**International Workshop on A Strategic Initiative of Computing: Systems and Applications (SISA): Integrating HPC, Big Data, AI and Beyond**

# Integrated Development of Parallelizing and Power Reducing Compiler and Multicore Architecture for HPC to Embedded Applications

# Hironori Kasahara

**Professor, Dept. of Computer Science & Engineering**

**Director, Advanced Multicore Processor Research Institute**

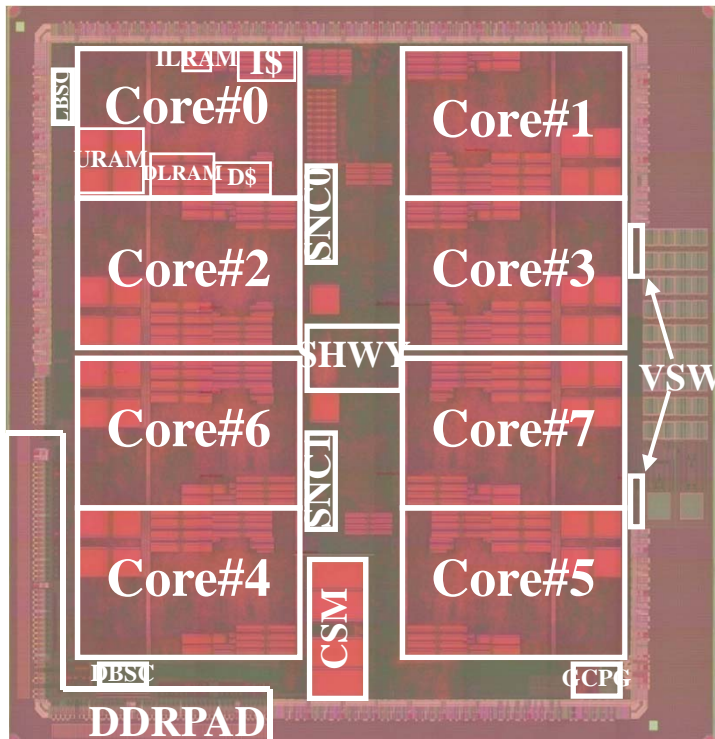**Waseda University (早稲田大学), Tokyo, Japan**

**IEEE Computer Society**
**President Elect 2017, President 2018**

**URL: http://www.kasahara.cs.waseda.ac.jp/**

**Waseda Univ. Green Computing R&D Center, Jan. 18, 2017**

# Multicores for Performance and Low Power

**Power consumption is one of the biggest problems for performance scaling from smartphones to cloud servers and supercomputers ("K" more than 10MW) .**



IEEE ISSCC08: Paper No. 4.5, M.ITO, … and  H. Kasahara, "An 8640 MIPS SoC with Independent Power-off Control of 8 CPUs and 8 RAMs by an Automatic Parallelizing Compiler"

**Power $\propto$ Frequency * Voltage$^2$**
**(Voltage $\propto$ Frequency)**
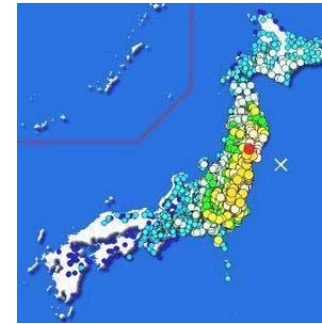
➡ **Power $\propto$ Frequency$^3$**

**If Frequency is reduced to 1/4 (Ex. 4GHz➔1GHz), Power is reduced to 1/64 and Performance falls down to 1/4 .**
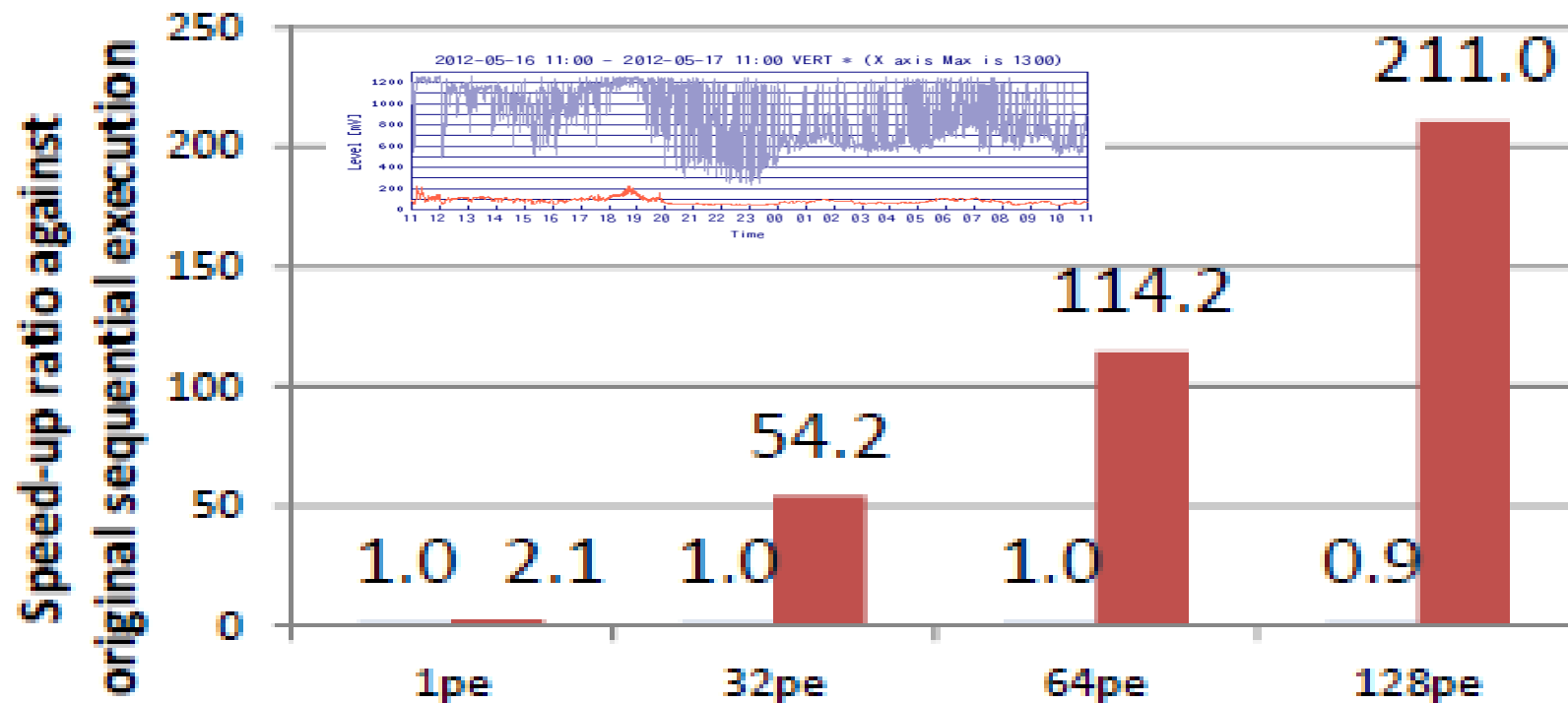**<Multicores>**
**If 8cores are integrated on a chip, Power is still 1/8 and Performance becomes 2 times.**

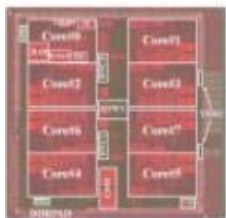# Earthquake Simulation "GMS" on Fujitsu M9000 Sparc CC-NUMA Server



Legend: ■ original (sun studio)  ■ proposed method



Speed-up ratio against original sequential execution:

| | 1pe | 32pe | 64pe | 128pe |
|---|---|---|---|---|
| original (sun studio) | 1.0 | 1.0 | 1.0 | 0.9 |
| proposed method | 2.1 | 54.2 | 114.2 | 211.0 |

**With 128 cores, OSCAR compiler gave us 100 times speedup against 1 core execution and 211 times speedup against 1 core using Sun (Oracle) Studio compiler.**
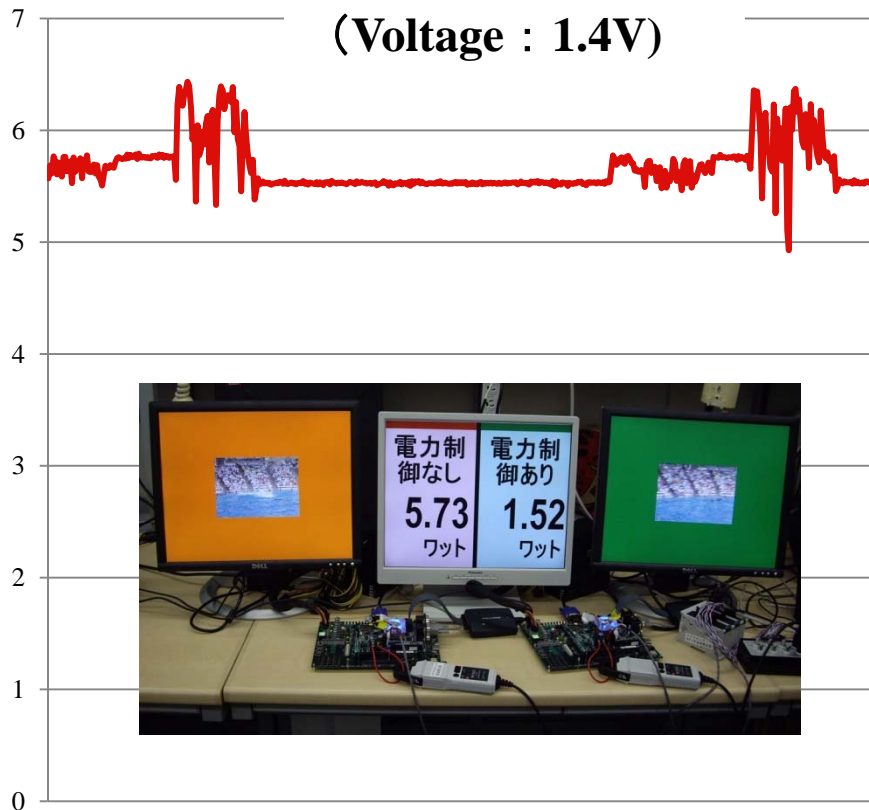
# Power Reduction of MPEG2 Decoding to 1/4 on 8 Core Homogeneous Multicore RP-2 by OSCAR Parallelizing Compiler

## MPEG2 Decoding with 8 CPU cores
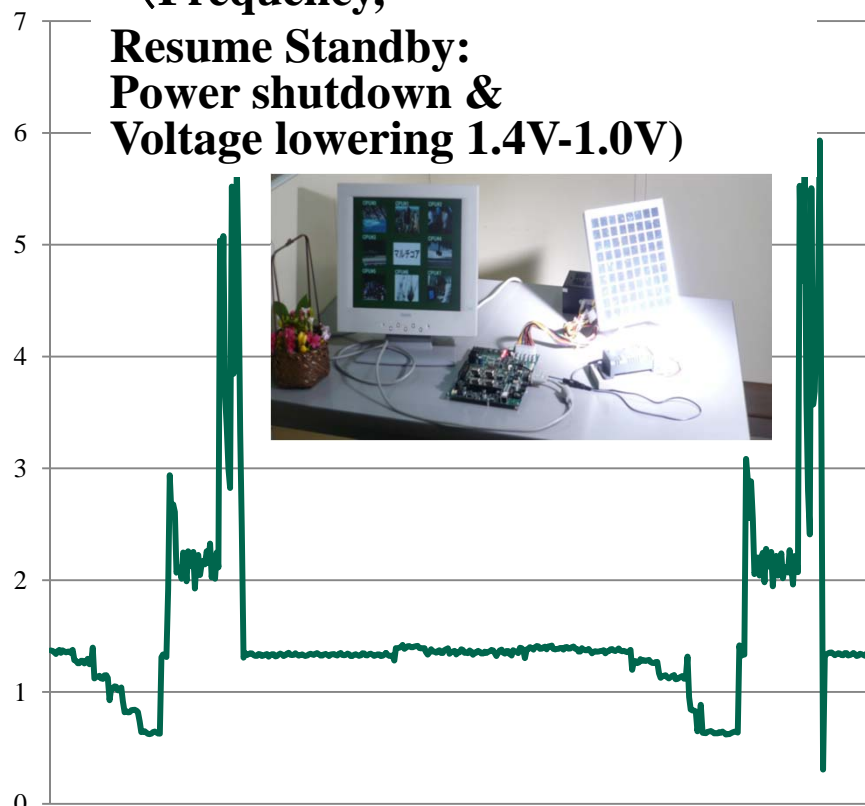
**Without Power Control （Voltage：1.4V)**

**With Power Control （Frequency, Resume Standby: Power shutdown & Voltage lowering 1.4V-1.0V)**



**Avg. Power 5.73 [W]**

**73.5% Power Reduction**

**Avg. Power 1.52 [W]**

# Renesas-Hitachi-Waseda Low Power 8 core RP2 Developed in 2007 in METI/NEDO project



| Process Technology | 90nm, 8-layer, triple-Vth, CMOS |
|---|---|
| Chip Size | 104.8mm$^2$ (10.61mm x 9.88mm) |
| CPU Core Size | 6.6mm$^2$ (3.36mm x 1.96mm) |
| Supply Voltage | 1.0V–1.4V (internal), 1.8/3.3V (I/O) |
| Power Domains | 17 (8 CPUs, 8 URAMs, common) |

**IEEE ISSCC08: Paper No. 4.5, M.ITO, … and H. Kasahara, "An 8640 MIPS SoC with Independent Power-off Control of 8 CPUs and 8 RAMs by an Automatic Parallelizing Compiler"**

# Demo of NEDO Multicore for Real Time Consumer Electronics at the Council of Science and Engineering Policy on April 10, 2008

第74回総合科学技術会議 【平成20年4月10日】

第74回総合科学技術会議の様子(1)

第74回総合科学技術会議の様子(2)

第74回総合科学技術会議の様子(3)

第74回総合科学技術会議の様子(4)

