# TERA**F**LUX.EU

## Exploiting Dataflow Parallelism in Teradevice Computing

University of Siena

Barcelona
Supercomputing Center

CAPS

[LABS]^hp

INRIA

Microsoft

THALES

University of Augsburg

University of Cyprus

Roberto Giorgi – University of Siena (coordinator)

Cagliari, Italy – Computing Frontiers

16/05/2012

MANCHESTER 1824
University of Manchester
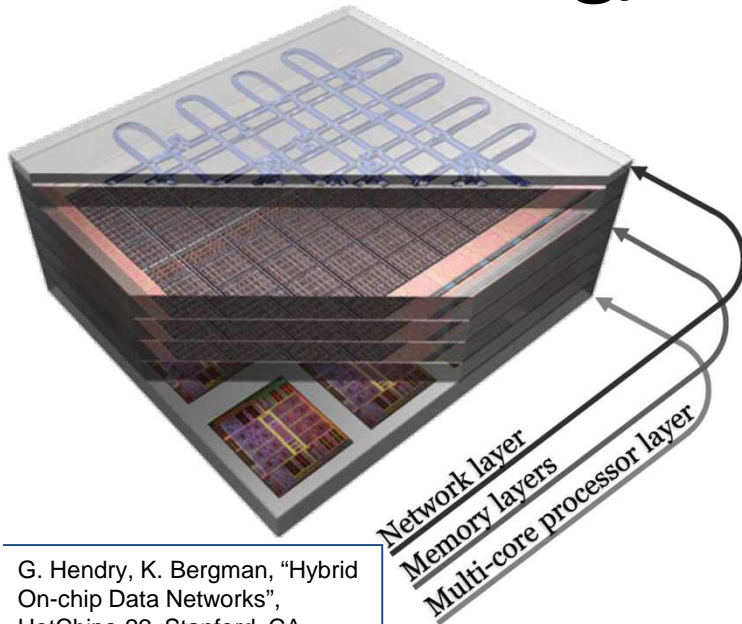
1

# What is **TERAFLUX** about

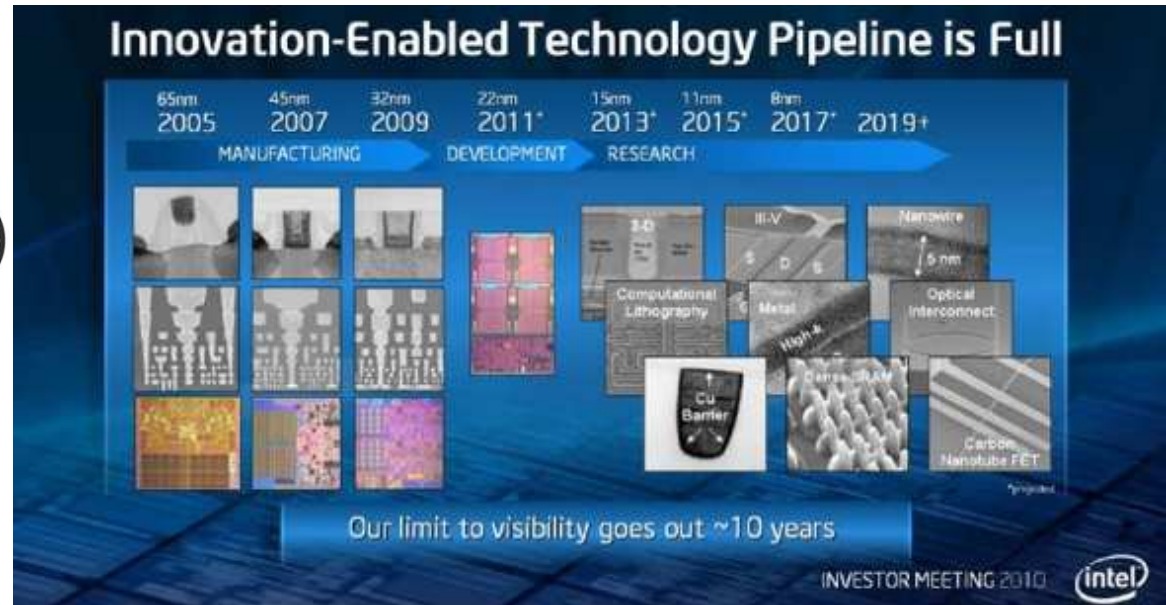# Architecture+Programmability+Reliability
of
Future (single chip)
Many-cores
(targeting 1000+ cores)

**TERA**<sup>**F**</sup>**LUX**

# Future Scenarios
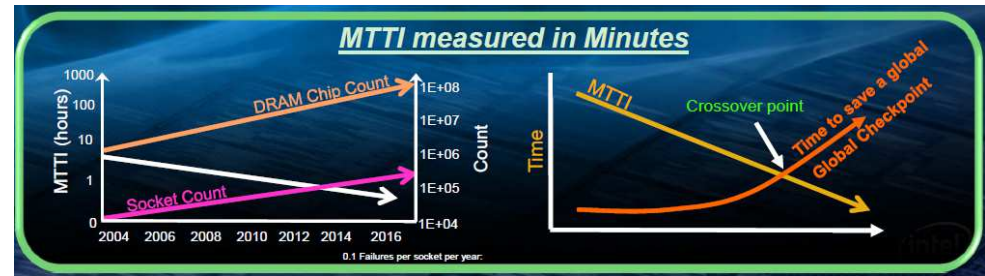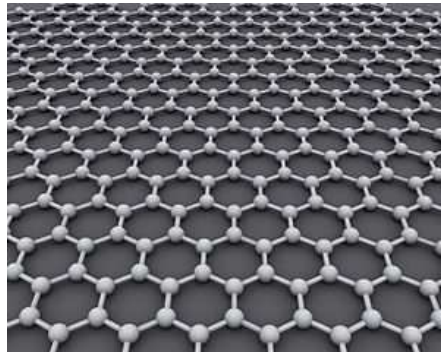## == 3D stacking, 8nm, 3D transistors, Graphene



G. Hendry, K. Bergman, "Hybrid On-chip Data Networks", HotChips-22, Stanford, CA – Aug. 2010



Innovation-Enabled Technology Pipeline is Full

| 65nm 2005 | 45nm 2007 | 32nm 2009 | 22nm 2011* | 15nm 2013* | 11nm 2015* | 8nm 2017* 2019+ |

MANUFACTURING — DEVELOPMENT — RESEARCH

Our limit to visibility goes out ~10 years

INVESTOR MEETING 2010  (intel)

Fab D1X (OR), 42 (AZ) starting the 14nm node in 2013



22 nm 3-D Tri-Gate Transistor

3-D Tri-Gate transistors form conducting channels on three sides of a vertical fin structure, providing "fully depleted" operation
Transistors have now entered the third dimension!





MTTI measured in Minutes

Pawloski, May 2011, Exascale Seminar, Ghent

**TERA$^F$LUX**

# Fundamental approach: DATAFLOW

A Scheme of Computation in which an activity is initiated by presence of the data it needs to perform its function
(Jack Dennis)

# Recent Projects/Efforts towards DATAFLOW

- Maxeler (UK) selling "dataflow computer" to J.P. Morgan → about 350x speedup vs. standard x86 cores

- DARPA funding 25M$ for UPHC program, encompassing:
  - Gao's dataflow execution model (codelet based) – SWARM by ETI
  - Intel's Runnamede project

**HPC** wire

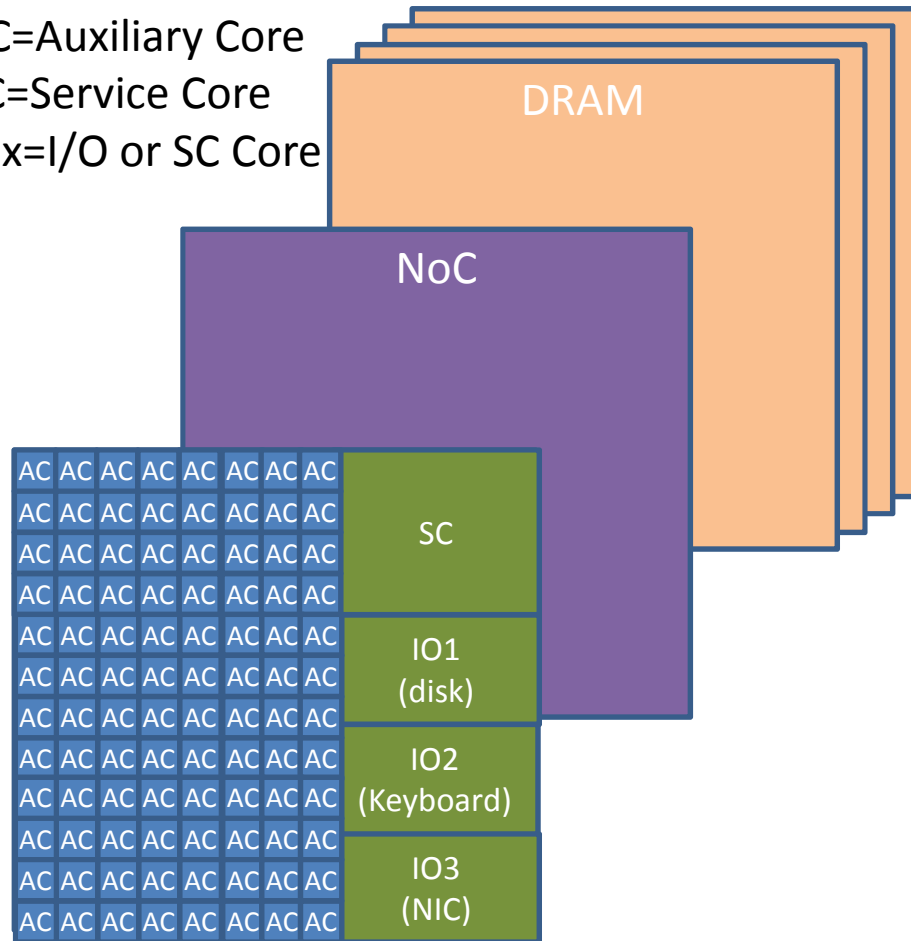J.P. Morgan Deploys Maxeler Dataflow Supercomputer for Fixed Income Trading
December 15, 2011

UPHC=Ubiquitus
High-Performance Computing

The Intel-lead UHPC team intends to develop new circuit topologies, new chip and system architectures, and new programming techniques to reduce the amount of energy required per computation by between 100x and 1000x compared to today's computing systems. Such dramatic reduction in energy consumption will allow these future systems to take full advantage of the increasing transistor budgets afforded by the steady advances in Moore's Law.

**TERAFLUX**

# TERAᶠLUX – Future Many-cores
## Key Challenges: Architecture+Programmability+Reliability

AC=Auxiliary Core
SC=Service Core
IOx=I/O or SC Core

DRAM

NoC

AC AC AC AC AC AC AC AC
AC AC AC AC AC AC AC AC
AC AC AC AC AC AC AC AC
AC AC AC AC AC AC AC AC
AC AC AC AC AC AC AC AC
AC AC AC AC AC AC AC AC
AC AC AC AC AC AC AC AC
AC AC AC AC AC AC AC AC
AC AC AC AC AC AC AC AC
AC AC AC AC AC AC AC AC
AC AC AC AC AC AC AC AC
AC AC AC AC AC AC AC AC

SC

IO1 (disk)

