

OSCAR Compiler for Automatic Multigrain Parallelization, Memory Optimization and Power Systems



Hironori Kasahara, Ph.D., IEEE Fellow

IEEE Computer Society President 2018

Professor, Dept. of Computer Science & Engineering

Director, Advanced Multicore Processor Research Institute

Waseda University, Tokyo, Japan

URL: <http://www.kasahara.cs.waseda.ac.jp/Reduction for>

1980 BS, 82 MS, 85 Ph.D. , Dept. EE, Waseda Univ.
1985 Visiting Scholar: U. of California, Berkeley
1986 Assistant Prof., 1988 Associate Prof., 1997,
Waseda Univ., Now Dept. of Computer Sci. & Eng.
1989-90 Research Scholar: U. of Illinois, Urbana-
Champaign, Center for Supercomputing R&D
2004 Director, Advanced Multicore Research
Institute, 2017 member: the Engineering Academy
of Japan and the Science Council of Japan

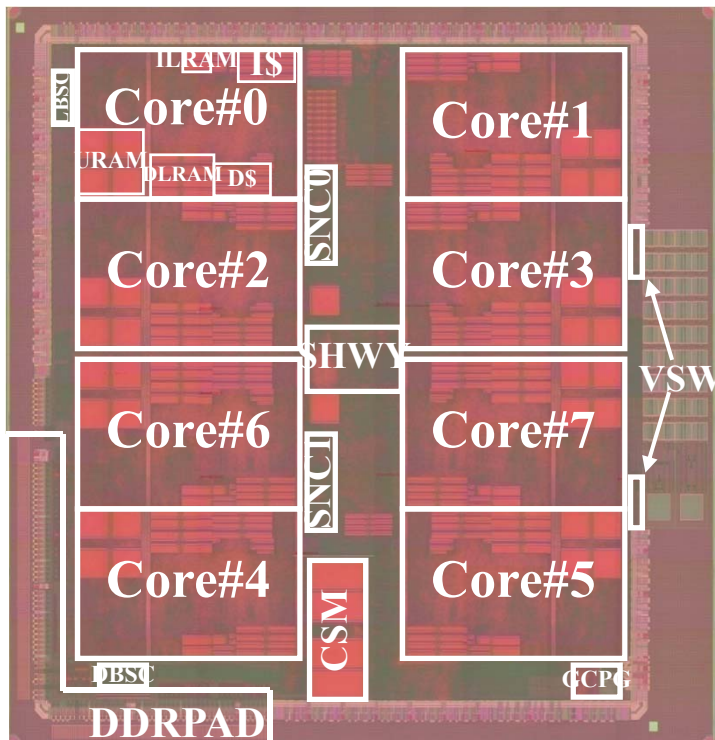
2005 STARC Academia-Industry Research Award
2008 LSI of the Year Second Prize
2008 Intel AsiaAcademic Forum Best Research Award
2010 IEEE CS Golden Core Member Award
2014 Minister of Edu., Sci. & Tech. Research Prize
2015 IPSJ Fellow
2017 IEEE Fellow, IEEE Eta Kappa Nu

Reviewed Papers: 216, Invited Talks: 162, Granted
Patents: 43 (Japan, US, GB, China), Articles in News
Papers, Web News, Medias incl. TV etc.: 584

Committees in Societies and Government 255
IEEE Computer Society President 2018,
IEEE CS: BoG(2009-14), Executive Committee(2017-
Multicore STC Chair (2012-), Japan Chair (2005-07),
IPSJ Chair: HG for Magazine. & J. Edit, Sig. on ARC.
【METI/NEDO】 Project Leaders: Multicore for
Consumer Electronics, Advanced Parallelizing
Compiler, Chair: Computer Strategy Committee
【Cabinet Office】 CSTP Supercomputer Strategic
ICT PT, Japan Prize Selection Committees, etc.
【MEXT】 Info. Sci. & Tech. Committee,
Supercomputers (Earth Simulator, HPCI Promo.,
Next Gen. Supercomputer K) Committees, etc.

