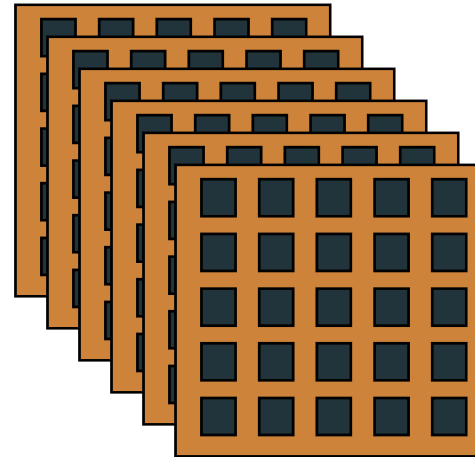# Multicore:
## Why is it happening now?
### eller
### Hur Mår Moore's Lag?

Erik Hagersten

Uppsala Universitet
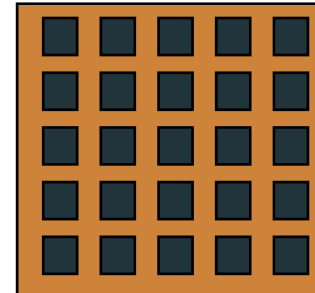
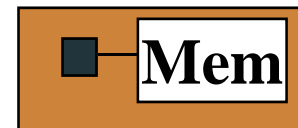# Darling, I shrunk the computer

**Mainframes**

**Super Minis:**

**Microprocessor:**

**Multicore: Many CPUs on a chip!**

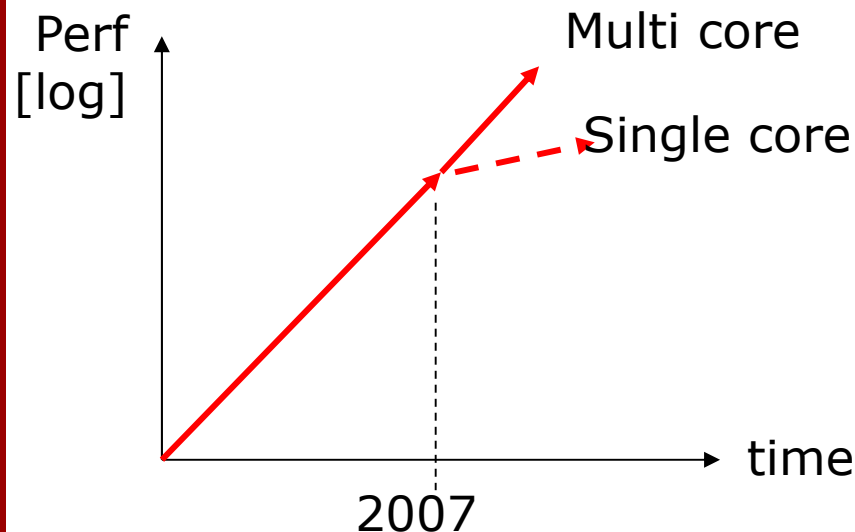Parallel
Comp
2012

# Multi-core CPUs:

- **Ageia** PhysX, a multi-core physics processing unit.
- **Ambric** Am2045, a 336-core Massively Parallel Processor Array (MPPA)
- **AMD**
  - Athlon 64, Athlon 64 FX and Athlon 64 X2 family, dual-core desktop processors.
  - Opteron, dual- and quad-core server/workstation processors.
  - Phenom, triple- and quad-core desktop processors.
  - Sempron X2, dual-core entry level processors.
  - Turion 64 X2, dual-core laptop processors.
  - Radeon and FireStream multi-core GPU/GPGPU (10 cores, 16 5-issue wide superscalar stream processors per core)
- **ARM** MPCore is a fully synthesizable multicore container for ARM9 and ARM11 processor cores, intended for high-performance embedded and entertainment **applications**.
- **Azul** Systems Vega 2, a 48-core processor.
- **Broadcom** SiByte SB1250, SB1255 and SB1455.
- **Cradle** Technologies CT3400 and CT3600, both multi-core DSPs.
- **Cavium** Networks Octeon, a 16-core MIPS MPU.
- **HP** PA-8800 and PA-8900, dual core PA-RISC processors.
- **IBM**
  - POWER4, the world's first dual-core processor, released in 2001.
  - POWER5, a dual-core processor, released in 2004.
  - POWER6, a dual-core processor, released in 2007.
  - PowerPC 970MP, a dual-core processor, used in the Apple Power Mac G5.
  - Xenon, a triple-core, SMT-capable, PowerPC microprocessor used in the Microsoft Xbox 360 game console.
- **IBM**, Sony, and Toshiba Cell processor, a nine-core processor with one general purpose PowerPC core and eight specialized SPUs (Synergystic Processing Unit) **optimized** for vector operations used in the Sony PlayStation 3.
- **Infineon** Danube, a dual-core, MIPS-based, home gateway processor.
- **Intel**
  - Celeron Dual Core, the first dual-core processor for the budget/entry-level market.
  - Core Duo, a dual-core processor.
  - Core 2 Duo, a dual-core processor.
  - Core 2 Quad, a quad-core processor.
  - Core i7, a quad-core processor, the successor of the Core 2 Duo and the Core 2 Quad.
  - Itanium 2, a dual-core processor.
  - Pentium D, a dual-core processor.
  - Teraflops Research Chip (Polaris), an 3.16 GHz, 80-core processor prototype, which the company says will be released within the next five years[6].
  - Xeon dual-, quad- and hexa-core processors.
- **IntellaSys** seaForth24, a 24-core processor.
- **Nvidia**
  - GeForce 9 multi-core GPU (8 cores, 16 scalar stream processors per core)
  - GeForce 200 multi-core GPU (10 cores, 24 scalar stream processors per core)
  - Tesla multi-core GPGPU (8 cores, 16 scalar stream processors per core)
- Parallax Propeller P8X32, an eight-core microcontroller.
- picoChip PC200 series 200-300 cores per device for DSP & wireless
- Rapport Kilocore KC256, a 257-core microcontroller with a PowerPC core and 256 8-bit "processing elements".
- Raza Microelectronics XLR, an eight-core MIPS MPU
- **Sun Microsystems**
  - UltraSPARC IV and UltraSPARC IV+, dual-core processors.
  - UltraSPARC T1, an eight-core, 32-thread processor.
  - UltraSPARC T2, an eight-core, 64-concurrent-thread processor.
- Texas Instruments TMS320C80 MVP, a five-core multimedia video processor.
- Tilera TILE64, a 64-core processor
- XMOS Software Defined Silicon quad-core XS1-G4