IO2 (Keyboard)
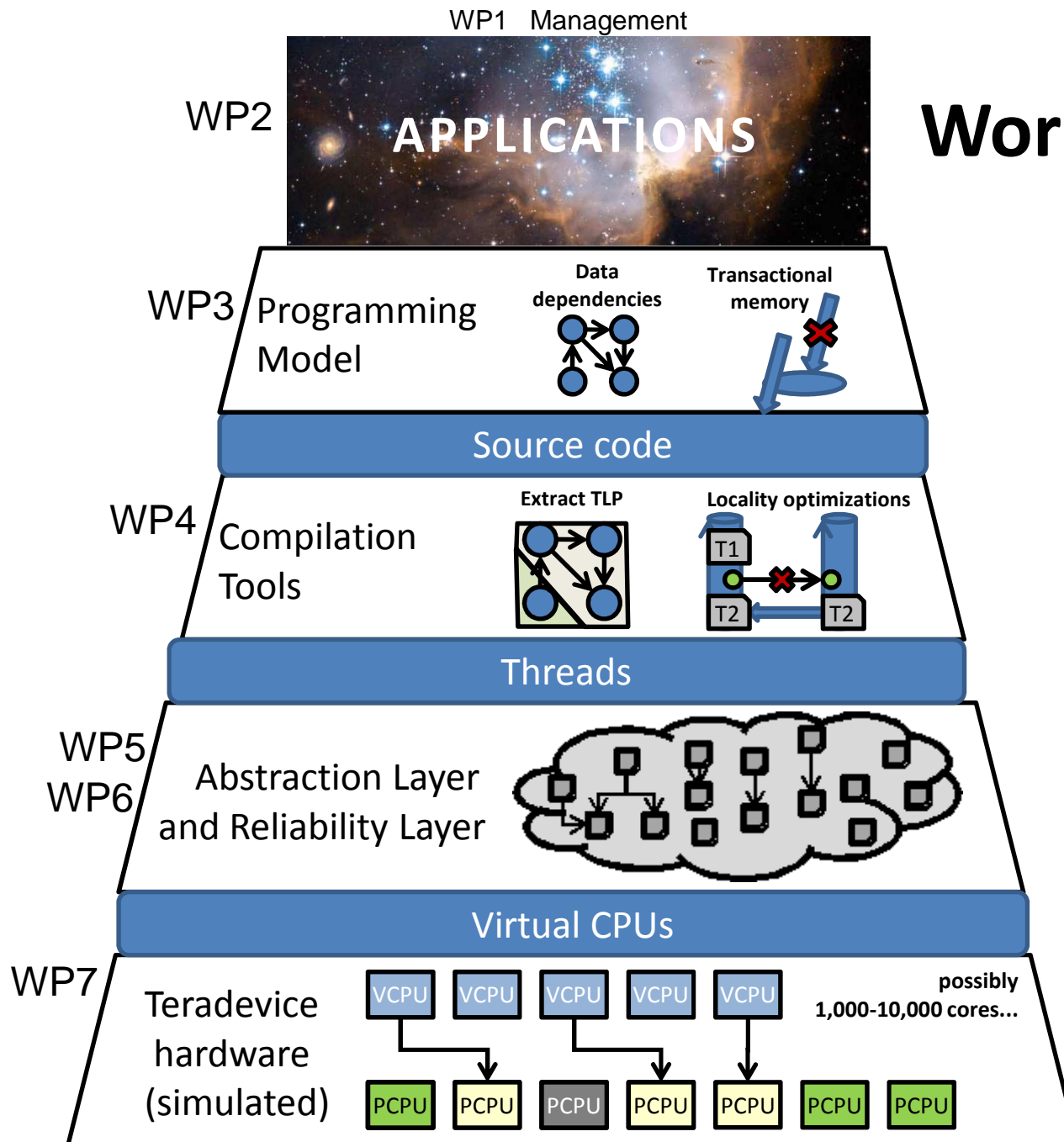
IO3 (NIC)

TERAFLUX Key Facts:
- FET – Integrated Project – 10 Partners – 4 years – 2010-2013
- 7.5 Meuro Total cost (5.7 Meuro EU funded)
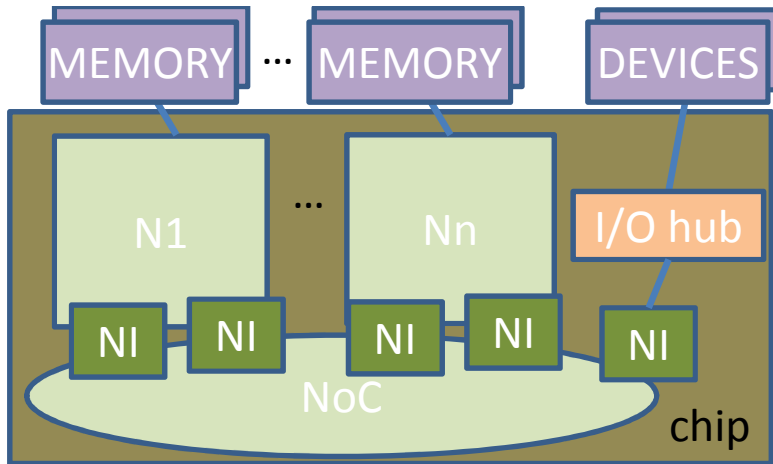
TERAFLUX-INCO Key Facts:
- FET – EU Worldwide Cooperation – Extends TERAFLUX to an additional USA partner for 21 months
- 546 Keuro Total cost (420 Keuro EU funded)

TERAᶠLUX

Roberto Giorgi – giorgi@unisi.it --- http://teraflux.eu

# TERA$^F$LUX.EU
# Working Hypothesis

- 1000 Billion- or 1 TERA-device computing platforms pose new challenges:
    - (at least) programmability, complexity of design, reliability

- TERAFLUX context:
    - High performance computing and applications (not necessarily embedded)

- TERAFLUX scope:
    - Exploiting a less exploited path (DATAFLOW) at each level of abstraction

Roberto Giorgi – giorgi@unisi.it --- http://teraflux.eu

# TERAFLUX Architectural template



LEGENDA:
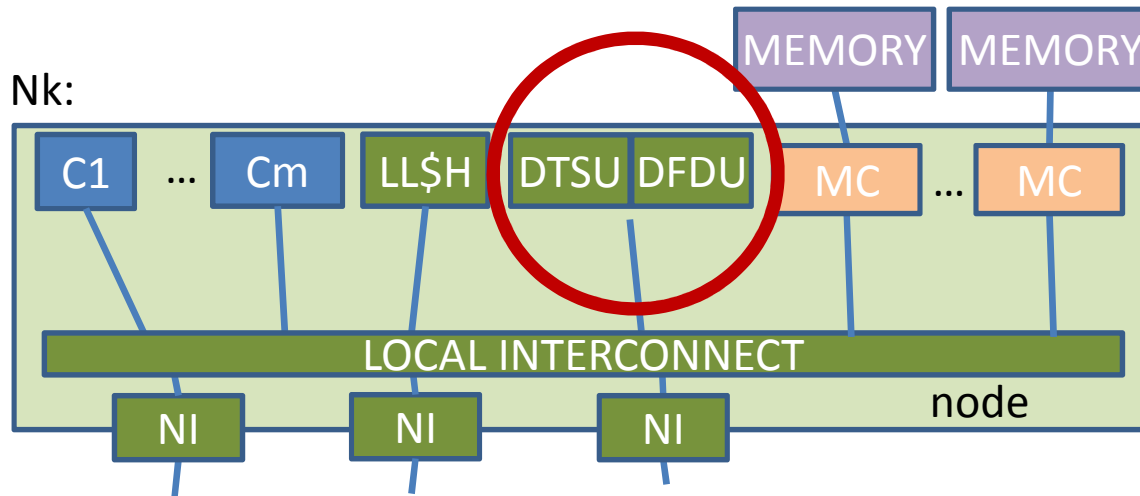n = # of nodes
m = # of cores per node
u = # of DRAM controllers insisting on the
        Unified Physical Address Space
z = # of I/O Hubs

Nk = k-th Node   (k=1..n)
NI = Network Interface
NoC = Network on Chip
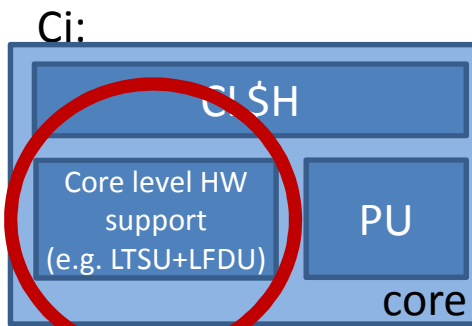
$C_j$ = j-th core   (j=1..m)
MC = Memory Controller
DTSU = Distributed Thread-Scheduler Unit
DFDU = Distributed Fault-Detection Unit
LL$H = Last Level Cache Hierarchy

CL$H = Core Level Cache Hierarchy
PU = Processing Unit
LTSU = Local Thread-Scheduler Unit
LFDU = Local Fault-Detection Unit

| CHIP LEVEL | EXTERNAL or OTHER LAYER |
| NODE LEVEL | NODE OPTIONAL |
| CORE LEVEL | |

**TERAFLUX**

# Our pillars

- FIXED and MOST-USED ISA (**x86**)

- MANYCORE FULL SYSTEM SIMULATOR (**COTSon**)

- REAL WORLD APPLICATIONS (e.g. GROMACS)

- SYNCHRONIZATION: **TRANSACTIONAL MEMORY**

- **GCC** based TOOL-CHAIN

- OFF-THE-SHELF COMPONENTS FOR CORES, OS, NOC,MEMORY HIERARCHY

- **FDU** AND **TSU** (Fault Detection Unit and Thread Scheduling Unit)

# A REVIEW OF RECENT MANY-CORES

TERA<sup>F</sup>LUX

Roberto Giorgi – giorgi@unisi.it --- http://teraflux.eu

# HotChips 2011

- Hot Chips papers suggest that the rest of the world is moving in a different direction: large numbers of relatively simple CPUs. But the trend is reinforcing a long-appreciated set of questions—as the number of cores grows, **how do you deal scalability with interconnect, memory hierarchy, coherency, and intra-thread synchronization?** Answers to these questions depend on the size of the design, the application space, and the heritage of the design team.
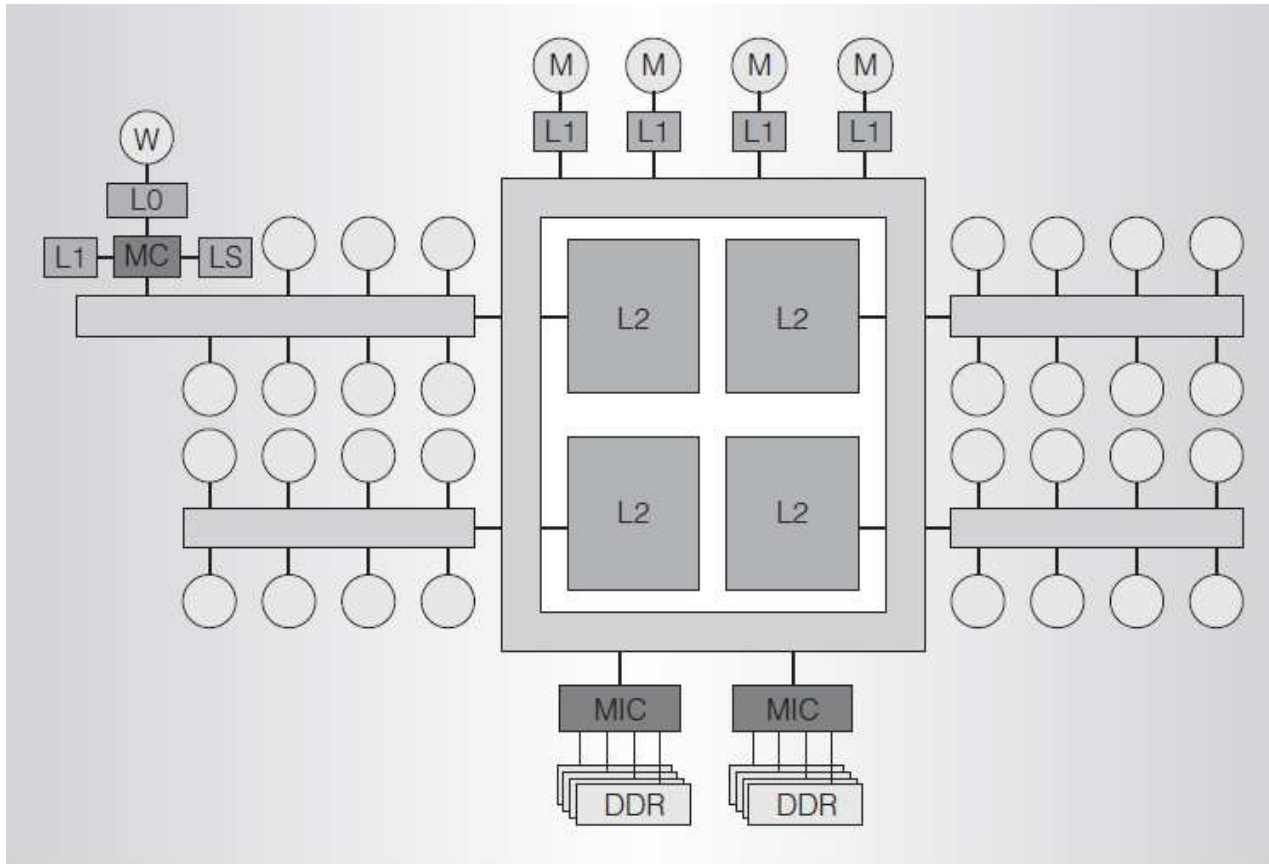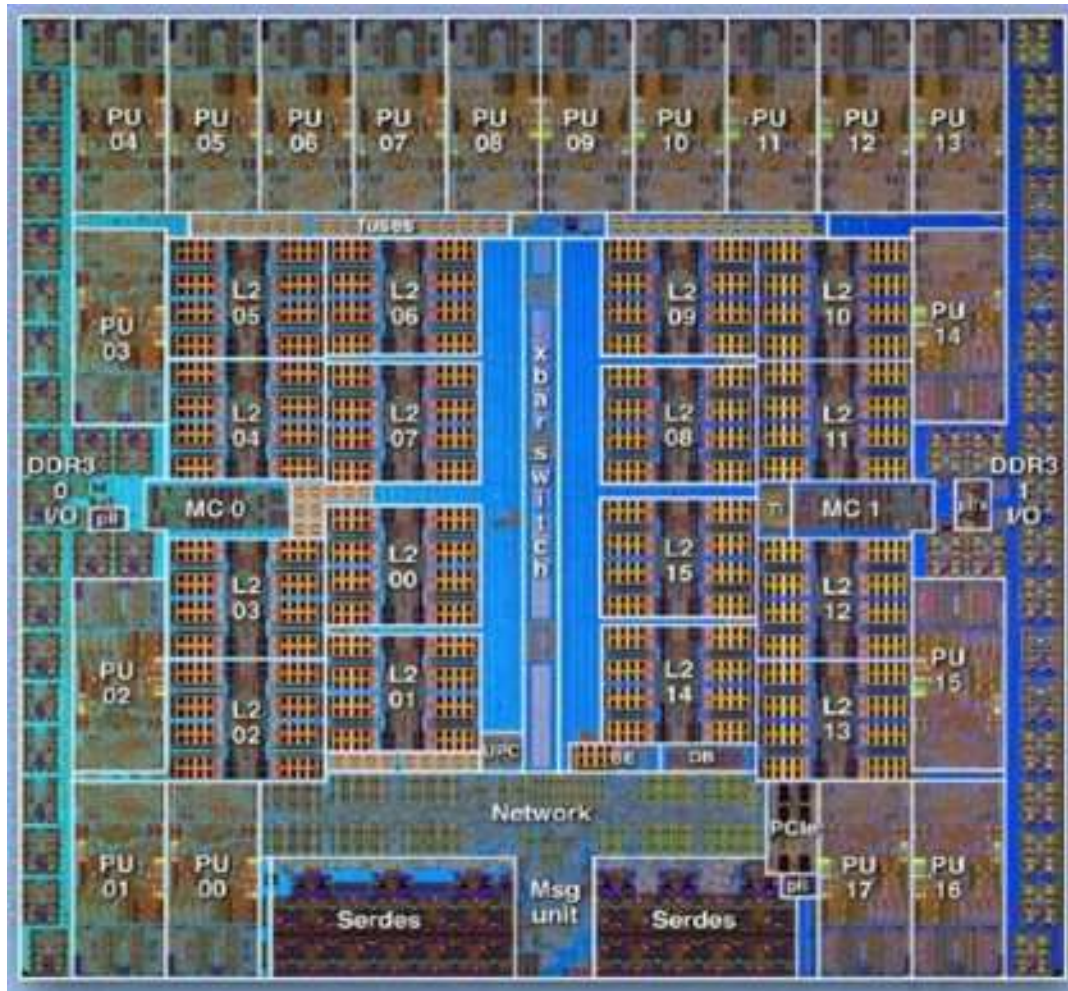
# SARC Architecture