**CSTP Members**

**Prime Minister:**
Mr. Y. FUKUDA

**Minister of State for Science, Technology and Innovation Policy:**
Mr. F. KISHIDA

**Chief Cabinet Secretary:**
Mr. N. MACHIMURA

**Minister of Internal Affairs and Communications :**
Mr. H. MASUDA
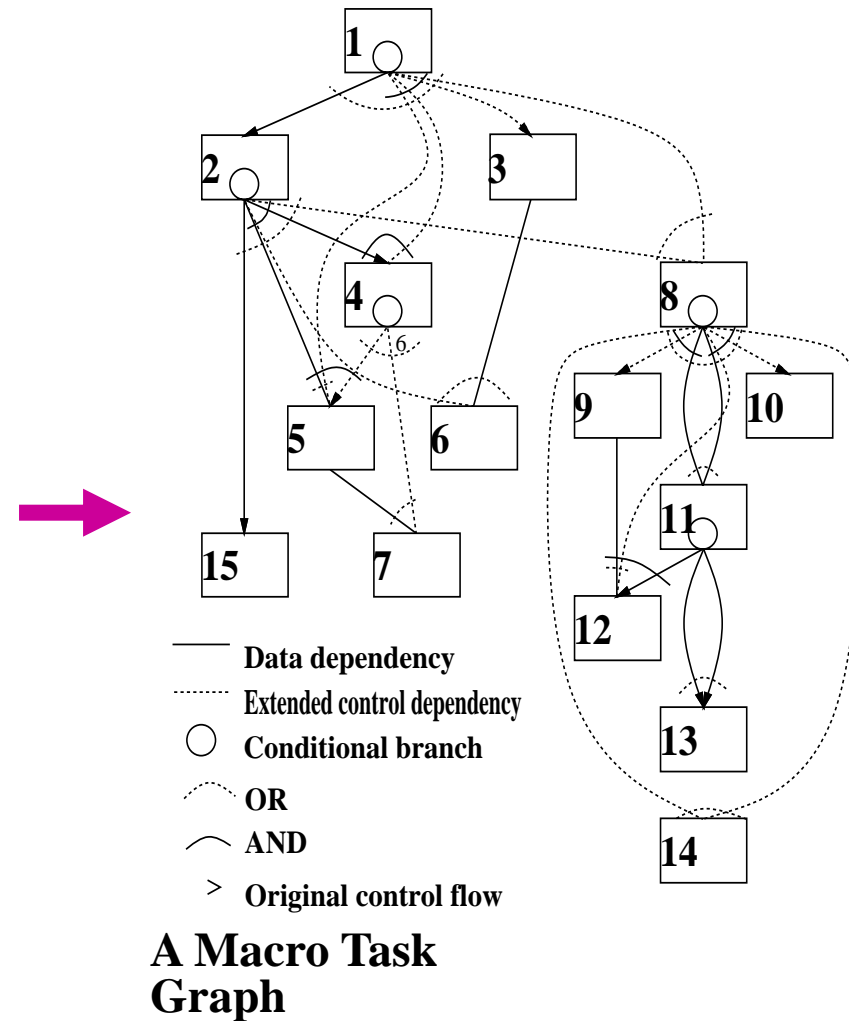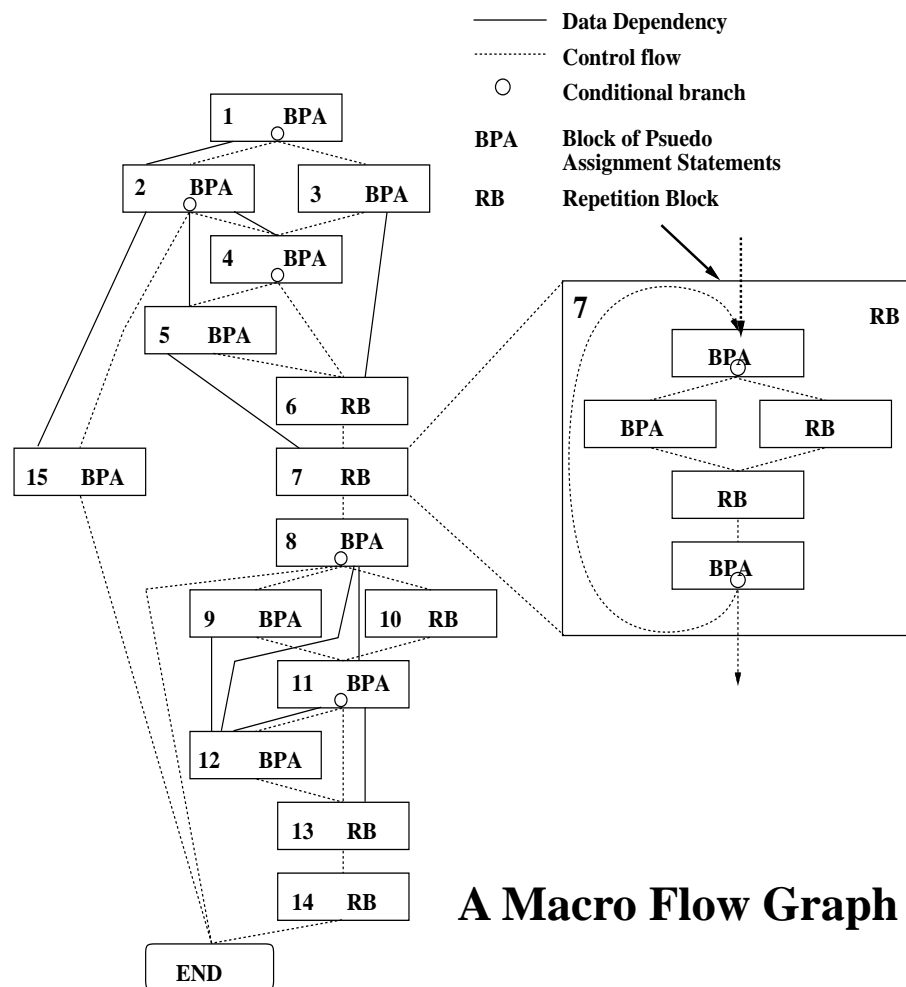
**Minister of Finance :**
Mr. F. NUKAGA

**Minister of Education, Culture, Sports, Science and Technology:**
Mr. K. TOKAI

**Minister of Economy,Trade and Industry:**
Mr. A. AMARI

# Earliest Executable Condition Analysis for Coarse Grain Tasks (Macro-tasks)



A Macro Flow Graph

A Macro Task Graph

# OSCAR Parallelizing Compiler

## To improve effective performance, cost-performance and software productivity and reduce power

### Multigrain Parallelization

coarse-grain parallelism among loops and subroutines, near fine grain parallelism among statements in addition to loop parallelism
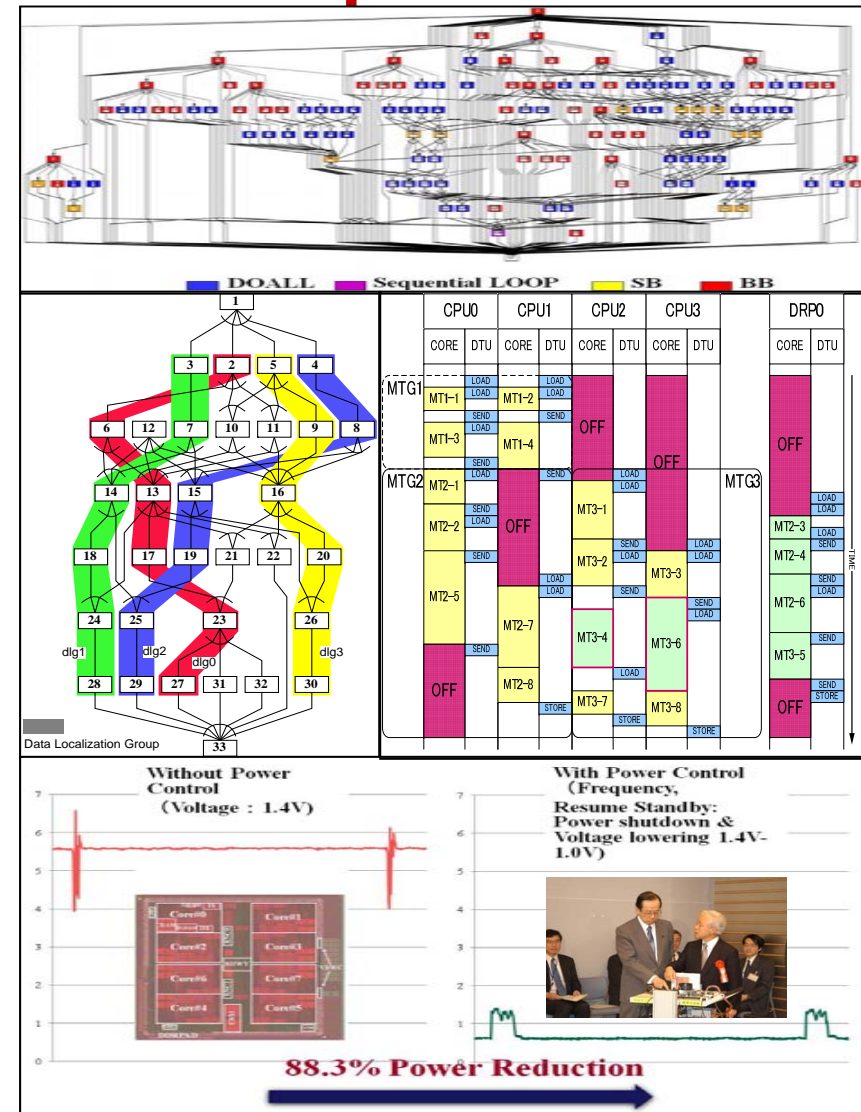
### Data Localization

Automatic data management for distributed shared memory, cache and local memory

### Data Transfer Overlapping

Data transfer overlapping using Data Transfer Controllers (DMAs)

### Power Reduction

Reduction of consumed power by compiler control DVFS and Power gating with hardware supports.



DOALL  Sequential LOOP  SB  BB

Data Localization Group

Without Power Control (Voltage : 1.4V)

With Power Control (Frequency, Resume Standby: Power shutdown & Voltage lowering 1.4V-1.0V)

88.3% Power Reduction

# Data Localization: Loop Aligned Decomposition

- **Decomposed loop into LRs and CARs**
  - LR ( Localizable Region): Data can be passed through LDM
  - CAR (Commonly Accessed Region): Data transfers are required among processors
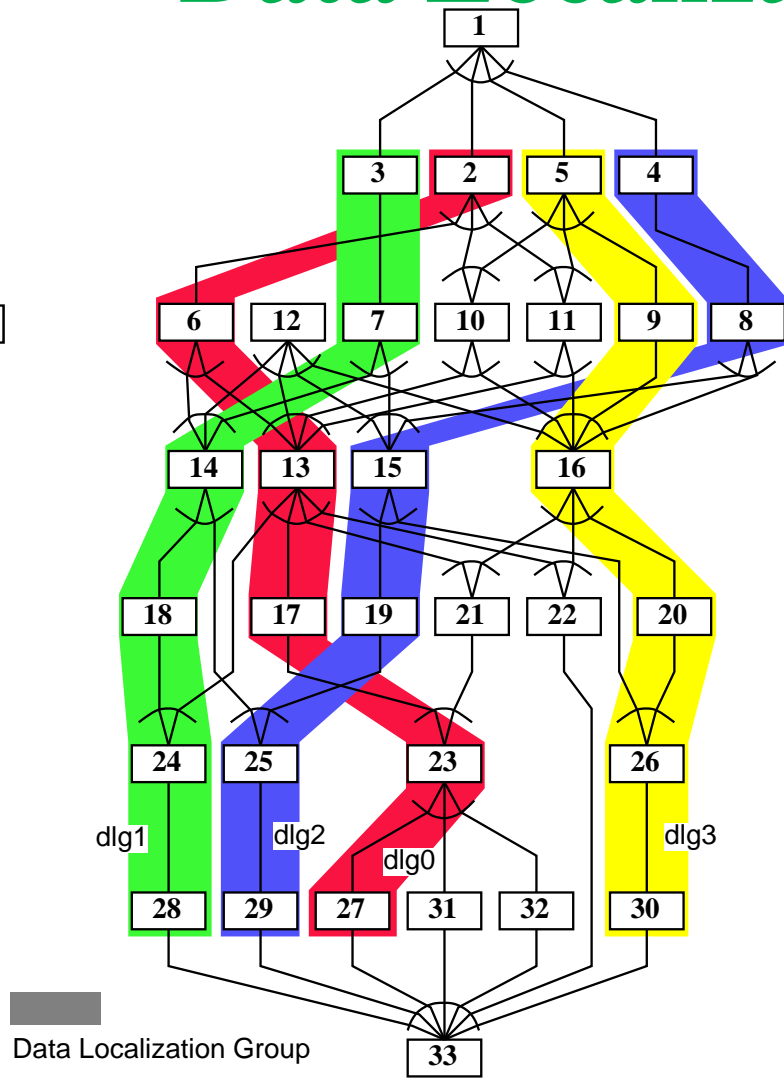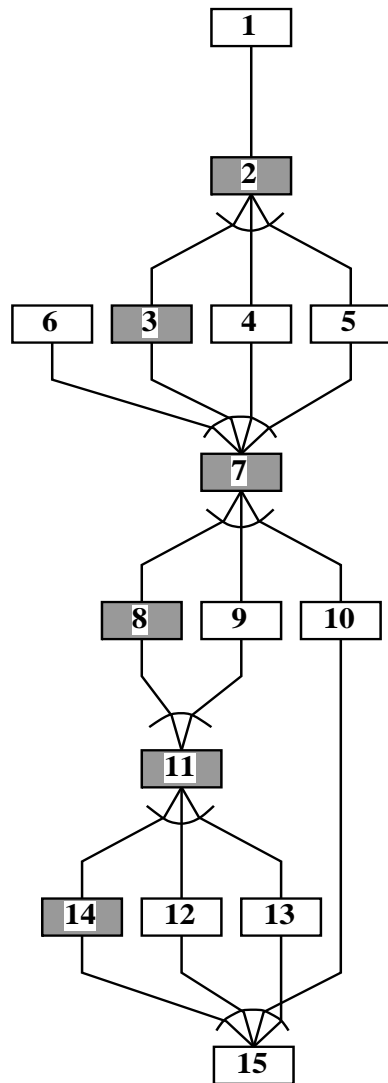
**Single dimension Decomposition**
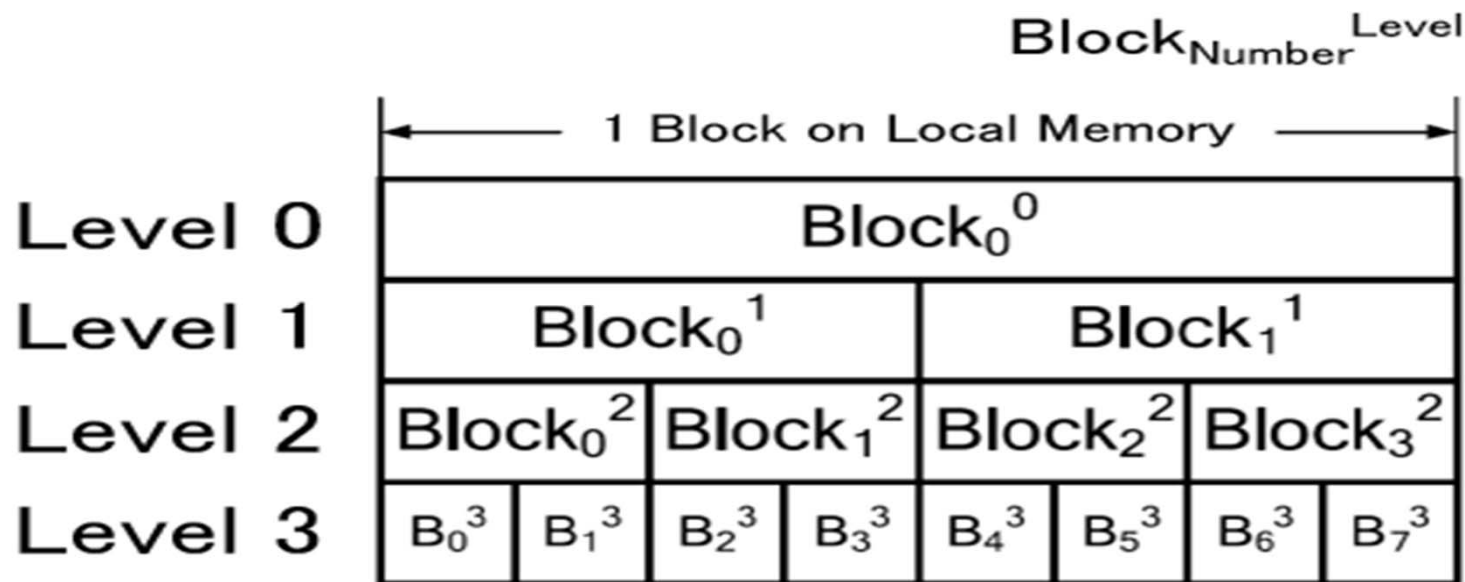
**Multi-dimension Decomposition**

