# OSCAR Automatic Parallelization and Power Reduction Compiler for Homogeneous and Heterogeneous Multicores

## Hironori Kasahara

**Professor, Dept. of Computer Science & Engineering**
**Director, Advanced Multicore Processor Research Institute**
**Waseda University, Tokyo, Japan**
**IEEE Computer Society Multicore STC Chair**
**IEEE Computer Society 2015 Election President Candidate**
URL: http://www.kasahara.cs.waseda.ac.jp/

# Green Computing Systems R&D Center
## Waseda University
### Supported by METI (Mar. 2011 Completion)

<R & D Target>

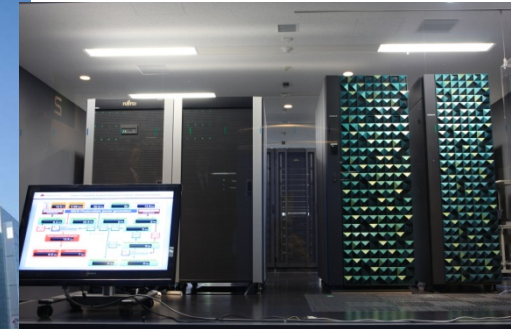**Hardware, Software, Application for Super Low-Power Manycore Processors**

➢**More than 64 cores**

➢**Natural air cooling (No fan)**
 **Cool, Compact, Clear, Quiet**

➢**Operational by Solar Panel**

<Industry, Government, Academia>

**Hitachi, Fujitsu, NEC, Renesas, Olympus, Toyota, Denso, Mitsubishi, OSCAR Tech.**

<Ripple Effect>

➢**Low $CO_2$ (Carbon Dioxide) Emissions**

➢**Creation Value Added Products**

  ➢**Consumer Electronics, Automobiles, Servers**

**Hitachi SR16000:**
  **Power7 128coreSMP**
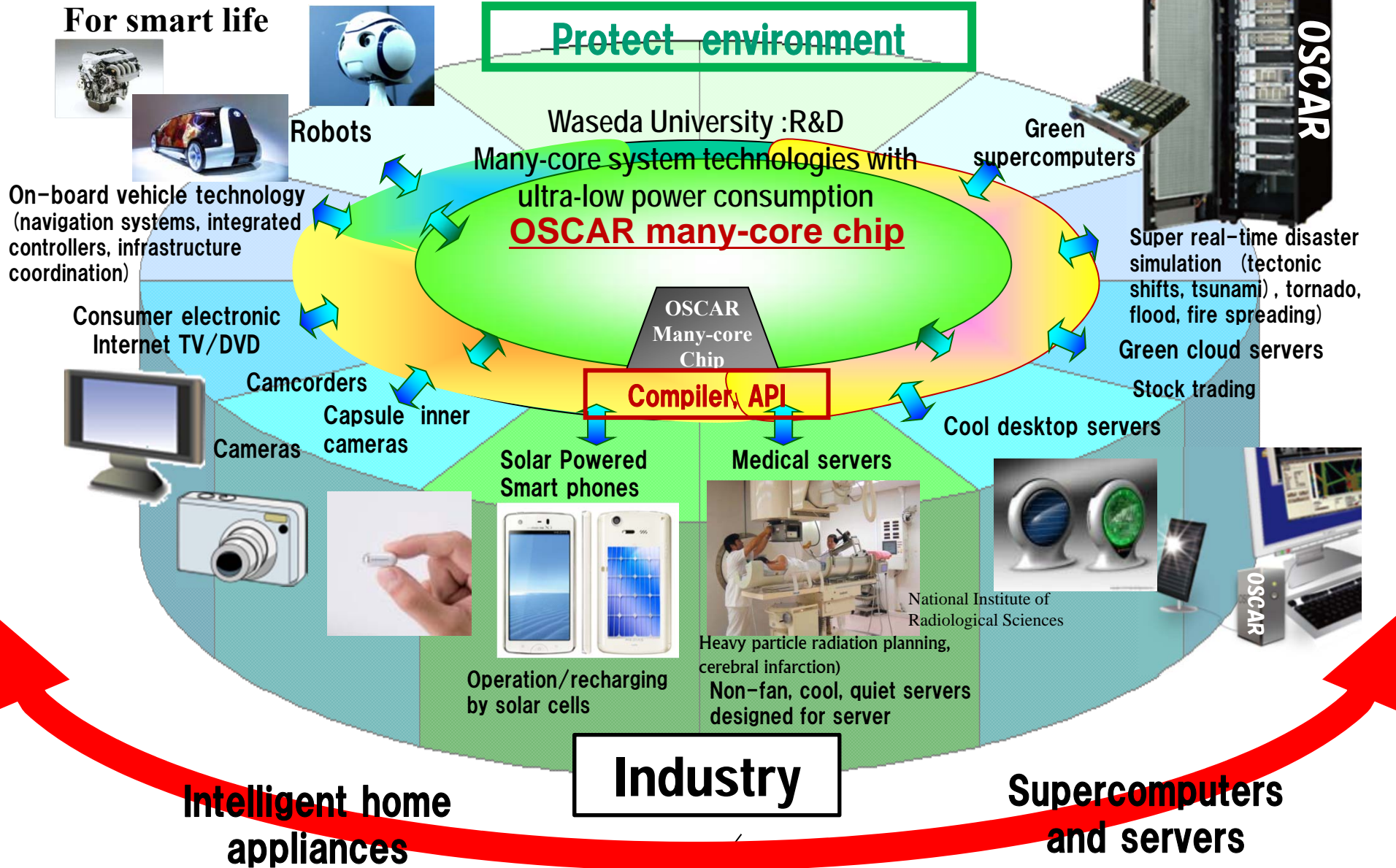**Fujitsu M9000**
  **SPARC VII 256 core SMP**

**Beside Subway Waseda Station, Near Waseda Univ. Main Campus**

**Industry-government-academia collaboration in R&D and target practical applications**

**Protect lives**

For smart life

**Protect environment**

OSCAR

Robots

Waseda University :R&D
Many-core system technologies with
ultra-low power consumption
**OSCAR many-core chip**

Green supercomputers

On-board vehicle technology
（navigation systems, integrated
controllers, infrastructure
coordination）

Super real-time disaster
simulation （tectonic
shifts, tsunami）, tornado,
flood, fire spreading）

Consumer electronic
Internet TV/DVD

OSCAR
Many-core
Chip

Green cloud servers

Stock trading

Camcorders

Capsule inner
cameras

**Compiler API**

Cool desktop servers

Cameras

Solar Powered
Smart phones

Medical servers

National Institute of
Radiological Sciences

OSCAR

Operation/recharging
by solar cells

Heavy particle radiation planning,
cerebral infarction）
**Non-fan, cool, quiet servers
designed for server**

