**2494-17**

**Workshop on High Performance Computing (HPC) Architecture
and Applications in the ICTP**

*14 – 25 October 2013*

**Eurotech company overview and HPC division**

G. Tecchiolli
*Eurotech, Udine*

# Workshop on HPC Architecture and Applications
# ICTP - Trieste

**Giampietro Tecchiolli**

**Trieste - 2013 10 21**

# Agenda

- **Introduction**
- **Beyond the Moore's Law**
    - MIC Computing
    - GPU Computing
- **Energy-Efficient Computing**
    - RISC Architectures (ARM)
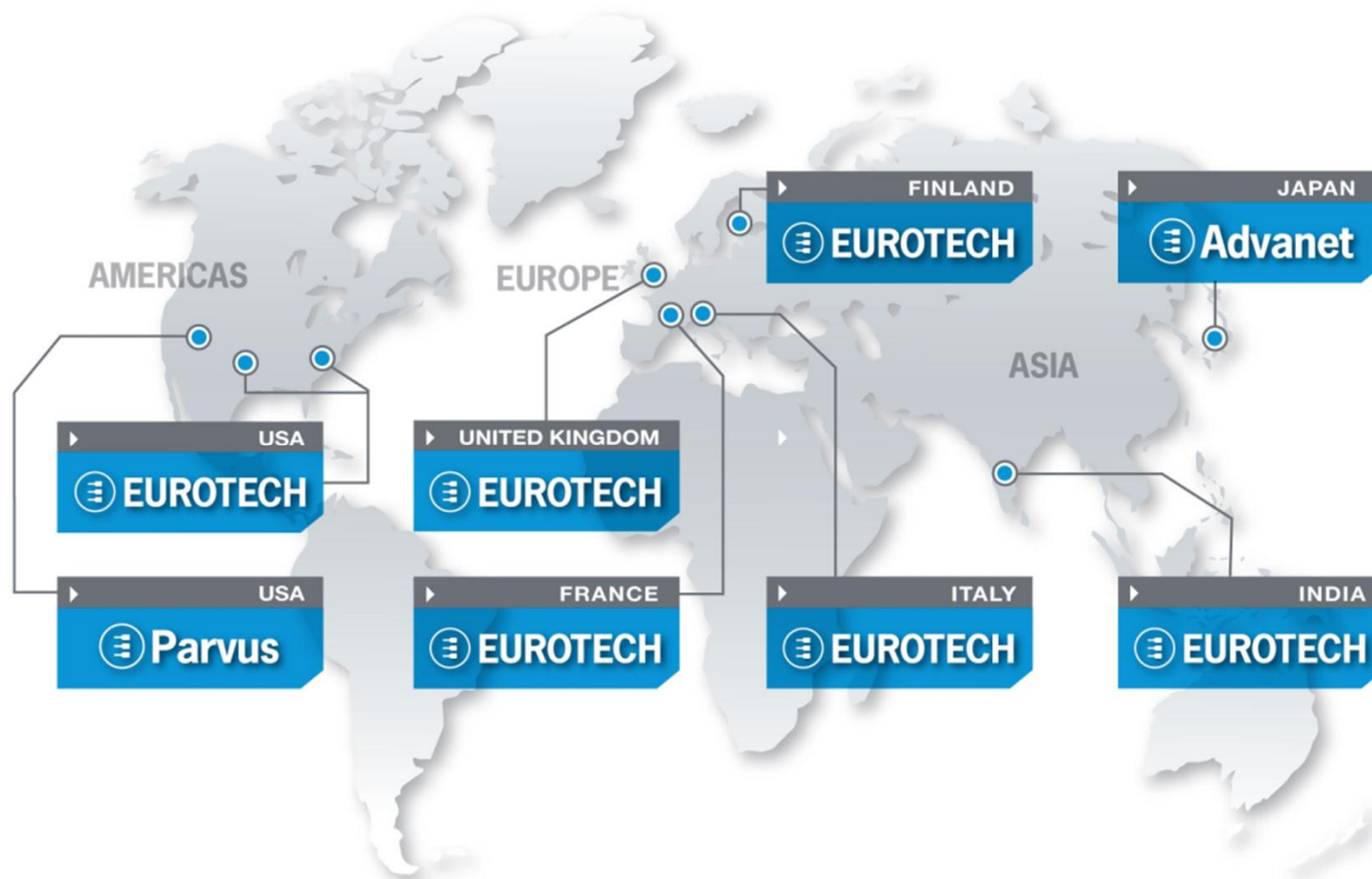    - Reconfigurable Architectures (FPGA)

⊜ EUROTECH

# Eurotech company overview and HPC division

# EUROTECH

# Eurotech Introduction
## Group Global Footprint

# HPC division highlights

- The Eurotech HPC division focuses on **designing, manufacturing, delivering and supporting high performance computing solutions**

- More than 14 years of history of delivering supercomputing systems and solutions to industry and academia

- First worldwide company to market hot water cooled high performance computers. First hot water cooled HPC in the market delivered in 2009.

- R&D capabilities nurtured in house and through collaboration with the best universities and research centres in Europe: INFN, Julich, Revensburg, Daisy…

- Funding member of ETP for HPC

ETP 4 HPC

AURORA

# Eurotech HPC project examples

Ape mille, 1999-2002
Ape next, 2002-2005

Q-Pace, 2007-2009

Janus, 2006-2008

**Universität Regensburg**

**Universidad Zaragoza**
1542

Aurora Science, 2008-2010

Eurora 2012

**CINECA**

Selex Elsag (E-security) 2011-2012

Deep project 2012

**JÜLICH** FORSCHUNGSZENTRUM

SELEX ELSAG
Secure Networking Solutions
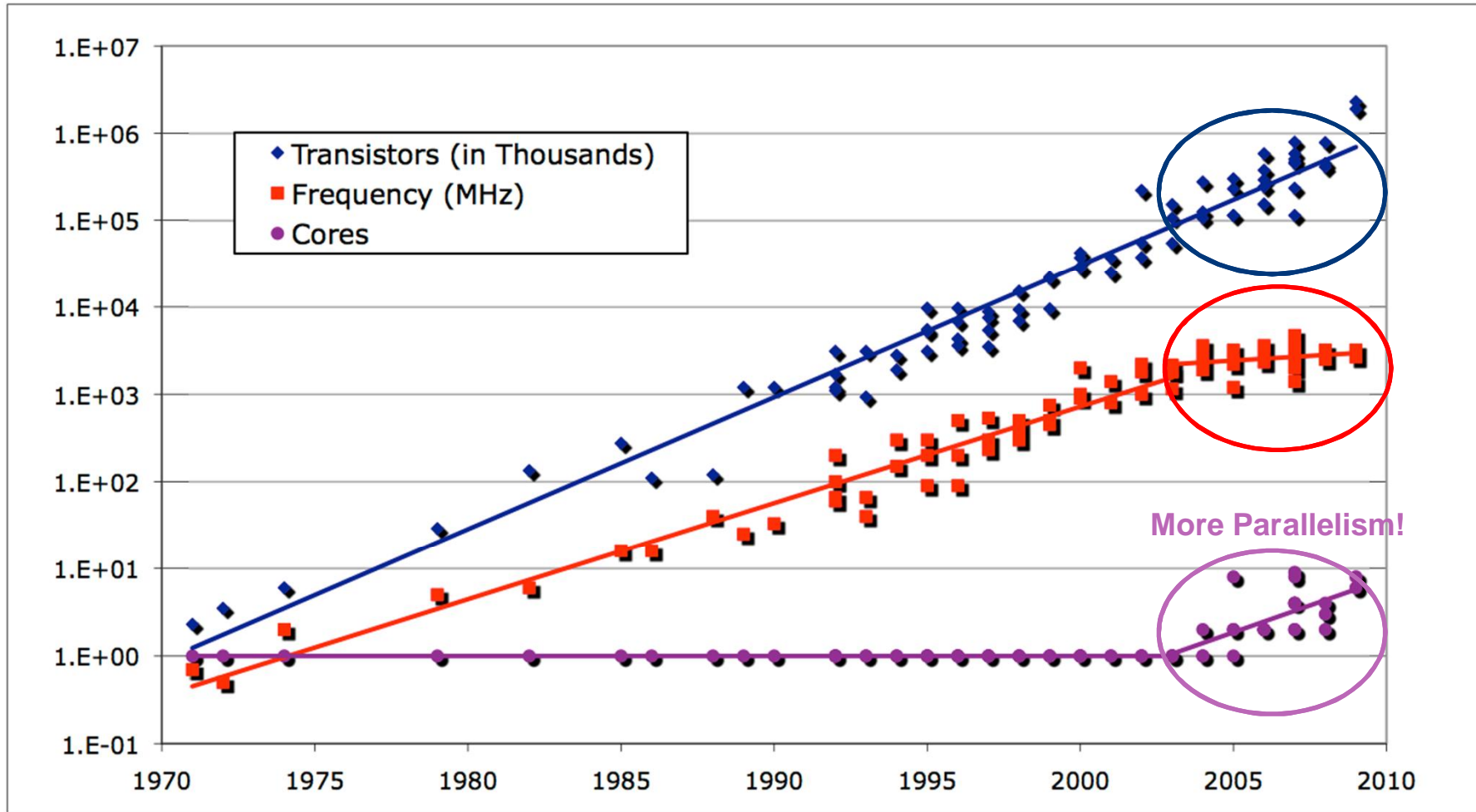A Finmeccanica Company