# Data Localization



MTG

MTG after Division

A schedule for two processors

Data Localization Group

dlg0  dlg1  dlg2  dlg3

PE0  PE1

10

# Local Memory Management Using Adjustable Blocks

- **Decide a suitable block size for each application**
  - **different from fixed block sizes like in cache**
  - **each block can be divided into smaller blocks with integer divisible size to handle small arrays and scalar variables**

$$\text{Block}_{\text{Number}}^{\text{Level}}$$

| | 1 Block on Local Memory | | | | | | |
|---|---|---|---|---|---|---|---|
| Level 0 | $\text{Block}_0^0$ | | | | | | |
| Level 1 | $\text{Block}_0^1$ | | | | $\text{Block}_1^1$ | | |
| Level 2 | $\text{Block}_0^2$ | | $\text{Block}_1^2$ | | $\text{Block}_2^2$ | | $\text{Block}_3^2$ |
| Level 3 | $B_0^3$ | $B_1^3$ | $B_2^3$ | $B_3^3$ | $B_4^3$ | $B_5^3$ | $B_6^3$ | $B_7^3$ |

11

# Multi-dimensional Template Arrays for Improving Readability

- **a mapping technique for arrays with varying dimensions**
  - **each block on LDM corresponds to multiple empty arrays with varying dimensions**
  - **these arrays have an additional dimension to store the corresponding block number**
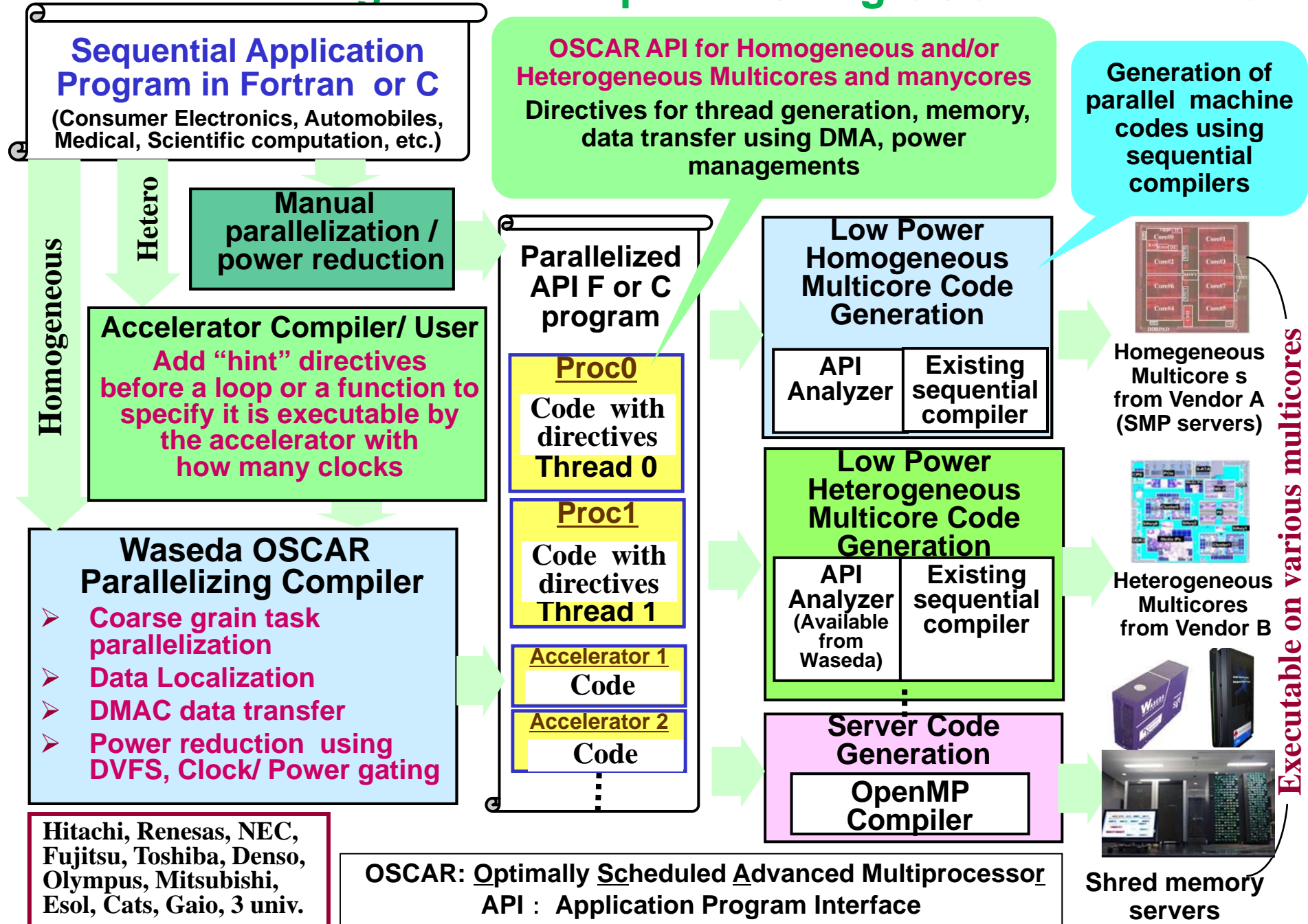    - **TA[Block#][] for single dimension**
    - **TA[Block#][][] for double dimension**
    - **TA[Block#][][][] for triple dimension**
    - **...**
- **LDM are represented as a one dimensional array**
  - **without Template Arrays, multi-dimensional arrays have complex index calculations**
    - **A[i][j][k] -> TA[offset + i' * L + j' * M + k']**
  - **Template Arrays provide readability**
    - **A[i][j][k] -> TA[Block#][i'][j'][k']**



TEMPLATE ARRAY FOR 1-DIMENSIONAL ARRAY

TEMPLATE ARRAY FOR 2-DIMENSIONAL ARRAY

TEMPLATE ARRAY FOR 3-DIMENSIONAL ARRAY

Block0
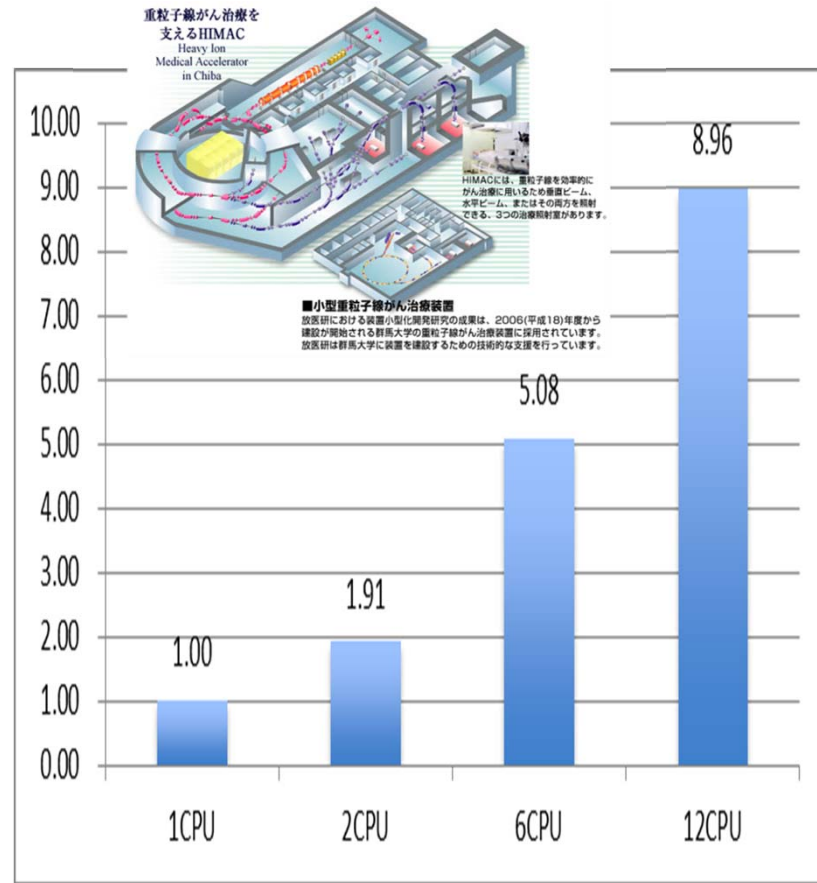Block1
Block2
⋮
Block7

LDM

# Block Replacement

➤ Appropriate memory blocks considering schedules are replaced

  – Dead, live and reuse information of each block is used.

  – Different from LRU using past information, this method uses future information available from static schedule.

➤ Block Replacement Priority

1. Dead Block (Variables) that will not be accessed later.

2. Live Blocks that are accessed only by the other cores.

3. Live Block that will be accessed by the current core latest.

4. Live Block that will be accessed by the current core soon and data transfer overhead can be hidden by DMA overlapped transfer.
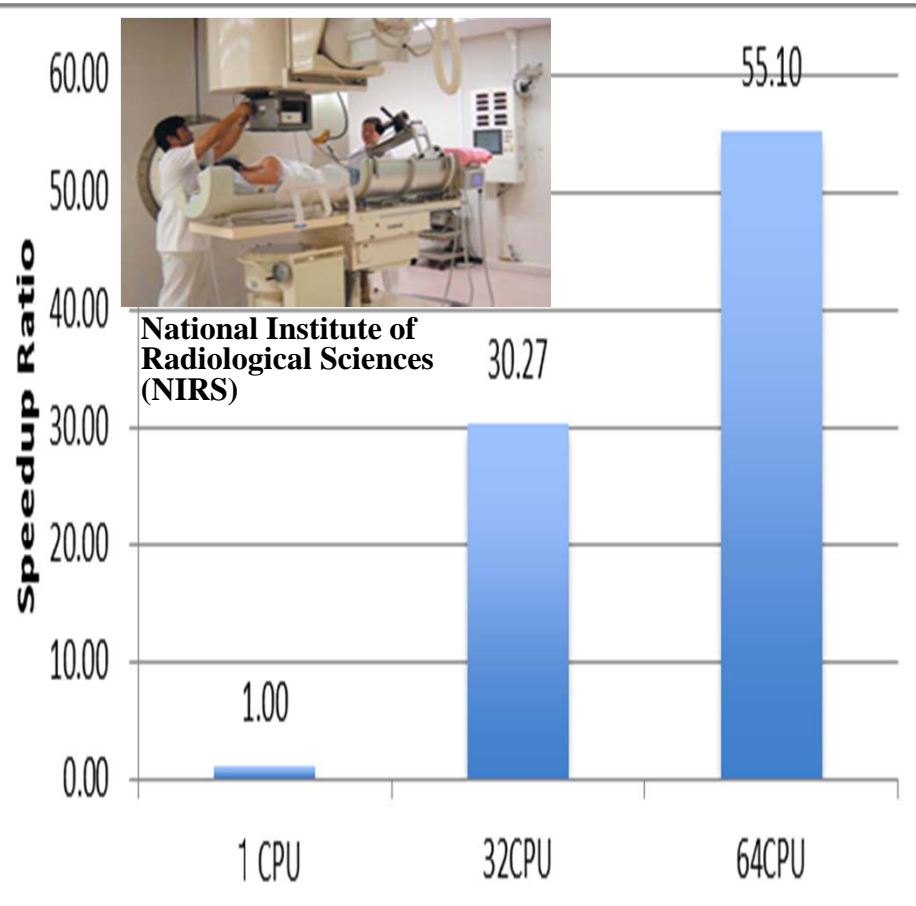
13

# Multicore Program Development Using OSCAR API V2.0

**Sequential Application Program in Fortran or C**
(Consumer Electronics, Automobiles, Medical, Scientific computation, etc.)