[source: Wikipedia]

MC 3

Parallel Comp 2012

UPPSALA UNIVERSITET

# Outline

- Why multicore now?
- Performance bottlenecks in MCs
- Commercial offerings
- Reflection for the future

Parallel
Comp
2012

Institutionen för informationsteknologi | www.it.uu.se

MC 4

© Erik Hagersten| user.it.uu.se/~eh

# Everybody is doing is!
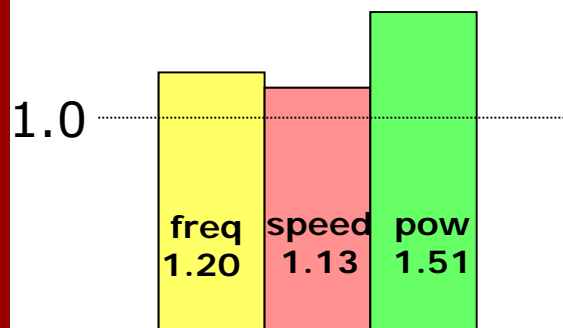## But, why now?

Perf [log]

Multi core

Single core

time

2007

1. Not enough ILP to get payoff from using more transistors

2. Signal propagation delay » transistor delay

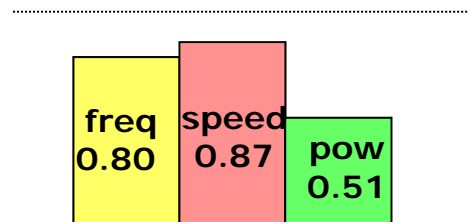3. Power consumption $P_{dyn} \sim C \bullet f \bullet V^2$

MC 5

Parallel Comp 2012

# Example: Freq. Scaling

$$P_{dyn} = C * f * V^2 \approx area * freq * voltage^2$$



1.0

| | | |
|---|---|---|
| freq 1.20 | speed 1.13 | pow 1.51 |

**20% higher freq.**

| | | |
|---|---|---|
| freq 0.80 | speed 0.87 | pow 0.51 |

**20% lower freq.**

| | | |
|---|---|---|
| freq 0.80 | speed 0.87 / speed 0.87 | pow 0.51 / pow 0.51 |

**20% lower freq. Two cores**

Parallel
Comp
2012

UPPSALA
UNIVERSITET

# Darling, I shrunk the computer

**Mainframes**

**Super Minis:**

**Microprocessor:** ■—|Mem|

**Chip Multiprocessor (CMP):**
**A multiprocessor on a chip!** ▦—|Mem|

Sequential execution (≈ one program)

