads (each):  $4 * 16 * 8 = 512$

- **Now, or in the near future: Hundreds of CPU slots**

- **And, by the time new software is ready: Thousands !!**

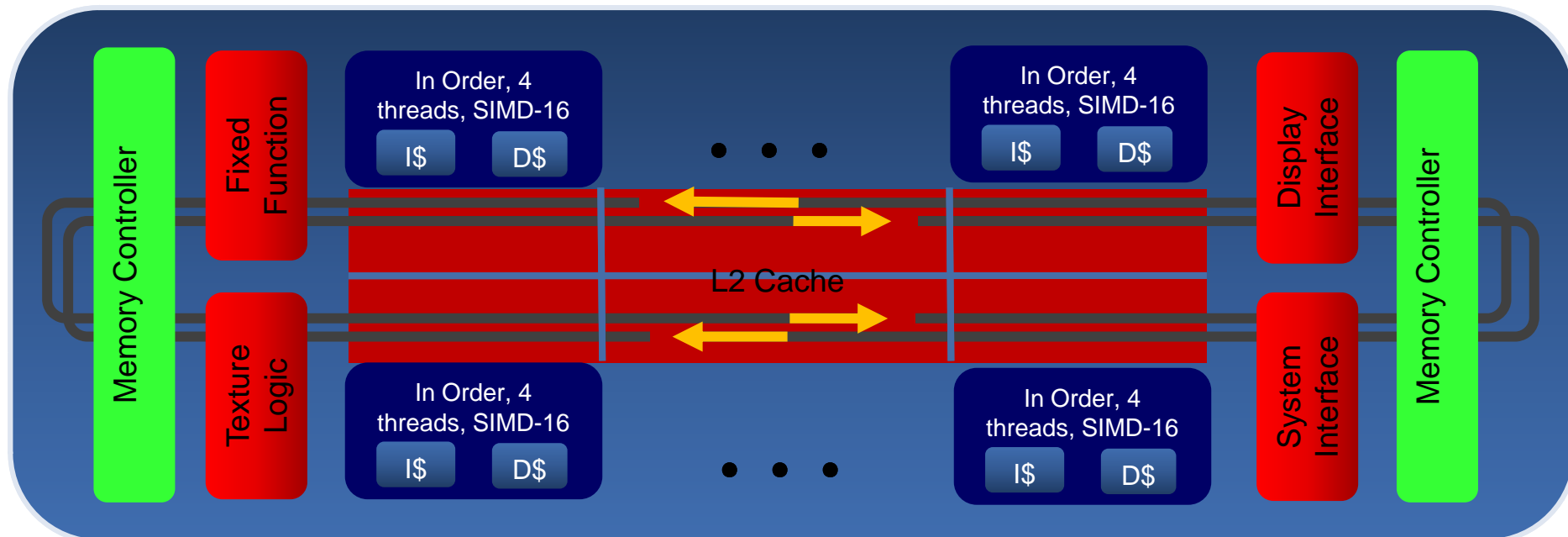
# Nvidia Fermi design

- **Streaming Multiprocessing Architecture**
- **32 CUDA cores per SM (512 total)**
- **8× peak double precision floating point performance**
  - 50% of peak single precision
- **Dual Thread Scheduler**
- **64 KB of RAM for shared memory and L1 cache (configurable)**



# Many-core accelerator

- **Many Inter Core Architecture:**
  - Announced at ISC10 (June 2010)
  - Based on the x86 architecture, 22nm (2012?)
  - Many-core (> 50 cores) + 4-way multithreaded + **512-bit vector unit**



# Conclusion

- **The parallel (CPU, GPU) hardware is here to stay!**
  - Realistic benchmarking becomes paramount !
- **Different teams will take different parallelisation approaches, based on:**
  - Political desire to change
  - The code size and the percentage that is performance sensitive
  - Loop constructs and vector, matrix availability in the code
  - Available effort for a substantial re-write
- **The potential is huge but everybody needs to understand the ROI**
  - Investment in human effort, against better use of available resources!