**Industry**

**Intelligent home
appliances**

**Supercomputers
and servers**

# 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
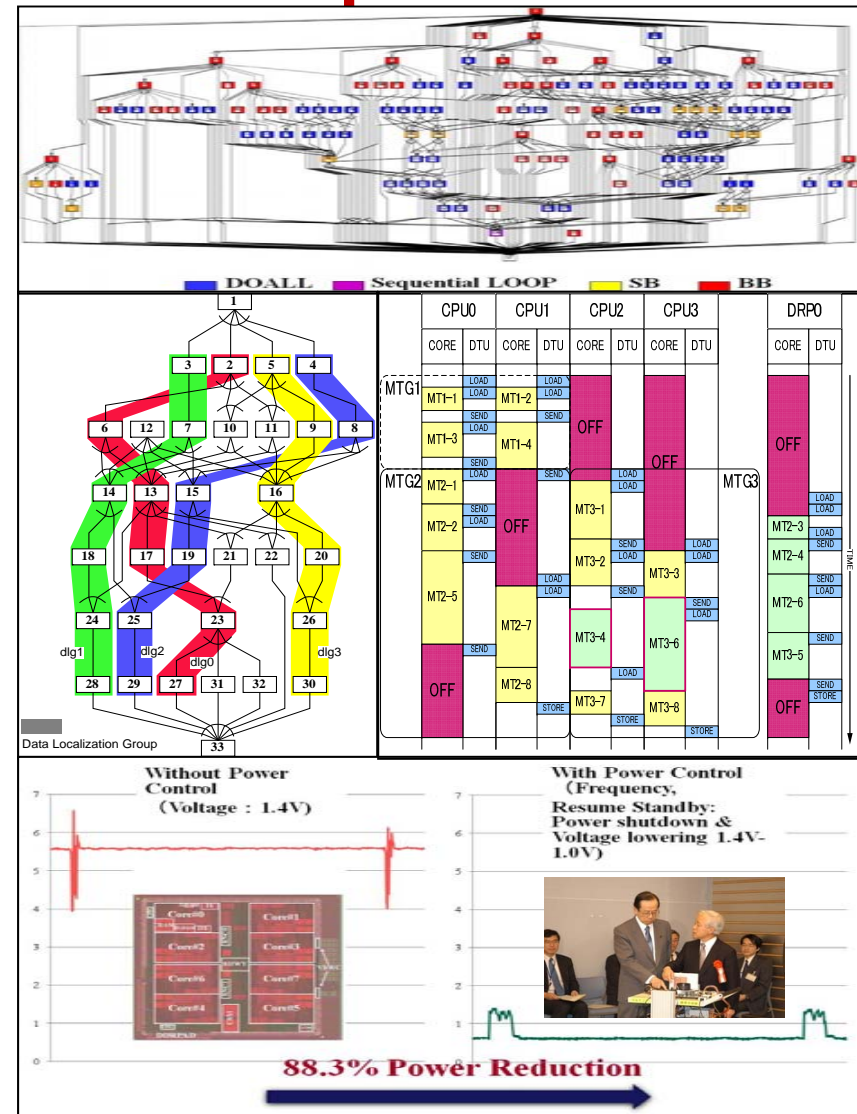
### Data Localization

Automatic data management for distributed shared memory, cache and local memory
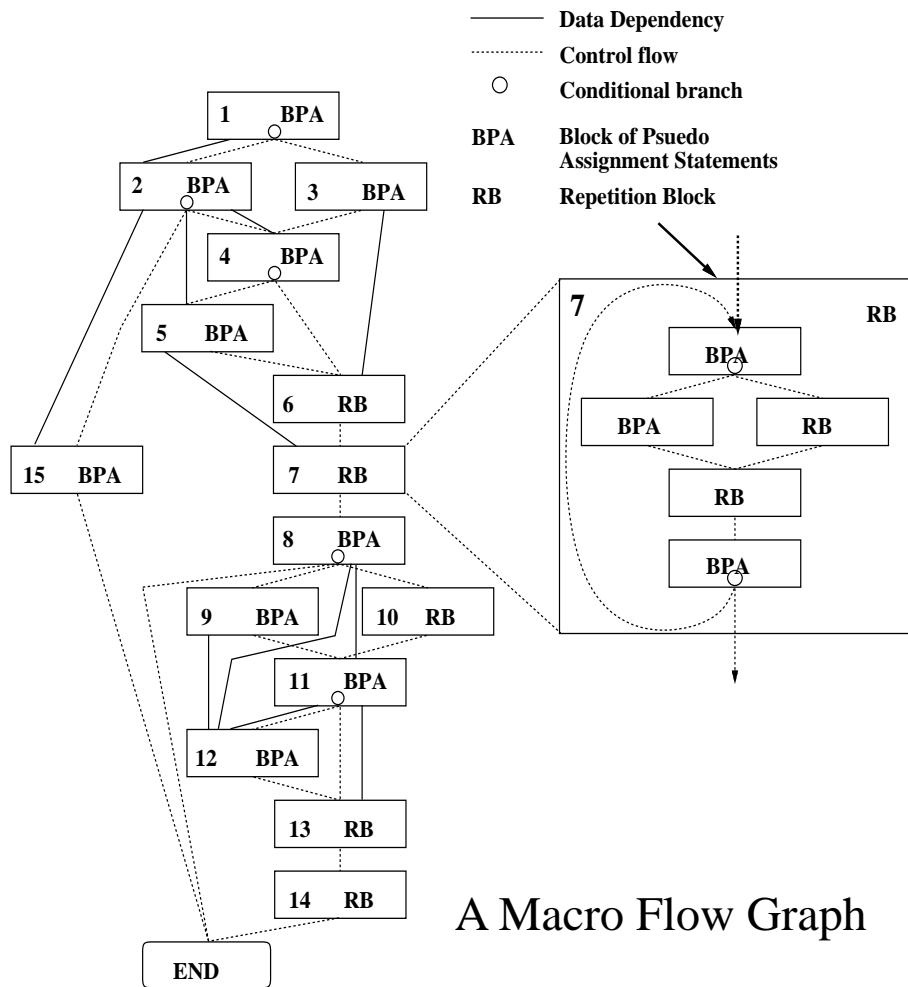
### Data Transfer Overlapping

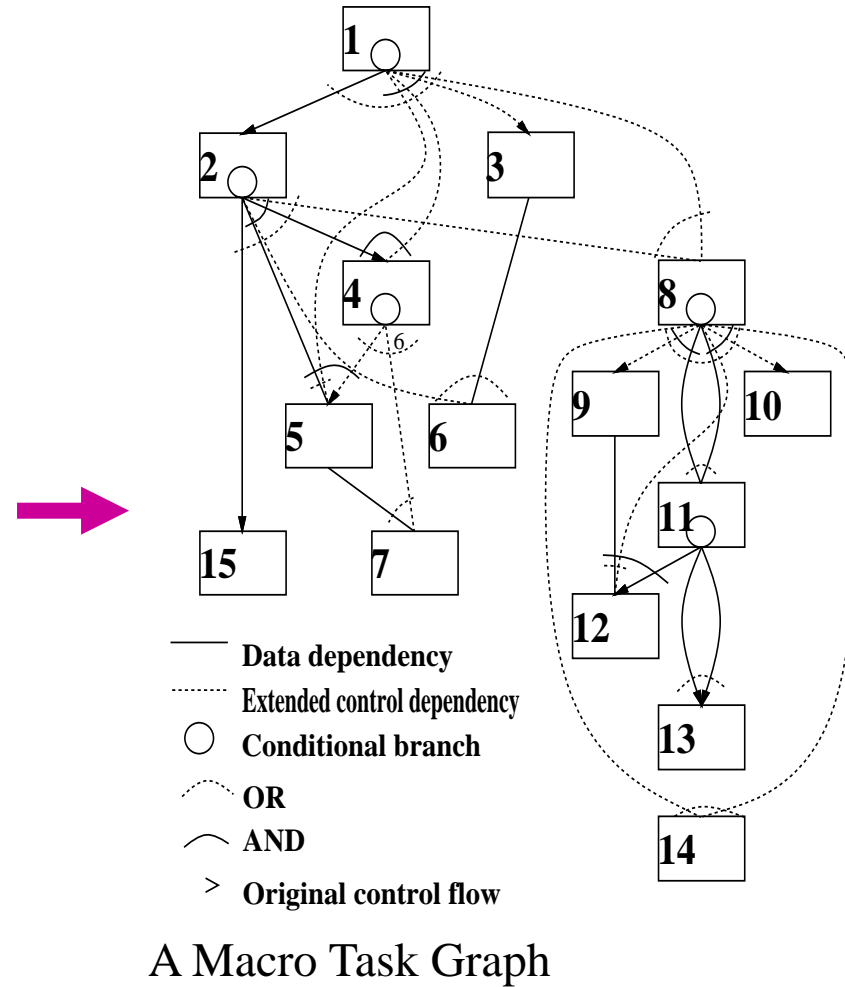Data transfer overlapping using Data Transfer Controllers (DMAs)

### Power Reduction

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



88.3% Power Reduction

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



Data Dependency
Control flow
○ Conditional branch
BPA — Block of Psuedo Assignment Statements
RB — Repetition Block

A Macro Flow Graph

Data dependency
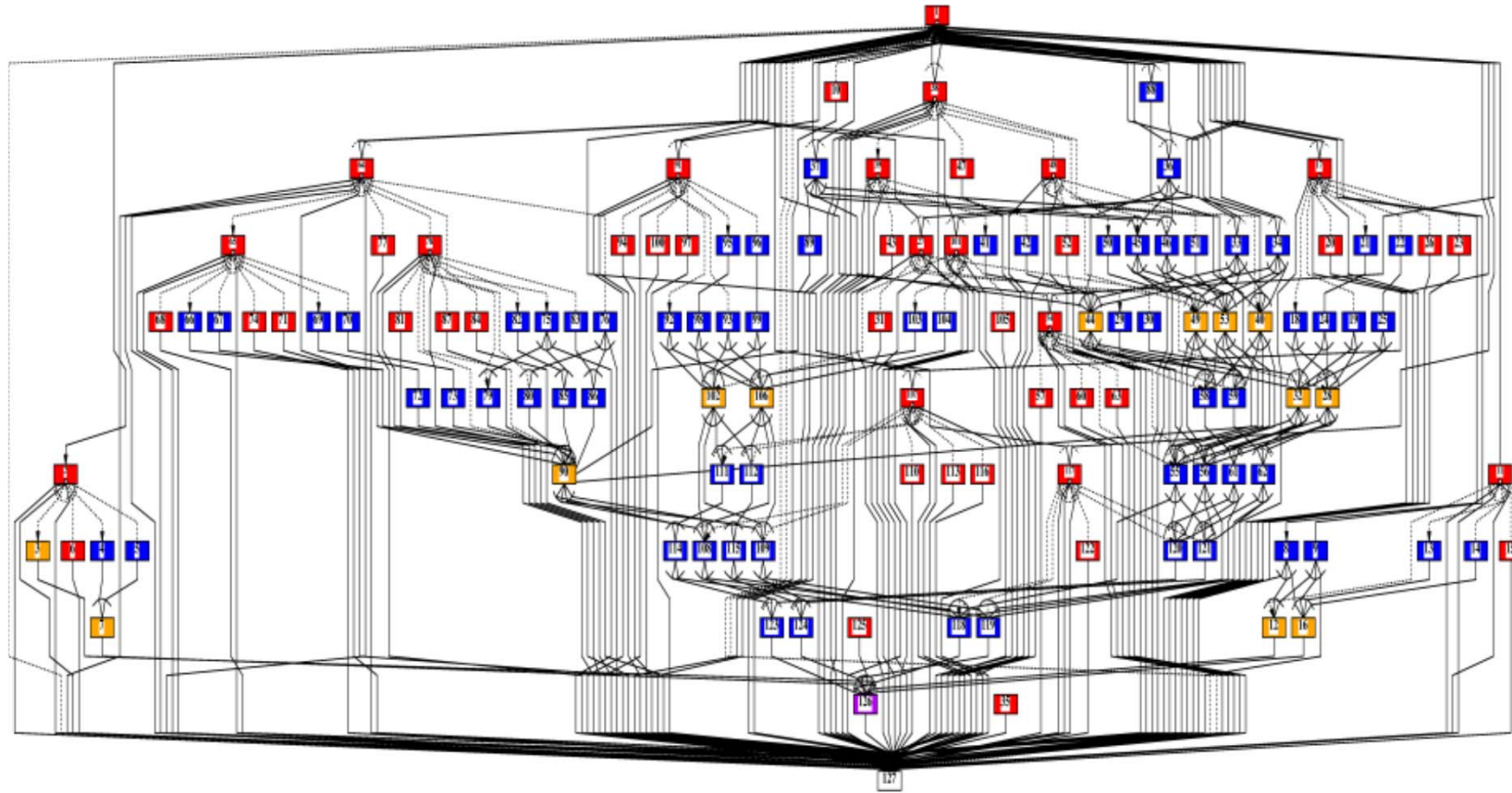Extended control dependency
○ Conditional branch
OR
AND
> Original control flow