**Parallel Apps (TLP)**

MC 7

# Shared Bottlenecks
## (the MCs on these slides are generalized)

CPU | CPU | CPU | CPU

L1 | L1 | L1 | L1

**L2 Cache**

L1 | L1 | L1 | L1

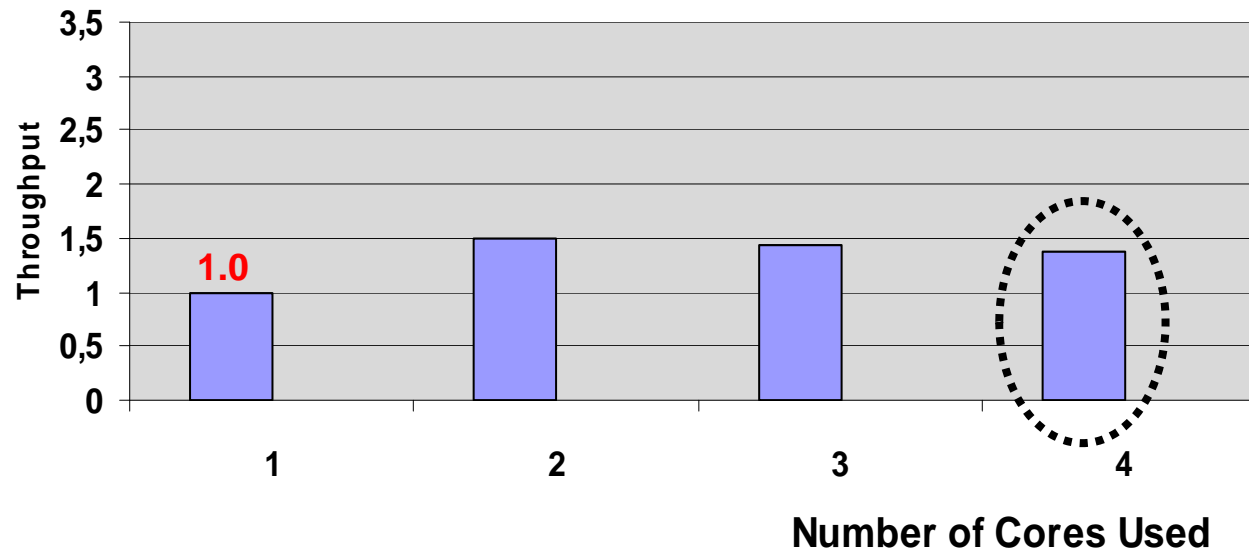CPU | CPU | CPU | CPU

Bandwidth

- Cache sharing effects
- Cache utilization
    - Loop order
    - Data allocation
    - Data usage...
- Data reuse
    - Tiling (aka Blocking)
    - Fusion...

**Shared Resources**

MC 8

Institutionen för informationsteknologi | www.it.uu.se

© Erik Hagersten | user.it.uu.se/~eh

# Example: Poor Throughput Scaling!

Example: 470.LBM
"Lattice Boltzmann Method" to simulate incompressible fluids in 3D



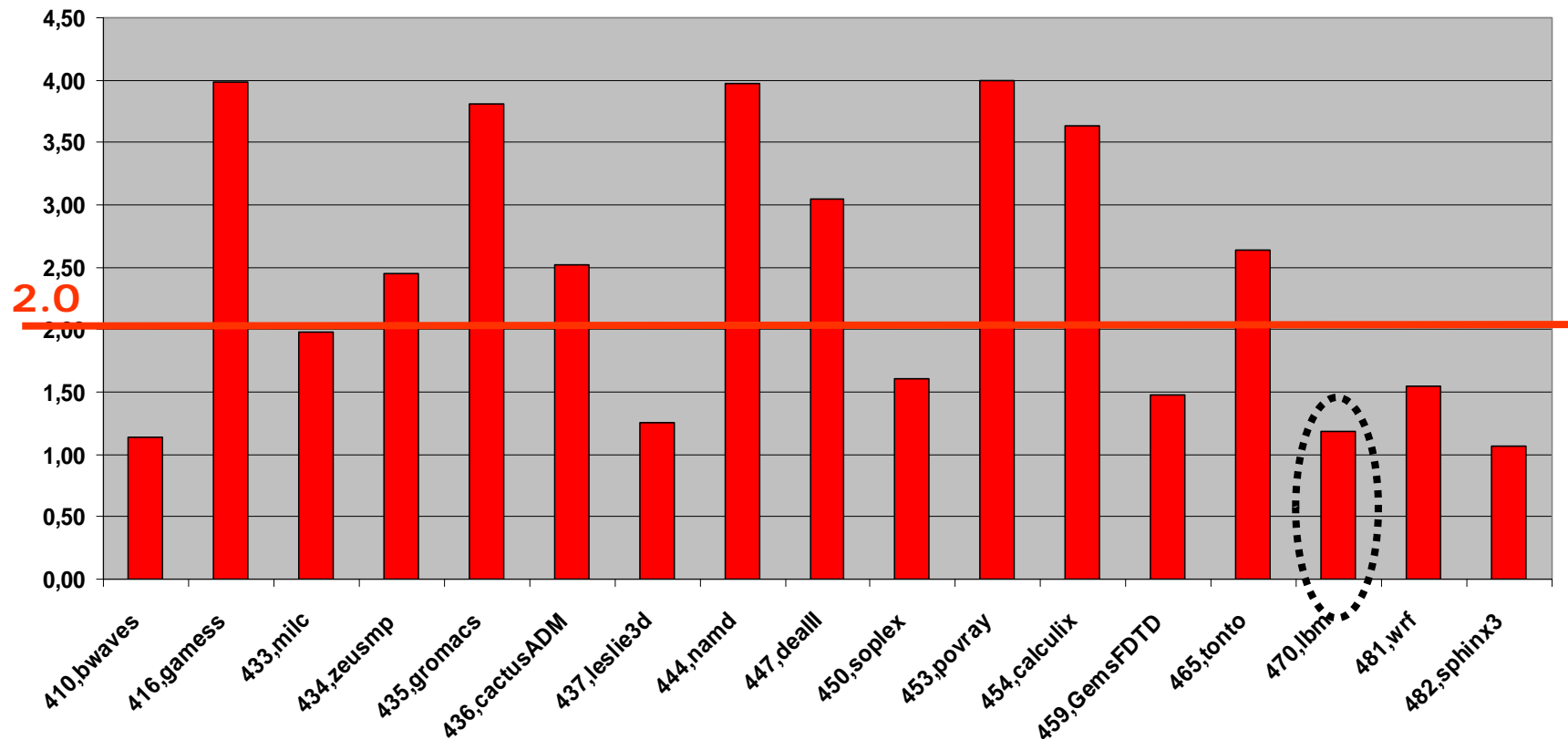**Throughput (as defined by SPEC):**
Amount of work performed per time unit when <u>several instances</u> of the application is executed simultaneously.
<u>Our TP study</u>: compare TP improvement when you go from 1 core to 4 cores
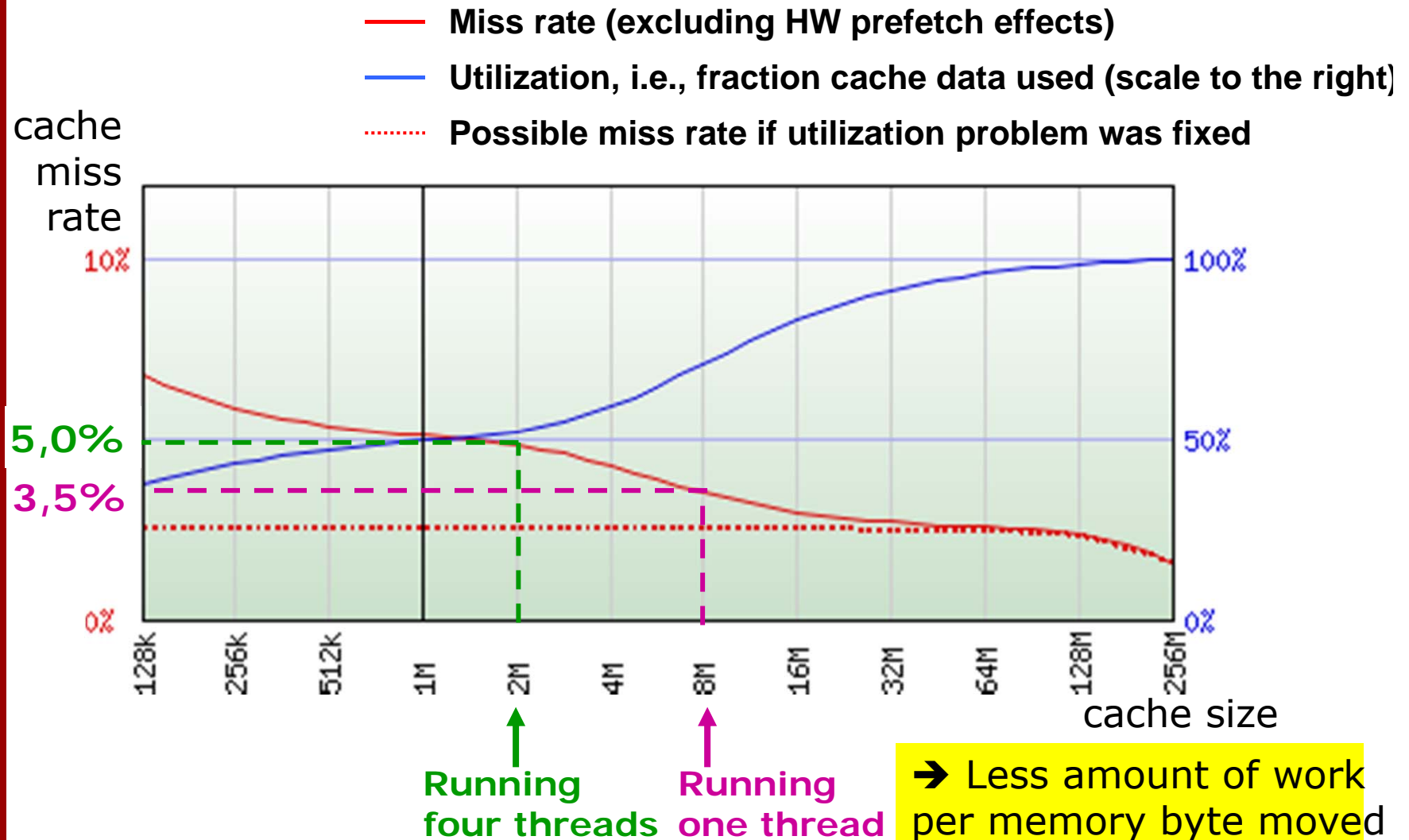
MC 9