Figure 1. Schematic of the SARC architecture. The number of masters, workers, level-2 (L2) blocks, and memory interface controllers is implementation dependent, as is their on-chip layout.

Table 1. Baseline SARC simulation parameters.

| Parameter | Value |
|---|---|
| Clock frequency | 3.2 GHz |
| Memory controllers | 4 × 2 DDR3 channels |
| Channel bandwidth | 12.8 Gbytes per second (GBps) (DDR3-1600) |
| Memory latency | Real DDR3-1600 |
| Memory interface controllers (MICs) policy | Closed-page, in-order processing |
| Shared L2 cache | 128 Mbytes (32 blocks æ 4 Mbytes), 4-way associative |
| L2 cache latency | 40 cycles |
| Local store | 256 Kbytes, 6 cycles |
| L0 cache | 32 Kbytes, 3 cycles |
| Interconnection links | 8 bytes/cycle (25.6 GBps) |
| Intracluster network on chip (NoC) | 2-bus (51.2 GBps) |
| Global NoC | 16-bus (409.6 GBps) |

More recently (20110908), Dimitris Nikoloupos confirmed me that there won't be anymore the LS as it will be integrated in the L2.

The SARC architecture, IEEE micro, Oct. 2010, vol. 30, n. 5, pp. 16-29

**TERA** **LUX**

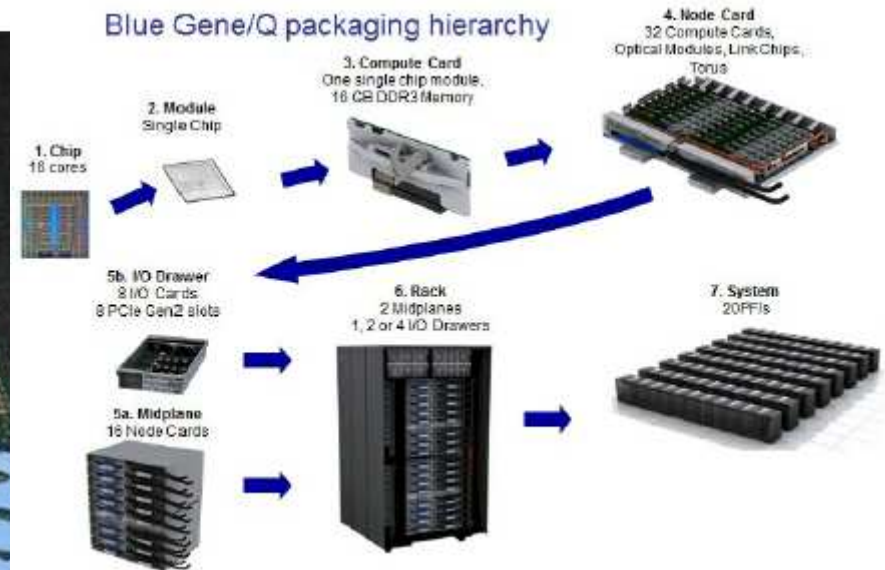Roberto Giorgi – giorgi@unisi.it --- http://teraflux.eu

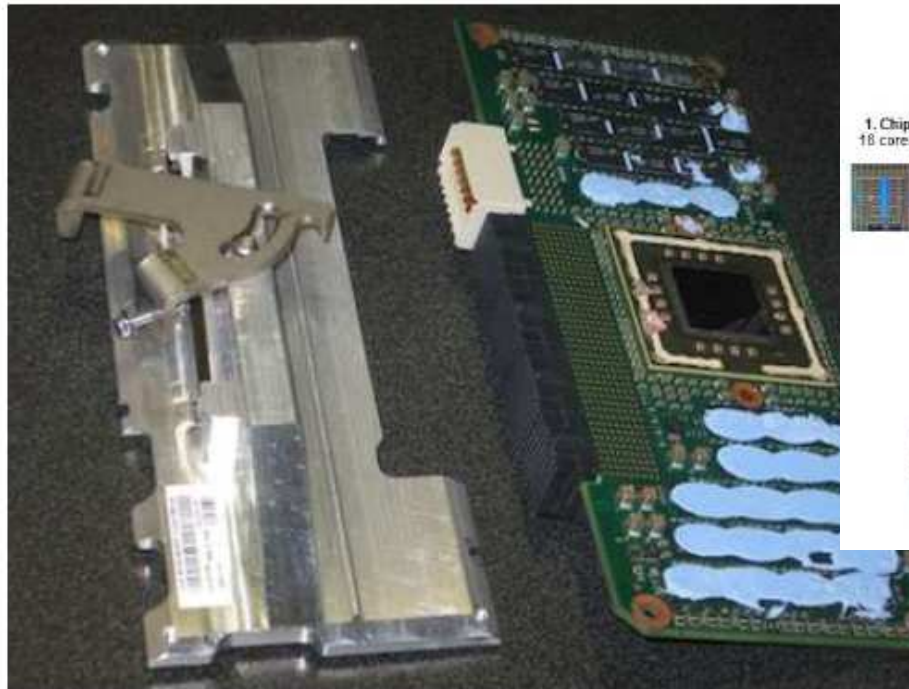# IBM BlueGene/Q



Rather than replicate more than 18 of these large cores, IBM chose to give each core hardware support for four concurrent threads, so under ideal circumstances the chip can behave almost as a 64-CPU system

- Ruud Haring, The IBM Blue Gene/Q Compute chip+SIMD floating-point unit, HotChips Symposium, Aug 2011.

**TERAFLUX**

Roberto Giorgi – giorgi@unisi.it --- http://teraflux.eu

# BlueGene/Q module and system



Blue Gene/Q packaging hierarchy

1. Chip
18 cores

2. Module
Single Chip

3. Compute Card
One single chip module.
16 GB DDR3 Memory

4. Node Card
32 Compute Cards,
Optical Modules, Link Chips,
Torus

5b. I/O Drawer
8 I/O Cards
8 PCIe Gen2 slots

5a. Midplane
16 Node Cards

6. Rack
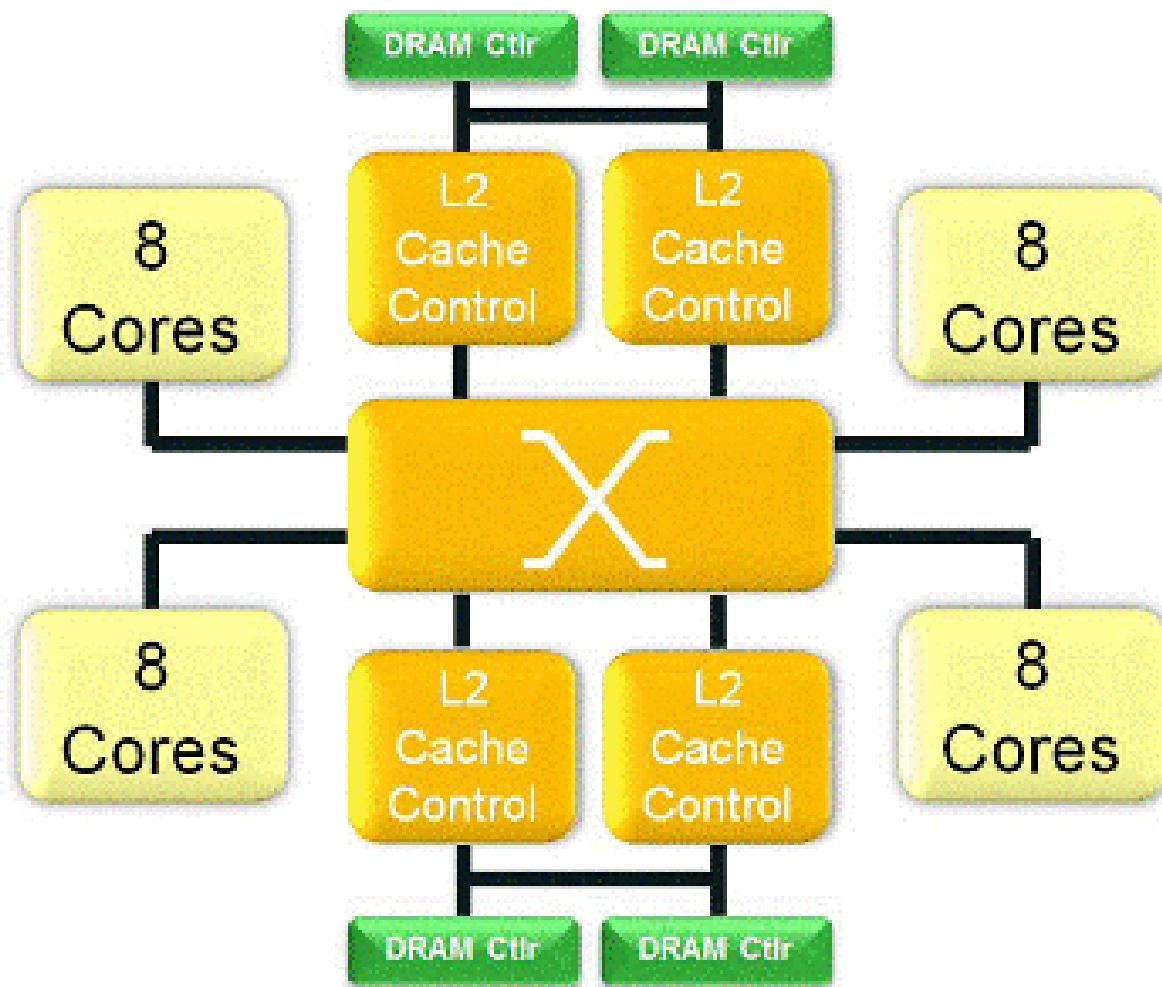2 Midplanes
1, 2 or 4 I/O Drawers

7. System
20PFIs

The BlueGene/Q module with DDR3 memory, five links and a water cooling system.

Thanks to the 64-bit support, the modules can now run 8 or 16 GB of DDR3 memory. Five links (2 GB/s per direction) connect each module to its neighbors, making it possible to create different 5D topologies. Half a rack with 8192 BlueGene/Q cores has already proven its capabilities in the Linpack benchmark. With 65.3 Tflops, the test system from the Thomas J. Watson Research Center scored 115th place in the new Top500 list. Its power consumption of 38.8 kW represented a new record value for energy efficiency at close to 1700 Mflops/watt. The Sequoia is supposed to get 96 fully equipped racks, which are supposed to deliver 20 Pflops of theoretical peak performance at the end of 2012.