A Macro Task Graph

# MTG of Su2cor-LOOPS-DO400

■ Coarse grain parallelism PARA_ALD = 4.3



| | DOALL | | Sequential LOOP | | SB | | BB |

# Data-Localization: Loop Aligned Decomposition

- **Decompose multiple loop (Doall and Seq) into CARs and LRs considering inter-loop data dependence.**

  - **Most data in LR can be passed through LM.**

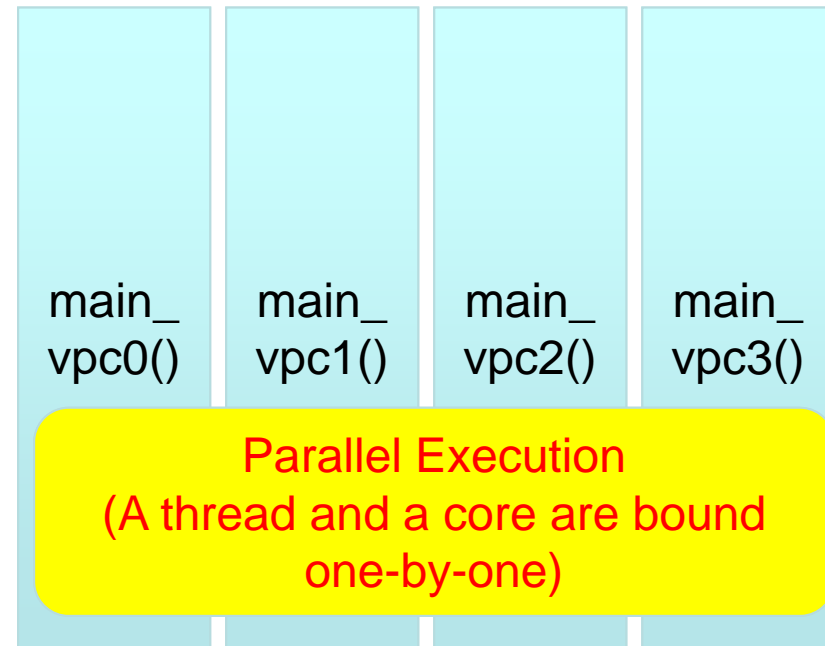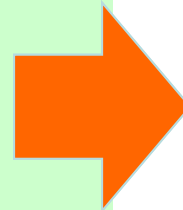  - **LR: Localizable Region, CAR: Commonly Accessed Region**



| | | | | |
|---|---|---|---|---|
| C RB1(Doall)<br>DO I=1,101<br>A(I)=2*I<br>ENDDO | | | | |
| C RB2(Doseq)<br>DO I=1,100<br>B(I)=B(I-1)<br>+A(I)+A(I+1)<br>ENDDO | | | | |
| RB3(Doall)<br>DO I=2,100<br>C(I)=B(I)+B(I-1)<br>ENDDO<br>C | | | | |

| LR | CAR | LR | CAR | LR |
|---|---|---|---|---|
| DO I=1,33 | DO I=34,35 | DO I=36,66 | DO I=67,68 | DO I=69,101 |
| DO I=1,33 | DO I=34,34 | DO I=35,66 | DO I=67,67 | DO I=68,100 |
| DO I=2,34 | | DO I=35,67 | | DO I=68,100 |

7

# 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**

# Thread Execution Model

```
#pragma omp parallel sections
    {
#pragma omp section
        main_vpc0();
#pragma omp section
        main_vpc1();
#pragma omp section
        main_vpc2();
#pragma omp section
        main_vpc3();
    }
```

| main_<br>vpc0() | main_<br>vpc1() | main_<br>vpc2() | main_<br>vpc3() |
|---|---|---|---|

**Parallel Execution**
**(A thread and a core are bound**
**one-by-one)**

VPC: Virtual Processor Core

# **Memory Mapping**

- **Placing variables on an onchip centralized shared memory (onchipCSM)**
  - #pragma oscar onchipshared (C)
  - !$oscar onchipshared (Fortran)
- **Placing variables on a local data memory (LDM)**
  - #pragma omp threadprivate (C)
  - !$omp threadprivate (Fortran)
  - This directive is an extension to OpenMP
- **Placing variables on a distributed shared memory (DSM)**
  - #pragma oscar distributedshared (C)
  - !$oscar distributedshared (Fortran)

# Data Transfer

- **Specifying data transfer lists**
  - #pragma oscar dma_transfer (C)
  - !$oscar dma_transfer (Fortran)
  - Containing following parameter directives
- **Specifying a contiguous data transfer**
  - #pragma oscar dma_contiguous_parameter (C)
  - !$oscar dma_contiguous_parameter (Fortran)
- **Specifying a stride data transfer**
  - #pragma oscar dma_stride_parameter
  - !$oscar dma_stride_parameter
  - This can be used for scatter/gather data transfer
- **Data transfer synchronization**
  - #pragma oscsar dma_flag_check
  - !$oscar dma_flag_check

# Hierarchical Barrier Synchronization

- **Specifying a hierarchical group barrier**
  - #pragma oscar group_barrier  (C)
  - !$oscar group_barrier (Fortran)

| 1st layer → | 8cores | | | | | |
|---|---|---|---|---|---|---|
| 2nd layer → | PG0(4cores) | | | | PG1(4cores) | |
| 3rd layer → | PG0-0 | PG0-1 | PG0-2 | PG0-3 | PG1-0(2cores) | PG1-1(2cores) |

# 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)**

# 110 Times Speedup against the Sequential Processing for GMS Earthquake Wave Propagation Simulation on Hitachi SR16000
## （Power7 Based 128 Core Linux SMP）

# 9.6 Times Speedup on 12 cores Power 8 against the Sequential Processing for GMS Earthquake Wave Propagation Simulation

- IBM S812L
  - CPU:POWER8
    - 12 cores
    - Clock Frequency 3.026GHz
  - Memory:60GB
  - OS:Redhat Linux 7.1
  - Backend Fortran compiler: IBM xlf 15.1.1
  - Evaluation using medium size input data(12GB)

**S812L(POWER8)**

# Performance of OSCAR Compiler on IBM p6 595 Power6 (4.2GHz) based 32-core SMP Server



**OpenMP codes generated by OSCAR compiler accelerate IBM XL Fortran for AIX Ver.12.1 about 3.3 times on the average**

Compile Option:
(*1) Sequential: -O3 –qarch=pwr6, XLF: -O3 –qarch=pwr6 –qsmp=auto, OSCAR: -O3 –qarch=pwr6 –qsmp=noauto
(*2) Sequential: -O5 -q64 –qarch=pwr6, XLF: -O5 –q64 –qarch=pwr6 –qsmp=auto, OSCAR: -O5 –q64 –qarch=pwr6 –qsmp=noauto
(Others) Sequential: -O5 –qarch=pwr6, XLF: -O5 –qarch=pwr6 –qsmp=auto, OSCAR: -O5 –qarch=pwr6 –qsmp=noauto

# Performance of OSCAR Compiler for Fortran Programs on a IBM p550q 8core Deskside Server

- **2.7 times speedup against loop parallelizing compiler on 8 cores**



■ Loop parallelization
■ Multigrain parallelizition



speedup ratio

tomcatv swim su2cor hydro2d mgrid applu turb3d apsi fpppp wave5 swim mgrid applu apsi

spec95                                    spec2000

# Performance for C programs on IBM p5 550Q



**5.8 times speedup against one processor on average**

# Performance of OSCAR Compiler on IBM pSeries690 (Power 4) RegattaH 16 Processors High-end Server



IBM XL Fortran for AIX Version 8.1

**3.5 times speedup in average**

Legend:
- XL Fortran(max)
- APC(max)