Parallel
Comp
2012

# Throughput Scaling, More Apps

**SPEC CPU 20006 FP Throughput improvements on 4 cores**



*Intel X5365 3GHz "Core2", 1333 MHz FSB, 8MB L2.*
*(Based on data from the SPEC web)*

MC 10

Institutionen för informationsteknologi | www.it.uu.se

© Erik Hagersten| user.it.uu.se/~eh

# Nerd Curve: 470.LBM

—— **Miss rate (excluding HW prefetch effects)**

—— **Utilization, i.e., fraction cache data used (scale to the right)**

········ **Possible miss rate if utilization problem was fixed**

cache miss rate

10%

5,0%

3,5%

0%

100%

50%

0%

128k  256k  512k  1M  2M  4M  8M  16M  32M  64M  128M  256M

cache size

↑ **Running four threads**    ↑ **Running one thread**

➔ Less amount of work per memory byte moved @ four threads

MC 11

# Nerd Curve (again)

**—— Miss rate (excluding HW prefetch effects)**

**—— Utilization, i.e., fraction cache data used (scale to the right)**

**········ Possible miss rate if utilization problem was fixed**

cache
miss
rate

orig application

optimized application

**Running
four threads**

cache size

→ Twice the amount of work
per memory byte moved

MC 12

© **Erik Hagersten**| **user.it.uu.se/~eh**

Parallel
Comp
2012

UPPSALA
UNIVERSITET

# BW in the Future?

#Cores ~ #Transistors

## Computation vs Bandwidth

| CPU | CPU |
|-----|-----|
| CPU | CPU |

Mem

. . .