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”

$$\text{Power} \propto \text{Frequency} * \text{Voltage}^2$$

(Voltage \propto Frequency)

➔ Power \propto Frequency³

If Frequency is reduced to 1/4
(Ex. 4GHz \rightarrow 1GHz),
Power is reduced to 1/64 and
Performance falls down to 1/4 .

<Multicores>

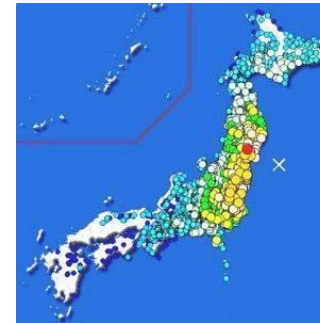
If 8cores are integrated on a chip,
Power is still 1/8 and
Performance becomes 2 times .

Parallel Soft is important for scalable performance of multicore (LCPC2015)

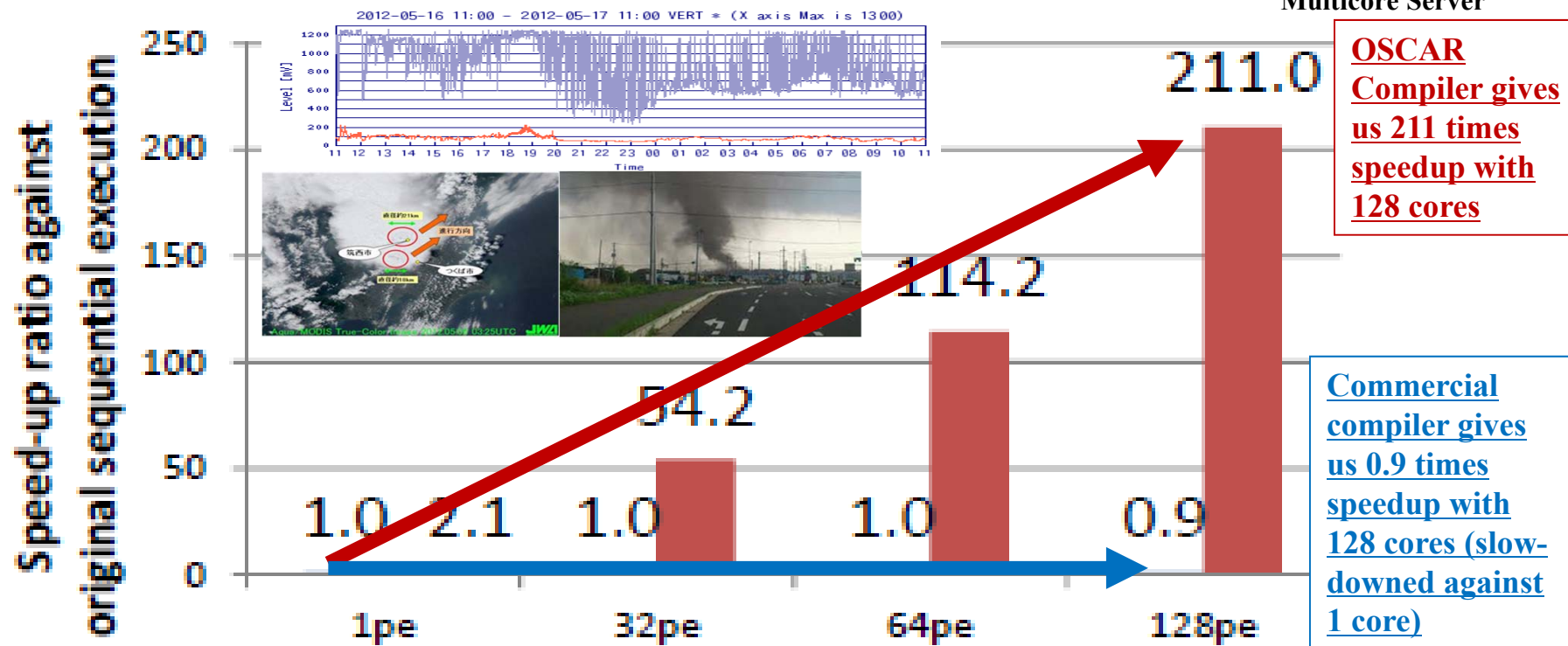
- Just more cores don't give us speedup
- Development cost and period of parallel software are getting a bottleneck of development of embedded systems, eg. IoT, Automobile