**OSCAR API for Homogeneous and/or Heterogeneous Multicores and manycores**
Directives for thread generation, memory, data transfer using DMA, power managements

**Generation of parallel machine codes using sequential compilers**

Hetero

Homogeneous

**Manual parallelization / power reduction**

**Accelerator Compiler/ User**
Add "hint" directives before a loop or a function to specify it is executable by the accelerator with how many clocks

**Waseda OSCAR Parallelizing Compiler**
- Coarse grain task parallelization
- Data Localization
- DMAC data transfer
- Power reduction using DVFS, Clock/ Power gating

**Parallelized API F or C program**

**Proc0**
Code with directives
Thread 0

**Proc1**
Code with directives
Thread 1

**Accelerator 1**
Code

**Accelerator 2**
Code

**Low Power Homogeneous Multicore Code Generation**

| API Analyzer | Existing sequential compiler |
|---|---|

**Low Power Heterogeneous Multicore Code Generation**

| API Analyzer (Available from Waseda) | Existing sequential compiler |
|---|---|

**Server Code Generation**
OpenMP Compiler

Homegeneous Multicore s from Vendor A (SMP servers)

Heterogeneous Multicores from Vendor B

Executable on various multicores

Shred memory servers

Hitachi, Renesas, NEC, Fujitsu, Toshiba, Denso, Olympus, Mitsubishi, Esol, Cats, Gaio, 3 univ.

OSCAR: Optimally Scheduled Advanced Multiprocessor
API : Application Program Interface

# Cancer Treatment
# Carbon Ion Radiotherapy
## (Previous best was 2.5 times speedup on 16 processors with hand optimization)



National Institute of Radiological Sciences (NIRS)

**8.9times speedup by 12 processors**

**Intel Xeon X5670 2.93GHz 12 core SMP (Hitachi HA8000)**

**55 times speedup by 64 processors**

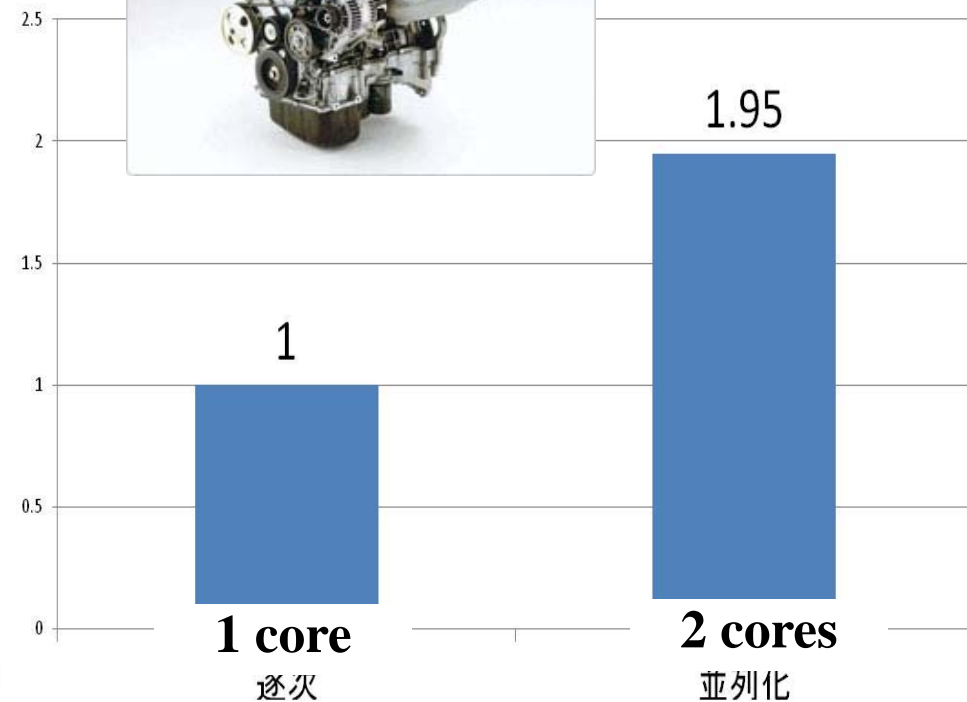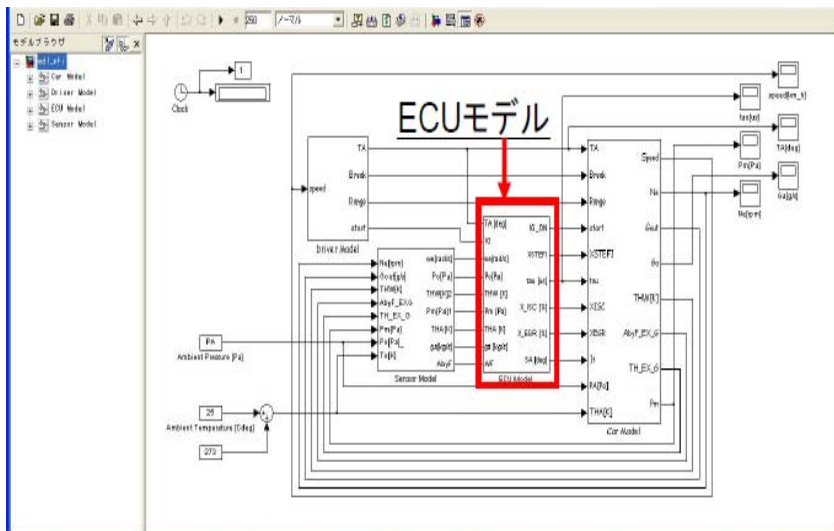**IBM Power 7    64 core SMP (Hitachi SR16000)**
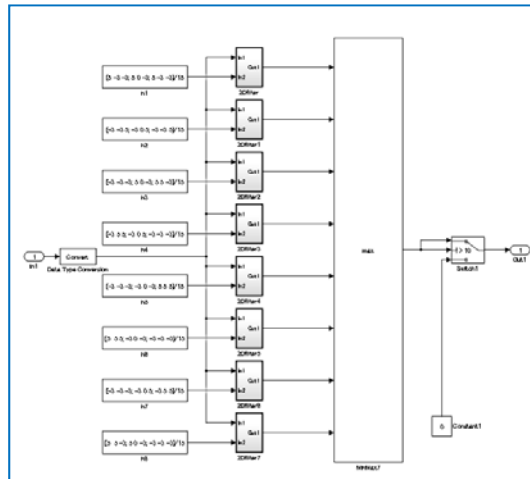
# Engine Control by multicore with Denso

**Though so far parallel processing of the engine control on multicore has been very difficult, Denso and Waseda succeeded 1.95 times speedup on 2core V850 multicore processor.**

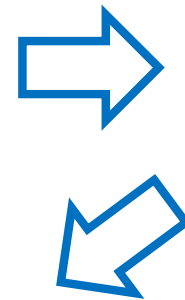**Hard real-time automobile engine control by multicore**

# OSCAR Compile Flow for Simulink Applications



**Simulink model**

**Generate C code using Embedded Coder**

```
/* Model step function */
void VesselExtraction_step(void)
{
  int32_T i;
  real_T u0;

  /* DataTypeConversion: '<S1>/Data Type Conversion' incorporates:
   *  Inport: '<Root>/In1'
   */
  for (i = 0; i < 16384; i++) {
    VesselExtraction_B.DataTypeConversion[i] = VesselExtraction_U.In1[i];
  }

  /* End of DataTypeConversion: '<S1>/Data Type Conversion' */

  /* Outputs for Atomic SubSystem: '<S1>/2Dfilter' */

  /* Constant: '<S1>/h1' */
  VesselExtraction_Dfilter(VesselExtraction_B.DataTypeConversion,
    VesselExtraction_P.h1_Value, &VesselExtraction_B.Dfilter,
    (P_Dfilter_VesselExtraction_T *)&VesselExtraction_P.Dfilter);

  /* End of Outputs for SubSystem: '<S1>/2Dfilter' */

  /* Outputs for Atomic SubSystem: '<S1>/2Dfilter1' */

  /* Constant: '<S1>/h2' */
  VesselExtraction_Dfilter(VesselExtraction_B.DataTypeConversion,
    VesselExtraction_P.h2_Value, &VesselExtraction_B.Dfilter1,
    (P_Dfilter_VesselExtraction_T *)&VesselExtraction_P.Dfilter1);
```
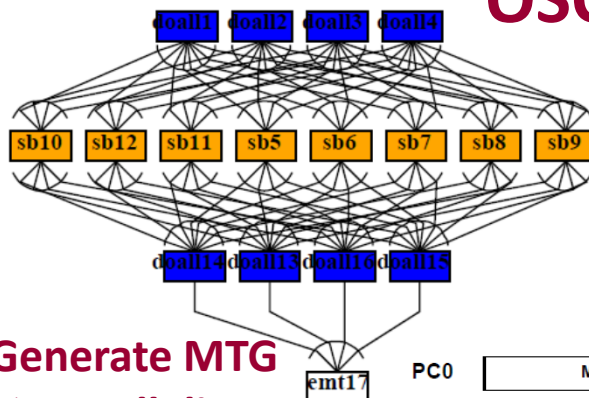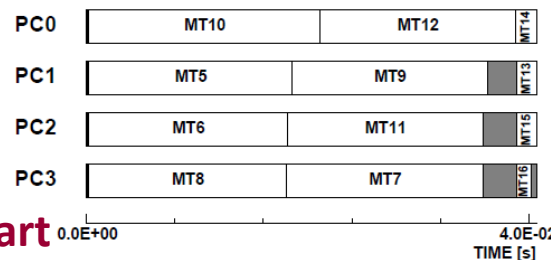
**C code**

## OSCAR Compiler



**(1) Generate MTG**
   **→ Parallelism**

**(2) Generate gantt chart**
   **→ Scheduling in a multicore**

```
void VesselExtraction_step ( )
{
    int thr1 ;
    int thr2 ;                    void thread_function_001 ( void )
    int thr3 ;                    {
    {                                 VesselExtraction_step_PE1 ( ) ;
        oscar_thread_create ( & thr1 ,
            thread_function_001 , (void*)1 ) ;
        oscar_thread_create ( & thr2 ,
            thread_function_002 , (void*)2 ) ;
        oscar_thread_create ( & thr3 ,
            thread_function_003 , (void*)3 ) ;

        VesselExtraction_step_PE0 ( ) ;

        oscar_thread_join ( thr1 ) ;
        oscar_thread_join ( thr2 ) ;
        oscar_thread_join ( thr3 ) ;
```
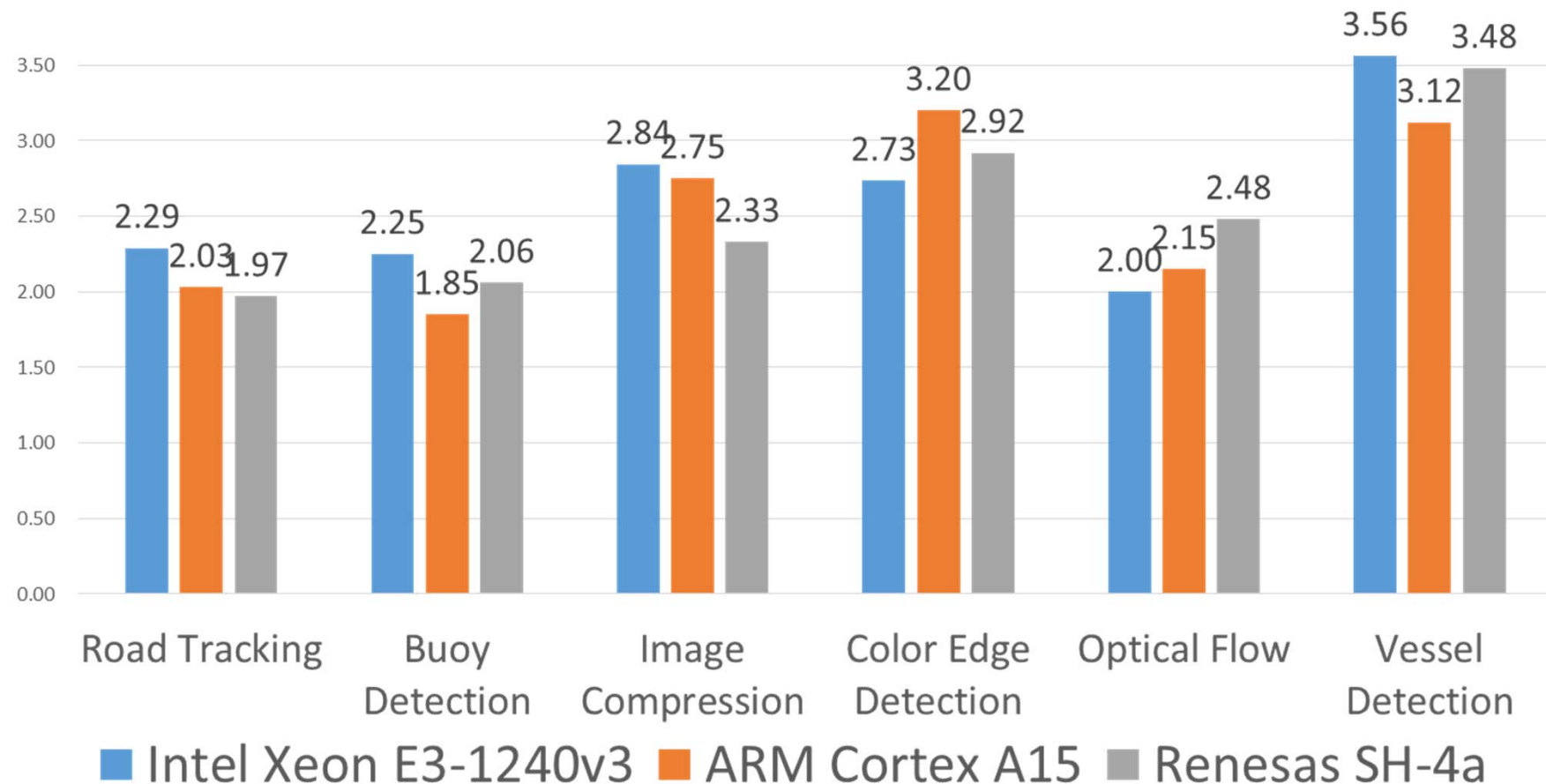
**(3) Generate parallelized C code using the OSCAR API**
   **→ Multiplatform execution (Intel, ARM and SH etc)**