Chart axis: 0 to 6 (y-axis), 2007 to 2015 (x-axis)

Legend: #T * T_freq / #P * P_freq

**Y e a r**

HPCWire.com this morning:
**Up Against the Memory Wall**
"Nevermind the cores. Just hand over the cache"

HPCWire December 07:
**More Than 16 Cores May Well Be Pointless**
[by Sandia Labs]

# Commercial x86 snapshot

## (I may have miss-quoted some details, get architecture details from vendors)
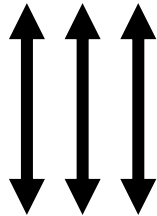
Erik Hagersten

Uppsala Universitet

# Intel Core2 Quad, 2006

South Bridge ⟷ I/O

North Bridge ⟷ DRAM

Front-side Bus (FSB)

Module

Another Module…

**Die 1**

L2$
6MB

| D$ 64kB | I$ 64kB | D$ 64kB | I$ 64kB |

CPU    CPU

**Die 2**

L2$
6MB

| D$ 64kB | I$ 64kB | D$ 64kB | I$ 64kB |

CPU    CPU

© Erik Hagersten| user.it.uu.se/~eh

Parallel
Comp
2012

# AMD Shanghai, 2007

Hyper Transport

DDR-2, DRAM

| L3   8MB |
|---|

| X-bar |
|---|

| L2$ 512kB | L2$ 512kB | L2$ 512kB | L2$ 512kB |
|---|---|---|---|

| D$ 64kB | I$ 64kB | D$ 64kB | I$ 64kB | D$ 64kB | I$ 64kB | D$ 64kB | I$ 64kB |
|---|---|---|---|---|---|---|---|
| CPU | | CPU | | CPU | | CPU | |

© Erik Hagersten| user.it.uu.se/~eh

Parallel Comp 2012

UPPSALA UNIVERSITET

# AMD MC System Architecture

# Intel: Nehalem, Core i7 Q1 2009 (4 cores)

QuickPath Interconnect

3x DDR-3 DRAM

| L3   8MB |
|----------|

| X-bar |
|-------|

| L2$ 256kB | L2$ 256kB | L2$ 256kB | L2$ 256B |
|-----------|-----------|-----------|----------|

| D$ 64kB | I$ 64kB | D$ 64kB | I$ 64kB | D$ 64kB | I$ 64kB | D$ 64kB | I$ 64kB |
|---------|---------|---------|---------|---------|---------|---------|---------|
| CPU, 2 thr | | CPU, 2 thr. | | CPU, 2 thr. | | CPU, 2 thr. | |

Up to 4 cores x 2 threads

MC 19

© Erik Hagersten| user.it.uu.se/~eh

Parallel Comp 2012

UPPSALA UNIVERSITET

# Nehalem "Core i7"

Parallel
Comp
2012

Institutionen för informationsteknologi | www.it.uu.se

MC 20

© Erik Hagersten| user.it.uu.se/~eh

# Intel: "Nehalem-Ex" (i7)

QuickPath Interconnect                                    4 x DDR-3

| L3   24MB |
|---|

| X-bar |
|---|

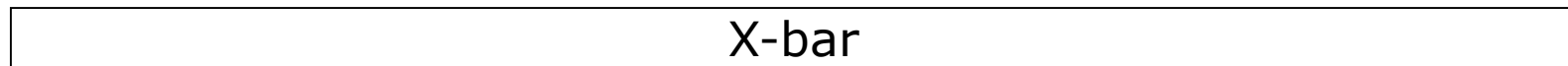| L2$ 256kB | L2$ 256kB | L2$ 256kB | L2$ 256B | | L2$ 256kB |
|---|---|---|---|---|---|
| D$ 64kB / I$ 64kB | D$ 64kB / I$ 64kB | D$ 64kB / I$ 64kB | D$ 64kB / I$ 64kB | ... | D$ 64kB / I$ 64kB |
| CPU, 2 thr | CPU, 2 thr. | CPU, 2 thr. | CPU, 2 thr. | | CPU |

8 cores x 2 threads

MC 21

Parallel
Comp
2012

UPPSALA
UNIVERSITET

# How is the silicon used (i7-Ex)?



| QPI0 | QPI1 | QPI2 | QPI3 |

Core3 Core4
Core2 Core5
System Interface
Core1 Core6
Core0 3 Mbyte Core7
MI MI



Execution Units | L1 Data Cache | L2 Cache & Interrupt Servicing
Memory Ordering & Execution | Paging
Instruction Reordering Scheduling & Retirement | Instruction Decode & Microcode | Branch Prediction
Instruction Fetch & L1 Cache