Earthquake wave propagation simulation GMS developed by National Research Institute for Earth Science and Disaster Resilience (NIED)

■ original (sun studio) ■ proposed method



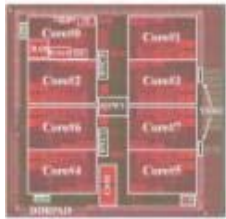
Fujitsu M9000 SPARC Multicore Server



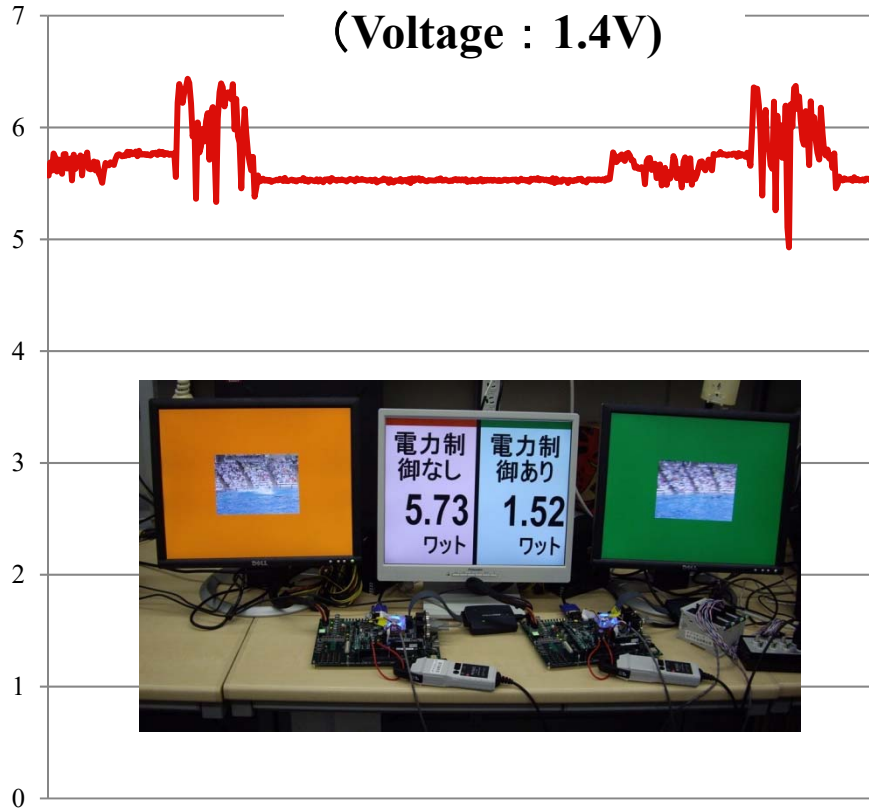
- Automatic parallelizing compiler available on the market gave us no speedup against execution time on 1 core on 64 cores
 - Execution time with 128 cores was slower than 1 core (0.9 times speedup)
- Advanced OSCAR parallelizing compiler gave us 211 times speedup with 128cores against execution time with 1 core using commercial compiler
 - OSCAR compiler gave us 2.1 times speedup on 1 core against commercial compiler by global cache optimization