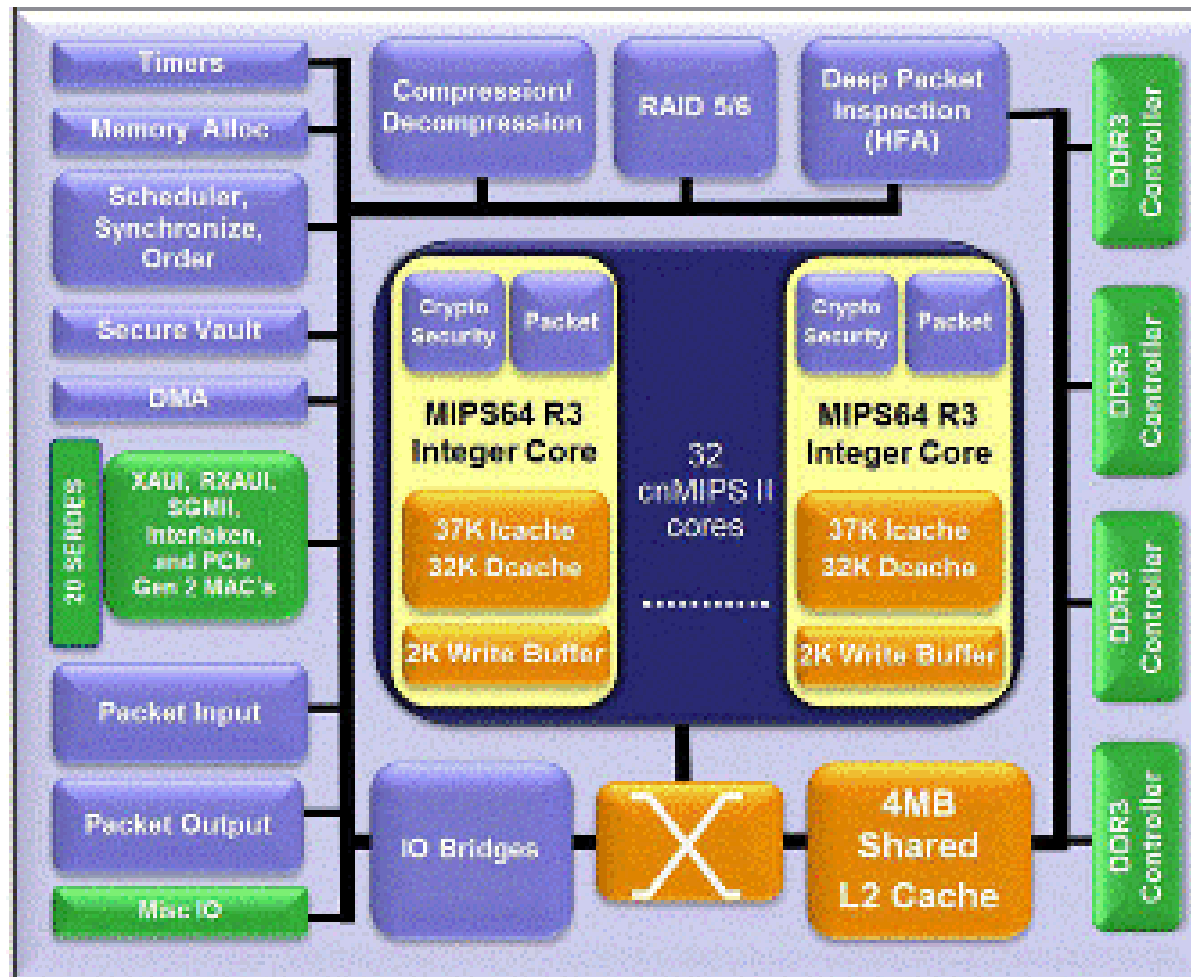# Cavium – Octeon II CN6880



Cavium relies primarily on locks for synchronization, assisted by facilities in the scheduler hardware

Cavium OCTEON II CN6880 Multi-Core MIPS64 Processor, HotChips Symposium, Aug. 2011.
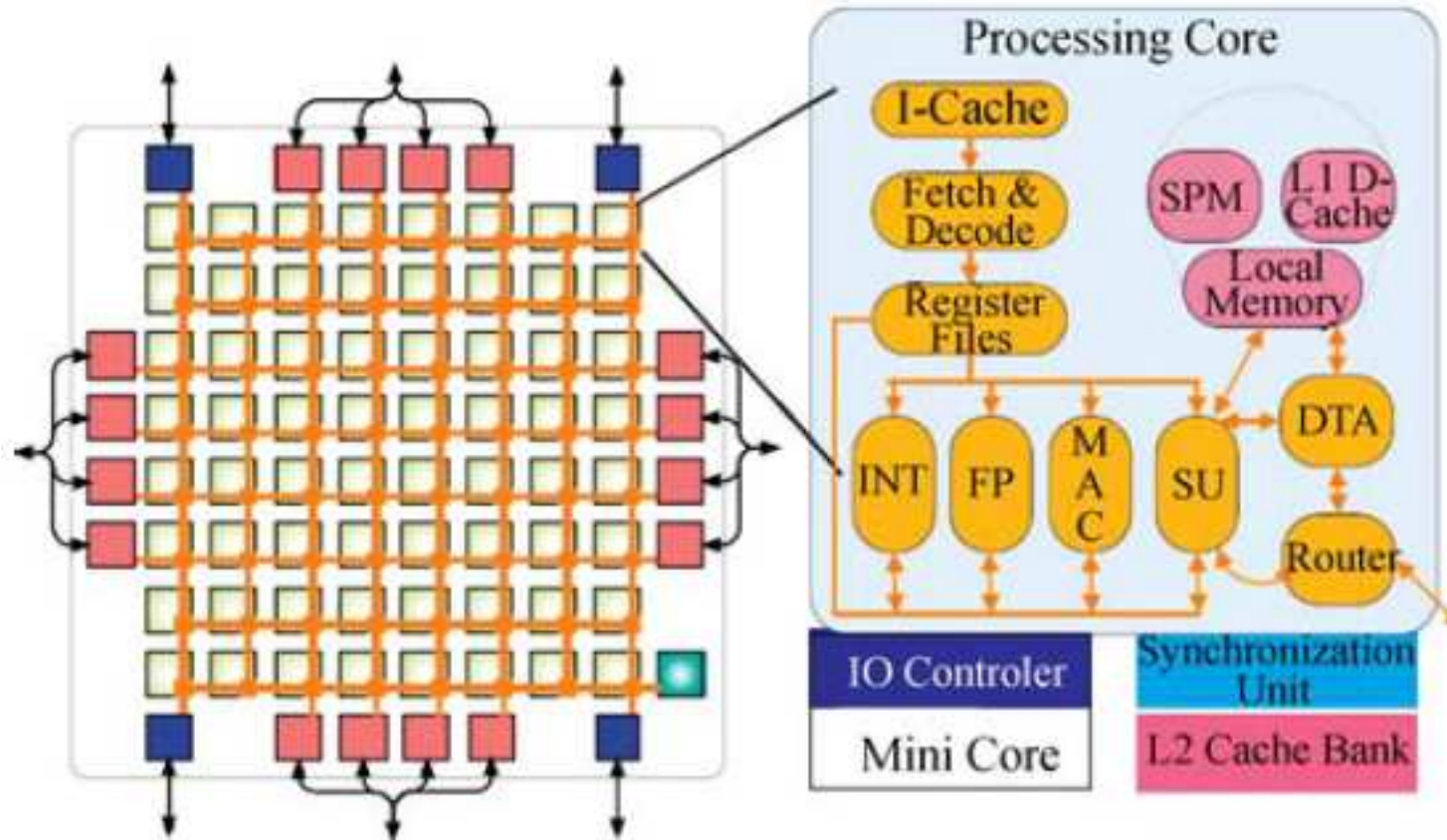
# Octeon II (2)

# Godson-T



Fig.1. Overview of Godson-T.

Cui HM, Wang L, Fang DR *et al.* Landing stencil code on Godson-T. JOURNAL OF COMPUTER SCIENCE ANDTECHNOLOGY 25(4): 886–894 July 2010.

**TERA<sup>F</sup>LUX**

# SPARC64™ VIIIfx Chip Overview



- **Architecture Features**
  - 8 cores
  - Shared 5 MB L2$
  - Embedded Memory Controller
  - 2 GHz
- **Fujitsu 45nm CMOS**
  - 22.7mm x 22.6mm
  - 760M transistors
  - 1271 signal pins
- **Performance (peak)**
  - 128GFlops
  - 64GB/s memory throughput
- **Power**
  - 58W (TYP, 30°C)
  - Water Cooling – Low leakage power and High reliability

**TER**

# K-Computer

## Compute nodes and Network
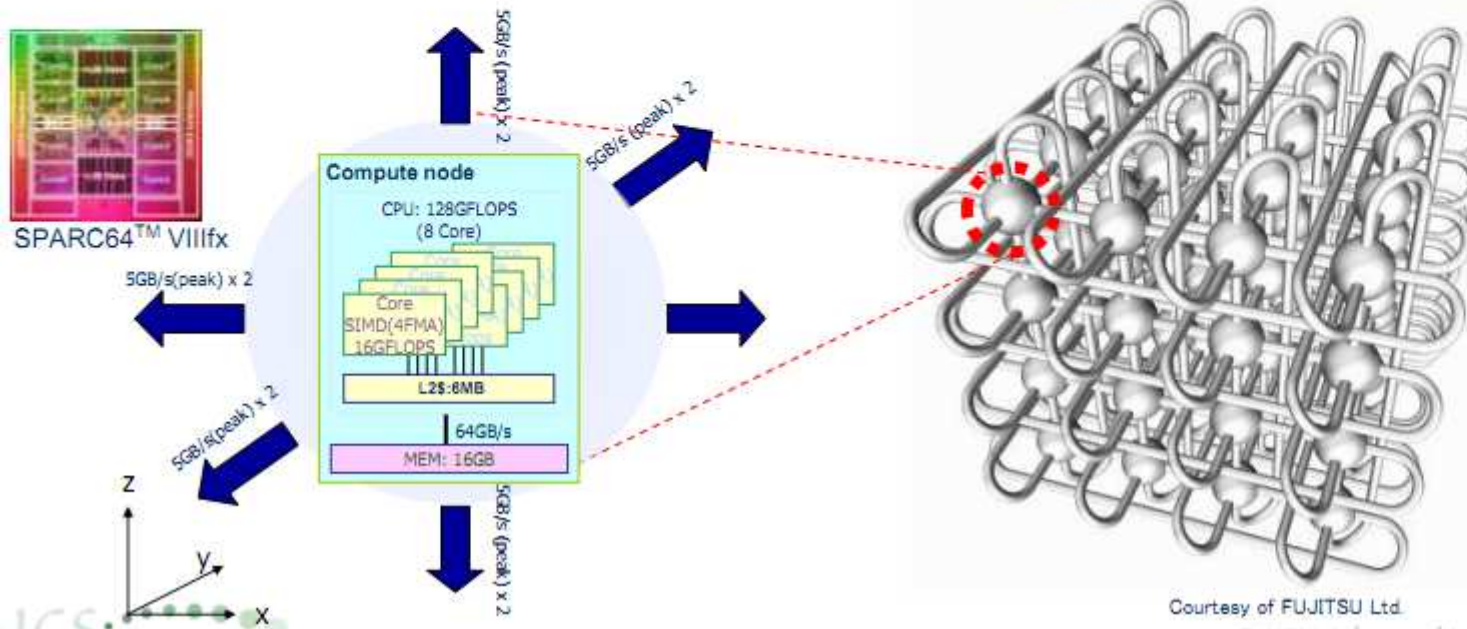
- Compute nodes (CPUs): > 80,000
    - Number of cores: > 640,000
- Peak performance: > 10PFLOPS
- Memory: > 1PB (16GB/node)

- 6-dimensional mesh/torus network: Tofu
    - 10 connections to each adjacent node
- Peak bandwidth: 5GB/s x 2 for each connection
- Logically 3-dimensional torus network



Courtesy of FUJITSU Ltd.

Mitsuo Yokokawa, Fumiyoshi Shoji, Atsuya Uno, Motoyoshi Kurokawa, and Tadashi Watanabe. 2011. *The K computer*: Japanese next-generation supercomputer development project. In*Proceedings of the 17th IEEE/ACM international symposium on Low-power electronics and design* (ISLPED '11).

**TERA$^F$LUX**

Roberto Giorgi – giorgi@unisi.it --- http://teraflux.eu

# K-Computer (2)



## Packaging of the system

✓ A rack consists of 24 system boards, 6 IO boards, power supply units, system storages, and diagnostic processors.

   ✓ A hose pipe is connected to the water loop under the floor.