# Speedups of MATLAB/Simulink Image Processing on Various 4core Multicores
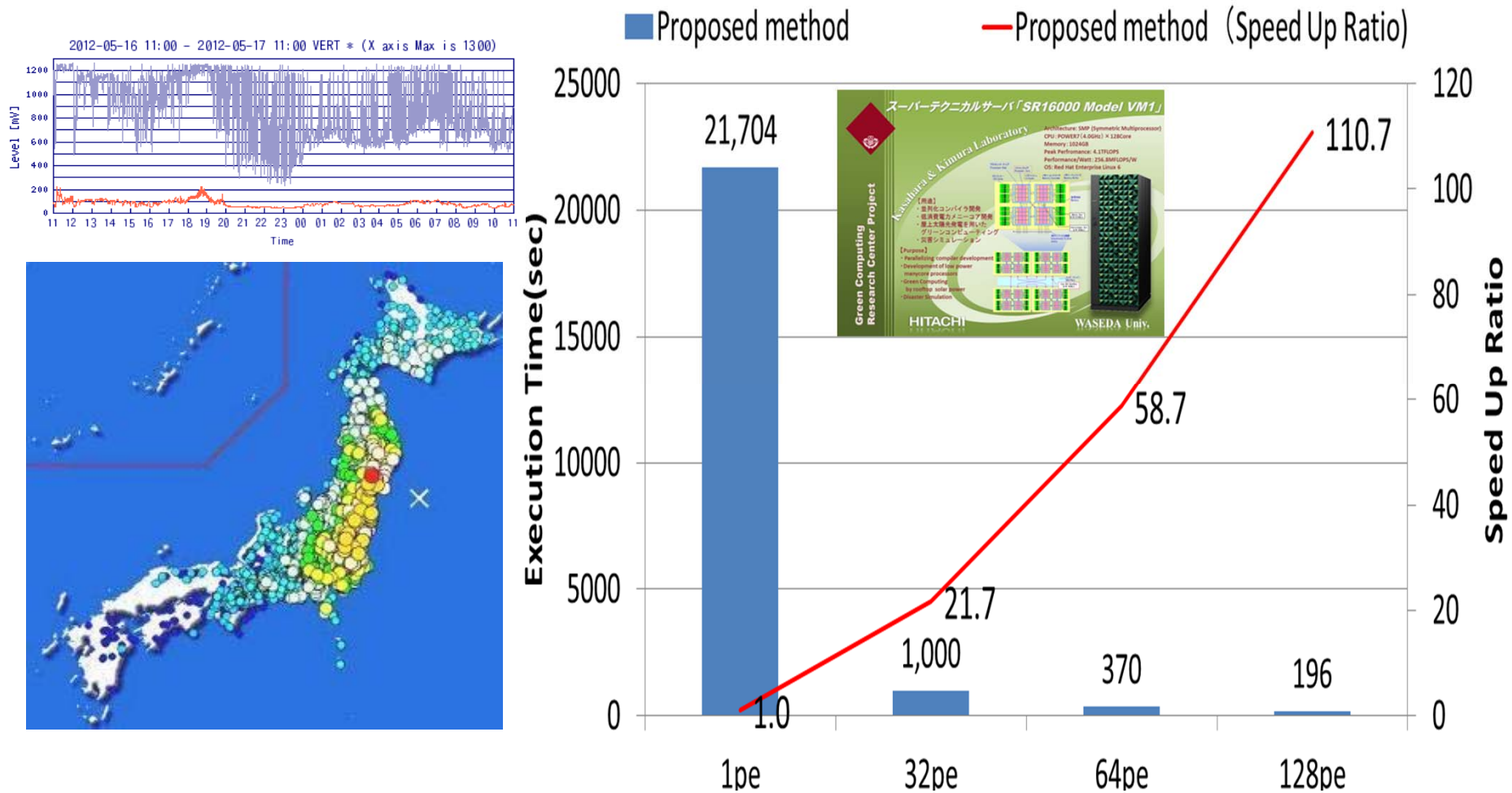## (Intel Xeon, ARM Cortex A15 and Renesas SH4A)



Road Tracking, Image Compression : http://www.mathworks.co.jp/jp/help/vision/examples
Buoy Detection : http://www.mathworks.co.jp/matlabcentral/fileexchange/44706-buoy-detection-using-simulink
Color Edge Detection : http://www.mathworks.co.jp/matlabcentral/fileexchange/28114-fast-edges-of-a-color-image--actual-color--not-converting-to-grayscale-/
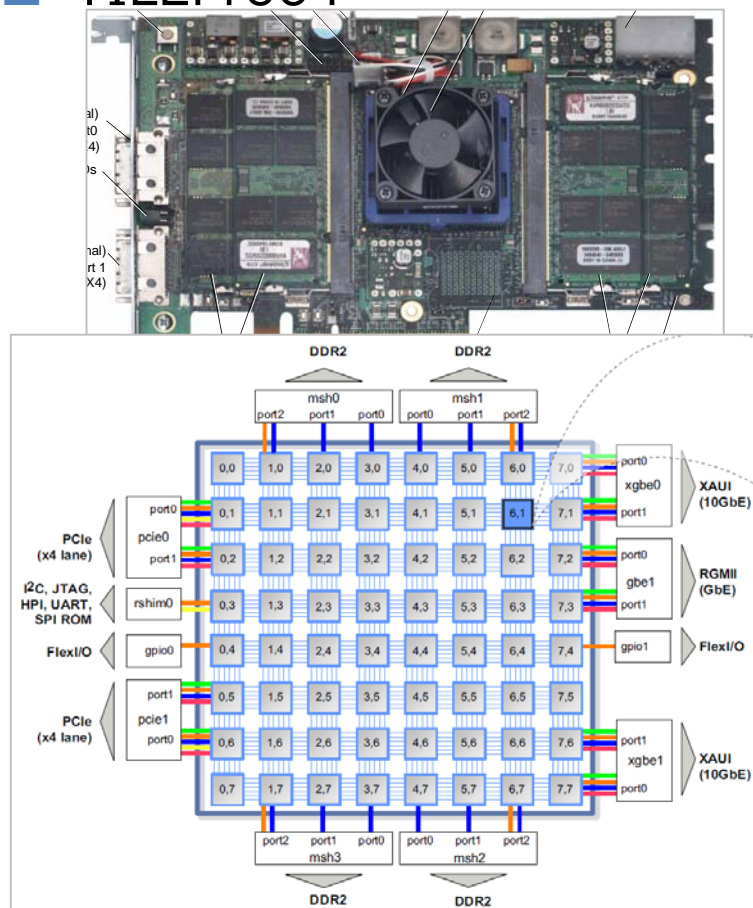Vessel Detection : http://www.mathworks.co.jp/matlabcentral/fileexchange/24990-retinal-blood-vessel-extraction/

# 110 Times Speedup against the Sequential Processing for GMS Earthquake Wave Propagation Simulation on Hitachi SR16000

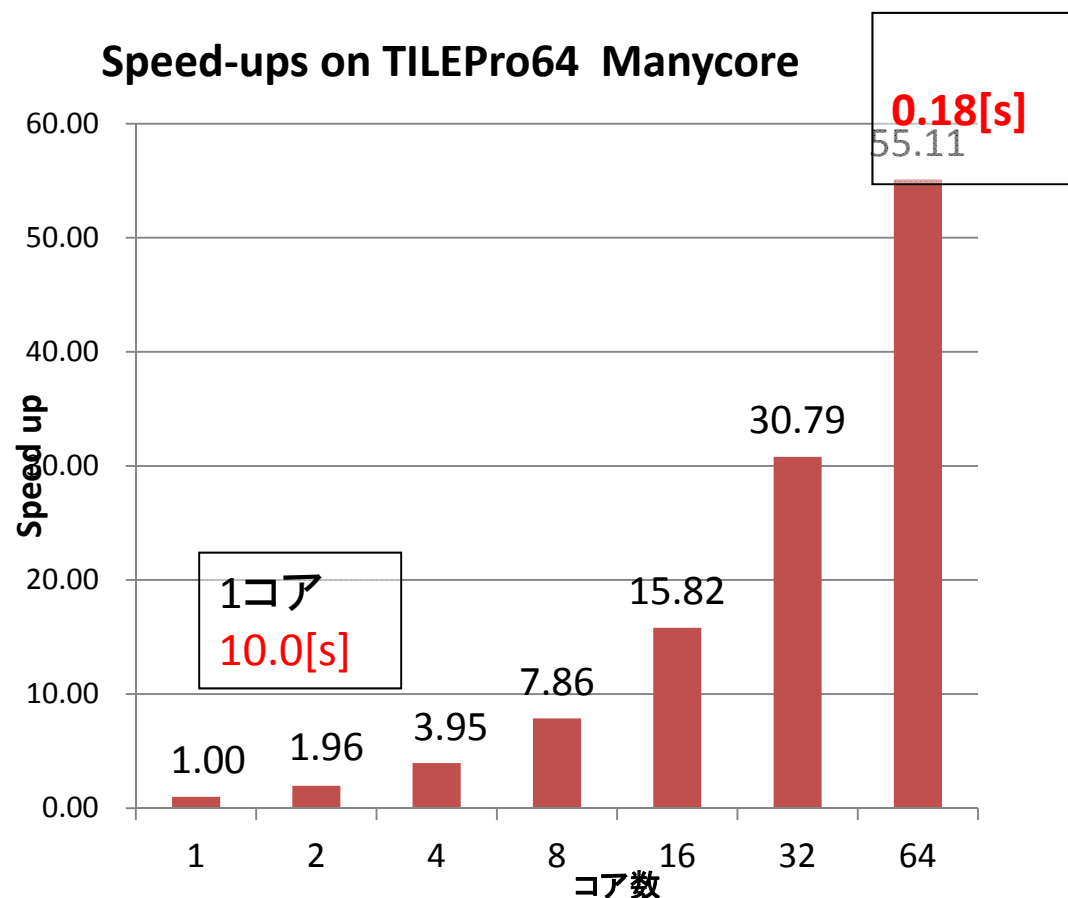## （Power7 Based 128 Core Linux SMP）

# Automatic Parallelization of Still Image Encoding Using JPEG-XR for the Next Generation Cameras and Drinkable Inner Camera
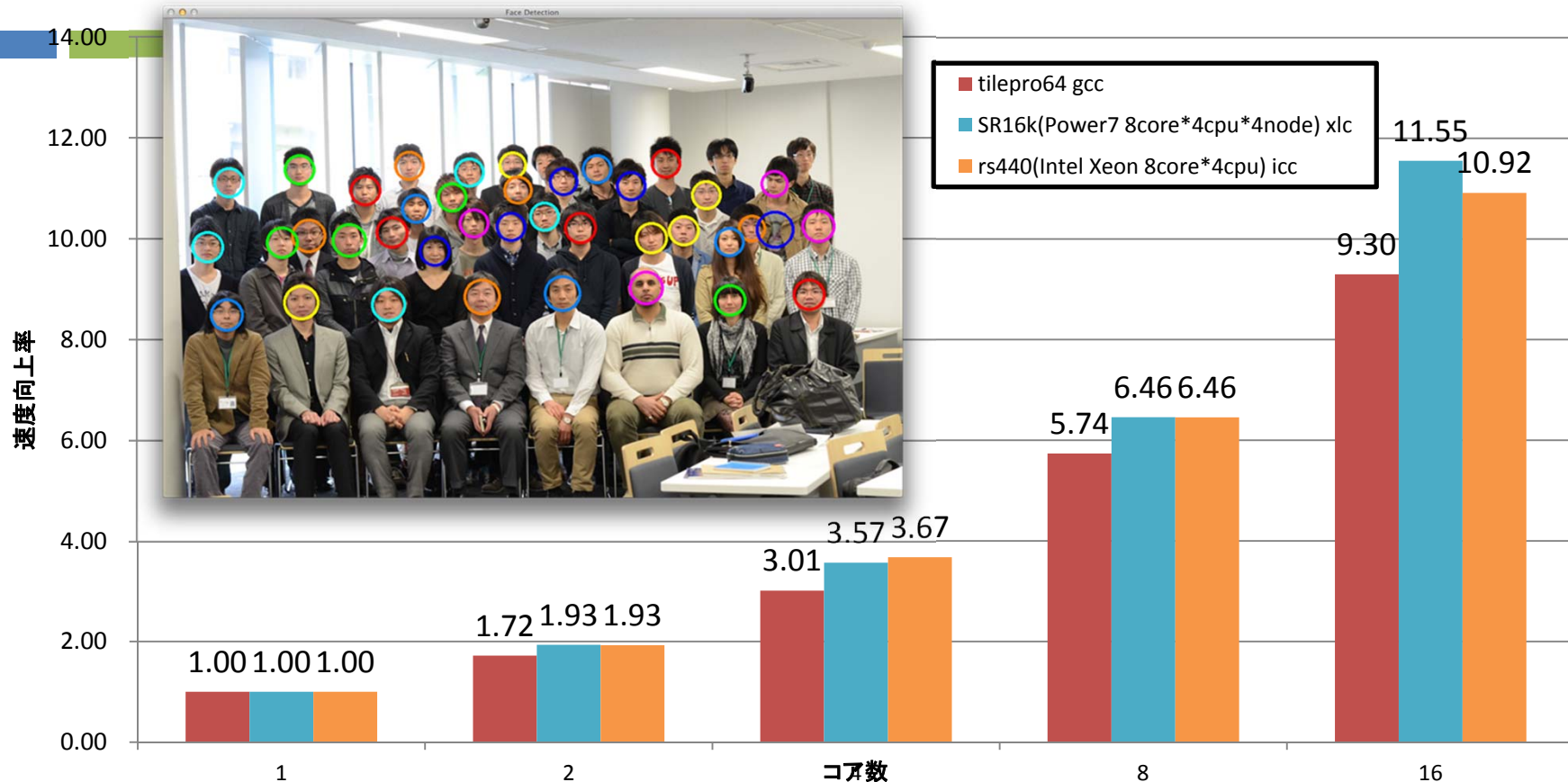
☐ TILEPro64