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

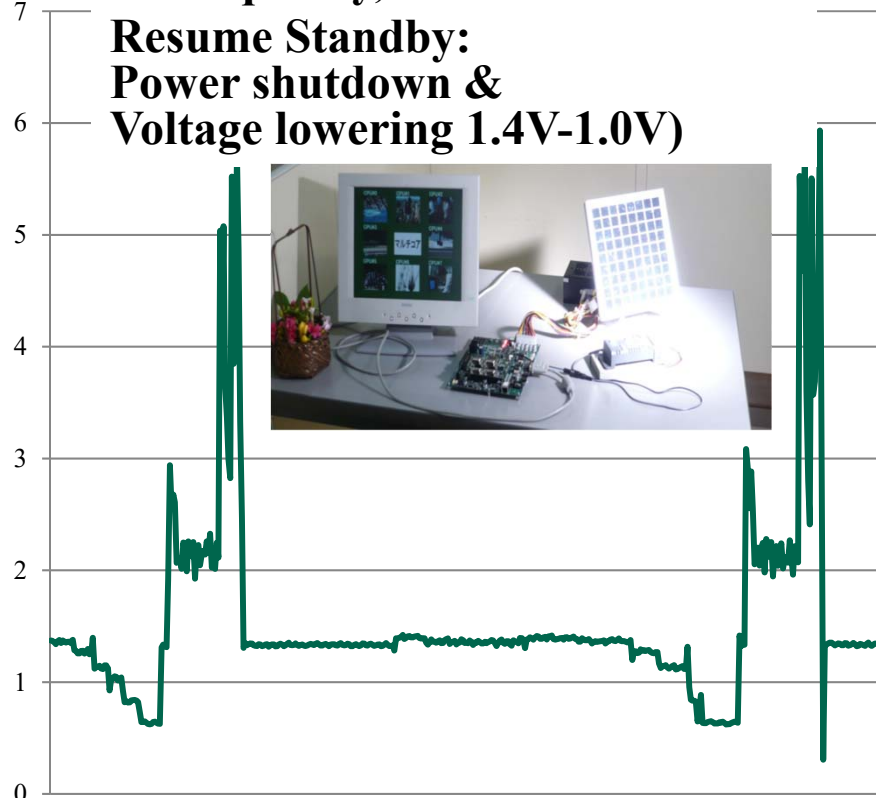
MPEG2 Decoding with 8 CPU cores



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]

OSCAR Parallelizing Compiler

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

Multigrain Parallelization (LCPC1991,2001,04)
 coarse-grain parallelism among loops and subroutines (2000 on SMP), near fine grain parallelism among statements (1992) in addition to loop parallelism