EUROTECH

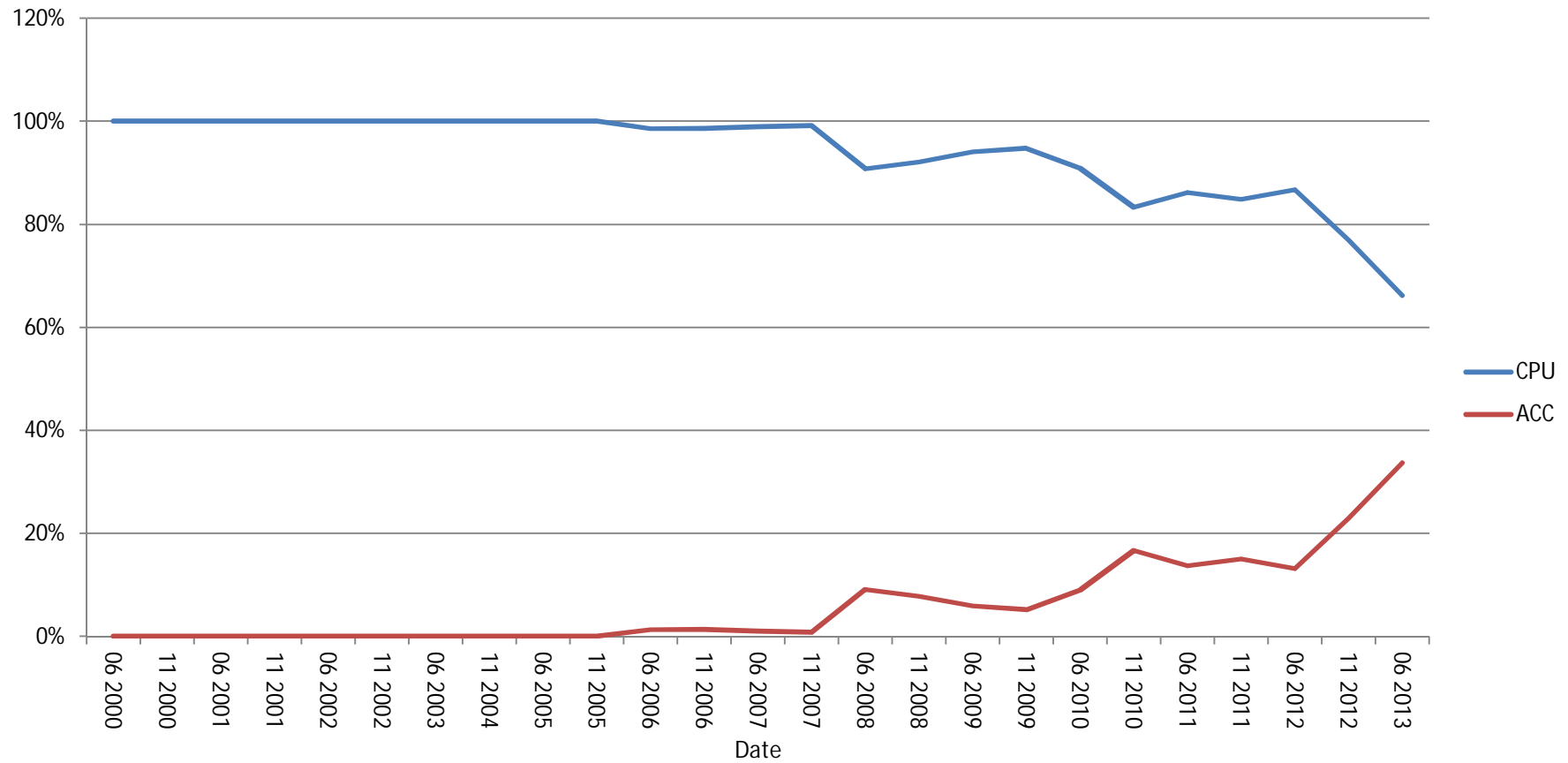# TRENDS IN HPC

# Top500 trend

# HPC Trend and Moore's Law
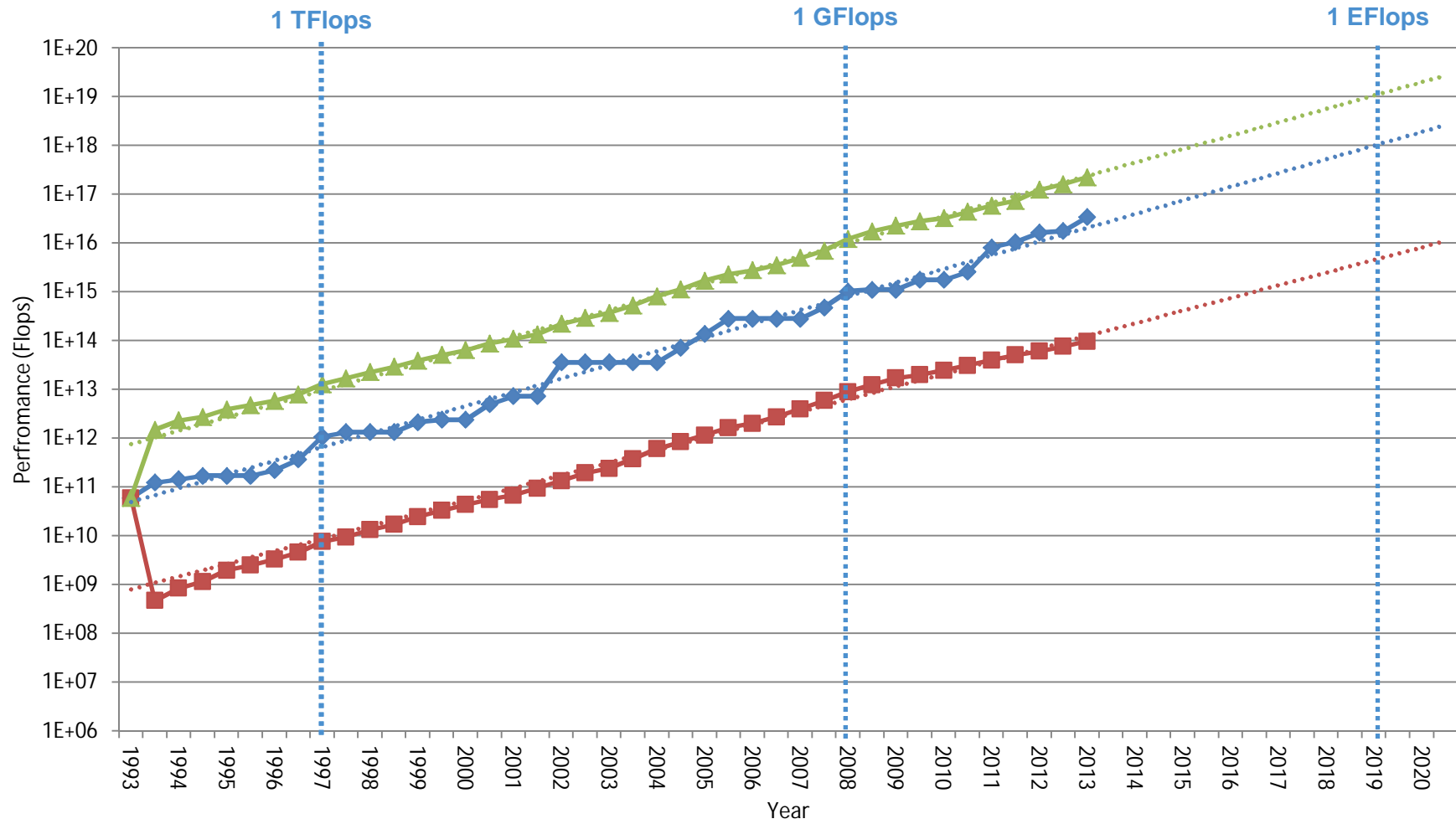


EUROTECH

# Scalar (CPU) vs Accelerated (ACC) computing
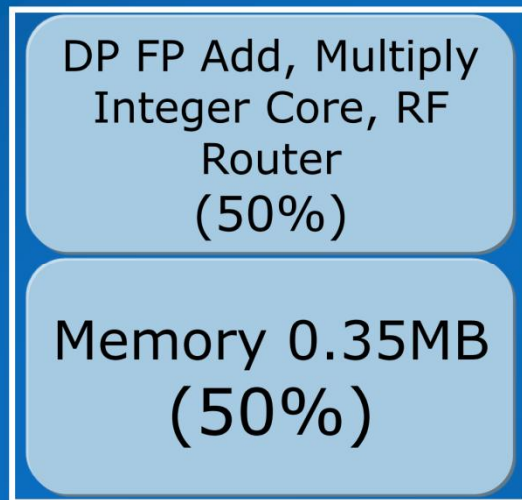
## (i.e. Beyond the Moore's Law)
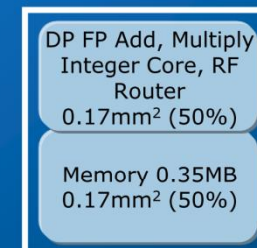


Total Performance

# Top500 trend

# Technology Outlook

| Technology (High Volume) | 45nm (2008) | 32nm (2010) | 22nm (2012) | 14nm (2014) | 10nm (2016) | 7nm (2018) | 5nm (2020) |
|---|---|---|---|---|---|---|---|
| Transistor density | 1.75 | 1.75 | 1.75 | 1.75 | 1.75 | 1.75 | 1.75 |
| Frequency scaling | 15% | 10% | 8% | 5% | 4% | 3% | 2% |
| Vdd scaling | -10% | -7.5% | -5% | -2.5% | -1.5% | -1% | -0.5% |
| Dimension & Capacitance | 0.75 | 0.75 | 0.75 | 0.75 | 0.75 | 0.75 | 0.75 |
| SD Leakage scaling/micron | 1X Optimistic to 1.43X Pessimistic | | | | | | |

45nm Core + Local Memory

| DP FP Add, Multiply Integer Core, RF Router (50%) |
|---|
| Memory 0.35MB (50%) |

7nm Core + Local Memory

| DP FP Add, Multiply Integer Core, RF Router 0.17mm$^2$ (50%) |
|---|
| Memory 0.35MB 0.17mm$^2$ (50%) |