Chart data — Speedup ratio (y-axis, 0.0 to 12.0):

| Benchmark | APC label | XL label | bottom time |
|---|---|---|---|
| tomcatv | 3.8s | 19.2s | 23.1s |
| swim | 3.4s | 16.7s | 38.3s |
| su2cor | 7.1s | 21.5s | 21.5s |
| hydro2d | 3.0s | 21.0s | 30.3s |
| mgrid | 3.1s | 16.4s | 35.1s |
| applu | 13.0s | 23.2s | 27.8s |
| turb3d | 3.5s | 37.4s | 39.5s |
| apsi | 22.5s | | 22.5s |
| fpppp | 85.8s | | 85.8s |
| wave5 | 18.8s | | 18.8s |
| wupwise | 28.9s | 126.5s | 126.5s |
| swim2k | 38.5s | 115.2s | 307.6s |
| mgrid2k | 28.8s | 107.4s | 291.2s |
| applu2k | 105.0s | 184.8s | 279.1s |
| sixtrack | 282.4s | | 282.4s |
| apsi2k | 321.4s | | 321.4s |

# 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.

# Speedup with 2cores for Engine Crankshaft Handwritten Program on Renesas RPX Multi-core Processor



**Macrotask graph with a lot of conditional branches**

**Macrotask graph after task fusion**

**1.6 times Speed up by 2 cores against 1core**

Branches are fused to macrotasks for static scheduling

Grain is too fine (us) for dynamic scheduling.

1core   2core

# OSCAR Compile Flow for Simulink Applications



**Simulink model**

**Generate C code using Embedded Coder**

**C code**

**OSCAR Compiler**

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

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

(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
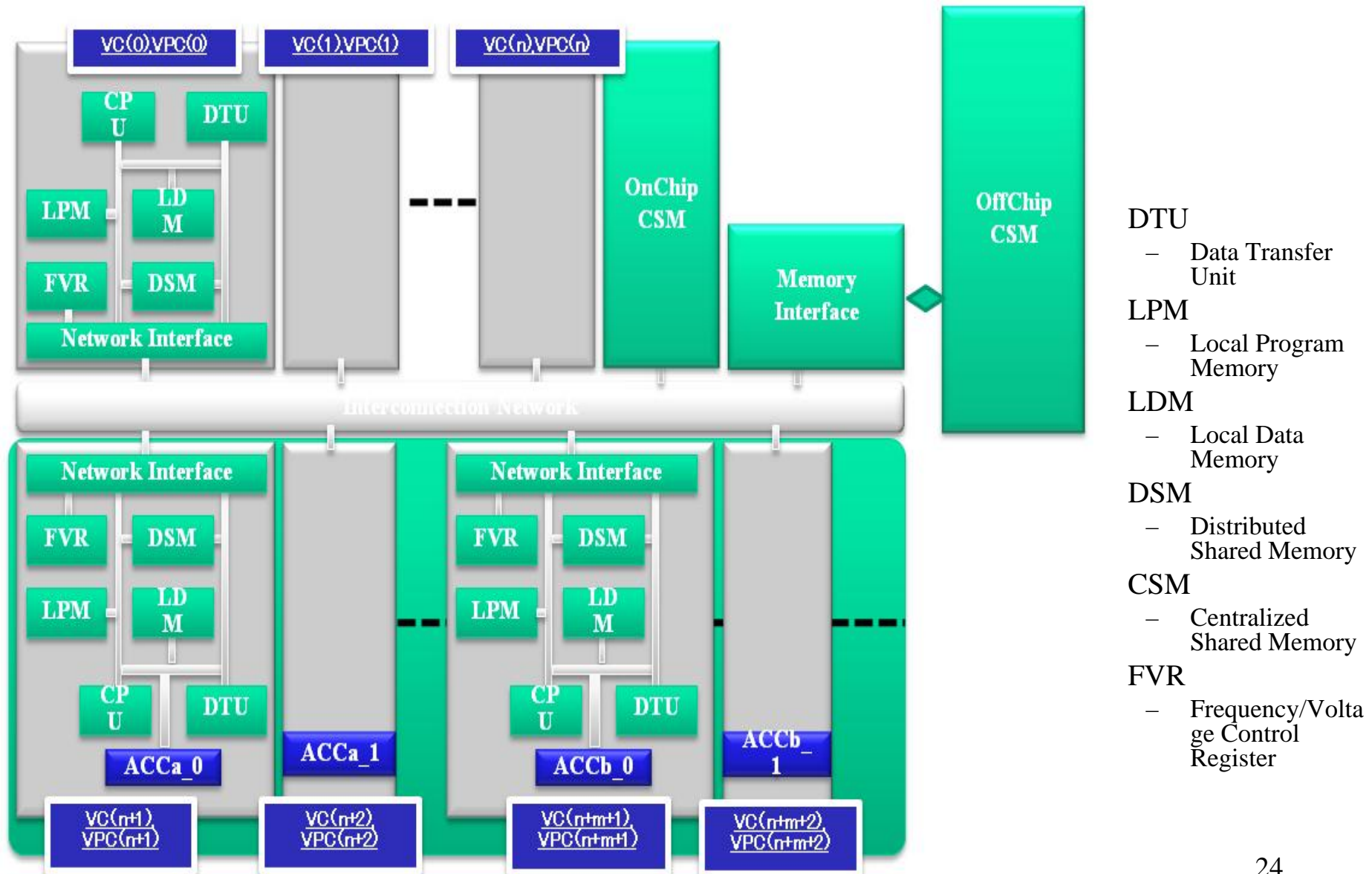## (Intel Xeon, ARM Cortex A15 and Renesas SH4A)



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-/
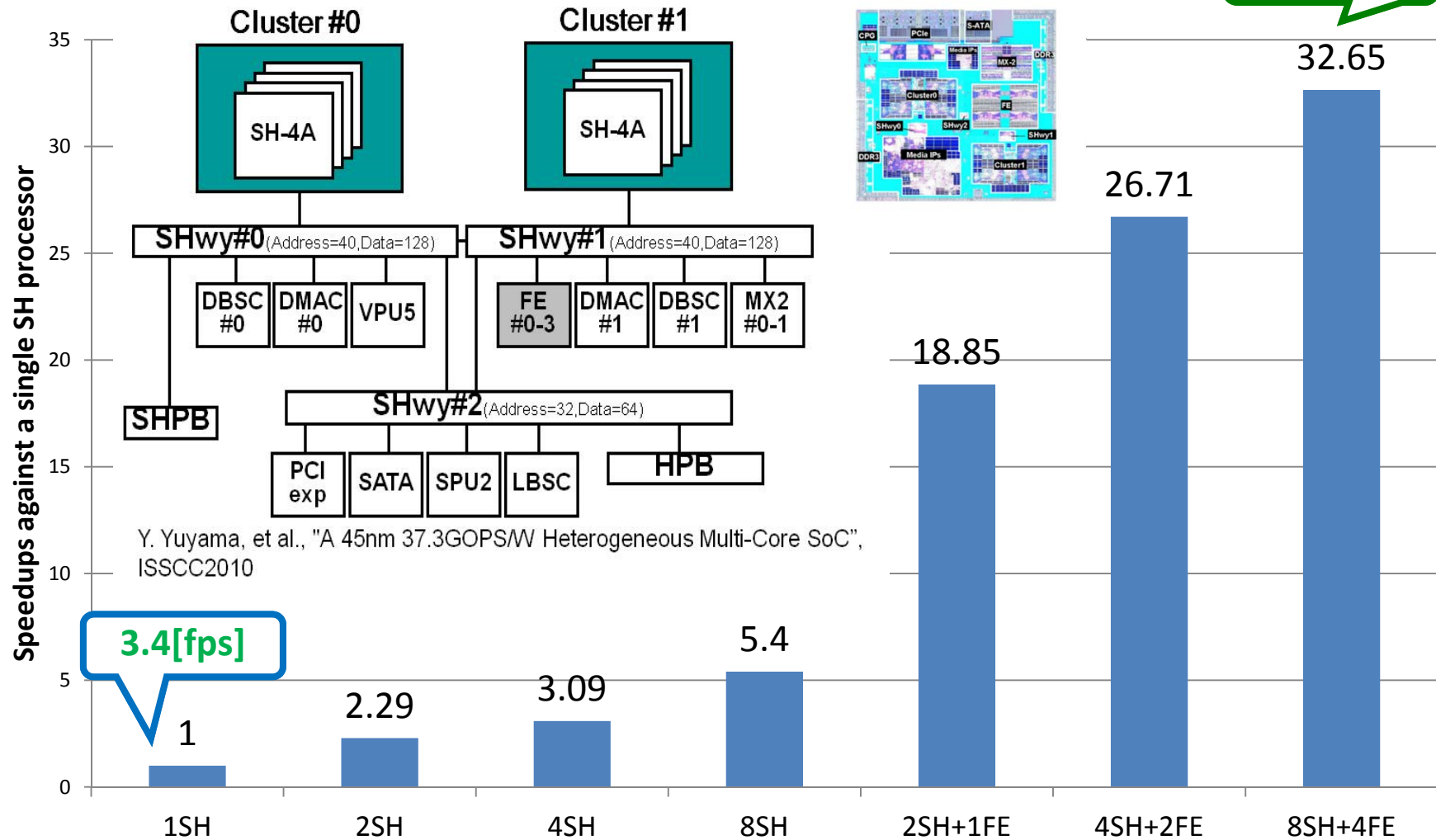Vessel Detection : http://www.mathworks.co.jp/matlabcentral/fileexchange/24990-retinal-blood-vessel-extraction/