Data Localization

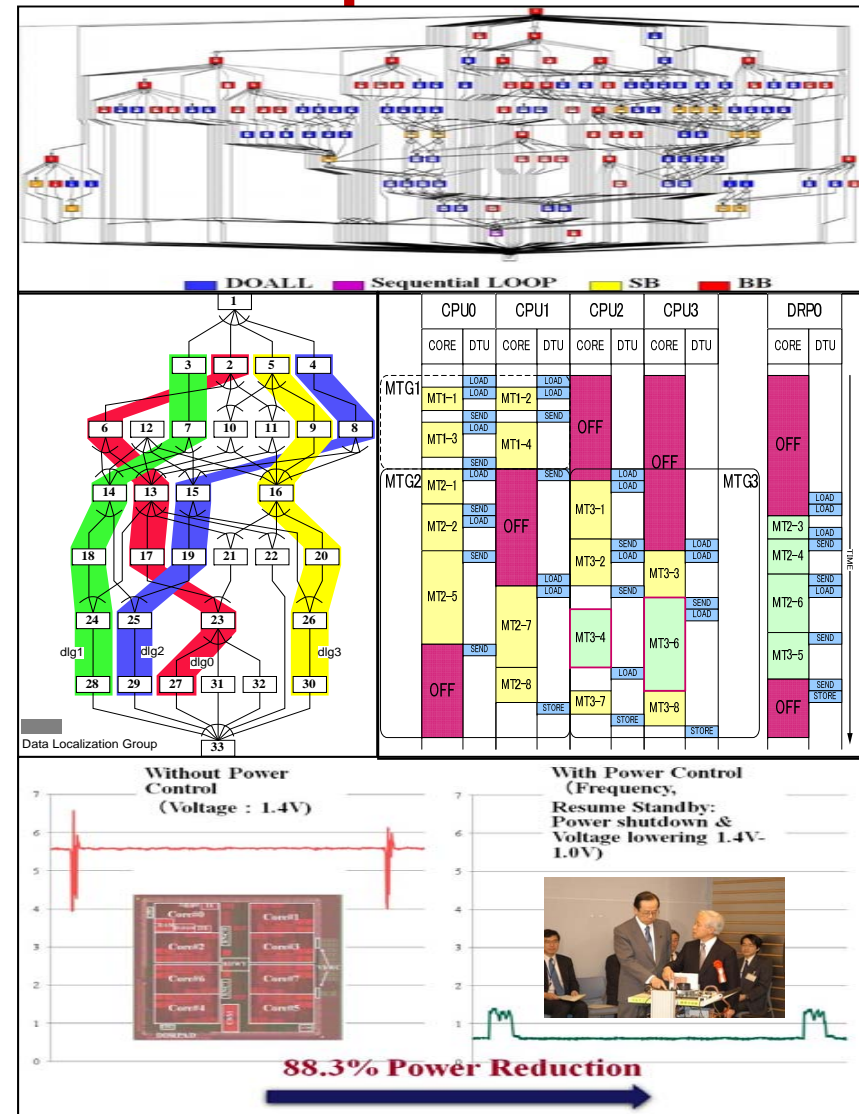
Automatic data management for distributed shared memory, cache and local memory (Local Memory 1995, 2016 on RP2, Cache2001,03)
 Software Coherent Control (2017)

Data Transfer Overlapping (2016 partially)