CPU

ICC

LSI for interconnect

System Board

~560mm

~460mm

796mm

750mm

2060mm

- S. Fumiyoshi, The K Computer: Project Overview

**TERA<sup>F</sup>LUX**

# IBM Cyclops-64



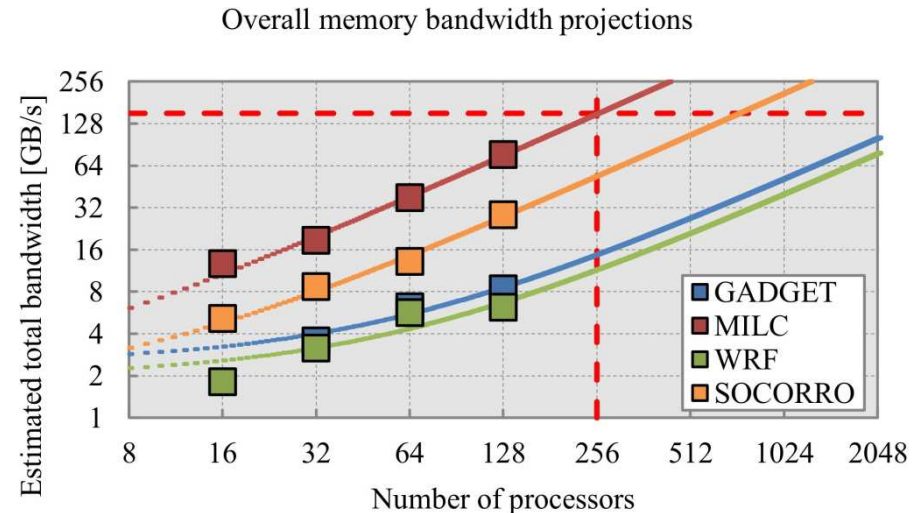Figure 1: IBM Cyclops-64 (C64) Many-Core Architecture: The architecture consists of 80 processors (Processor 0 -79). Each processor has two Thread Units (TUs) called TU 0 and TU 1. Both share one Floating-Point Unit (FPU) and one crossbar port (MPG). Each TU is connected to a SRAM bank, which can be accessed by all other TUs via the crossbar. Ten TUs share one Instruction Cache (IC). The system has four on-chip DDR2 memory controllers to access off-chip memory. The A-Switch is used to connect to the six surrounding neighbors in a 3D-mesh network.

Jürgen Ributzka, Yuhei Hayashi, Joseph B. Manzano, and Guang R. Gao. 2011. The elephant and the mice: the role of non-strict fine-grain synchronization for modern many-core architectures. In *Proceedings of the international conference on Supercomputing* (ICS '11). ACM, New York, NY, USA, 338-347.

TERA**F**LUX

# TERAFLUX RESEARCH OVERVIEW

# WP2 - APPLICATIONS

- MPI applications

- Measures with executions with 16, 32, 64, 128 procs.
  - Linear regression projections

- Off-chip Memory Bandwidth
  - 1K cores will need more than 256GB/s sustained bandwidth
  - 3x than current DDR3

- Total memory footprint for MPI applications
  - Total memory footprint increases with the number of processors
  - Manycores with more than 100 cores will require a few dozens GBs of main memory

- Alternative programming models are required to deal with memory requirements for scalability



Overall memory bandwidth projections

M. Pavlovic, et al. "On the Memory System Requirements of Scientific Applications: Four Case Studies", In IEEE Intl. Symp. On Workload Characterization (IISWC), Nov 2011.

**TERA**F**LUX**

# WP3 – PROGRAMMING MODEL

Transactions and Dataflow additions to Scala

- – Modified to the Scala compiler to include transactional constructs – surveyed other possibilities using closures

- – Runtime STM support

- – Statically Typed Dataflow Library

- – Reimplementation of the Scala parallel collection using dataflow plus transactions

- – Analysis for Lee-TM of benefits of Dataflow plus transactions

http://apt.cs.man.ac.uk/projects/TERAFLUX/MUTS/

Daniel Goodman, Behram Khan, Salman Khan, Chris Kirkham, Mikel Lujan and Ian Watson.
MUTS: Native Scala Constructs for Software Transactional Memory. In: Scala Days Workshop,
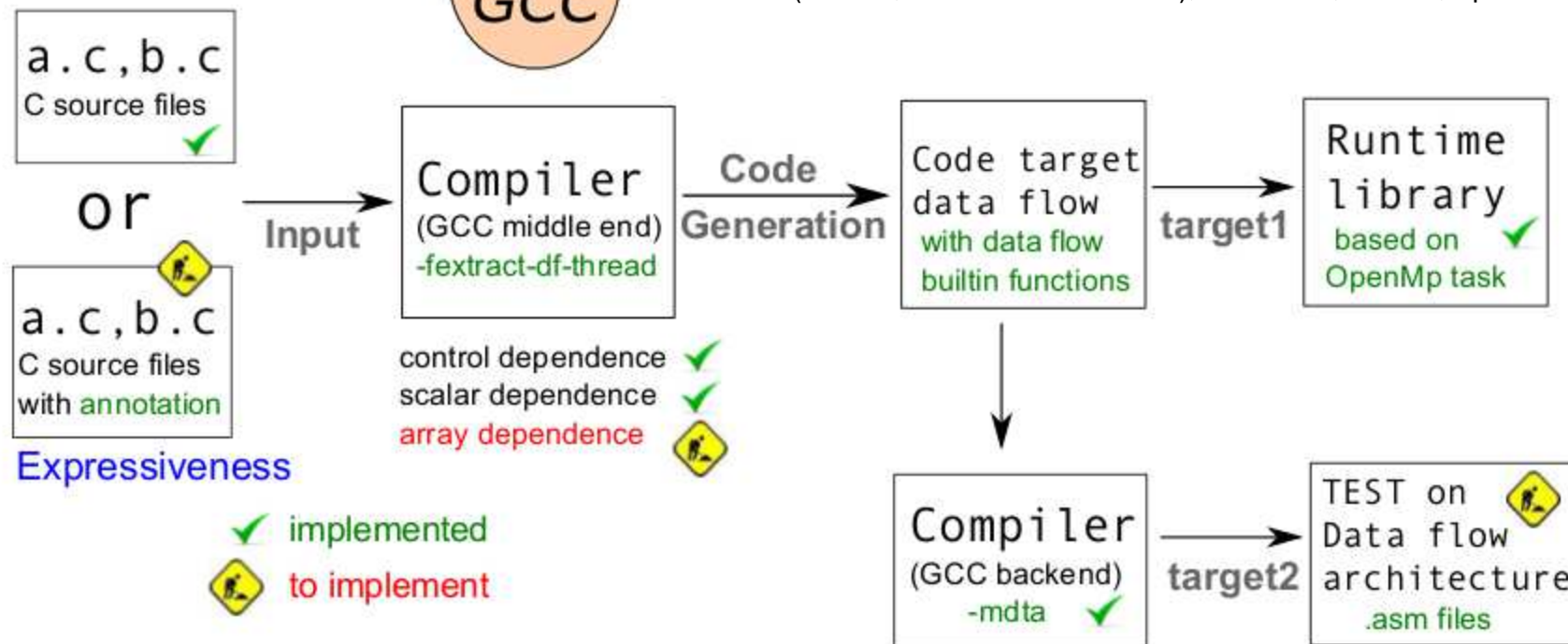Stanford, California, June 2-3, 2011.

**TERA<sup>F</sup>LUX**

# WP4 – COMPILATION TOOLS

## Compilation for Dataflow Threads
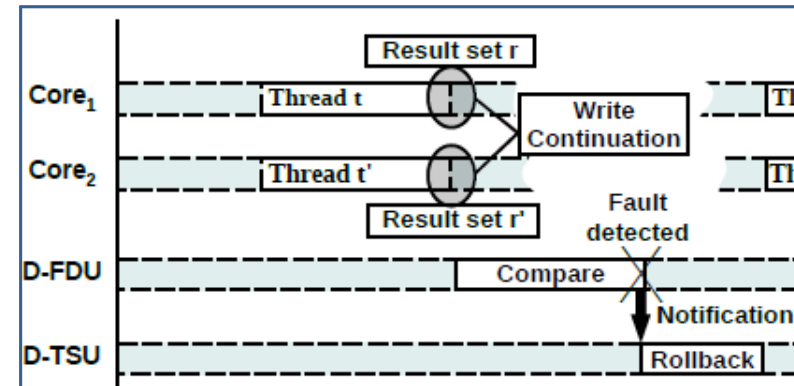### Automatic DF Thread Extraction

GCC

Feng Li, Antoniu Pop, and Albert Cohen. Extending loop distribution to ps-dswp. In: 1st Workshop on Intermediate Representations (WIR'11, associated with CGO), Chamonix, France, April 2011.

a.c,b.c
C source files ✓

or

Input

a.c,b.c
C source files
with annotation

**Expressiveness**

Compiler
(GCC middle end)
-fextract-df-thread

Code
Generation

control dependence ✓
scalar dependence ✓
array dependence

Code target
data flow
with data flow
builtin functions

target1

Runtime
library
based on ✓
OpenMp task

✓ implemented
to implement

Compiler
(GCC backend)
-mdta ✓

target2

TEST on
Data flow
architecture
.asm files

**TERA**FLUX

# WP5 - RELIABILITY

- DOUBLE EXECUTION detects control flow AND data errors

- Runs each thread twice, once
  as a leading thread $t$ and
  second time as a trailing thread $t'$



- The duplicated threads can run on
  the same core or on different cores
  of the same node/cluster

- Each execution generates signature of output results

- At completion compare the two signatures, if consistent,
  the D-TSU writes its results to subsequent thread frames
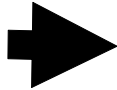
- If not, no commitment
  and recovery

Sebastian Weis, Arne Garbade, Julian Wolf, Bernhard Fechner, Avi Mendelson, Roberto Giorgi, and Theo Ungerer. A Fault Detection and Recovery Architecture for a Teradevice Dataflow System. In Data-Flow Execution Models for Extreme Scale Computing (DFM) 2011 Workshop Proceedings. IEEE Computer Society, 2011

**TERA<sup>F</sup>LUX**

# WP6 - ARCHITECTURE
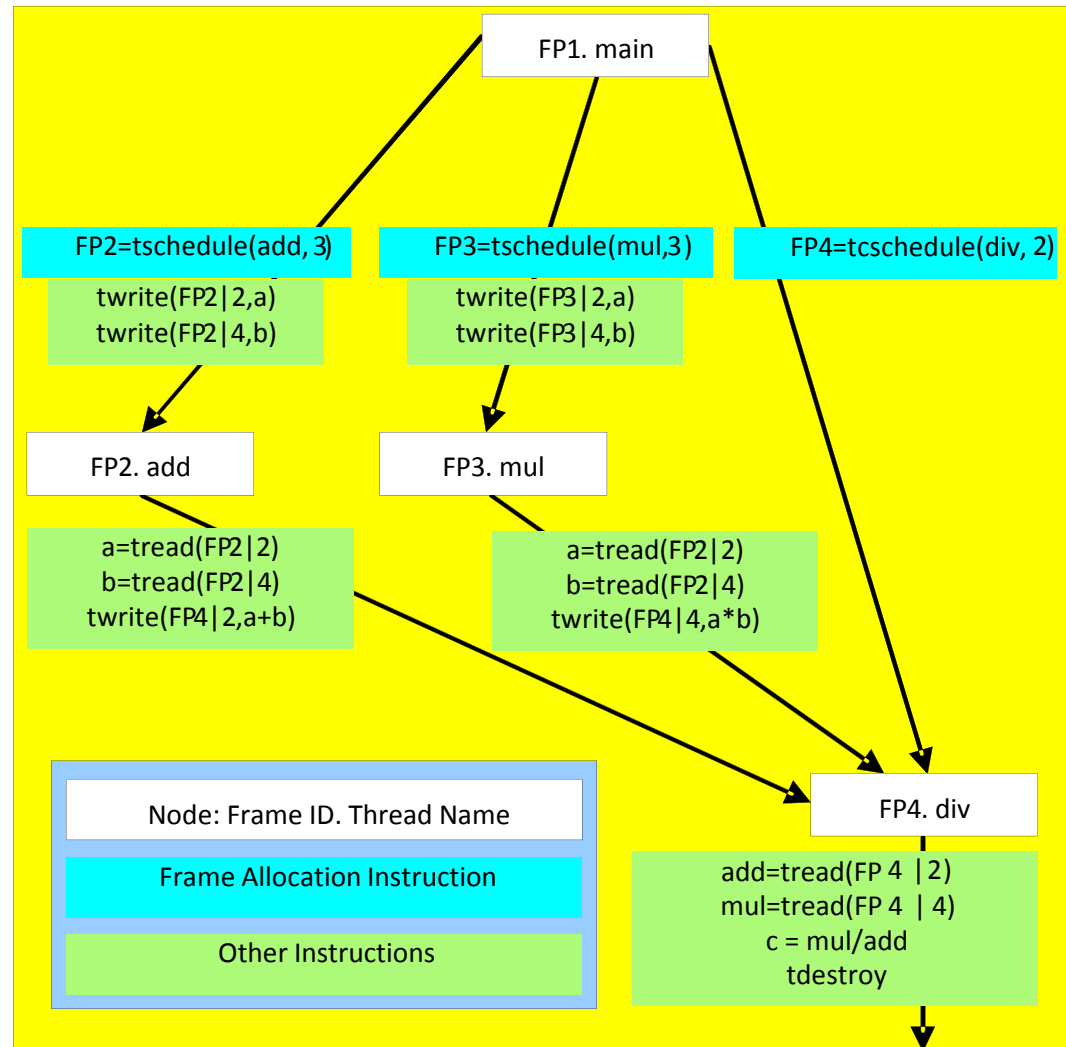
.c (sequential)

```
void main() {
    int a = 4;
    int b = 4;
    int add, mul, c;

    add = a + b;
    mul = a * b;
    c = mul/add;
}
```