~0.6mm

0.34mm2, 4.6GHz, 9.2GF, 0.24 to 0.46W

6mm2, 3.5GHz, 7GF, 1.2W

3

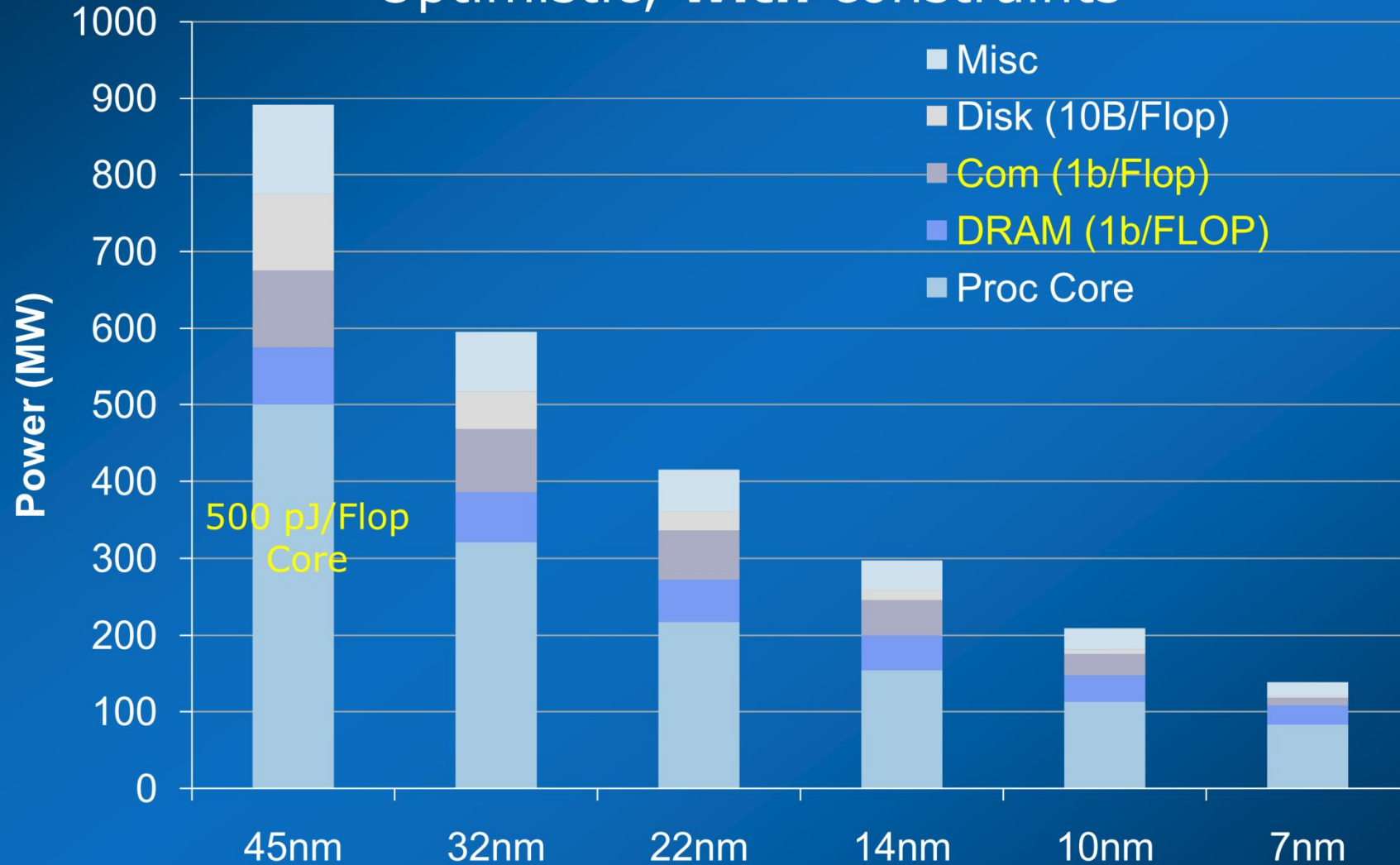# Energy per Operation



Source: Stekhar Borkar, Intel Corp., Aug 2011

# Peak Exascale System Power
## Optimistic, **with** constraints

Legend:
- Misc
- Disk (10B/Flop)
- Com (1b/Flop)
- DRAM (1b/FLOP)
- Proc Core

500 pJ/Flop Core

Y-axis: Power (MW), scale 0 to 1000

X-axis: 45nm, 32nm, 22nm, 14nm, 10nm, 7nm

6

# *Darpa Study 2008*
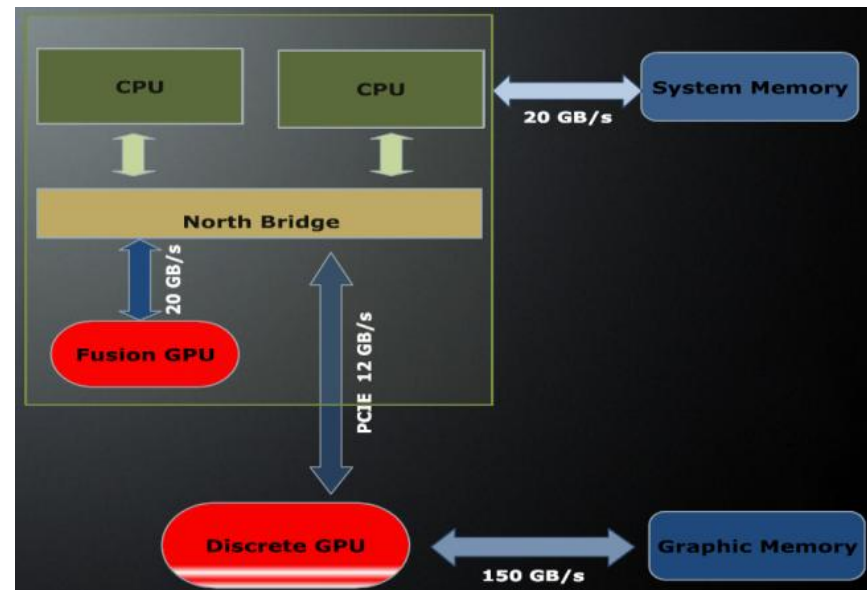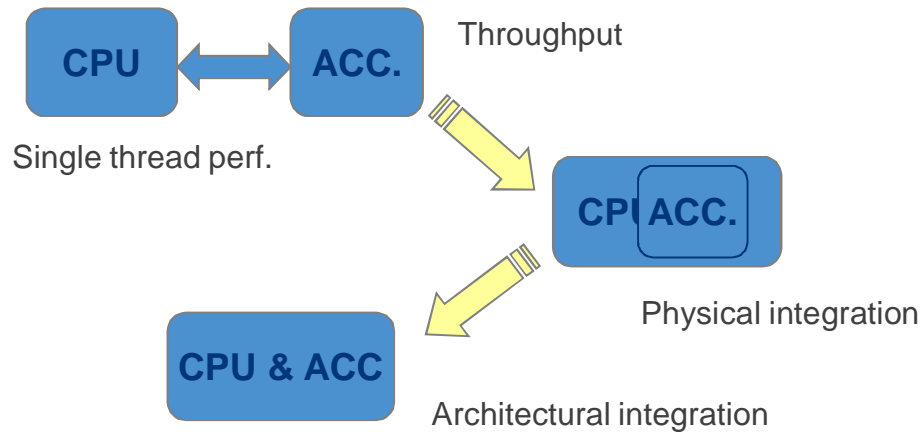## (i.e. more Energy-Efficient Computing)



From Peter Kogge,
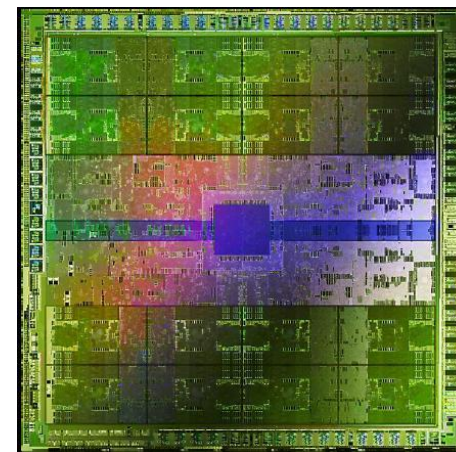DARPA Exascale Study

# ACCELERATED COMPUTING

# What is an Accelerator.

- A set (one or more) of very simple execution units that can perform few operations (with respect to standard CPU) with very high efficiency. When combined with full featured CPU (CISC or RISC) can accelerate the "nominal" speed of a system.
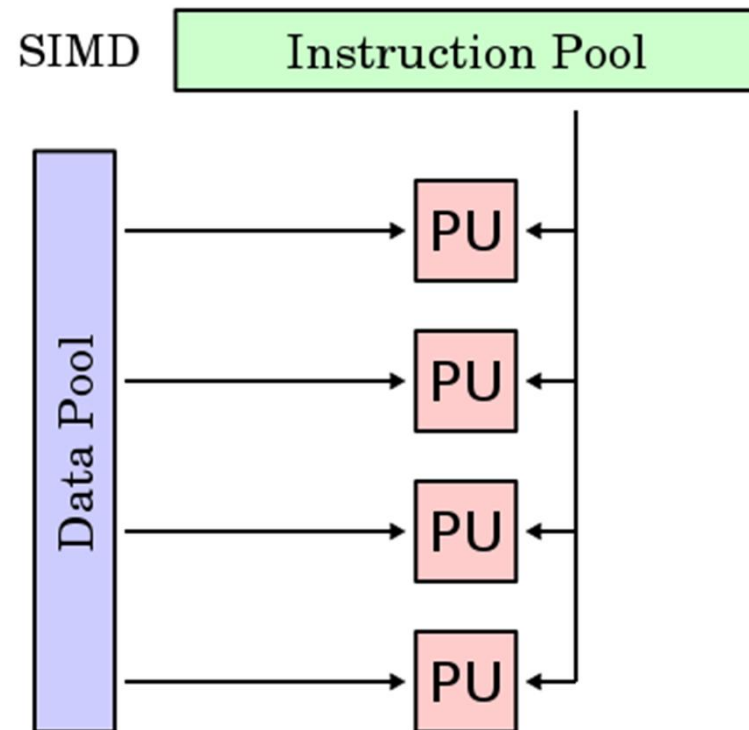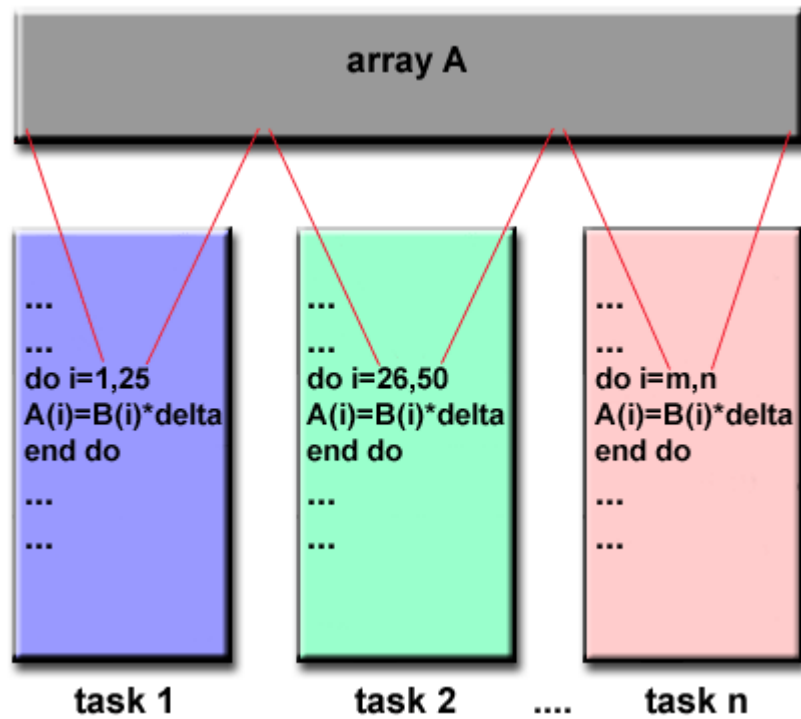