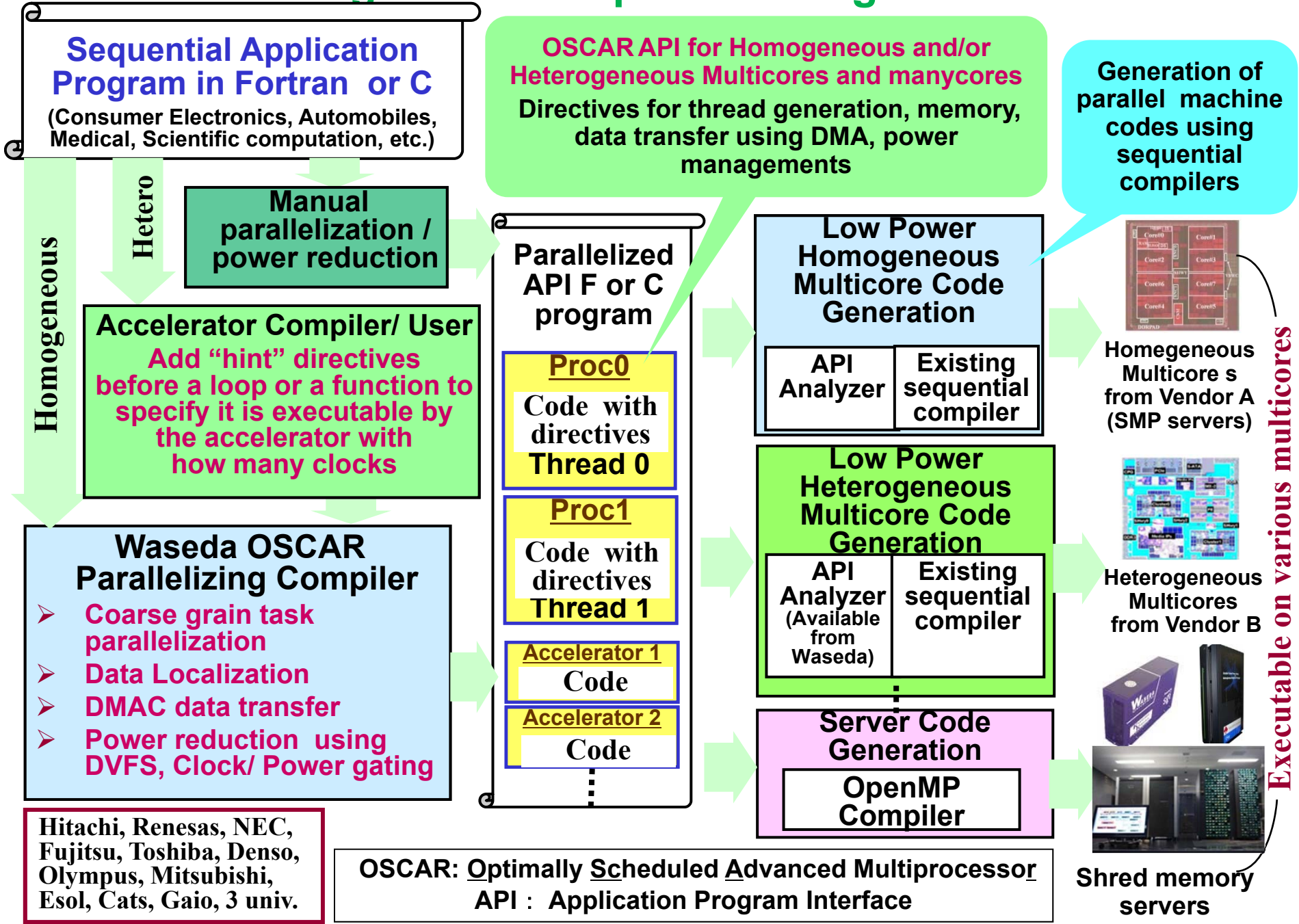
Data transfer overlapping using Data Transfer Controllers (DMAs)

Power Reduction

(2005 for Multicore, 2011 Multi-processes, 2013 on ARM)
 Reduction of consumed power by compiler control DVFS and Power gating with hardware supports.



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

Manual parallelization / power reduction

Homogeneous

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

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

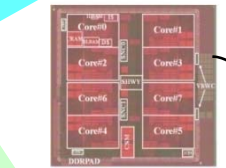
OpenMP Compiler

Waseda OSCAR Parallelizing Compiler

- Coarse grain task parallelization
- Data Localization
- DMAC data transfer
- Power reduction using DVFS, Clock/ Power gating

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

OSCAR: Optimally Scheduled Advanced Multiprocessor API : Application Program Interface



Homogeneous Multicores from Vendor A (SMP servers)



Heterogeneous Multicores from Vendor B

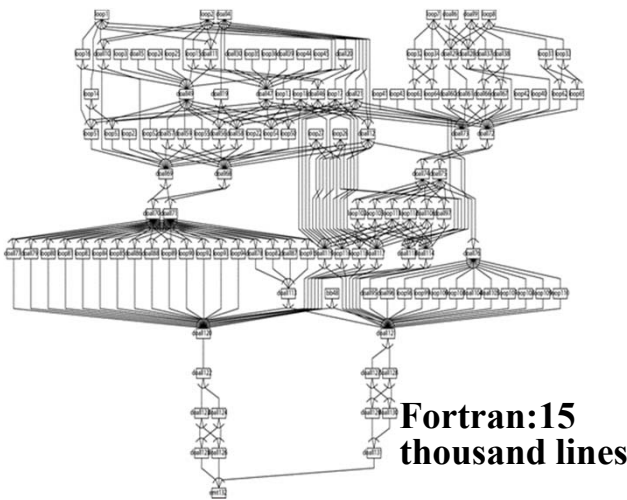
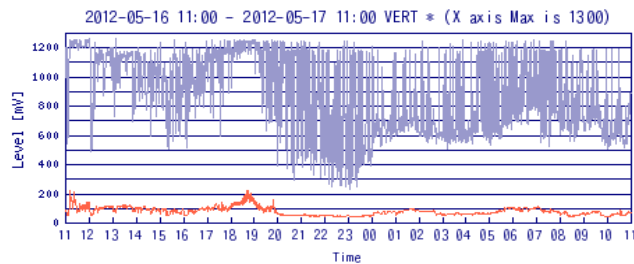