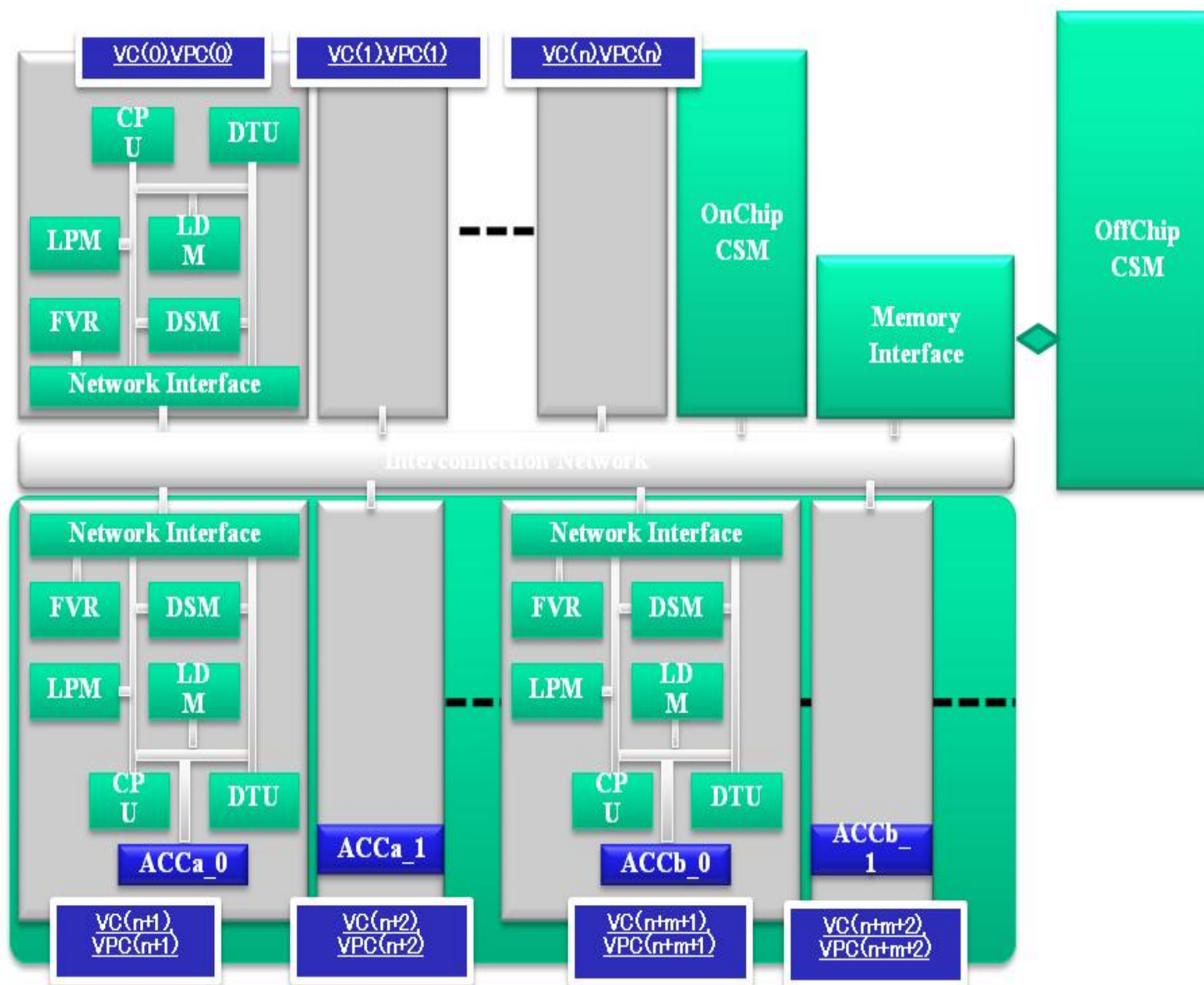**Speed-ups on TILEPro64 Manycore**

0.18[s]

1コア
10.0[s]



55 times speedup with 64 cores against 1 core

# Parallel Processing of Face Detection on Manycore, Highend and PC Server



□ OSCAR compiler gives us 11.55 times speedup for 16 cores against 1 core on SR16000 Power7 highend server.

# OSCAR Heterogeneous Multicore



DTU
- Data Transfer Unit

LPM
- Local Program Memory

LDM
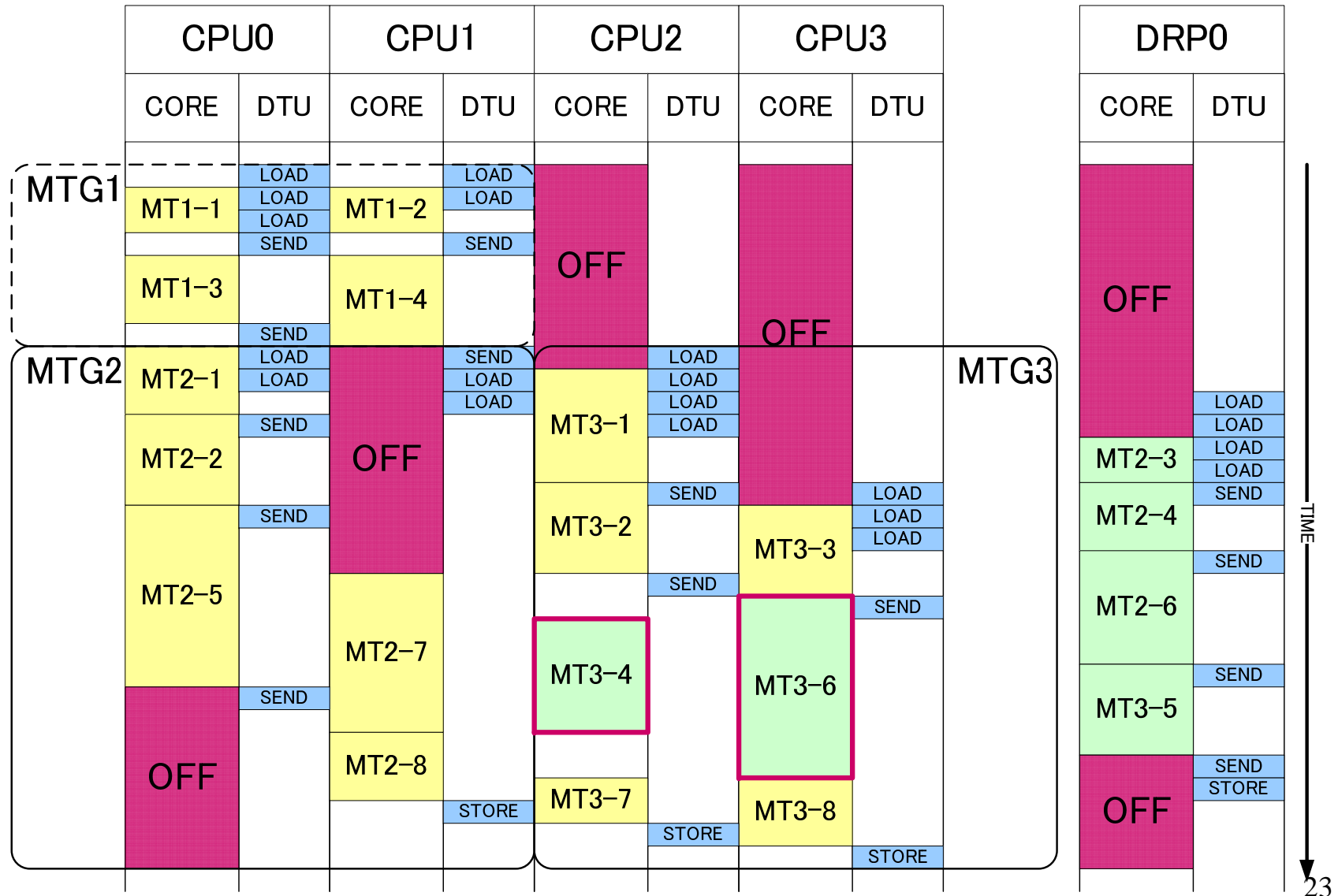- Local Data Memory

DSM
- Distributed Shared Memory
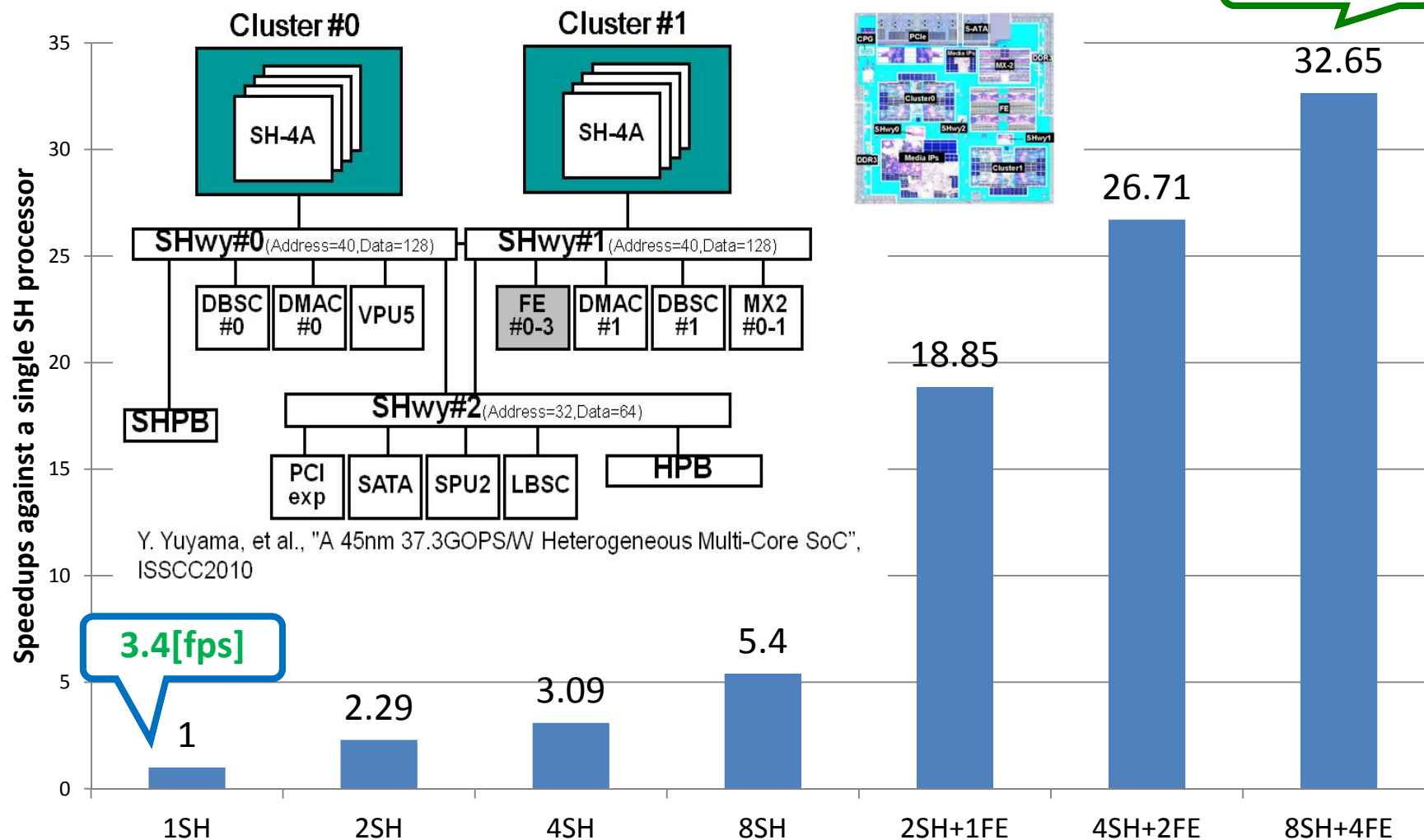
CSM
- Centralized Shared Memory

FVR
- Frequency/Voltage Control Register

22

# An Image of Static Schedule for Heterogeneous Multi-core with Data Transfer Overlapping and Power Control

# 33 Times Speedup Using
# OSCAR Compiler and OSCAR API on RP-X
## (Optical Flow with a hand-tuned library)



Y. Yuyama, et al., "A 45nm 37.3GOPS/W Heterogeneous Multi-Core SoC", ISSCC2010

111[fps]

3.4[fps]

| Configuration | Speedup |
|---|---|
| 1SH | 1 |
| 2SH | 2.29 |
| 4SH | 3.09 |
| 8SH | 5.4 |
| 2SH+1FE | 18.85 |
| 4SH+2FE | 26.71 |
| 8SH+4FE | 32.65 |

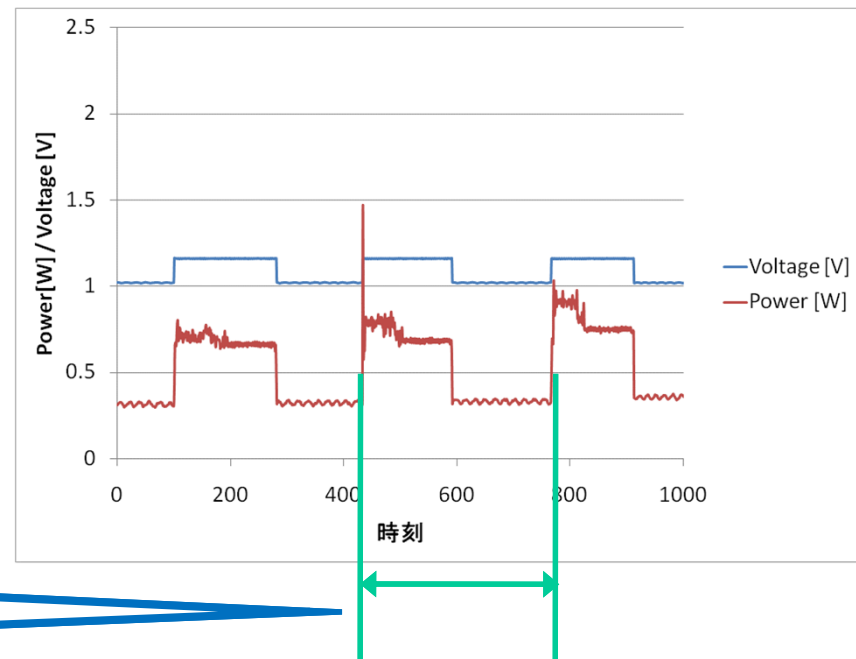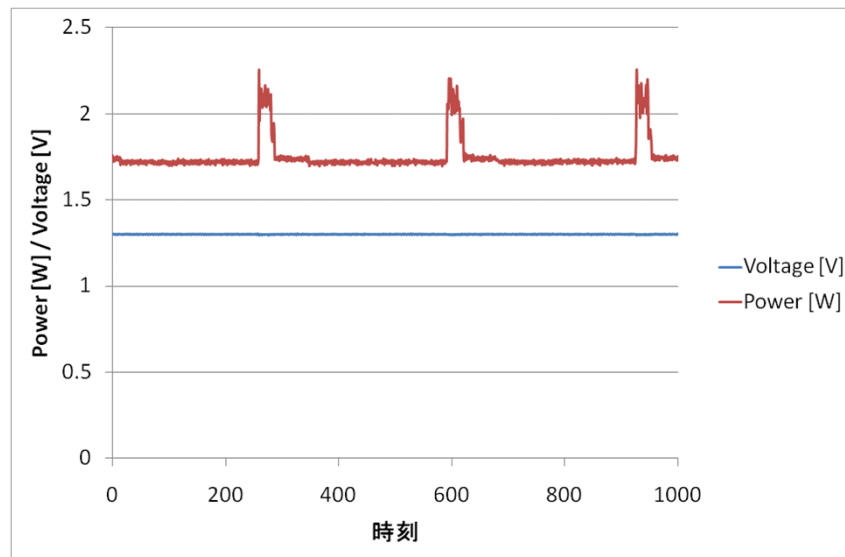Speedups against a single SH processor