- Main approaches to accelerators:
  - ➢ Task Parallelism (MIMD) → MIC
  - ➢ Data Parallelism (SIMD) → GPU
  - ➢ Reconfigurable Devices

# GPU

# DATA Parallelism (SIMD)
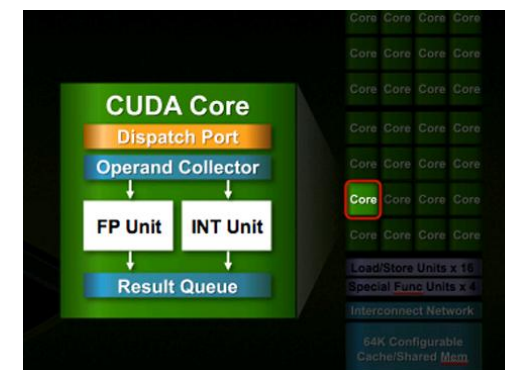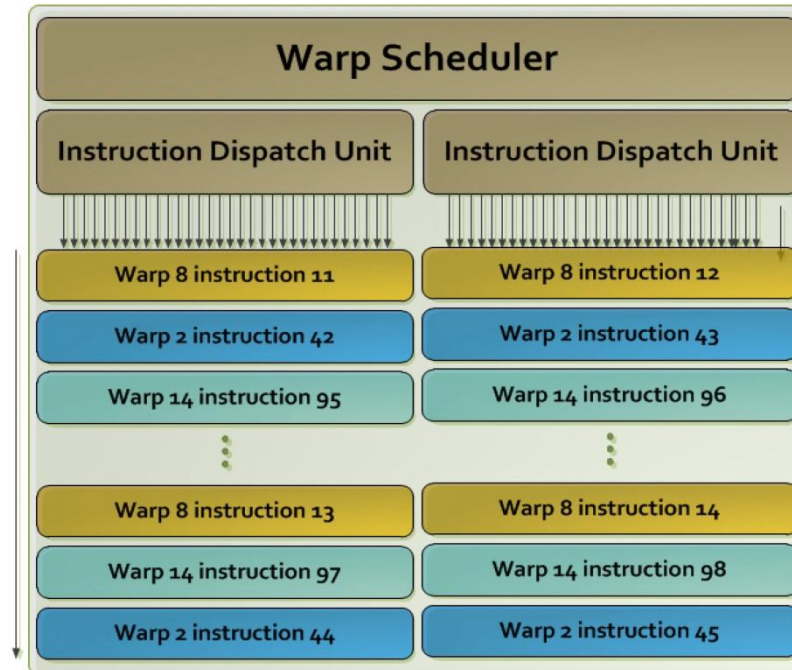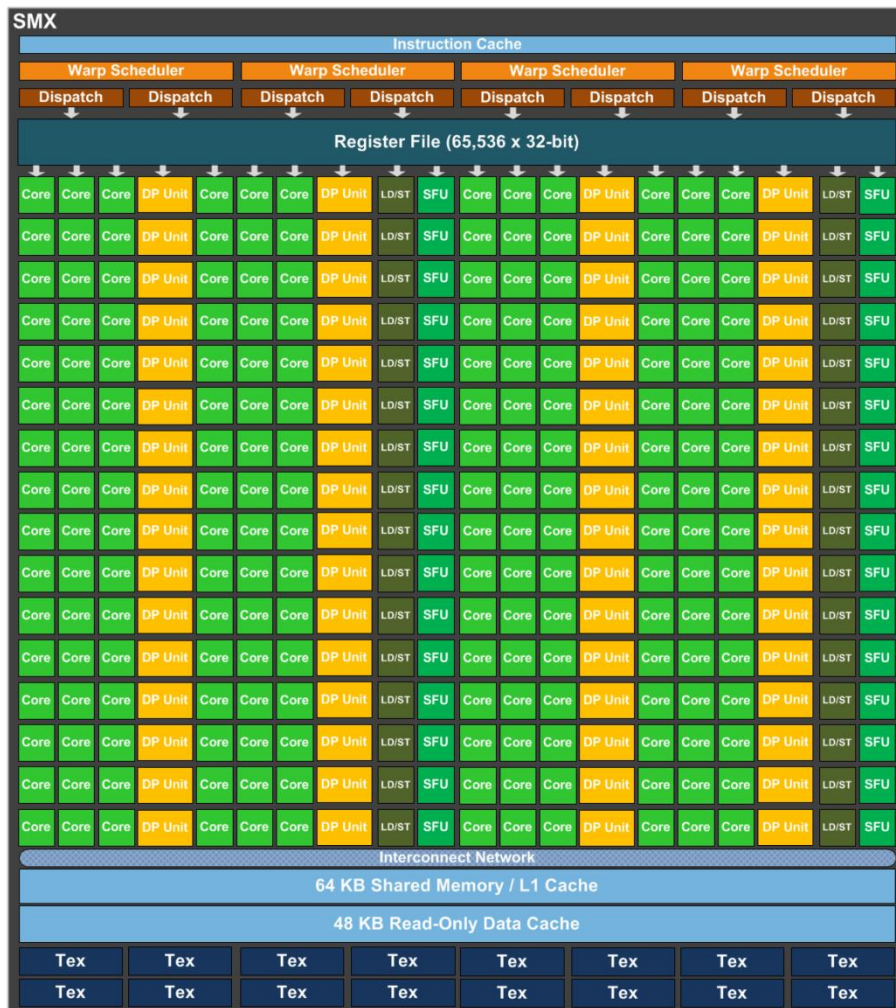
```
do i=1,10000
          A(i) = B(i) * delta
end do
```

# GigaThread Engine (Grid)

# SMX Processor & Warp Scheduler & Core

# DAXPY (Y ← a * X + Y) in action

```
// Invoke DAXPY
daxpy(n, 2.0, x, y);
// DAXPY in C
void daxpy(int n, double a, double *x, double *y)
{
    for (int i = 0; i < n; ++i)
        y[i] = a*x[i] + y[i];

}
```

1. A vectorizable (i.e. data parallel) loop is managed at **Grid** Level (pool of threads)
2. The body of a vectorizable loop is managed by a **thread block** an managed by the SMX processors
3. Each **thread** is managed by the warp scheduler that dynamically issues the PTX instruxctions in a pool of excution units (cores)

y[18] = a*x[18] + y[18]

| Type | More descriptive name | Closest old term outside of GPUs | Official CUDA/ NVIDIA GPU term | Book definition |
|---|---|---|---|---|
| **Program abstractions** | Vectorizable Loop | Vectorizable Loop | Grid | A vectorizable loop, executed on the GPU, made up of one or more Thread Blocks (bodies of vectorized loop) that can execute in parallel. |
| | Body of Vectorized Loop | Body of a (Strip-Mined) Vectorized Loop | Thread Block | A vectorized loop executed on a multithreaded SIMD Processor, made up of one or more threads of SIMD instructions. They can communicate via Local Memory. |
| | Sequence of SIMD Lane Operations | One iteration of a Scalar Loop | CUDA Thread | A vertical cut of a thread of SIMD instructions corresponding to one element executed by one SIMD Lane. Result is stored depending on mask and predicate register. |
| **Machine object** | A Thread of SIMD Instructions | Thread of Vector Instructions | Warp | A traditional thread, but it contains just SIMD instructions that are executed on a multithreaded SIMD Processor. Results stored depending on a per-element mask. |
| | SIMD Instruction | Vector Instruction | PTX Instruction | A single SIMD instruction executed across SIMD Lanes. |
| **Processing hardware** | Multithreaded SIMD Processor | (Multithreaded) Vector Processor | Streaming Multiprocessor | A multithreaded SIMD Processor executes threads of SIMD instructions, independent of other SIMD Processors. |
| | Thread Block Scheduler | Scalar Processor | Giga Thread Engine | Assigns multiple Thread Blocks (bodies of vectorized loop) to multithreaded SIMD Processors. |
| | SIMD Thread Scheduler | Thread scheduler in a Multithreaded CPU | Warp Scheduler | Hardware unit that schedules and issues threads of SIMD instructions when they are ready to execute; includes a scoreboard to track SIMD Thread execution. |
| | SIMD Lane | Vector Lane | Thread Processor | A SIMD Lane executes the operations in a thread of SIMD instructions on a single element. Results stored depending on mask. |