Shred memory servers

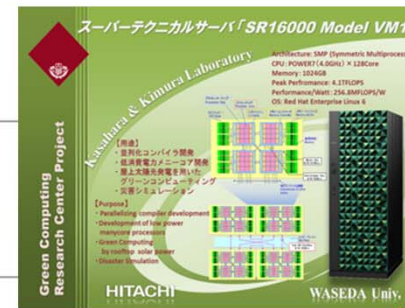
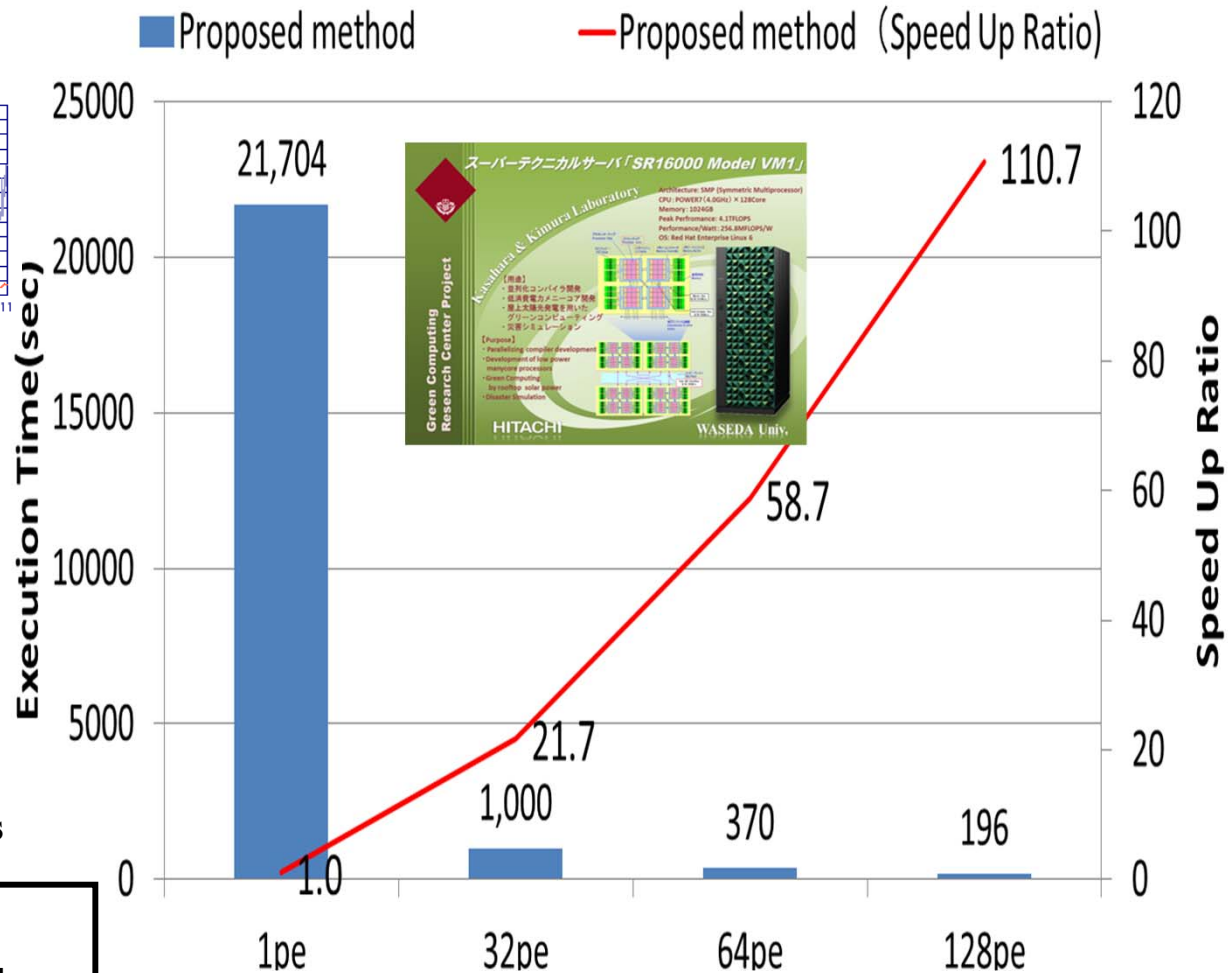
Executable on various multicores