# 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

# 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

# Power Reduction in a real-time execution controlled by OSCAR Compiler and OSCAR API on RP-X (Optical Flow with a hand-tuned library)
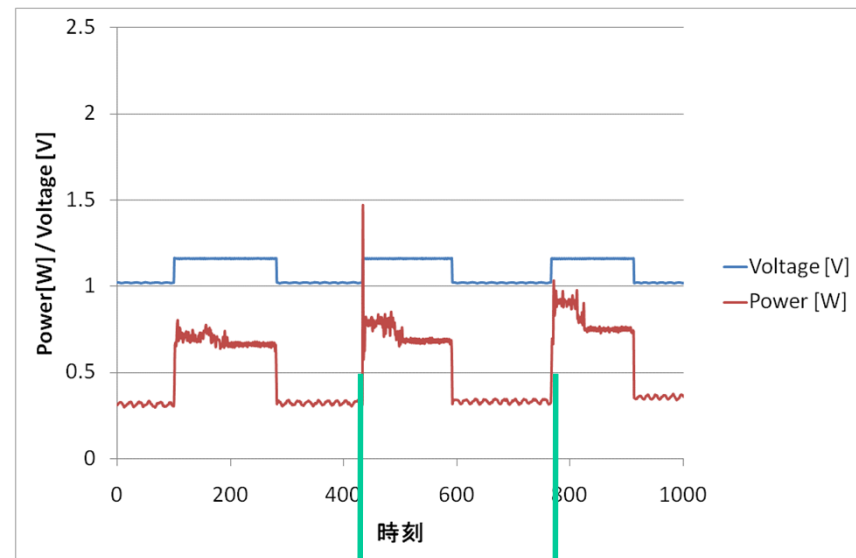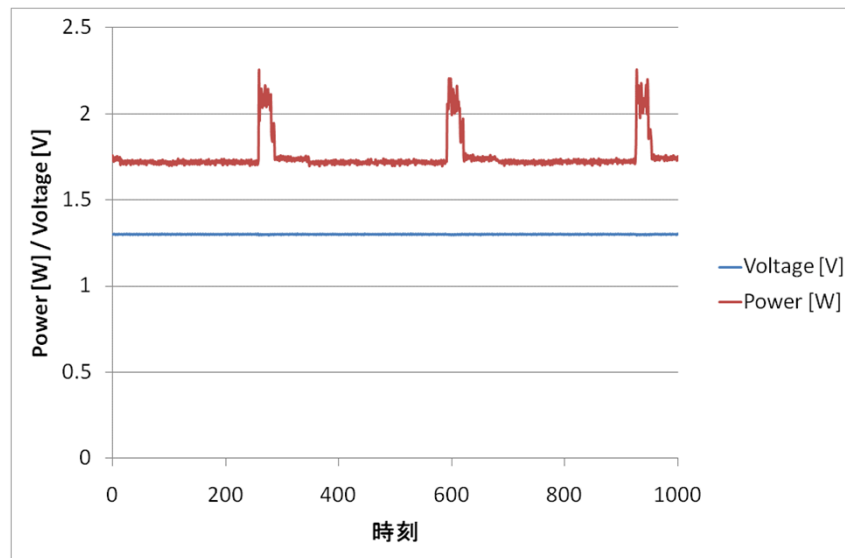
**Without Power Reduction**

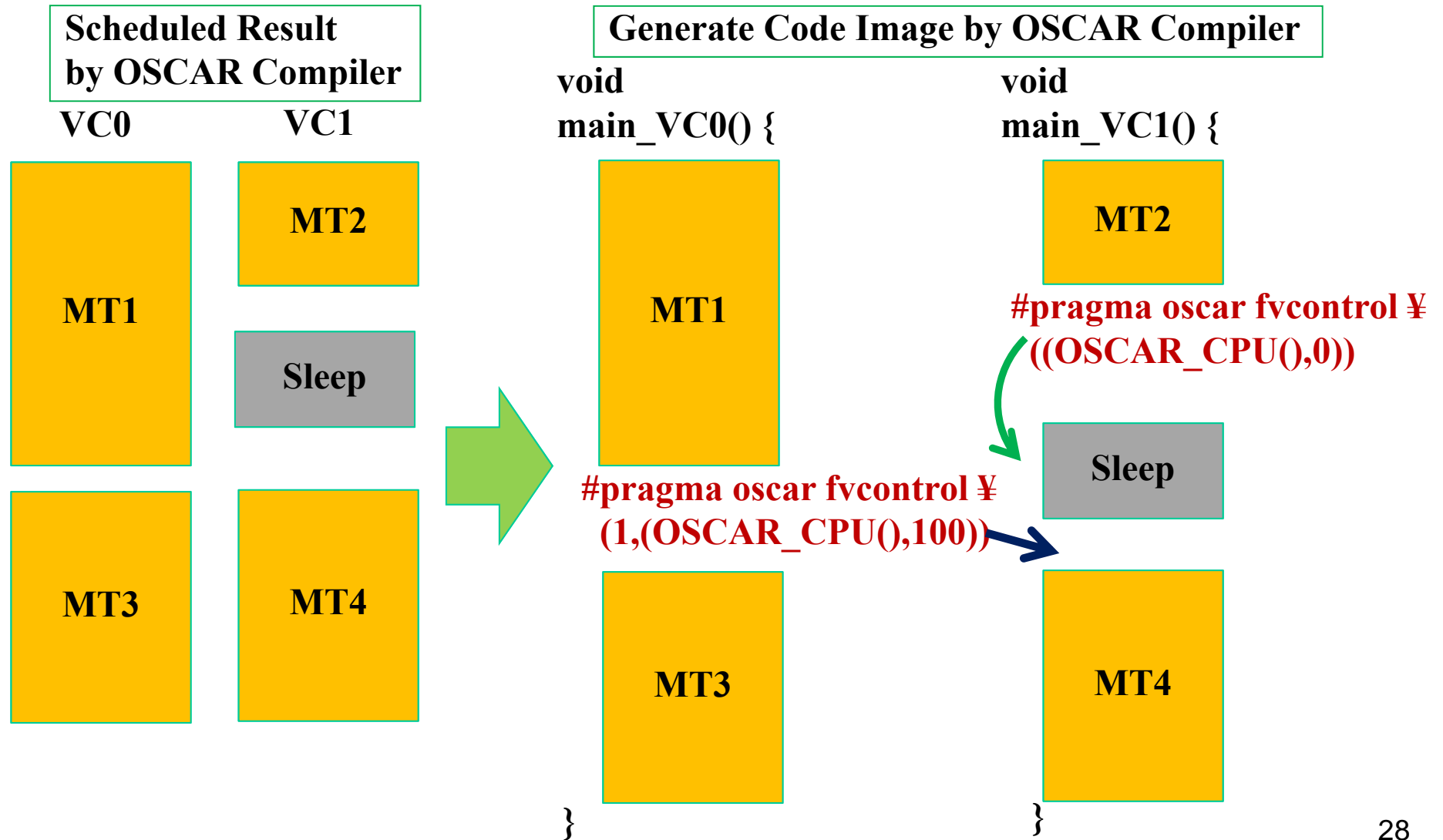**With Power Reduction by OSCAR Compiler**

**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

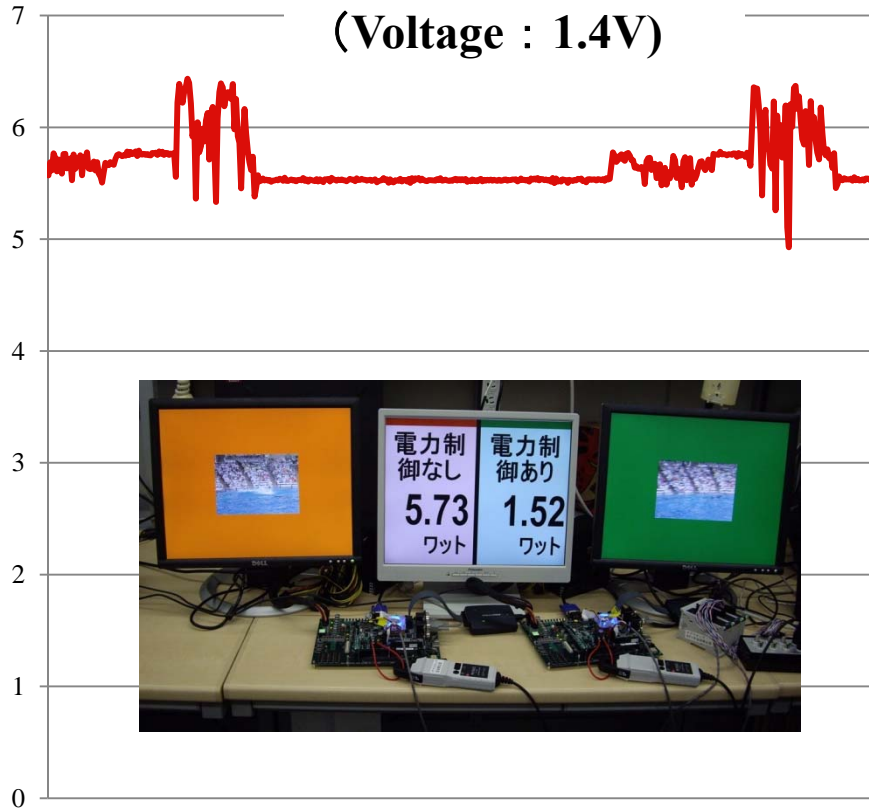#pragma oscar fvcontrol ¥
((OSCAR_CPU(),0))