.x86 assembly (parallel)

FP1. main

| FP2=tschedule(add, 3) | FP3=tschedule(mul,3) | FP4=tcschedule(div, 2) |

twrite(FP2|2,a)
twrite(FP2|4,b)

twrite(FP3|2,a)
twrite(FP3|4,b)

FP2. add

FP3. mul

a=tread(FP2|2)
b=tread(FP2|4)
twrite(FP4|2,a+b)

a=tread(FP2|2)
b=tread(FP2|4)
twrite(FP4|4,a*b)

Node: Frame ID. Thread Name

Frame Allocation Instruction

Other Instructions

FP4. div

add=tread(FP 4 |2)
mul=tread(FP 4 | 4)
c = mul/add
tdestroy

**TERA<sup>F</sup>LUX**

Roberto Giorgi – giorgi@unisi.it  --- http://teraflux.eu

.x86 assembly (parallel)

```
main:    movq         $4,       %R8
         movq         $4,       %R9
         movq         $1,       %RAX
         cmpq         $1,       %RAX
         TSCHEDULE    $add,     $3,       %R10
         TSCHEDULE    $mult,    $3,       %R11
         TSCHEDULE    $div,     $2,       %R12
         TWRITE       %R8,      %R10,     $2
         TWRITE       %R9,      %R10,     $3
         TWRITE       %R12,     %R10,     $4
         TWRITE       %R8,      %R11,     $2
         TWRITE       %R9,      %R11,     $3
         TWRITE       %R12,     %R11,     $4
         TDESTROY
```

```
add:     TREAD     $2,       %R8
         TREAD     $3,       %R9
         TREAD     $4,       %R10
         movq      %R8,      %R11
         addq      %R9,      %R11
         TWRITE    %R11,     %R10,     $2
         TDESTROY
```

```
mult:    TREAD     $2,       %R8
         TREAD     $3,       %R9
         TREAD     $4,       %R10
         movq      %R8,      %RAX
         mulq      %R9
         movq      %RAX,     %R11
         TWRITE    %R11,     %R10,     $3
         TDESTROY
```

```
div:     TREAD     $2,       %R8
         TREAD     $3,       %R9
         movq      $0,       %RDX
         movq      %R9,      %RAX
         divq      %R8
         movq      %RAX,     %R10
         TDESTROY
```

TERA<sup>F</sup>LUX

# T* (or T86) ISE: TSCHEDULE/TDESTROY

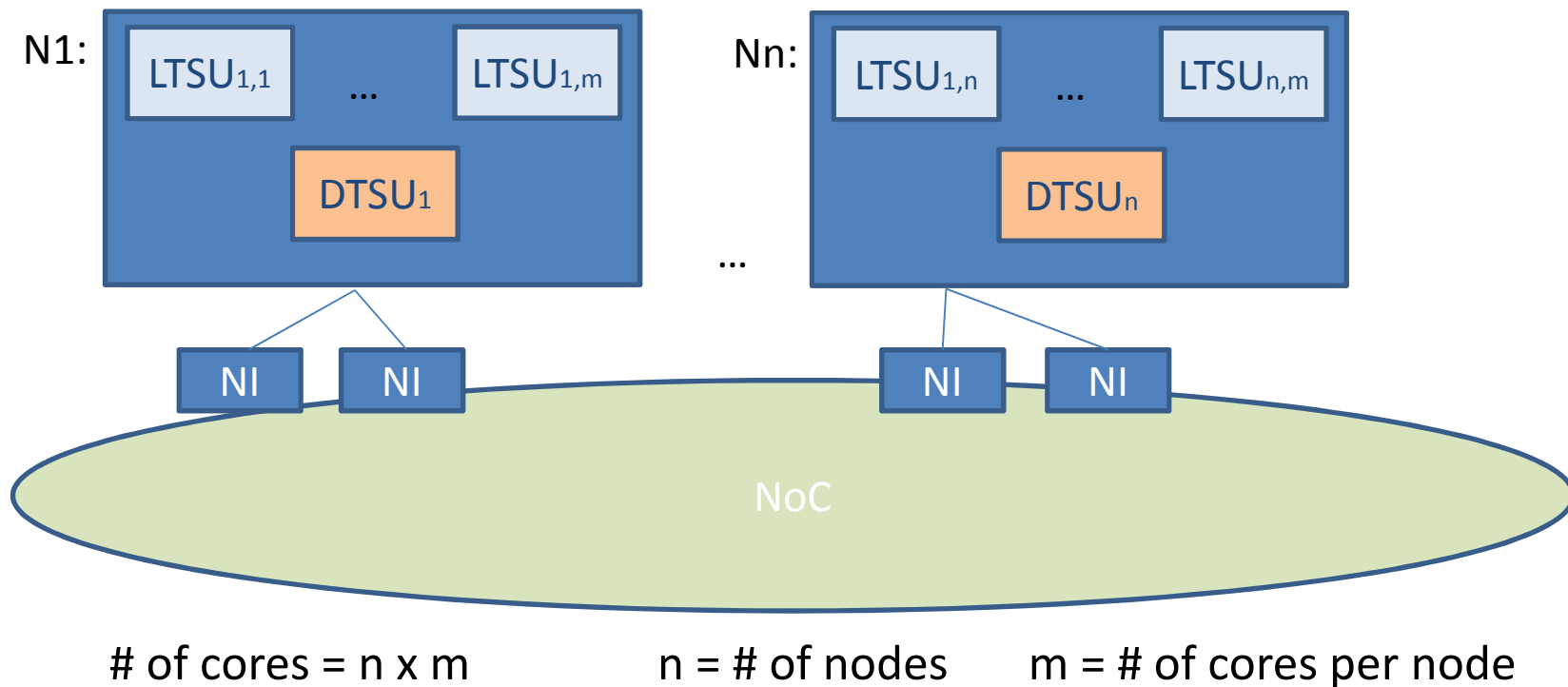|  | T* INSTRUCTIONS | IMPLIED COMPILER TARGET |
|---|---|---|
| Synopsis | TSCHEDULE RS1, RS2, RD | TSCHEDULE(<IP>, <SC>, &<frame_pointer>) |
| Description | This instruction allocates the resources (a DF-frame of size RS2 words and a corresponding entry in the Distributed Thread Scheduler – or DTS) for a new DF-thread and returns its Frame Pointer (FP) in RD. RS1 specifies the Instruction Pointer (IP) of the first instruction of the code of this DF-thread and RS2 specifies the Synchronization Count (SC). Finally the zero flag is logically negated. | |
| Notes | The allocated DF-thread is not executed until its SC reaches 0. The TSCHEDULE can be conditional or non-conditional based on the value stored in the zero flag. If the zero flag is set to 1 then the TSCHEDULE will take effect, otherwise it is ignored. Two subsequent TSCHEDULE implement an if-then-else. | |
| Synopsis | TDESTROY | TDESTROY |
| Description | The thread that invokes TDESTROY finishes and its DF-frame is freed, (the corresponding entry in the Thread Scheduling Unit is also freed). | |
| Notes | - | |

Antoni Portero, Zhibin Yu, Roberto Giorgi, "T-Star (T*): An x86-64 ISA Extension to support thread execution on many cores", *HiPEAC ACACES-2011*, ISBN:978 90 382 17987, Fiuggi, Italy, July 2011, pp. 277-280.

TERAFLUX

# T* (or T86) ISE: TWRITE/TREAD

| | T* INSTRUCTIONS | IMPLIED COMPILER TARGET |
|---|---|---|
| **Synopsis** | **TWRITE RS, RD, offset** | *(<frame_pointer> + <offset>) = (<source_register>) |
| **Description** | The data in RS is stored into the DF-frame pointed to by RD at the specified offset. | |
| **Notes** | Side Effect: The Distributed Thread Scheduler decrements the SC of the corresponding DF-thread entry (located through the FP):   $SC_{FP} = SC_{FP}-1$ | |
| **Synopsis** | **TREAD offset, RD** | (<destination_register>)  =  *(<self_frame_pointer> + <offset>) |
| **Description** | Loads the data indexed by 'offset' from the self (current thread) DF-frame into RD. | |
| **Notes** | Assumption: the DTS has to load into the register implicitly used by TREAD the value <self_frame_pointer>. In a x86-64 implementation, we can reserve RAX for this purpose. | |

**TERA<sup>F</sup>LUX**

# T* (or T86) ISE: TALLOC/TFREE

| | T* INSTRUCTIONS | IMPLIED COMPILER TARGET |
|---|---|---|
| **Synopsis** | **TALLOC RS1, RS2, RD** | **\<pointer\> = TALLOC (\<size\>, \<type\>)** |
| **Description** | Allocates a block of memory of RS1 words. The pointer to it is stored in RD. RS2 specifies the special purpose memory type. | |
| **Notes** | The Distributed Thread Scheduler tracks the memory allocated. An implementation can code \<type\> in the 2 LSB of \<size\> | |
| **Synopsis** | **TFREE (RS)** | **TFREE(\<pointer\>)** |
| **Description** | Frees memory pointed to by RS. | |
| **Notes** | The Thread Scheduling Unit tracks the memory deallocated. | |

TERA\<sup\>F\</sup\>LUX

Roberto Giorgi – giorgi@unisi.it  --- http://teraflux.eu

# DTS – Distributed Thread Scheduler
## (formerly called TSU)

N1:

| | |
|---|---|
| $LTSU_{1,1}$ | ... $LTSU_{1,m}$ |
| | $DTSU_1$ |

Nn:

| | |
|---|---|
| $LTSU_{1,n}$ | ... $LTSU_{n,m}$ |
| | $DTSU_n$ |

...

NI    NI

NI    NI

NoC

# of cores = n x m        n = # of nodes        m = # of cores per node

**TERAᶠLUX**

# Distributed Thread Scheduler
## A 2-node x 4-core example



- Core (Cjk)
  - Off-the-shelf cores (may include L1, L2slice)
  - Minimal ISA extension
- Local Thread Scheduling Unit (LTSU)
  - Manages threads on this Core
- Distributed Thread Scheduling Unit (DTSU)
  - Distributes threads among Nodes
- Node
  - Groups Cores+Resources

# Scheduling Example



- LTSU11 → DTSU1:
  - I need a new frame
- DTSU1 → LTSU13:
  - You're available, give one frame to LTSU11
- LTSU13 → LTSU11:
  - Here is a frame you can use

Fast communication → low overhead
What if every PE in the cluster is busy?

- LTSU14 → DTSU1:
  - "I need a new frame"
- DTSU1 → DTSU2:
  - LSE14 needs a frame, I don't have it
- DTSU2 → LTSU22:
  - Give a frame to LSE14
- LTSU22 → LTSU14:
  - Here is a frame you can use

# Fine-Grain Thread Scheduling

|  | Plurality | CUDA | TFlux | SSI | DTA |
|---|---|---|---|---|---|
| HW/SW | HW + SW | HW + SW | SW | HW | HW |
| Prog. Model | Custom C language extensions | Custom C language extensions | Custom C preprocessor macros (#pragma) | Standard thread-oriented C libraries (pthreads) | Standard thread-oriented C libraries (pthreads) |
| Exec. Model | - Pool of RISC processors<br>- Uniform shared memory<br>- Hardware scheduler , synchronizer and load balancer | - Each thread block is split into warps (thread block.<br>- Each thread block is executed by only one multiprocessor.<br>- A multi-processor can execute several blocks concurrently. | High-level threads are mapped to OS threads, using the standard OS programming interfaces as backend. | Subset static interleaved scheduling of fine-grained / coarse-grained threads performed by a hardware Task Scheduling Unit (TSU) | - Decoupling memory and execution activity of non-blocking threads.<br>- Threads are synchronized and communicate each other in a producer-consumer fashion. |
| Architecture | Complex hardware subsystem: synchronizer, scheduler, load balancer | SIMD, on-chip shared memory | Everything is implemented in software: can run on any general purpose architecture | Hardware scheduler (TSU) + extensions to a VLIW prototype | Thread management and frame memory management implemented in hardware |