# DAXPY (Y ← a * X + Y) in CUDA

```
// Invoke DAXPY
daxpy(n, 2.0, x, y);
// DAXPY in C
void daxpy(int n, double a, double *x, double *y)
{
        for (int i = 0; i < n; ++i)
                y[i] = a*x[i] + y[i];
}
```

```
// Invoke DAXPY with 256 threads per Thread Block
__host__
int nblocks = (n+ 255) / 256;
    daxpy<<<nblocks, 256>>>(n, 2.0, x, y);
// DAXPY in CUDA
__device__
void daxpy(int n, double a, double *x, double *y)
{
    int i = blockIdx.x*blockDim.x + threadIdx.x;
    if (i < n) y[i] = a*x[i] + y[i];
}
```
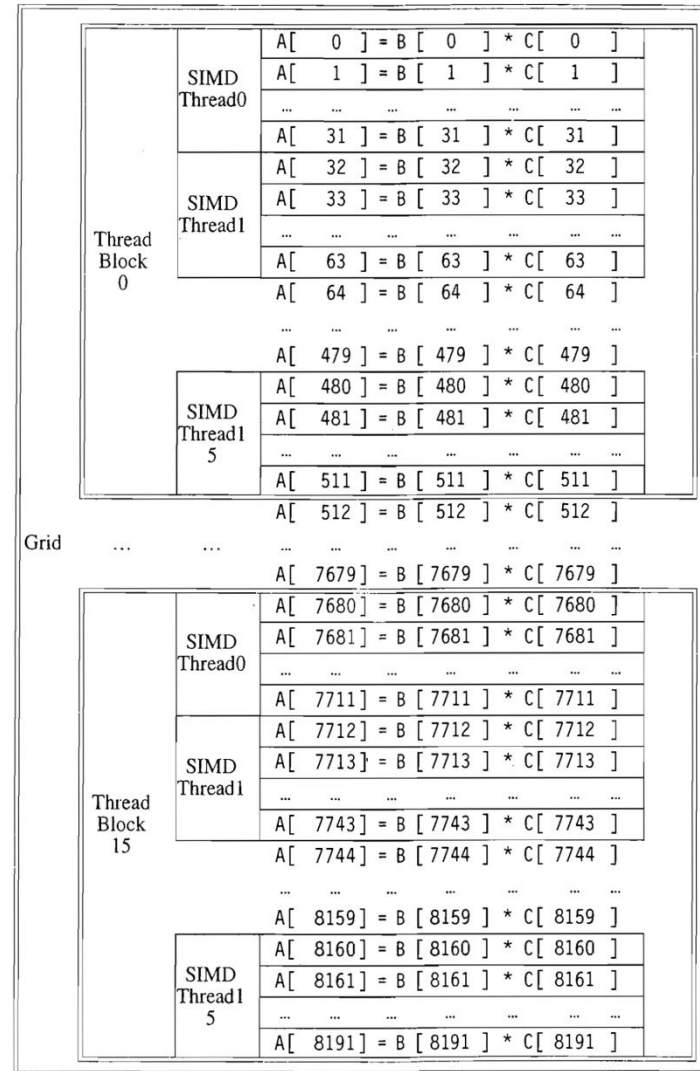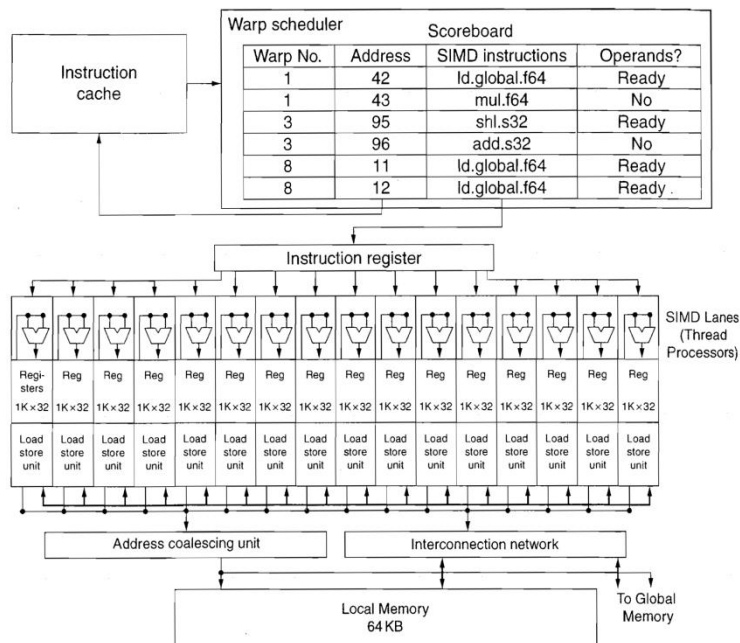
1. A vectorizable (i.e. data parallel) loop is managed at Grid Level (pool of threads)
2. The body of a vectorizable loop is managed by a thread block an managed at SMX level
3. Each thread is managed by the warp scheduler that dynamically issues the PTX instruxctions in a pool of excution units (cores)

- __*host*__ instructs the compiler to generate code and to allocate data in the CPU scope
- __*device*__ instructs the compiler to generate code and to allocate data in the GPU scope

# DAXPY Summary

```
// Invoke DAXPY
daxpy(n, 2.0, x, y);
// DAXPY in C
void daxpy(int n, double a, double *x, double *y)
{
        for (int i = 0; i < n; ++i)
                y[i] = a*x[i] + y[i];
}
```

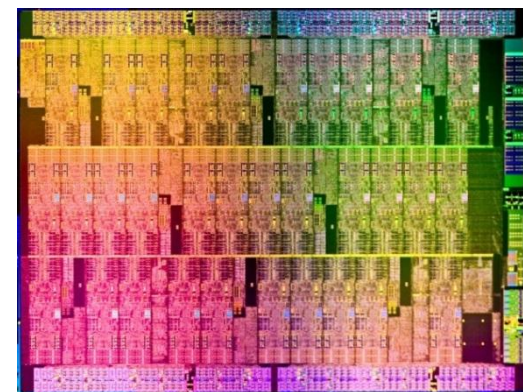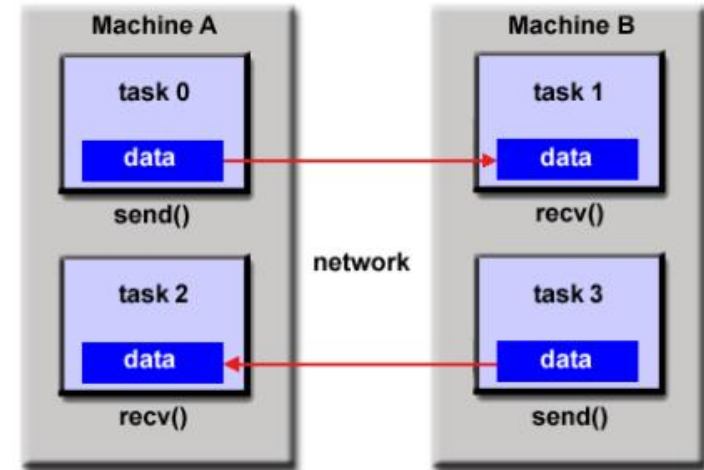# Execution Model



Processing flow on CUDA

# MIC

# TASK Parallelism (MIMD)

# Xeon PHI Architecture



Typical Platform consists of:
1 to 2 Intel Xeon processors (CPUs)
1 to 8 Intel Xeon Phi Coprocessors per host

EUROTECH

# Core Architecture



- Up to 32 in-order cores
- 4 hardware threads per core
- Two pipelines
  - Pentium® processor family-based scalar units
  - Fully-coherent L1 and L2 caches
  - 64-bit addressing
- All new vector unit
  - 512-bit SIMD Instructions – not Intel® SSE, MMX™, or Intel® AVX
  - 32 512-bit wide vector registers
    - Hold 16 singles or 8 doubles per register
  - Pipelined one-per-clock throughput
    - 4 clock latency, hidden by round-robin scheduling of threads
  - Dual issue with scalar instructions

# A versatile combination

# Execution Models

## Offload Execution Mode

1. Host system offloads part or all of the computation from one or multiple processes or threads running on host

2. The application starts execution on the host

3. As the computation proceeds it can decide to send data to the coprocessor and let that work on it and the host and the coprocessor may or may not work in parallel.

OpenMP 4.0 TR being proposed and implemented in Intel® Composer XE provides directives to perform offload computations. Composer XE also provides some custom directives to perform offload operations.

# Execution Models
## Coprocessor Native Execution Mode

- An Xeon Phi hosts a Linux micro OS in it and can appear as another machine connected to the host like another node in a cluster.

- This execution environment allows the users to view the coprocessor as another compute node.

- In order to run natively, an application has to be cross compiled for Xeon Phi operating environment. Intel® Composer XE provides simple switch to generate cross compiled code.

EUROTECH

# Execution Models
## Symmetric Execution

- The application processes run on both the host and the Phi coprocessor and communicate through some sort of message passing interface like MPI.

- This execution environment treats Xeon Phi card as another node in a cluster in a heterogeneous cluster environment.