Sleep

MT4

}

28

# Power Reduction of MPEG2 Decoding to 1/4 on 8 Core Homogeneous Multicore RP-2 by OSCAR Parallelizing Compiler
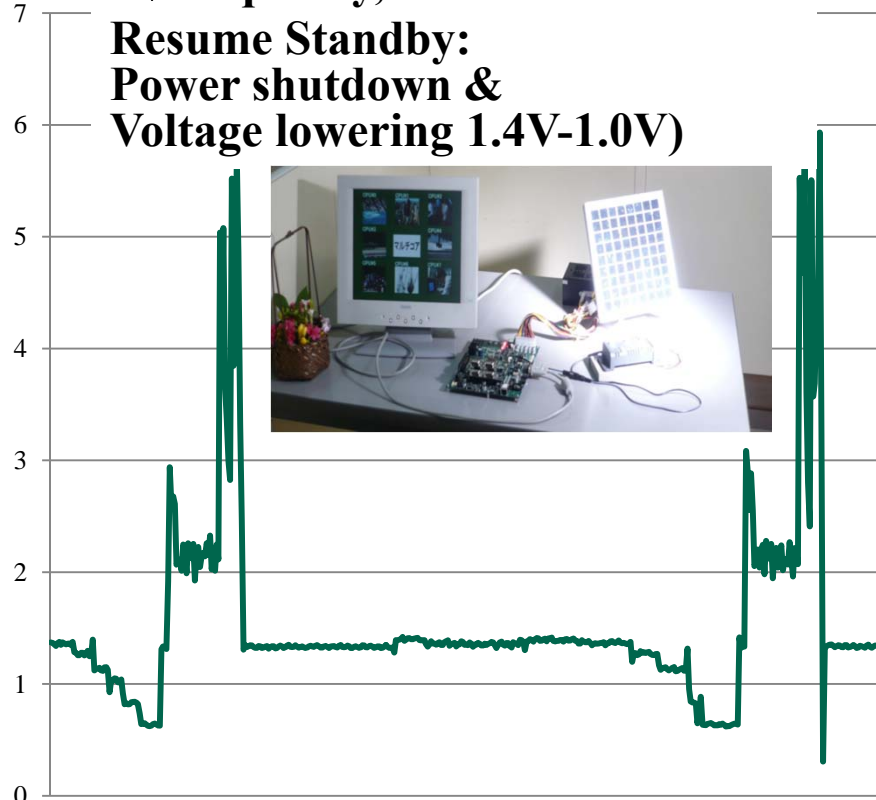
**MPEG2 Decoding with 8 CPU cores**



**Without Power Control（Voltage：1.4V)**

**With Power Control（Frequency, Resume Standby: Power shutdown & Voltage lowering 1.4V-1.0V)**
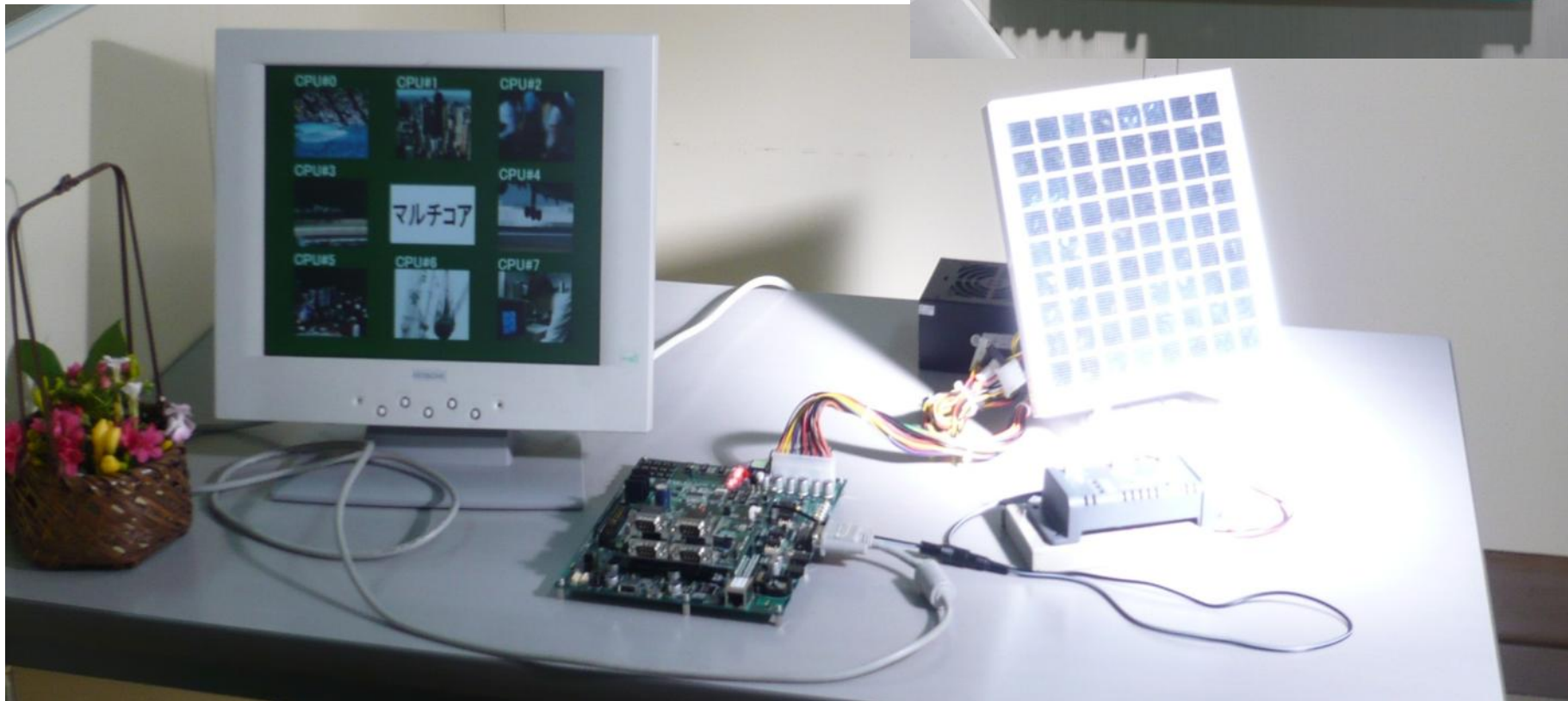
**Avg. Power 5.73 [W]**