SSI=Subset Static Interleaved

DTA=Decoupled Threaded Architecture

**TERA FLUX**

# fib(21): number of threads

# WP7: Evaluating a MANY-CORE chip of the future (2020), i.e., 1000+ cores on a chip

# AMD SIMnow and COTSon

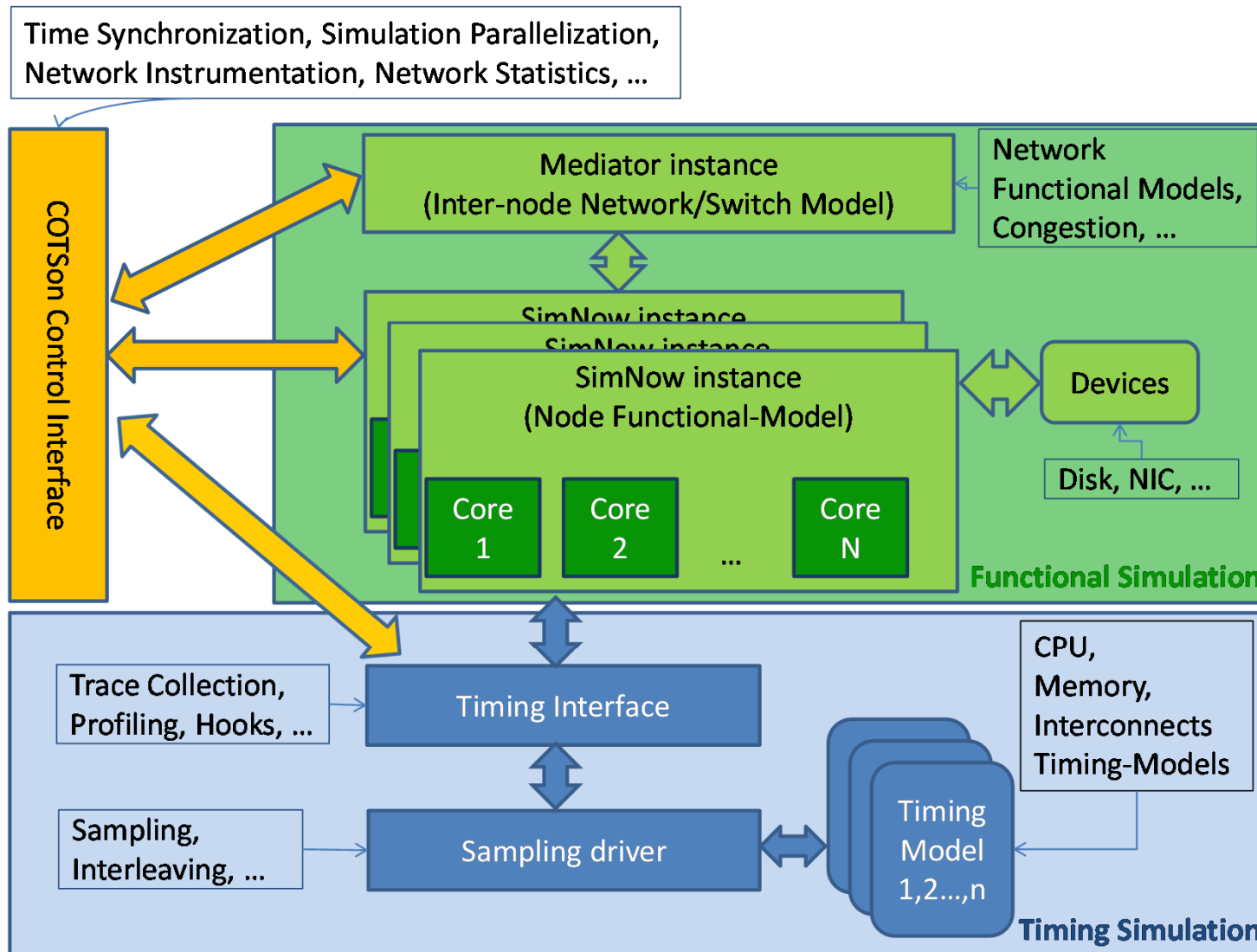TERA<sup>F</sup>LUX

# COTSon Overview



Time Synchronization, Simulation Parallelization, Network Instrumentation, Network Statistics, …

**COTSon Control Interface**

**Mediator instance**
**(Inter-node Network/Switch Model)**

Network Functional Models, Congestion, …

SimNow instance
SimNow instance
**SimNow instance**
**(Node Functional-Model)**

Devices

Disk, NIC, …

| Core 1 | Core 2 | … | Core N |

**Functional Simulation**

Trace Collection, Profiling, Hooks, …

**Timing Interface**

CPU, Memory, Interconnects Timing-Models

Sampling, Interleaving, …

**Sampling driver**

Timing Model 1,2…,n

**Timing Simulation**

**TERA**F**LUX**

**The ambition of TERAFLUX is however
to be able of *changing* such machine in a flexible way,
while tackling research challenges on programmability,
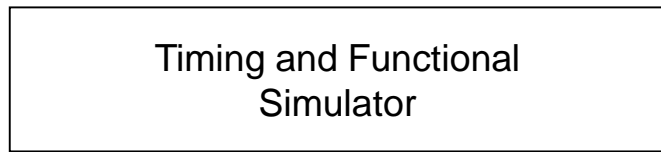architectural design and reliability.
Therefore, we have the need to stress the COTSon platform,
in order to being able to simulate 1000 cores.**

**TERA<sup>F</sup>LUX**

# Comparison among different approaches for doing research related to 1000-core computing system (Information revised from data of the RAMP project)

| | SMP | Cluster | FPGA | Emulator | Simulator |
|---|---|---|---|---|---|
| **Scalability (1K cores)** | C | A | A | A | A |
| **Cost (1K cores)** | F(€40M) | C | B(€0.1-0.2M) | A+ (€0.01M) | A+(€0.01M) |
| **Power/Space (Kw, racks)** | D (120 kw, 12 racks) | D (120 kw, 12 racks) | A (1.5 kw, 0.3 racks) | A+ (0.1 kw, 0.1racks) | A+ (0.1 kw, 0.1 racks ) |
| **Observability** | D | C | A+ | A+ | A+ |
| **Reproducibility** | B | D | A+ | A+ | A+ |
| **Reconfigurability** | D | C | A+ | A+ | A+ |
| **Credibility** | A+ | A+ | B+/A- | F/D | C |
| **Development time** | B | B | C | A+ | A+ |
| **Performance (clock)** | A (2GHz) | A(3GHz) | C (0.1 GHz) | B(≈ 0.9 of original) | C(1/10 to 1/1000 SMP) |
| **x86-64 ISA** | A+ | A+ | F | A+ | A+ |
| **Modifiable** | F | F | B | A | A |
| **GPA** | D | D | B+/A- | B | A |

TERA$^F$LUX

# FUNCTIONAL/TIMING SIMULATION

| Timing and Functional Simulator |
|---|

**Integrated (SimOS)**
- Complex, no reuse, very slow

| Timing Simulator | → ← | Functional Simulator |
|---|---|---|

- Complete Timing
- No? Function

- No Timing
- Complete Function

**Timing-Directed (Exec-driven)**
+ Timing feedback
- Tight Coupling
- Very slow

| Timing Simulator | → ⋯ | Functional Simulator |
|---|---|---|

- Complete Timing
- Partial Function

- No Timing
- Complete Function

**Timing-First (Multifacet)**
+ Timing feedback
+ Using existing simulators
+ Software development advantages
- Slow

| Functional Simulator | → | Timing Simulator |
|---|---|---|

- Complete Function

- Partial Timing

**Functional-First (Trace-driven)**
+Fast
- No timing feedback

speed

accuracy

*Source: Multifacet Project (www.cs.wisc.edu/multifacet) - [Mauer02-sigmetrics-Full_System Timing_First Simulation]*

**TERA**^F**LUX**

Roberto Giorgi – giorgi@unisi.it --- http://teraflux.eu

# COTSon: FUNCTIONAL-DIRECTED

- A variant of "functional first"

  - Adds timing feedback at coarse granularity
    (100s – 1000s of instructions)

- Applications see an approximation of time

  - May miss some fine-grain timing interaction

- Compatible with fast (caching) emulators and samplers

# OTHER RECENT X64 SIMULATORS

|  | Sniper | Graphite | Gem5 | COTSon | MARSSx86 |
|---|---|---|---|---|---|
| Timing-directed/integrated |  |  | X |  |  |
| Func-directed | X | X |  | X | X |
|  |  |  |  |  |  |
| User-level | X | X | X | X | X |
| Full-system |  |  | X | X | X |
|  |  |  |  |  |  |
| Archs Supported | x64 | x64 | x64 Alpha SPARC | x64 | x64 |
| Parallel (in-node) | X | X |  | Multi-node |  |
| Shared caches | X |  | X | X | X |

Heirman120401-ISPASS Tutorial
The SNIPER multi-core simulator

**TERA<sup>F</sup>LUX**

# Simulation booting up 1024 cores. (1) COTSon execution of 32 SimNow instances. (2) Each instance manages 32 cores. Host: 48 cores, 256 GB memory



**Note: the simulation is PARALLEL at GUEST NODE-LEVEL and it's also possible to distribute the simulation on several HOST NODES**

TERA<sup>F</sup>LUX

# INITIAL EXPERIMENTAL RESULTS

The proposed simulation framework has been validated running applications and benchmarks on a target machine with up to 1024 cores, operating in accordance with dataflow principle on standard cores

We run several applications and benchmarks based on well established programming models (mainly OpenMP and MPI):



- NAS Parallel Benchmark (NPB)
- Graph500 and HPL 2.0 Linpack
- Sequential Recursive Fibonacci

**TERA<sup>F</sup>LUX**

# Major Technical Innovations in TERAFLUX

- Fragmenting the Applications in Finer grained DF-threads:
  - DF-threads allow an easy way to decouple memory accesses, therefore hiding memory latencies, balancing the load, managing fault, temperature information without fine grain intervention of the software.
- Possibility to repeat the execution of a DF-thread in case this thread happened to be on a core later discovered as faulty
- Taking advantage of a "direct" dataflow communication of the data (through what we call DF-frames).
- Synchronizing threads while taking advantage of native dataflow mechanism (e.g. several threads can be synchronized at a barrier)
  - DF-threads allow (atomic ) Transactional semantics (DF meets TM)
- A Thread Scheduling Unit would allow fast thread switching and scheduling, besides the OS scheduler; scalable and distributed
- A Fault Detection Unit works in conjunction with TSU

**TERA<sup>F</sup>LUX**

# TERAFLUX SIMULATOR (COTSon)

# http://cotson.sf.net

## HP-Labs COTSon is **OPEN-SOURCE**

**FUTURE AND
EMERGING
TECHNOLOGIES
PROJECT N. 249013**

**SEVENTH FRAMEWORK
PROGRAMME THEME
FET proactive 1 (ICT-2009.8.1)
Concurrent Tera-Device Computing**

COOPERATION

# TERA<sup>F</sup>LUX

## Exploiting dataflow parallelism in Teradevice Computing

**PROJECT NUMBER: 249013**

**http://teraflux.eu**

# BACKUP SLIDES

**TERAFLUX**

# TERAFLUX TOOLCHAIN (Jan. 2011)



**TERA**<sup>F</sup>**LUX**

Roberto Giorgi – giorgi@unisi.it --- http://teraflux.eu

Author - Partner

Free frame table

| CID | FFN |
|-----|-----|
|     |     |
|     |     |
|     |     |
|     |     |

Control logic

TCRQ – simple FIFO queue

- FreeFrameTable (CID – Core ID, FFN – number of free frames)
  - Keeps track of the occupancy of processors inside a cluster
  - Updated on each TDestroy and accepted TSchedule
- TCRQ – ThreadScheduleRequestQueue
  - Holds unserved ThreadScheduleRequest messages
  - Message is pushed into queue when there are no free local resources
  - Message is popped from the queue when either TDestroyor BroadcastResponse arrives
- Control logic
  - Responsible for both inter and intra node communication and updating the messages inside a scheduler

## Map table

| V | M | ThID | Fptr |
|---|---|------|------|
|   |   |      |      |
|   |   |      |      |
|   |   |      |      |
|   |   |      |      |

## Store buffer

| V | ThID | Offset | DATA |
|---|------|--------|------|
|   |      |        |      |
|   |      |        |      |
|   |      |        |      |
|   |      |        |      |

- Map table (V – Valid, M – Mapped, ThID – thread ID, Fptr – Frame pointer)
    - Keeps track of the issued resource requests for the execution of new threads
    - a ThID is assigned to a thread when new Tschedule instruction occurs; it is used just inside that core
    - a Fptr is assigned when a TScheduleResponse message arrives; it is unique globally in the system
    - Can be cleared on the thread completion

- Store buffer (V – valid, ThID – thread ID, offset – for storing in a frame, DATA – data to store)
    - Keeps track of the TStores issued for the threads that didn't receive a TScheduleResponse yet (those kept in Map table and still not mapped)
    - On each TStore for the new thread that still doesn't have resources assigned, a new entry is created
    - When TScheduleResponse arrives, all entries are checked and TStore messages are sent (entry invalidated) if there is any matching
    - If TStore occurs for a thread that already has its resources assigned, there is no need to use the buffer

# Distributed Thread Scheduler Unit

- On new TScheduleRequestMessage checks the availability in local node
  - If yes – forwards it
  - If no – put the message in FRQ and send broadcast
- Message is removed from FRQ when FfreeMessage or BroadcastResponse arrive
- Other messages are just forwarded



DISTRIBUTED SCHEDULER ELEMENT

TERA**F**LUX

# Local Thread Scheduling Unit

- On TScheduleRequestMessage
  - Choose a free frame for execution of the new thread
  - Send TScheduleResponseMessage to the issuing processor
- On TScheduleResponseMessage simply update the continuation with the frame identifier
- On store send DataStore message (group them if the destination is the same)

# Non-blocking resource assignment (1)

- Avoid waiting from the distributed scheduler by introducing Virtual frame pointers
- Two additional structures – map table and store buffer



Even if we don't speed-up the starting time of new threads, execution time is shorter and processor becomes free earlier

**TERA<sup>F</sup>LUX**

# One Physical Machine running two Virtual Machine instances that communicate through the Virtual Network (Mediator).



ADVANTAGES
This setup can run both on a real machines (at least at smal scale for tests) AND on the COTSon simulator
It allows us to modify system parameters like e.g. number o. cores in each simulated instance.
It allows for a parallelization of the simulation (the several instances are running in parallel on the available cores – load balancing automatically provided by the Host OS scheduler).
Possible to avoid copying buffers among instances because they reside in the Host Shared Memory Network
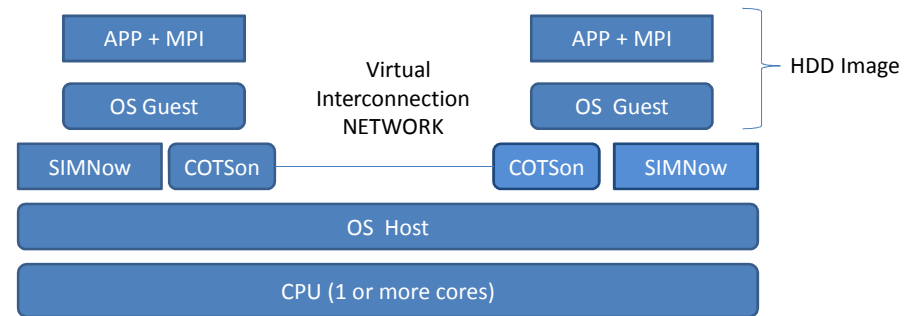Possibility to take advantage of RVI/VT-x virtualization mechanisms across different Physical Machines (under development).
The communication and synchronization among the simulation instances adds up to the Application traffic, but could bypass TCP/IP and avoid using the Physical Interconnection Network.
No need to use the Physical Network.