MC 22



Tag

Data array

Source: JSSC Jan 2010, Rusu et. al

Parallel Comp 2012

# How is the silicon used?

Institutionen för informationsteknologi | www.it.uu.se    MC 23

Source: JSSC Jan 2010, Rusu et. al

# AMD Istanbul, 6 cores

Institutionen fö

# AMD Magny-Cours



MCM

C C C C C C    C C C C C C
¢ ¢ ¢ ¢ ¢ ¢    ¢ ¢ ¢ ¢ ¢ ¢
$ $ $ $ $ $    $ $ $ $ $ $
€              €

~Istanbul

DDR3    DDR3    HT  HT    HT  HT    DDR3    DDR3

MC 25

# Some other multicores

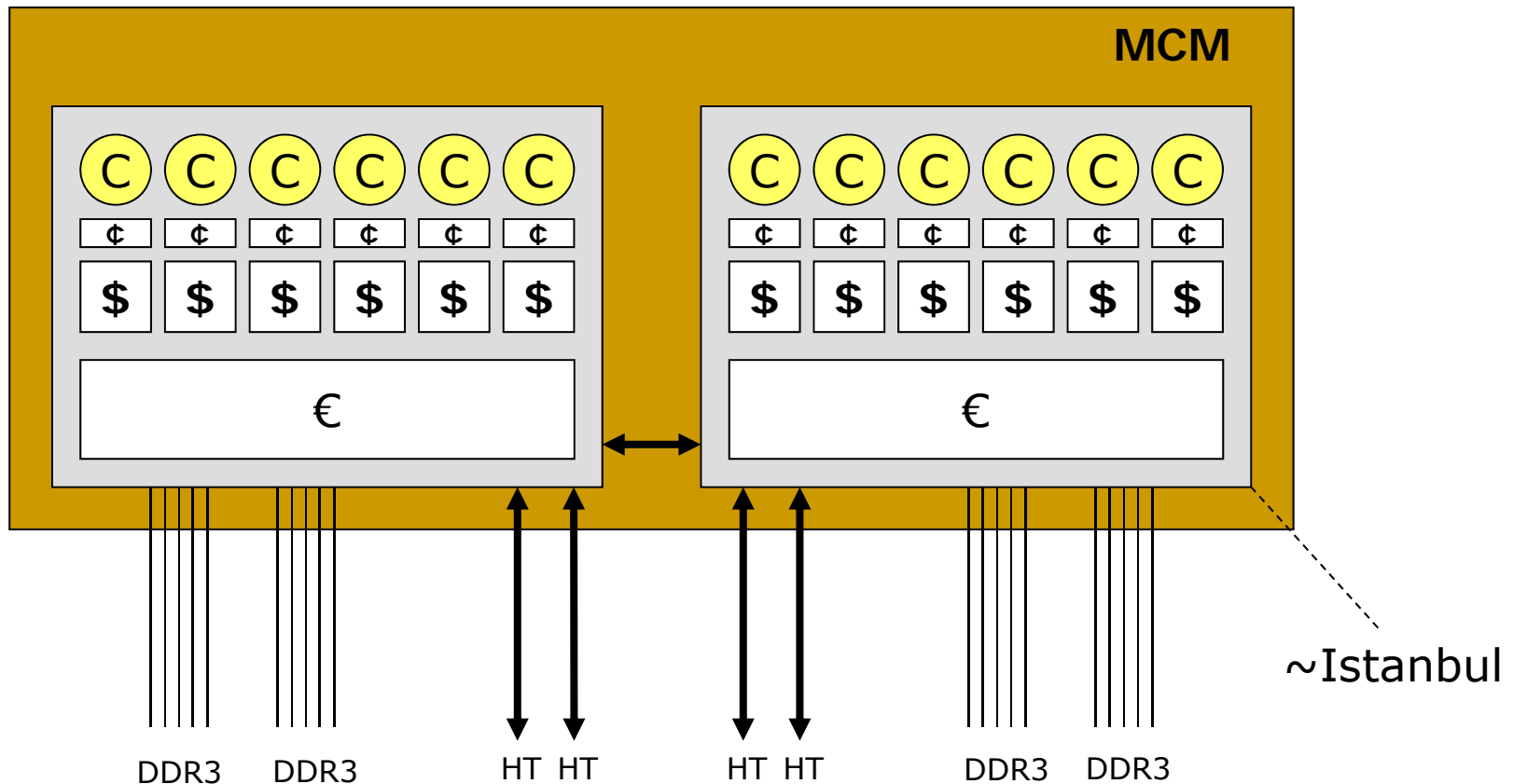## (I may have miss-quoted some details, get architecture details from vendors)

Erik Hagersten

Uppsala Universitet

# Sun Niagara, 2005!!

**4 x DDR-2 = 25GB/s (!)**

| Memory ctrl | Memory ctrl | Memory ctrl | Memory ctrl |
|---|---|---|---|
| L2 | L2 | L2 | L2 |

**Shared L2**

**Xbar = 134 GB/s**

...8...

| L1I | L1D (wt) | L1I | L1D (wt) |

CPU   ... 8 ...   CPU

**Four interleaved threads**

**Now: Victoria's falls: 16 core with 16 threads each**

Parallel Comp 2012

UPPSALA UNIVERSITET

# Niagara Chip



UltraSPARC-Core

Sun Microsystems

# TILERA Architecture



Tile Processor

Core + Switch = Tile

Core

S

The tile is the basic building block

**64 cores connected in a mesh**
**Local L1 + L2 caches**
**Shared distributed L3 cache**
**Linux + ANSI C**
**New Libraries**
**New IDE**
**Stream computing**
**...**

# IBM Cell Processor

Mem



**IBM Cell**

Parallel
Comp
2012

Institutionen för informationsteknologi | www.it.uu.se

MC 30

© Erik Hagersten| user.it.uu.se/~eh

# So-called accelerators

- Sits on the IO bus (!!)

- GP Graphics processors, aka GPGPU
  [e.g. NVIDIA, AMD/ATI]

- Specialized accelerators
  [e.g., FPGA/Mitrionics, ASIC/ClearSpeed]

- Specialized languages for the above
  [CUDA, Ct, Rapid Mind, Open-CL, ...]