# Execution Models
## Summary

# Programming PHI

```
1. Offloading a function call
#pragma offload target (mic)
    foo();
```

```
foo() { .... } // Compiled for mic
```

```
2. Calculating Pi with automatic offload
#pragma offload target (mic)
#pragma omp parallel for reduction(+:pi)
 for (i=0; i<count; i++)
 {
     float t = (float)((i+0.5)/count);
     pi += 4.0/(1.0+t*t);
 }
pi /= count
```

```
3. Using MKL with offload
void your_hook()
{
    float *A, *B, *C; /* Matrices */
    #pragma offload target(mic)
     in(transa, transb, N, alpha, beta) \
     in(A:length(matrix_elements)) \
     in(B:length(matrix_elements)) \
     in(C:length(matrix_elements)) \
     out(C:length(matrix_elements)alloc_if(0))
     sgemm(&transa, &transb, &N, &N,
            &N, &alpha, A, &N, B, &N, &beta, C,
     &N);
}
```

# Heterogeneous Compiler



**Source Code**

```
main()
{
f();
}
```

```
f()
{
    #pragma offload
    a = b + g();
}
```

```
__attribute__
((target(mic))) g()
{
}
```

**Linux* Host Program**

```
main()
{
copy_code_to_mic();
f();
unload_mic();
}
```

```
f() {
    if (mic_available()){
        send_data_to_mic();
        start f_part_mic();
        receive_data_from_mic();
    } else
        f_part_host();
}
```

```
f_part_host()
{a = b + g();}
```

```
g() {...}
```

**Intel ®MIC Program**

This all happens automatically when you issue a single compile command

```
f_part_mic()
{a = b + g_mic();}
```

```
g_mic() {...}
```

EUROTECH

# ENERGY-EFFICIENT COMPUTING

# ARM

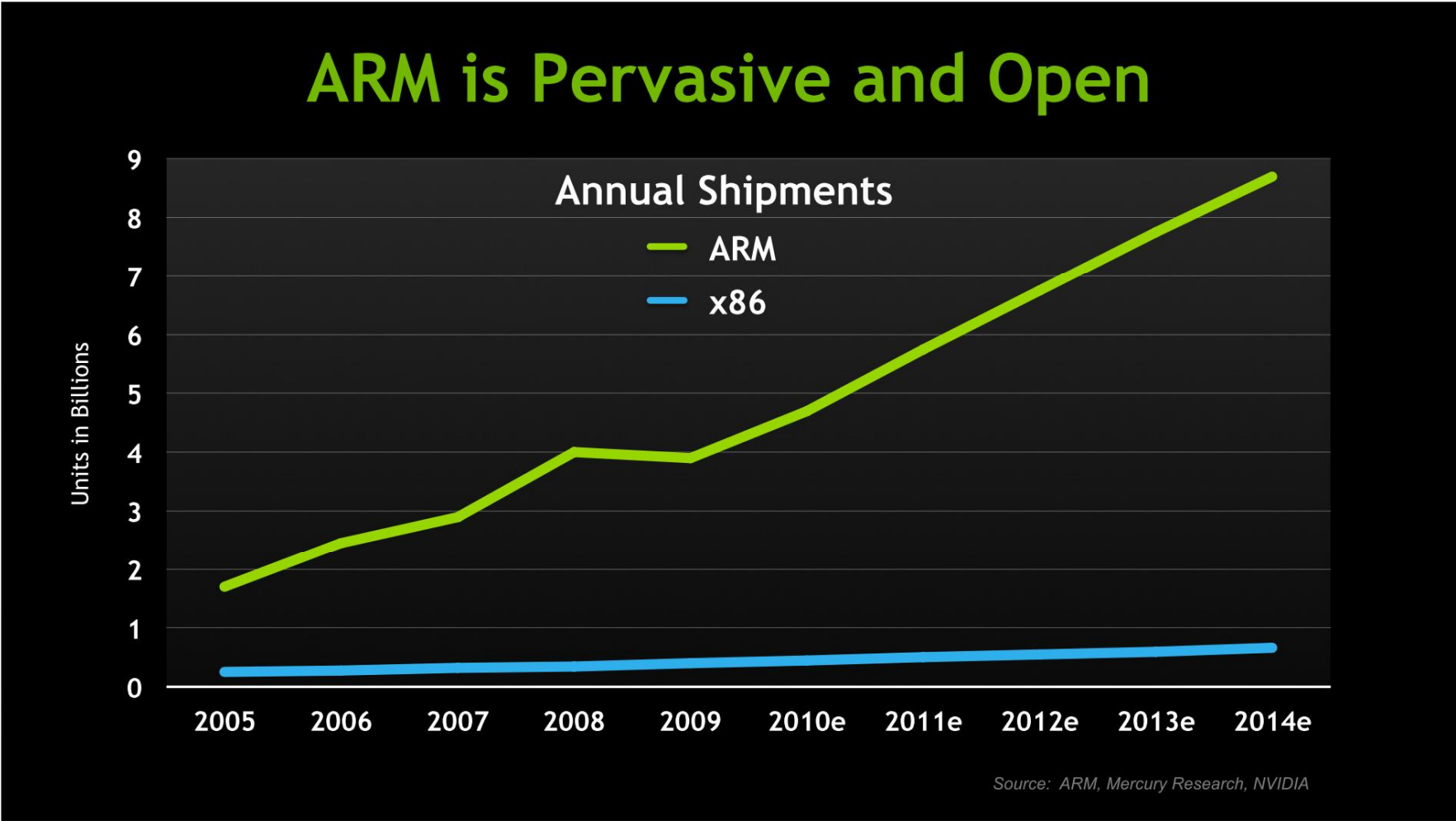# Next Step in the Comodity Chain

- Total cores in Jun'12 Top500
  - 13.5 Mcores
- Tablets sold in Q4 2011
  - 27 Mtablets
- Smartphones sold Q4 2011
  - > 100 Mphones

HPC

Servers

Desktop

Mobile

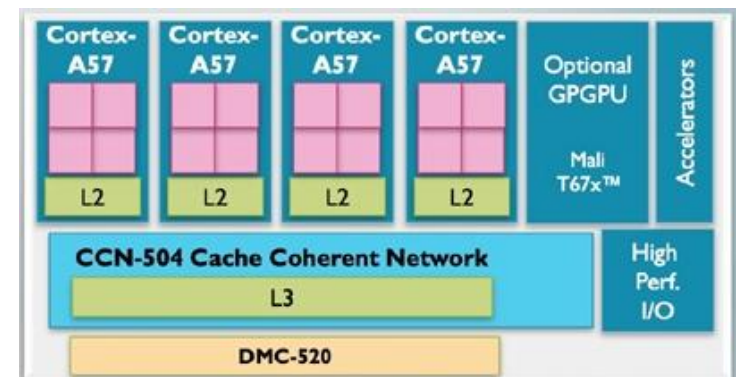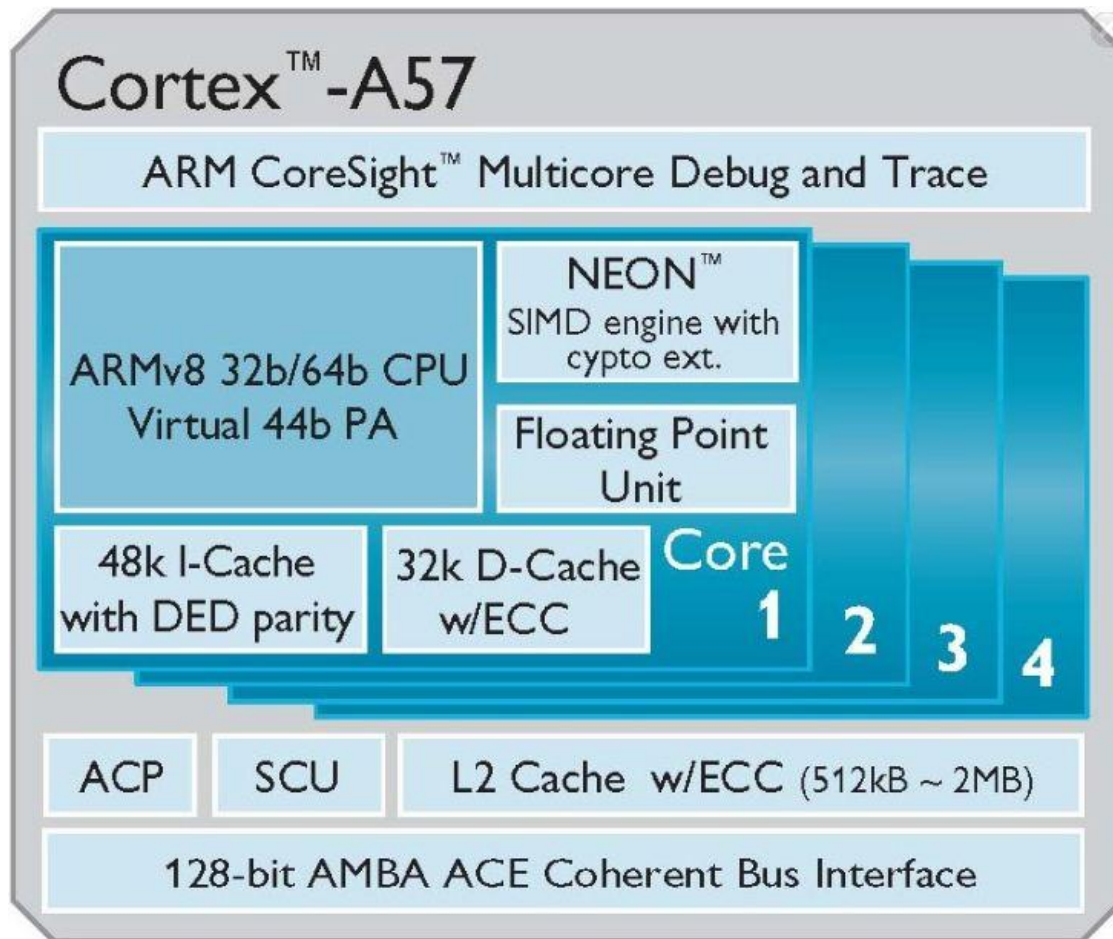# ARM is the most sold CPU Architecture