**73.5% Power Reduction**

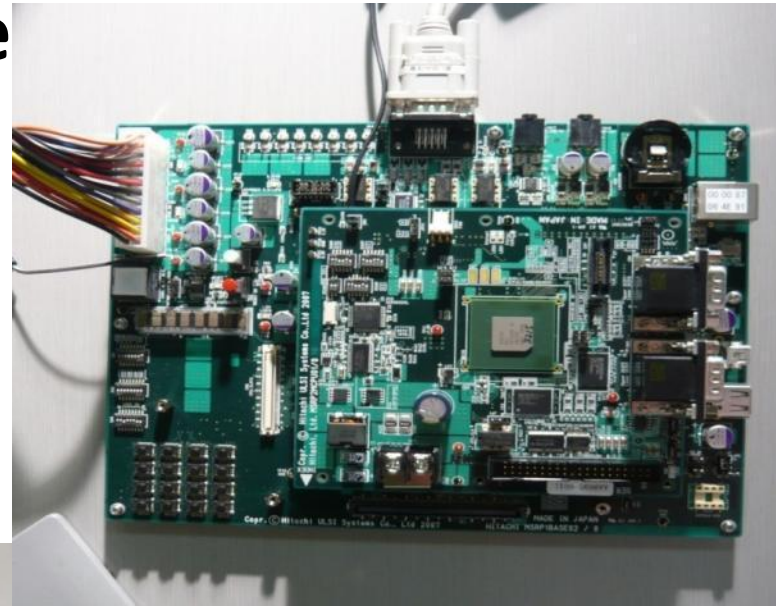**Avg. Power 1.52 [W]**

# Low Power High Performance Multicore Computer with Solar Panel

> **Clean Energy Autonomous**

> **Servers operational in deserts**
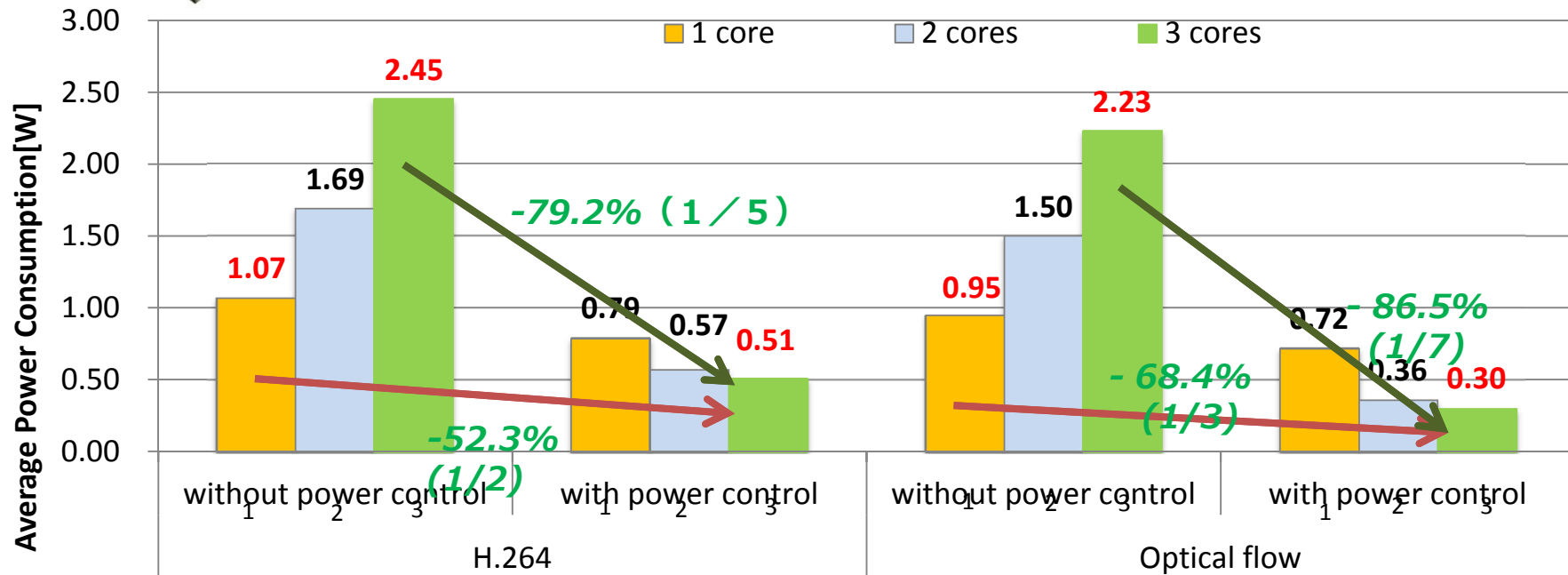
# Power on 4 cores ARM CortexA9 with Android

## H.264 decoder & Optical Flow  (Using 3 cores)
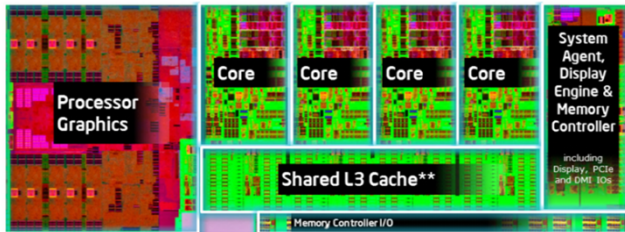
**ODROID X2**

**Samsung Exynos4412 Prime, ARM Cortex-A9 Quad core 1.7GHz～0.2GHz, used by Samsung's Galaxy S3**



- ➢ **On the same 3 cores, the power control reduced the power to 1/5～1/7 against no power control.**
- ➢ **The power control reduced the power to 1/2～1/3 compared with the ordinary sequential execution on 1 core without power control.**

31

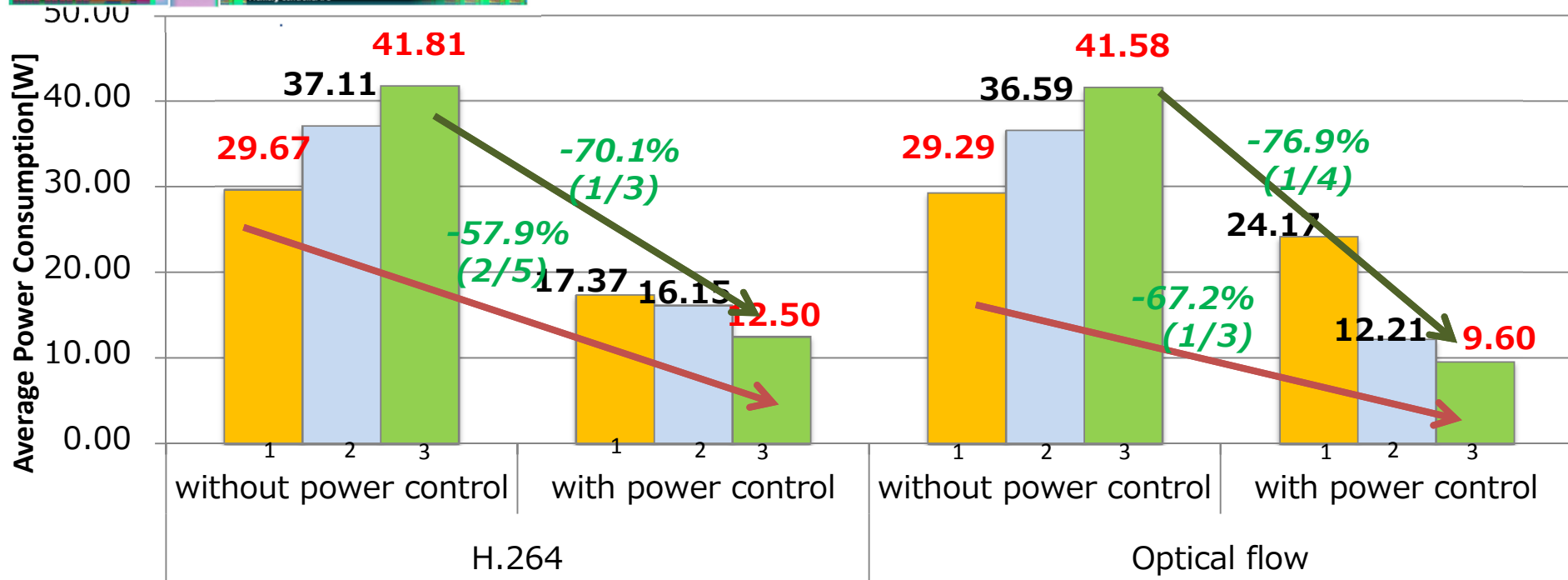# Power Reduction on Intel Haswell 3cores
## H.264 decoder & Optical Flow

**H81M-A, Intel Core i7 4770k**

**Quad core, 3.5GHz〜0.8GHz**



- The power consumption was reduced to **1/3〜1/4** by OSCAR power control against no power control on the same 3 processor cores.
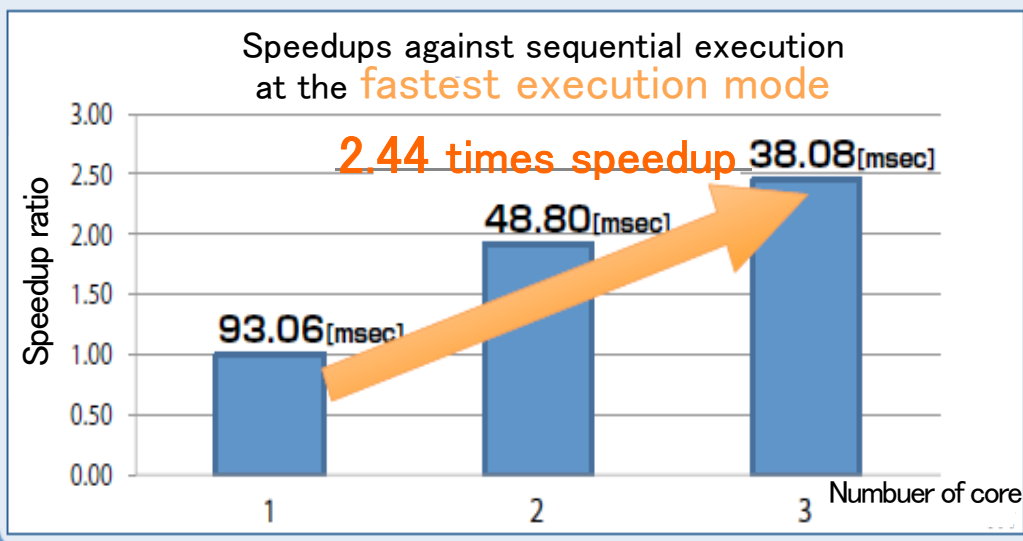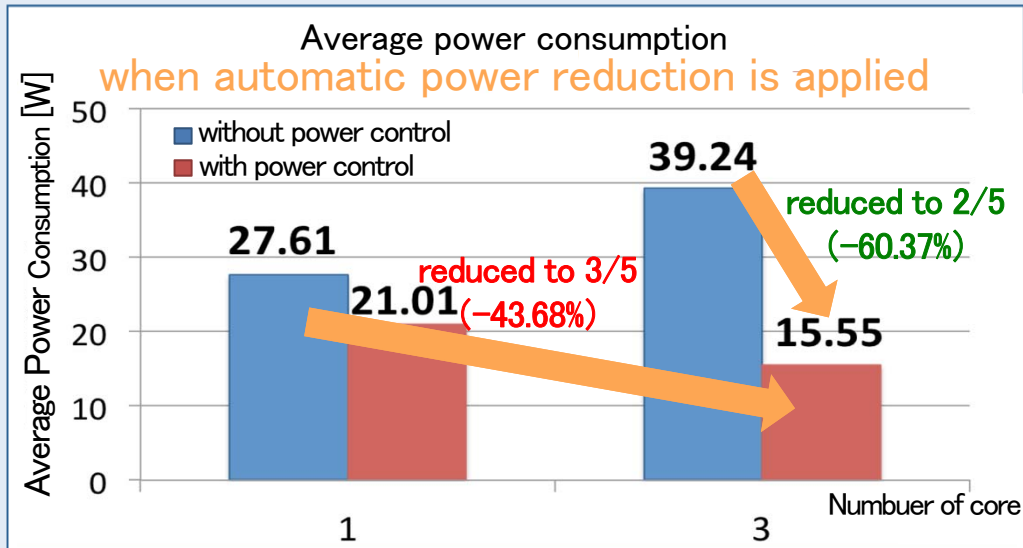- The power reduced to **2/5〜1/3 by the compiler power control against 1core no power control.**