DISADVANTAGES

- Taking into account that we aim to flexibly change the programming model and architecture (e.g. the dataflow based execution model and architecture), this setup may end up in poor performance when N (number of nodes) increases.

- Tightens the Application to the Machine, which is exactly the opposite direction that we follow globally in TERAFLUX: we aim to decouple the Application (WP2) from the Machine with appropriate Programming Models (WP3), Compilation Tools (WP4) and Execution Models (WP6).

- The MPI run-time is constantly involved to appropriately schedule the ready tasks/threads on the available nodes.

- The Physical architecture that is more natural to model is a Distributed Machine not like the general one we aim in TERAFLUX.

**TERA<sup>F</sup>LUX**

Roberto Giorgi – giorgi@unisi.it --- http://teraflux.eu

# VM instances governed by a Single Source Image (SSI) OS

ADVANTAGES

- Allows us to run Shared Memory applications like OpenMP ones (can still run MPI as if it was a single big node).

- Can run both on a real machines (at least at small scale for tests) and on the COTSon simulator.

- It allows us to modify system parameters like e.g. number of cores in each simulated instance.

- It allows for a parallelization of the simulation (the several instances are running in parallel on the available cores – load balancing automatically provided by the Host OS scheduler).

- Possible to avoid copying buffers among instances because they reside in the Host Shared Memory Network

- Possibility to take advantage of RVI/VT-x virtualization mechanisms across different Physical Machines (under development).

- The communication and synchronization among the simulation instances adds up to the Application traffic, but could bypass TCP/IP and avoid using the Physical Interconnection Network.
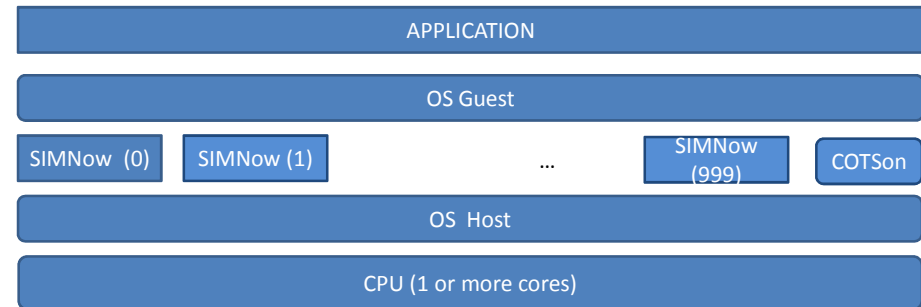
- Load Balancing for the Application is managed by the Guest OS

- No need to use the Physical Network.



| APPLICATION | | | | |
| OS Guest | | | | |
| SIMNow (0) | SIMNow (1) | ... | SIMNow (999) | COTSon |
| OS Host | | | | |
| CPU (1 or more cores) | | | | |

DISADVANTAGES
This setup requires the use of a Distributed OS as Guest OS (like e.g., Kerrighed [KERRIGHED10], which offers the view of a unique SMP machine on top of a cluster) or in general a SSI (Single System Image) OS.
Relatively poor performance when N (number of nodes) increases;
Partially tightens the Application to the Machine, which is in the opposite direction in respect to what we follow globally in TERAFLUX: we aim to decouple the Application (WP2) from the Machine with appropriate Programming Models (WP3), Compilation Tools (WP4) and Execution Models (WP6).
The underlying Guest Architecture is a "cluster", which is then more naturally mapped to a physical Distributed Machine not a generic one like we aim in TERAFLUX.

**TERA<sup>F</sup>LUX**

Roberto Giorgi – giorgi@unisi.it --- http://teraflux.eu

# One core aware of all the other cores

ADVANTAGES
- Allows us to run Shared Memory applications like OpenMP ones (can still run MPI as if it was a single big node).
- Can run both on a real machines (at least at small scale for tests) AND on the COTSon simulator as provided at the Month-1 of the TERAFLUX project.
- It allows us to modify system parameters like e.g. number of cores in each simulated instance.
- It allows for a parallelization of the simulation (the several instances are running in parallel on the available cores – load balancing automatically provided by the Host OS scheduler).
- Possible to avoid copying buffers among instances because they reside in the Host Shared Memory Network.
- Possibility to take advantage of RVI/VT-x virtualization mechanisms across different Physical Machines (under development).
- The communication and synchronization among the simulation instances adds up to the Application traffic, but could bypass TCP/IP and avoid using the Physical Interconnection Network.
- Load Balancing for the Application is managed by the Guest OS
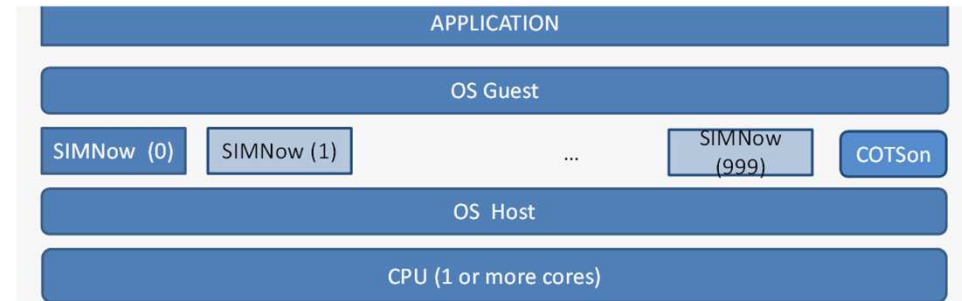- No need to use the Physical Network.
- No need to use a very different OS like an SSI OS.
- The underlying Guest Architecture is a shared memory machine, however thanks to the availability of a global address space, there is now full possibility of evolving the machine in a more "general one" like the one we aim to evolve during the TERAFLUX project. The TERAFLUX Execution Model can decouple completely the architecture of the machine.

| APPLICATION |
| OS Guest |
| SIMNow (0) | SIMNow (1) | ... | SIMNow (999) | COTSon |
| OS Host |
| CPU (1 or more cores) |

DISADVANTAGES

- Relatively poor performance when N (number of nodes) increases; however, as other simulator like COREMU [Wang11] already demonstrated a high speed up in simulations even with 255 cores, we have good confidence that we can improve much the simulation speed going in a similar direction.

- Requires some patches to the Linux OS; however we shall need to patch anyway the Memory Manager and the Scheduler in order to properly support the TERAFLUX threads

**TERA$^F$LUX**

# NAS benchmarks running in COTSon

- Machine 37nodes of 4 cores. One node Master and 36 Slaves
- Two examples from the set:

NAS Parallel Benchmarks 3.3 -- BT Benchmark
 No input file inputbt.data. Using compiled default
Class A: Size:  64x  64x  64
Iterations: 200   dt:  0.0008000
Number of active processes:    36

|  | Time in seconds | Mop/s total | Mop/s/process |
|---|---|---|---|
| BT 1-4 | 398.93 | 421.84 | 11.72 |
| BT 2-4 | 422.08 | 398.7 | 11.08 |
| BT 4-4 | 398.17 | 422.65 | 11.74 |

NAS Parallel Benchmarks 3.3 – CG

Size:     14000

Iterations:    15

Number of active processes:    32

Number of nonzeroes per row:       11

Eigenvalue shift: .200E+02

|  | Time in seconds | Mop/s total | Mop/s/process |
|---|---|---|---|
| CG 1-4 | 46.26 | 32.35 | 1.01 |
| CG 2-4 | 48.45 | 30.89 | 0.97 |
| CG 4-4 | 46.65 | 32.08 | 1 |

**TERA$^F$LUX**

# NAS benchmarks running in COTSon (cont.)

NAS Parallel Benchmarks 3.3 -- EP Benchmark
 Number of random numbers generated:      536870912
 Number of active processes:              32
EP Benchmark Results: CPU Time =    5.5777, N = 2^  28
─────────────────────────────────────
 NAS Parallel Benchmarks 3.3 -- FT Benchmark
 No input file inputft.data. Using compiled defaults
 Size         : 256x 256x 128 (Class A)
 Iterations     :        6, Number of processes :     32
 Processor array   :       1x 32, Layout type    :    1D
─────────────────────────────────
 NAS Parallel Benchmarks 3.3 -- IS Benchmark

 Size:  8388608  (class A), Iterations:   10
 Number of processes:    32, IS Benchmark Completed
 Class      =     A, Size       =    8388608
 Iterations   =       10
─────────────────────────────────────
 NAS Parallel Benchmarks 3.3 -- LU Benchmark
 Size:   64x  64x  64 (Class A),Iterations:  250
 Number of processes:    32
─────────────────────────────────
 NAS Parallel Benchmarks 3.3 -- MG Benchmark
 No input file. Using compiled defaults
 Size: 256x 256x 256  (class A)
 Iterations:   4,  Number of processes:     32

| | Time in seconds | Mop/s total | Mop/s/ process |
|---|---|---|---|
| | | 410.53 | |
| EP 1-4 | 4.9301 | 108.9 | 3.4 |
| EP 2-4 | 5.5777 | 96.25 | 3.01 |
| EP 4-4 | 4.9324 | 108.85 | 3.4 |
| FT 1 -4 | 187.25 | 38.11 | 1.19 |
| FT 2 -4 | 185.43 | 38.49 | 1.2 |
| FT 4 -4 | 201.85 | 35.36 | 1.1 |
| IS 1-4 | 2.59 | 32.39 | 1.01 |
| IS 2 -4 | 2.67 | 31.45 | 0.98 |
| IS 4 -4 | 2.57 | 32.64 | 1.02 |
| LU 1-4 | 188.96 | 631.33 | 19.73 |
| LU 2 -4 | 185.15 | 644.32 | 20.14 |
| LU 4 -4 | 183.93 | 648.6 | 20.27 |
| MG 1-4 | 171.15 | 22.74 | 0.71 |
| MG 2 -4 | 176.11 | 22.1 | 0.69 |

# Plurality

- Plurality: http://www.plurality.com/profile.html

    - Architecture: general-purpose accelerator for multicore/manycore system-on-chip (SoC)

    - Task-oriented programming model: the programmer  has to perform a partitioning of the program into specific tasks (task-map)

    - The body of each task is a traditional sequential code

    - Each core is a RISC processor

    - Scheduler, synchronization and load balancing among cores are done by a complex hardware subsystem that communicates with all the RISC processors

    - Uniform shared memory access

TERA<sup>F</sup>LUX

# CUDA

– Programming model: extensions to standard C language (CUDA libraries)

– DRAM memory addressing + on-chip shared memory

– However a single process must run spread across multiple disjoint memory spaces (???)

– Recursive functions are not supported (must be converted to loops)

– Bus bandwidth and latency between CPU and GPU may be a bottleneck (acceleratore esterno)

TERA<sup>F</sup>LUX

# TFlux

- TFlux:

  - Paralell processing system targeted to commodity, Linux-based shared-memory multiprocessor systems

  - Data-Driven multi-threading

  - Programming model: takes as input a C program, argumented with TFlux-specific compiler directives (#pragma's)

  - Everything implemented in software

TERA<sup>F</sup>LUX

# Subset Static Interleaved (SSI)

- Interleaved threads

  – Advantage: operations latencies become shorter in terms of executed instructions from the same thread

- Combination of blocked multithreading and static interleaved multithreading:

  – A set of background threads + a set of foreground threads. Foreground threads are interleaved until a stall occurs (e.g. cache miss). When a foreground thread stalls and a background thread is ready for execution we exchange them so that the foreground thread becomes a background thread and vice versa

- TSU: task scheduling unit (in hardware)

TERA<sup>F</sup>LUX