# Power Reduction in a real-time execution controlled by OSCAR Compiler and OSCAR API on RP-X (Optical Flow with a hand-tuned library)

**Without Power Reduction**

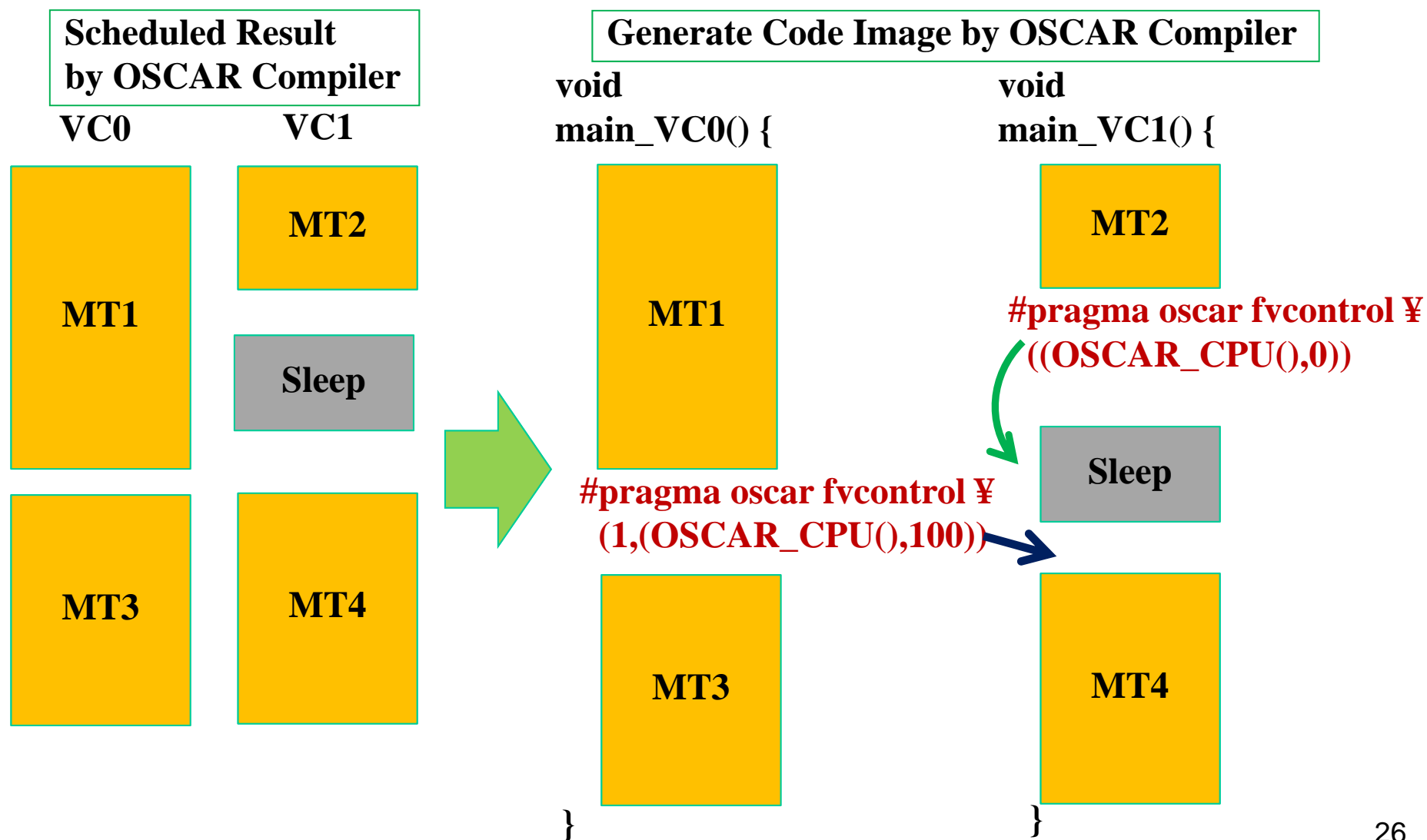**With Power Reduction by OSCAR Compiler**

**70% of power reduction**

**Average:1.76[W]**

**Average:0.54[W]**



**1cycle : 33[ms]
→30[fps]**

# Low-Power Optimization with OSCAR API

**Scheduled Result by OSCAR Compiler**

VC0 | VC1

MT1

MT2

Sleep

MT3 | MT4

**Generate Code Image by OSCAR Compiler**

void
main_VC0() {

MT1

#pragma oscar fvcontrol ¥
(1,(OSCAR_CPU(),100))

MT3

}

void
main_VC1() {

MT2

#pragma oscar fvcontrol ¥
((OSCAR_CPU(),0))

Sleep
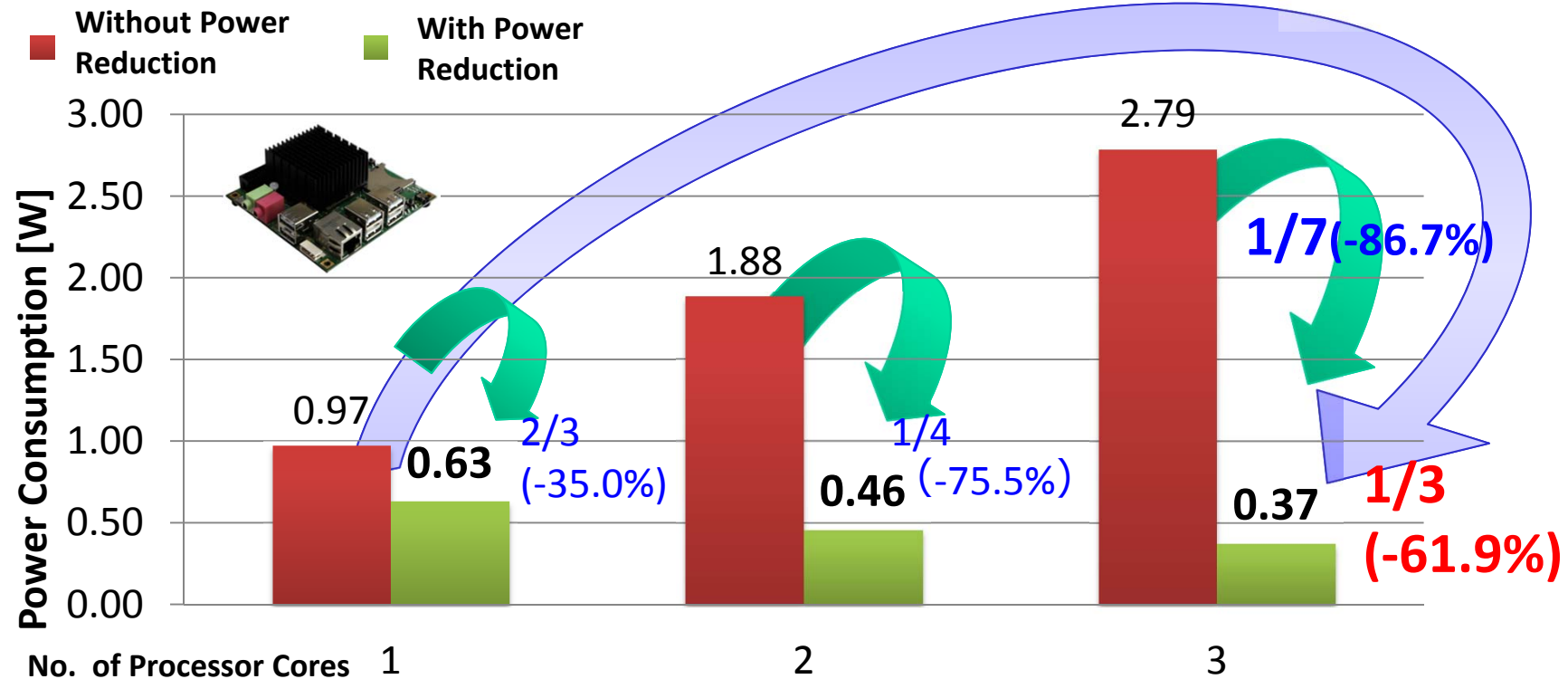
MT4

}

26

# Automatic Power Reduction for
# MPEG2 Decode on Android Multicore

## ODROID X2 ARM Cortex-A9 4 cores

http://www.youtube.com/channel/UCS43lNYEIkC8i_KIgFZYQBQ

**Without Power Reduction** (red)  **With Power Reduction** (green)

Power Consumption [W]

3.00
2.50
2.00
1.50
1.00
0.50
0.00

0.97
**0.63**
2/3 (-35.0%)

1.88
**0.46**
1/4 (-75.5%)

2.79
**0.37**
1/3 (-61.9%)

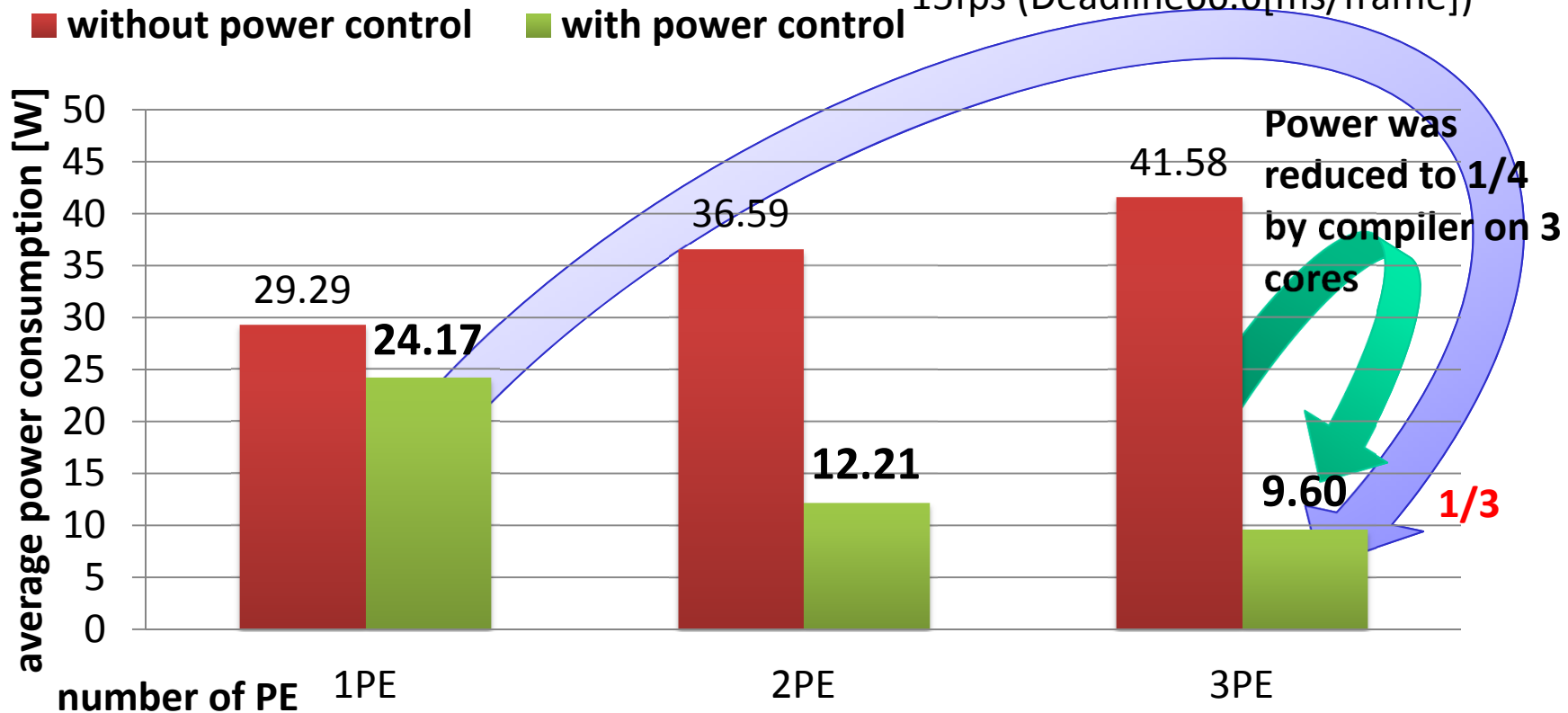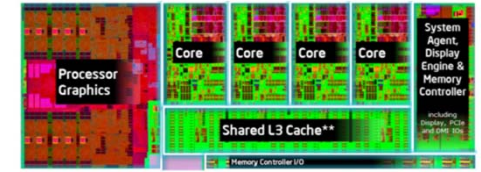**1/7(-86.7%)**

No. of Processor Cores  1  2  3

- **On 3 cores, Automatic Power Reduction control successfully reduced power to 1/7 against without Power Reduction control.**

- **3 cores with the compiler power reduction control reduced power to 1/3 against ordinary 1 core execution.**

27

# Power Reduction on Intel Haswell for Real-time Optical Flow

**Intel CPU Core i7 4770K**

For HD 720p(1280x720) moving pictures
15fps (Deadline66.6[ms/frame])



■ **without power control**  ■ **with power control**

Power was reduced to 1/4 by compiler on 3 cores

1/3

- 1PE: without 29.29, with **24.17**
- 2PE: without 36.59, with **12.21**
- 3PE: without 41.58, with **9.60**

average power consumption [W] — axis 0 to 50

number of PE: 1PE, 2PE, 3PE

Power was reduced to 1/4 (9.6W) by the compiler power optimization on the same 3 cores (41.6W).

Power with 3 core was reduced to 1/3 (9.6W) against 1 core (29.3W).