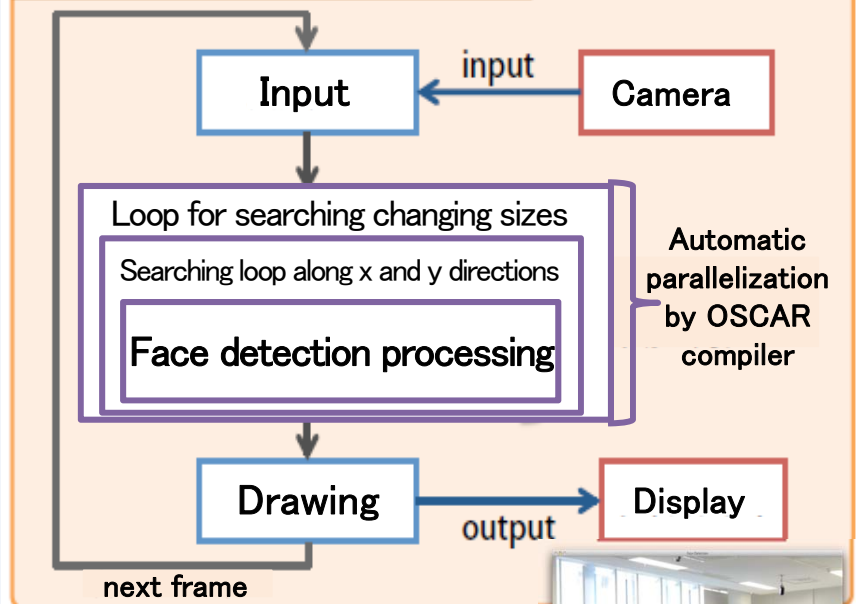# Automatic Power Reduction by OSCAR Compiler on Intel Haswell 4 Core Multicore
## - Power Consumption for real-time face detection was reduced to 2/5 -

WASEDA UNIVERSITY

- OSCAR Compiler
- Intel Haswell
- Power Reduction

## Parallel processing of face detection program on Intel Haswell 4cores

Average power consumption
when automatic power reduction is applied

- without power control
- with power control

27.61
21.01 (−43.68%)    reduced to 3/5
39.24
reduced to 2/5 (−60.37%)
15.55

Average Power Consumption [W]

Numbuer of core: 1, 3

Speedups against sequential execution
at the fastest execution mode

2.44 times speedup    38.08[msec]

93.06[msec]
48.80[msec]

Speedup ratio

Numbuer of core: 1, 2, 3

## Parallelization flow of OpenCV face detection program

Input ← input ← Camera

Loop for searching changing sizes

Searching loop along x and y directions

Face detection processing

Automatic parallelization by OSCAR compiler
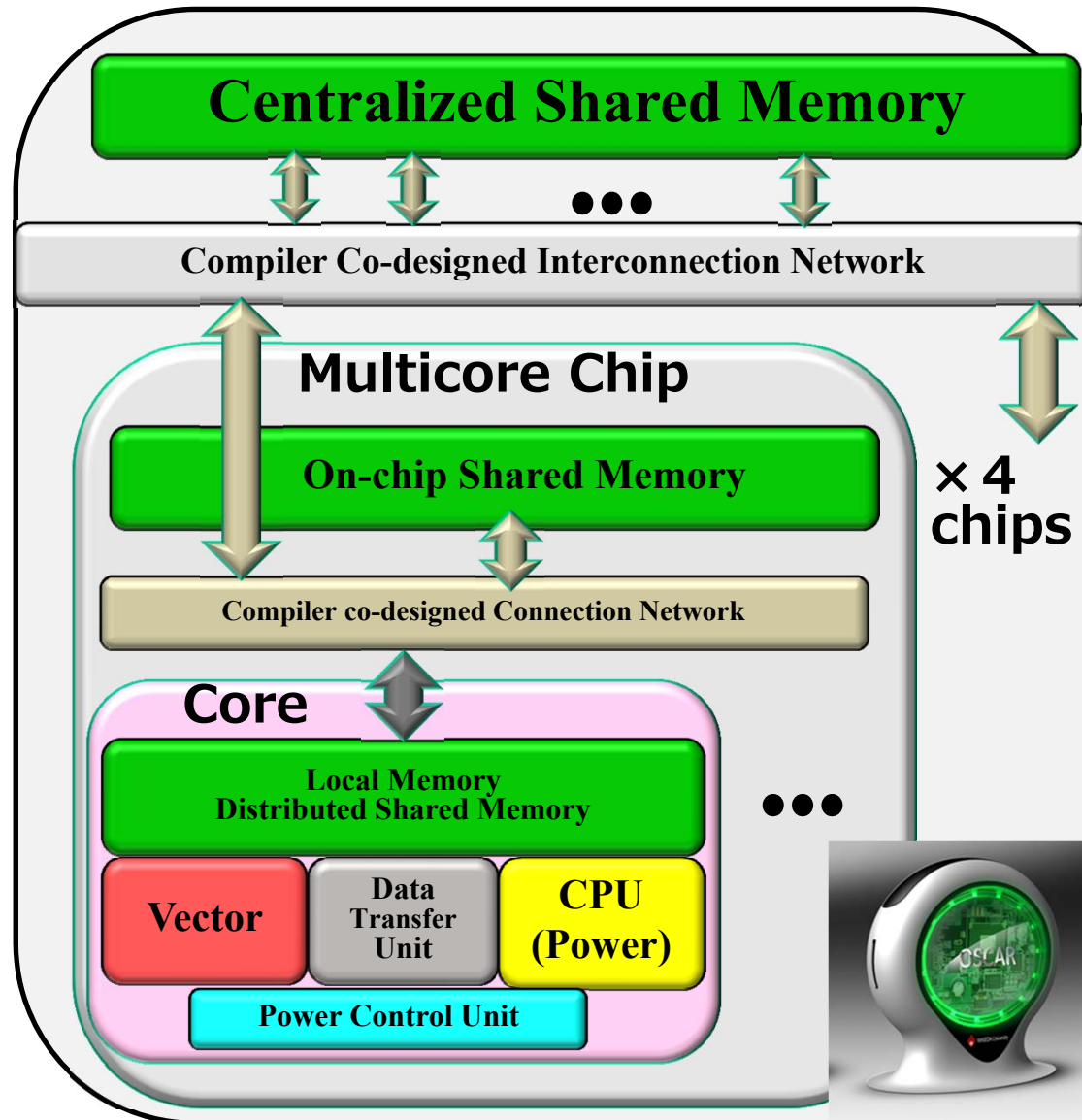
Drawing → output → Display

next frame

## Power measurement on Intel Haswell board

CPU : Inel Core i7 4770K
Number of core : 4
Clock frequency : 3.5GHz～0.8GHz
Mother board : ASUS H81M−A

Inserting power measurement circuit between PMIC and CPU

# OSCAR Vector Multicore with Power Processor Core 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)

Power Control Unit

**Target:**

➢ Solar Powered with compiler power reduction.

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

## Summary

- OSCAR Automatic Parallelizing and Power Reducing Compiler has succeeded speedup and/or power reduction of scientific applications including medical applications and natural disaster simulation, and real-time applications like Automobile Engine Control and MATLAB/Simulink, and media applications codec and face detection on various homogeneous and heterogeneous multicores.

- In automatic parallelization on Power Architectures:

    OSCAR Compiler has been developed on Power architectures like Power 4, 5, 5+, 6, 7 and 8.

    - for "Earthquake Wave Propagation Simulation", 110 times speedup on 128 cores of IBM Power 7 and 9.6 times speedup on 12 cores IBM Power8 against 1 core,

    - for "Cancer Treatment Using Carbon Ion", 55 times speedup on 64 cores of IBM Power 7 against 1 core

- In automatic power reduction, consumed powers were reduced to 1/2 or 1/3 using 3 cores on ARM Cortex A9 and Intel Haswell.

- We are planning to realize automatic power reduction using "On-Chip Regulator" on Power Processors with accelerators.