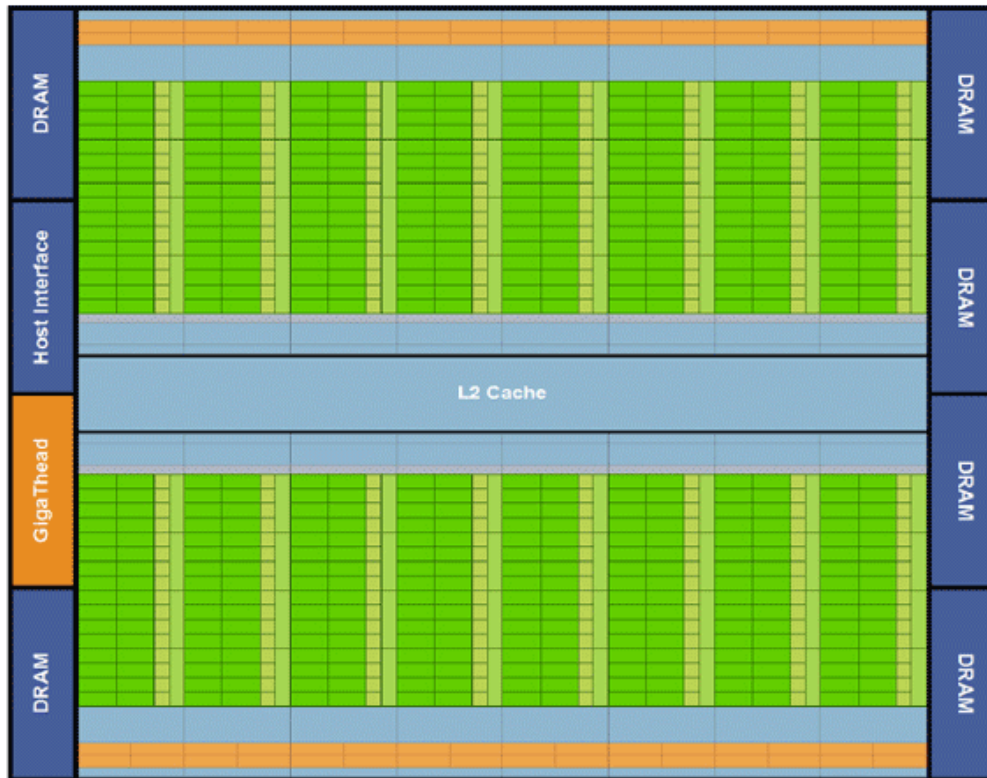# So-called accelerators

- Sits on the IO bus (!!)

- GP Graphics processors, aka GPGPU?
  [e.

- Spe
  [e.

**My view: Not very general purpose yet!**

•Fits well for a few **VERY IMPORTANT** app domains!

•Limited applicability?

•Programmer productivity?

•Application life time?

•New generation devices will be more useful...

- Spe
  [CUDA, Ct, Rapid Mind, Open-CL, ...]