# Power Reduction of Face Recognition on Intel Haswell 3 cores by OSCAR Compiler - Reduced Power to 2/5 on Intel-

Kasahara & Kimura Lab, Waseda University, TOKYO
http://www.kasahara.cs.waseda.ac.jp

WASEDA University

- OSCAR Compiler
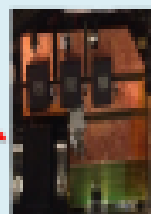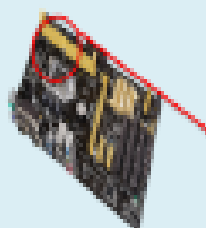- Intel Haswell
- Power Reduction

## Measuring Environment
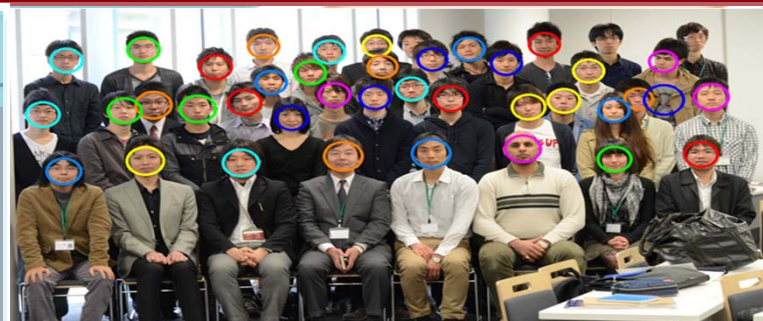
CPU : Intel Core i7 4770K
No. of Cores : 4
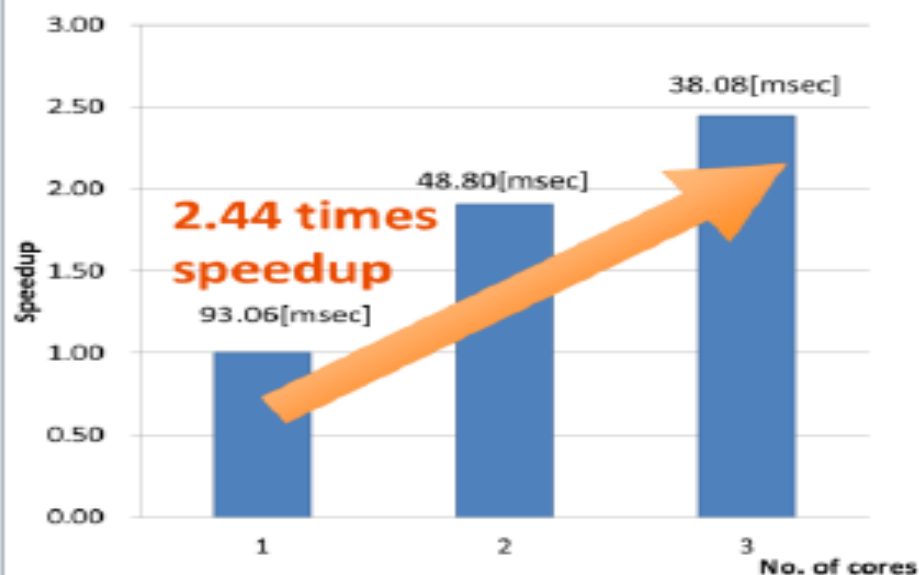Frequency : 3.5GHz～0.8GHz
Motherboard : ASUS H81M-A

Measuring current from CPU power source

## Speedup and Power reduction on Intel Haswell 3 Cores

### Speedup
#### at Fastest Execution Mode

2.44 times speedup

- 93.06[msec]
- 48.80[msec]
- 38.08[msec]

Speedup

No. of cores: 1, 2, 3

### Average Power Consumption
#### at Power Reduction Mode

- ■ Without power control
- ■ With power control

Average Power[W]

- 21.02
- 20.15
- Reduced to 5/7 (-23.22%)
- 41.72
- Reduced to 2/5 (-61.31%)
- 16.14

No. of cores: 1, 3

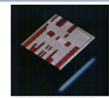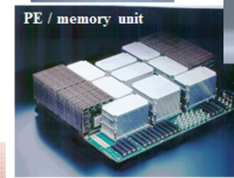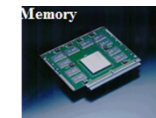# OSCAR Technology

**Started up on Feb.28, 2013:**
**Licensing the all patents and OSCAR compiler from Waseda Univ.**

**CEO:  Dr. T. Ono (Ex- CEO of First Section-listed Company,**
**VP of National Univ., Invited Prof. of Waseda U. )**
**Executives: Mr. T. Ito (Visiting Prof. Tokyo Agricult. and Eng. U.)**
**Prof. K. Shirai (Ex-President of Waseda U**
**Chairman of Japanese Open Univ. )**
**CTO: Mr. M. Takamura (Ex-Fellow Fujitsu Lab.,**
**Fujitsu VPP500, 5000 & NWT Development Leader )**
**Mr. K. Ashida(Ex-VP Sumitomo Trading,**
**Ashida Consult. CEO, A leader of Business World**
**Auditor: Dr. S. Matsuda ( Prof. Emeritus Waseda U.**
**Ex-President Ventures and Entrepreneurs Society )**
**Advisors: Dr. T. Sato ( Patent Attorney, Ex-President of**
**Patent Attorneys Assoc., Gov. IP Committee)**
**Ms. K. Ishiguro ( Lawyer, Supreme Court Trainer)**
**Mr.  A. Fukuda (Leader of Alumni Assoc.)**
**Prof. K. Kimura (Waseda Univ. )**
**Prof. H. Kasahara ( Waseda Univ. )**

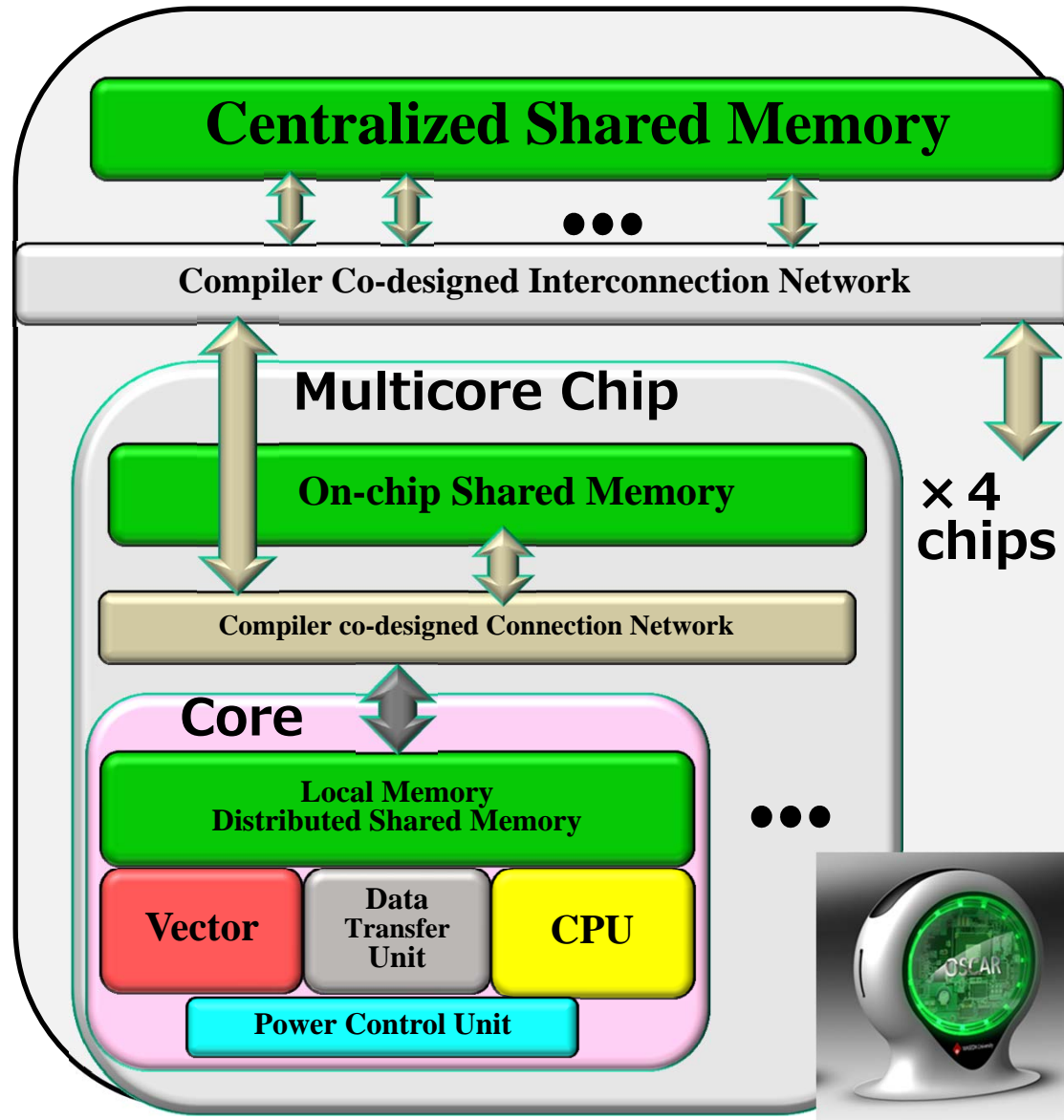**Fujitsu VPP5000**

Memory    Cabinet

PE / memory unit

CMOS LSI
Copyright 2008 FUJITSU LIMITED    51

OSCAR TECHNOLOGY CORPORATION

# OSCAR Vector Multicore and Compiler for Embedded to Severs with OSCAR Technology

**Centralized Shared Memory**

Compiler Co-designed Interconnection Network

**Multicore Chip**

**On-chip Shared Memory**

× 4 chips

Compiler co-designed Connection Network

**Core**

**Local Memory Distributed Shared Memory**

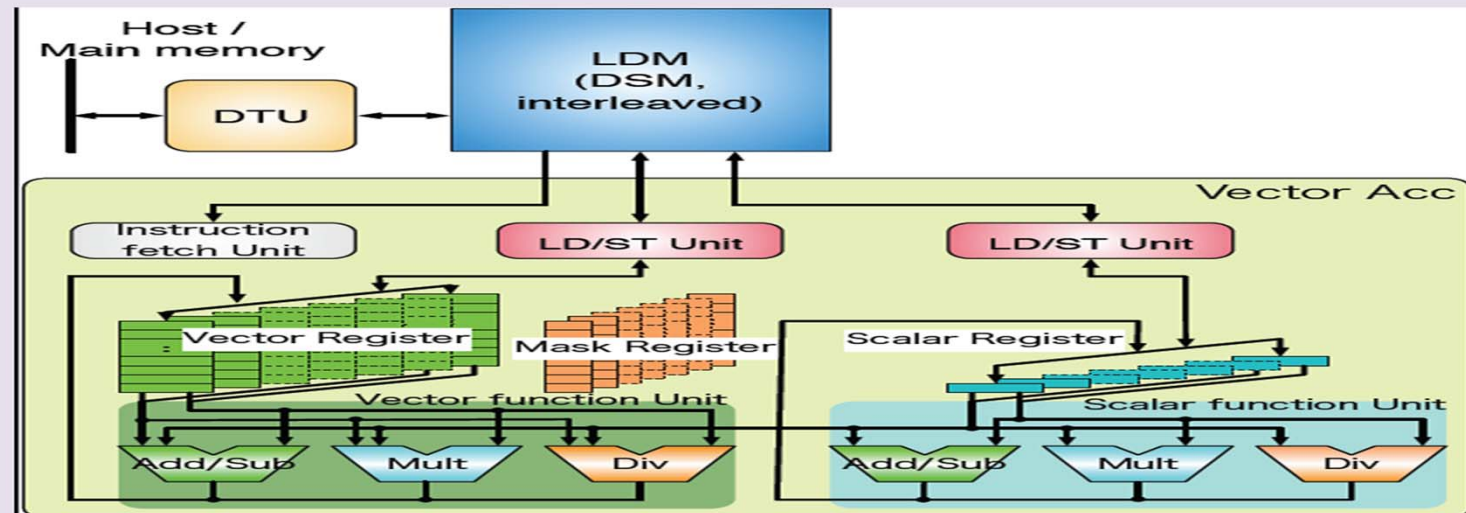| Vector | Data Transfer Unit | CPU |
|--------|--------------------|-----|

**Power Control Unit**

**Target:**

➢ Solar Powered with compiler power reduction.

➢ Fully automatic parallelization and vectorization including local memory management and data transfer.

# Vector Accelerator

## Features
- **Attachable for any CPUs (Intel, ARM, IBM)**
- **Data driven initiation by sync flags**



## Function Units [tentative]
- **Vector Function Unit**
  - **8 double precision ops/clock**
  - **64 characters ops/clock**
  - **Variable vector register length**
  - **Chaining LD/ST & Vector pipes**
- **Scalar Function Unit**

## Registers[tentative]
- **Vector Register  256Bytes/entry, 32entry**
- **Scalar Register   8Bytes/entry**
- **Floating Point Register  8Bytes/entry**
- **Mask Register  32Bytes/entry**

# Summary

➢ **Waseda University Green Computing Systems R&D Center supported by METI has been researching on low-power high performance Green Multicore hardware, software and application with government and industry including Hitachi, Fujitsu, NEC, Renesas, Denso, Toyota, Olympus and OSCAR Technology.**

➢ **OSCAR Automatic Parallelizing and Power Reducing Compiler has succeeded speedup and/or power reduction of scientific applications including "Earthquake Wave Propagation", medical applications including "Cancer Treatment Using Carbon Ion", and "Drinkable Inner Camera", industry application including "Automobile Engine Control", "Smartphone", and "Wireless communication Base Band Processing" on various multicores from different vendors including Intel, ARM, IBM, AMD, Qualcomm, Freescale, Renesas and Fujitsu.**

➢ **In automatic parallelization, 110 times speedup for "Earthquake Wave Propagation Simulation" on 128 cores of IBM Power 7 against 1 core, 55 times speedup for "Carbon Ion Radiotherapy Cancer Treatment" on 64cores IBM Power7, 1.95 times for "Automobile Engine Control" on Renesas 2 cores using SH4A or V850, 55 times for "JPEG-XR Encoding for Capsule Inner Cameras" on Tilera 64 cores Tile64 manycore.**

➢ **The compiler will be available on market from OSCAR Technology.**

➢ **In automatic power reduction, consumed powers for real-time multi-media applications like Human face detection, H.264, mpeg2 and optical flow were reduced to 1/2 or 1/3 using 3 cores of ARM Cortex A9 and Intel Haswell and 1/4 using Renesas SH4A 8 cores against ordinary single core execution.**

# Fujitsu VPP500/NWT: PE Unit