ARM is Pervasive and Open

Annual Shipments
— ARM
— x86

Units in Billions

2005 2006 2007 2008 2009 2010e 2011e 2012e 2013e 2014e

Source: ARM, Mercury Research, NVIDIA

# ARM vs x86

**Table 1. Summary of RISC and CISC Trends.**

| | Format | Operations | Operands |
|---|---|---|---|
| **RISC / ARM** | ○ Fixed length instructions<br>○ Relatively simple encoding<br>○ ARM: 4B, THUMB(2B, optional) | ○ Simple, single function operations<br>○ Single cycle | ○ Operands: registers, immediates<br>○ Few addressing modes<br>○ ARM: 16 general purpose registers |
| **CISC / x86** | ○ Variable length instructions<br>○ Common insts shorter/simpler<br>○ Special insts longer/complex<br>○ x86: from 1B to 16B long | ○ Complex, multi-cycle instructions<br>○ Transcendentals<br>○ Encryption<br>○ String manipulation | ○ Operands: memory, registers, immediates<br>○ Many addressing modes<br>○ x86: 8 32b & 6 16b registers |
| **Historical Contrasts** | ○ CISC decode latency prevents pipelining<br>○ CISC decoders slower/more area<br>○ Code density: RISC < CISC | ○ Even w/ $\mu$code, pipelining hard<br>○ CISC latency may be longer than compiler's RISC equivalent | ○ CISC decoder complexity higher<br>○ CISC has more per inst work, longer cycles<br>○ Static code size: RISC > CISC |
| **Convergence Trends** | ○ **$\mu$-op cache minimizes decoding overheads**<br>○ **x86 decode optimized for common insts**<br>○ **I-cache minimizes code density impact** | ○ **CISC insts split into RISC-like micro-ops; optimizations eliminated inefficiencies**<br>○ **Modern compilers pick mostly RISC insts; $\mu$-op counts similar for ARM and x86** | ○ **x86 decode optimized for common insts**<br>○ **CISC insts split into RISC-like micro-ops; x86 and ARM $\mu$-op latencies similar**<br>○ **Number of data cache accesses similar** |
| **Empirical Questions** | ○ How much variance in x86 inst length?<br>   Low variance $\Rightarrow$ common insts optimized<br>○ Are ARM and x86 code densities similar?<br>   Similar density $\Rightarrow$ No ISA effect<br>○ What are instruction cache miss rates?<br>   Low $\Rightarrow$ caches hide low code densities | ○ Are macro-op counts similar?<br>   Similar $\Rightarrow$ RISC-like on both<br>○ Are complex instructions used by x86 ISA?<br>   Few complex $\Rightarrow$ Compiler picks RISC-like<br>○ Are $\mu$-op counts similar?<br>   Similar $\Rightarrow$ CISC split into RISC-like $\mu$-ops | ○ Number of data accesses similar?<br>   Similar $\Rightarrow$ no data access inefficiencies |

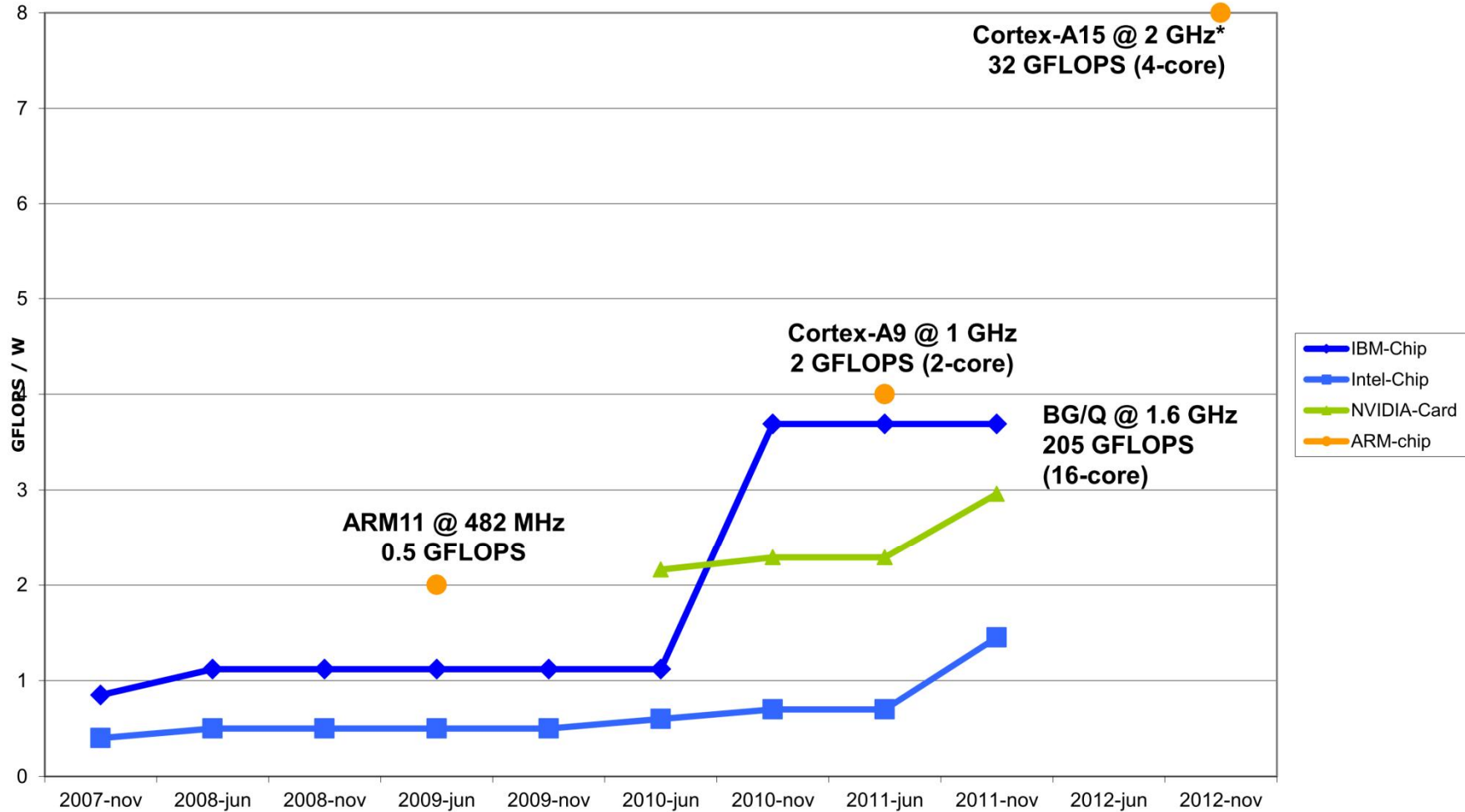# High Performance ARM CPU



EUROTECH

# ARM Improvement in DFP



- IBM BG/Q and Intel AVX implement DP in 256-bit SIMD
  - 8 DP ops / cycle
- ARM quickly moved from optional floating-point to state-of-the-art
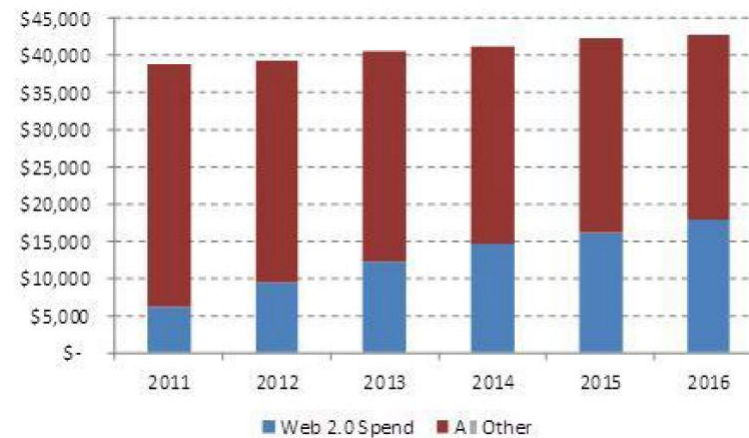  - ARMv8 ISA introduces DP in the NEON instruction set (128-bit SIMD)

# ARM Energy Efficiency

# SM - Summary

| Vendor | 2011 Revenue | Market Share | 2010 Revenue | Market Share | 2011/2010 Revenue Growth |
|--------|-------------|--------------|--------------|--------------|--------------------------|
| 1. IBM | $16,456 | 31.5% | $15,342 | 31.1% | 7.3% |
| 2. HP | $15,301 | 29.3% | $15,388 | 31.2% | -0.6% |
| 3. Dell | $7,814 | 15.0% | $7,106 | 14.4% | 10.0% |
| 4. Oracle | $3,215 | 6.2% | $3,204 | 6.5% | 0.3% |
| 5. Fujitsu | $2,516 | 4.8% | $2,197 | 4.4% | 14.5% |
| Others | $6,964 | 13.3% | $6,159 | 12.5% | 13.1% |

According to IDC, there were 8.05M x86 servers shipped in 2011. At a blended average price of $4,837, the x86 server TAM was $39B.