# Fermi from nVIDIA
## a huge step in the right direction
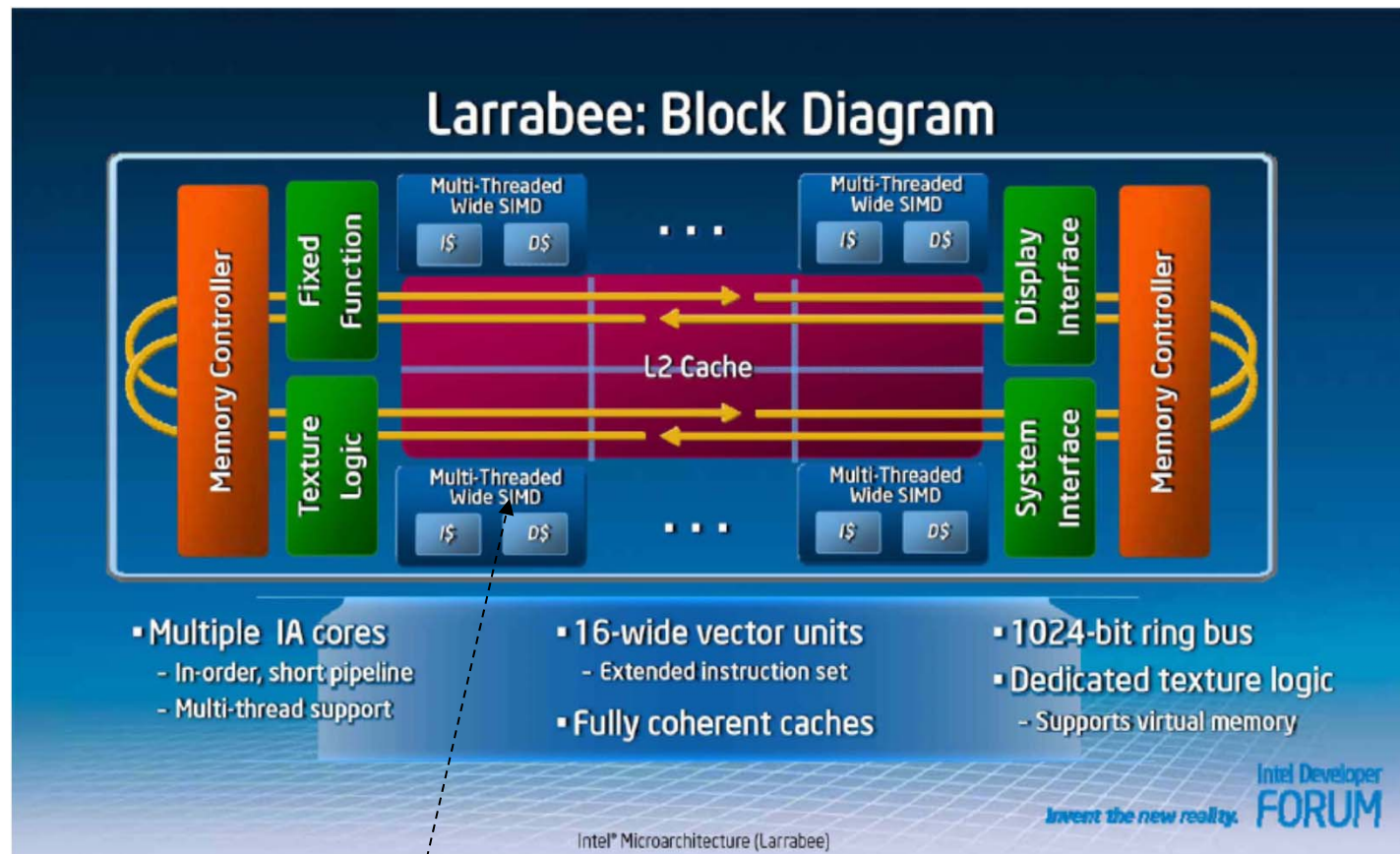


- 512 "processors" (**P**)
- 16 **P**/StreamProcessor (**SP**)
- Full DP-FP IEEE support
- 64kB L1 cache /**SP**
- 768kB global shared cache
- Atomic instructions
- ECC correction
- Debugging support
- …

**David Black-Schaffer to give you the full story**

Parallel
Comp
2012

Institutionen för informationsteknologi | www.it.uu.se

MC 33

© Erik Hagersten| user.it.uu.se/~eh

# Scaling the x86 Manycore computing
# Larrabee (now called MIC) from Intel 2012-2013??

"more than 50 cores"



Larrabee: Block Diagram

- **Multiple IA cores**
  - In-order, short pipeline
  - Multi-thread support
- **16-wide vector units**
  - Extended instruction set
- **Fully coherent caches**
- **1024-bit ring bus**
- **Dedicated texture logic**
  - Supports virtual memory

Intel® Microarchitecture (Larrabee)

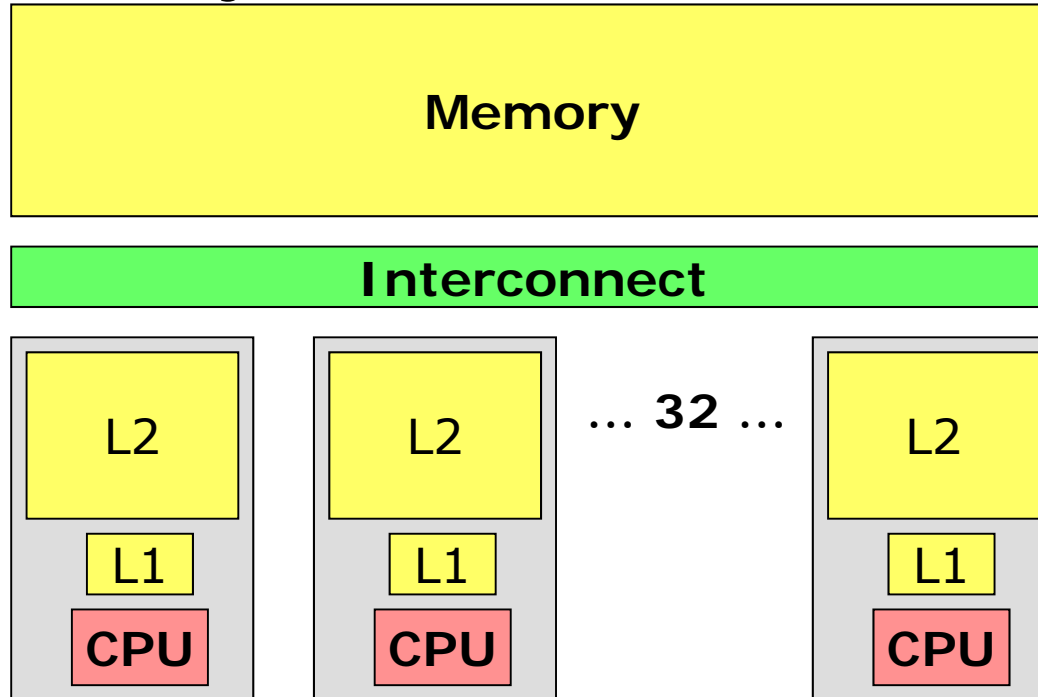SIMD instructions (16-way??)

MC 34

© Erik Hagersten| user.it.uu.se/~eh

# Wrapping up about multicores

Erik Hagersten

Uppsala Universitet

# Looks and Smells Like an SMP (aka UMA)?

**SMP system**

**Multicore system**

Memory

Memory

Interconnect

| L2 | L2 | ... **32** ... | L2 |

L1 | L1 | | L1

**CPU** | **CPU** | | **CPU**

L2

1 ...8... 1

T T

## Well, how about:

- Cost of parallelism?
- Cache capacity per thread?
- Memory bandwidth per thread?
- Cost of thread communication? …

MC 36

Parallel
Comp
2012

# What matters for multicore performance?

- Are we buying…
  - CPU frequency?
  - Number of cores?
  - MIPS and FLOPS?
  - Memory bandwidth?
  - Cache capacity?
  - Memory capacity?
  - Performance/Watt?
  - Dark Silicon is around the corner!

Parallel Comp 2012

# MC Questions for the Future

- How to get parallelism?
- How to get performance/data locality?
- How to debug?
- A case for new funky languages?
- A case for automatic parallelization?
- Are we buying:
  - compute power,
  - memory capacity, or
  - memory bandwidth?
- Will 128 cores be mainstream in 5 years?
- Will the CPU market diverge into desktop/capacity/capability/special-purpose CPUs again?
- **A non-question: will it happen?**

Parallel
Comp
2012

UPPSALA
UNIVERSITET