Small Socket Servers: 21% of Server Market by 2016 95% CAGR

Source: Oppenheimer, Cloudy with a Chance of ARM, March 30th 2012



Web 2.0 Spend    All Other

# FPGAs

# Altera Stratix IV



**Altera Stratix IV GX 230 / 1152 pins**

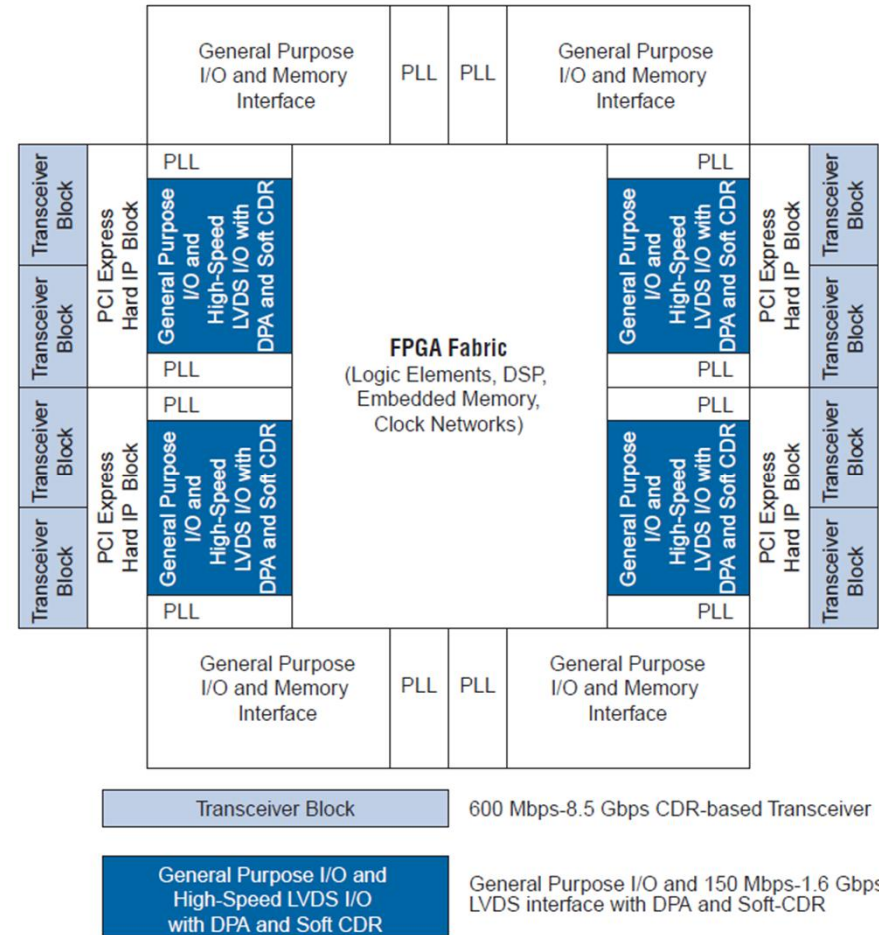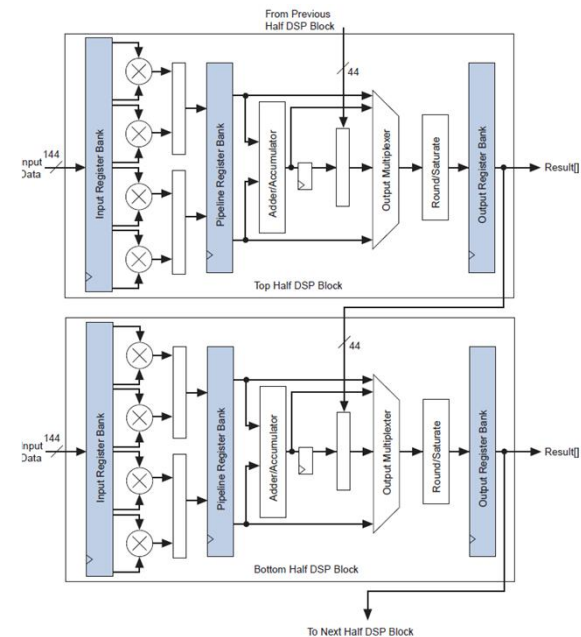# Stratix IV FPGA main features

- 91200 ALMs (Adaptive Logic Modules)
- 18 MB SRAM
- 228000 Logic Elements (~ gates)
- High speed transceivers
- 2 PCIe 2.0 hard IP blocks
- 161 DSP blocks,
  (322 36*36 mult 1288 18*18 mult. etc
  up to 55GFlops Matrix Inversion)
  All of them can be used in parallel

# The Aurora Tigon node card

| Node card |
| --- |
| 2 X Intel Xeon E5 series |
| 2 X Nvidia Kepler K20 or |
| 2 X Intel Xeon Phi 5120D |
| 2730 GFlops |
| 800 W power consumption |
| 64 GB soldered memory |



- The **node card** is the main processing unit of an Eurotech Aurora Tigon machine

- An aluminum cold plate that cools the board smoothing temperature distributions and assuring maximum heat extraction efficacy

- A large, high end FPGA allow implementation of a point to point 3D-Torus network

# Aurora Tigon – Node Card Architecture



DDR3 DDR3 DDR3 DDR3

**SDB-EP (R)**

**SBD-EP (R)**

DDR3 DDR3 DDR3 DDR3

**Patsburg**

USB, GbE, SATA

x16

X 8

X 8

X 8

x16

PCIe3 x16

**FPGA**   Acceleration Engine

3D Torus NIC/Router   Global Sync

ConnectX Infiniband

PCIe3 x16

**PCIe3 x16 Expansion slot**
for GPU and other external devises

X   Y   Z

3D Torus:
160(x)+160(y)+160(z) =480 Gbps

Sync Tree
Barriers, jitter sync, etc

FDR Infiniband
**56 Gbps**

**PCIe3 x16 Expansion slot**
for GPU and other external devises

QPI          PCIe3          DMI 4x + PCIe3 4x

**EUROTECH**

# The Aurora Tigon node card

Cooling Plate

Infiniband QDR

Stratix V - FPGA

2 x Nvidia K20

2 x Intel Xeon E5 sockets

Local Disk

Soldered RAM

# Aurora Accelerator core



non reconfigurable block

IP block

Aurora FPGA Device

Application

API

FPGA i/O Driver

Aurora OS

I/O Interface

Algorithm executed in hardware

FPGA Configuration Flash memory

Link registers: Specific to programming environments: key to easy integration with C code to VHDL (hardware instantiation) programming interface
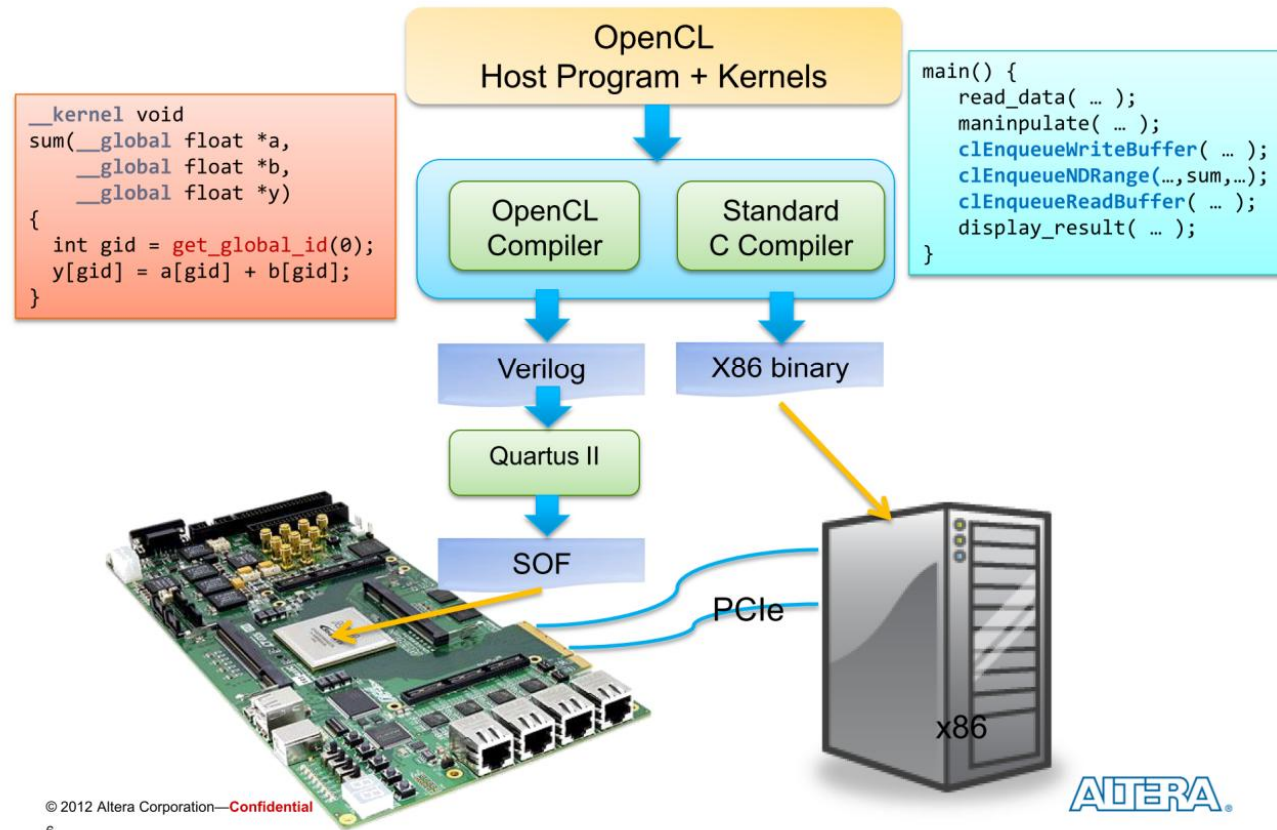
SW FPGA Configuration using Aurora node PCIe interfaces

EUROTECH

# OpenCL and FPGAs



Altera OpenCL Compile Flow

# Energy Efficient FP Computing

| Example | Device | Configuration (Channel Size/ Matrix Size/ Vector Size) | Throughput (kMatrices/ sec) | Fmax (MHz) | Performance (GFLOPS) | Total Power[1] (W) | GFLOPS/ W |
|---|---|---|---|---|---|---|---|
| Cholesky | Stratix V | 1 / 360×360 / 90 | 1.43 | 189 | 91 | 16 | 5.7 |
| | | 20 / 60×60 / 60 | 118.35 | 234 | 39 | 15 | 2.6 |
| | | 64 / 30×30 / 30 | 544.28 | 288 | 26 | 10 | 2.5 |
| | Arria V | 6 / 90×90 / 45 | 35.22 | 197 | 38 | 9.1 | 4.2 |
| | | 64 / 30×30 / 30 | 349.62 | 184 | 16 | 7.1 | 2.3 |
| QR | Stratix V | 1 / 400×400 / 100 | 0.315 | 203 | 162 | 26 | 6.2 |
| | | 1 / 200×100 / 100 | 8.76 | 207 | 141 | 23 | 6.1 |
| | | 1 / 200×100 / 50 | 6.17 | 260 | 99 | 16 | 6.2 |
| | | 1 / 100×50 / 50 | 32.82 | 259 | 66 | 13 | 5.1 |
| | Arria V | 1 / 200×100 / 50 | 4.05 | 171 | 65 | 9.1 | 7.1 |
| | | 1 / 100×50 / 50 | 21.54 | 170 | 44 | 8.1 | 5.4 |

Table 1. Power efficiency of Stratix V and Arria V FPGAs running Cholesky and QR solvers.

[1] Power values have an error margin of ±1 %.

EUROTECH

# Thank You

www.eurotech.com