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

(Power7 Based 128 Core Linux SMP) [\(LCPC2015\)](#)



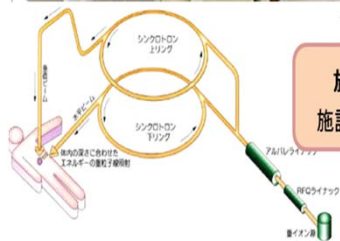
First touch for distributed shared memory and cache optimization over loops are important for scalable speedup



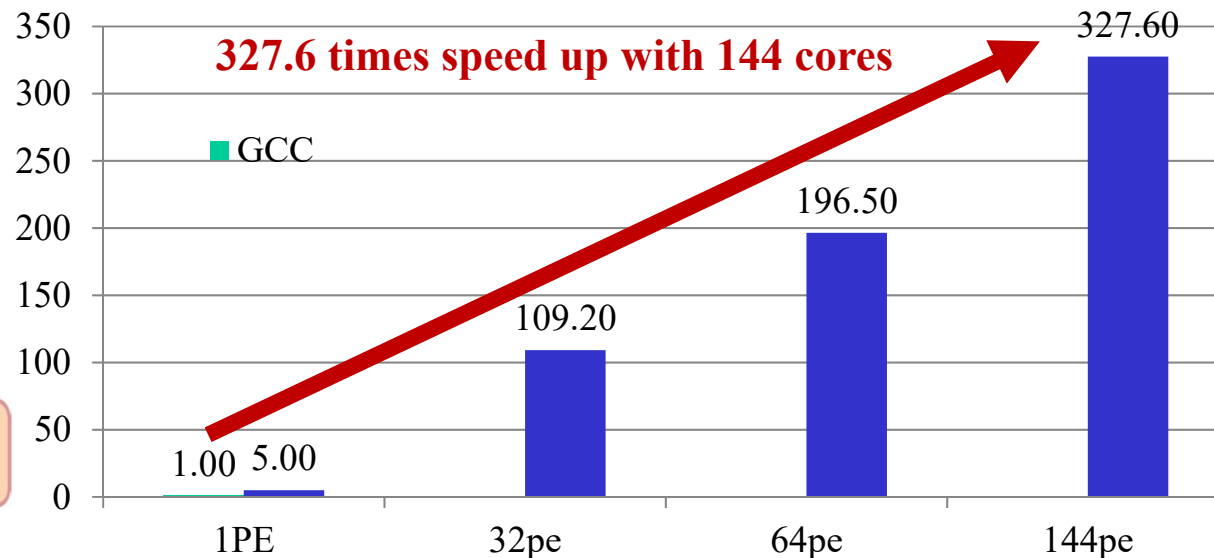
Performance on Multicore Server for Latest Cancer Treatment Using Heavy Particle (Proton, Carbon Ion)

327 times speedup on 144 cores

Hitachi 144cores SMP Blade Server BS500:
Xeon E7-8890 V3(2.5GHz 18core/chip) x8 chip

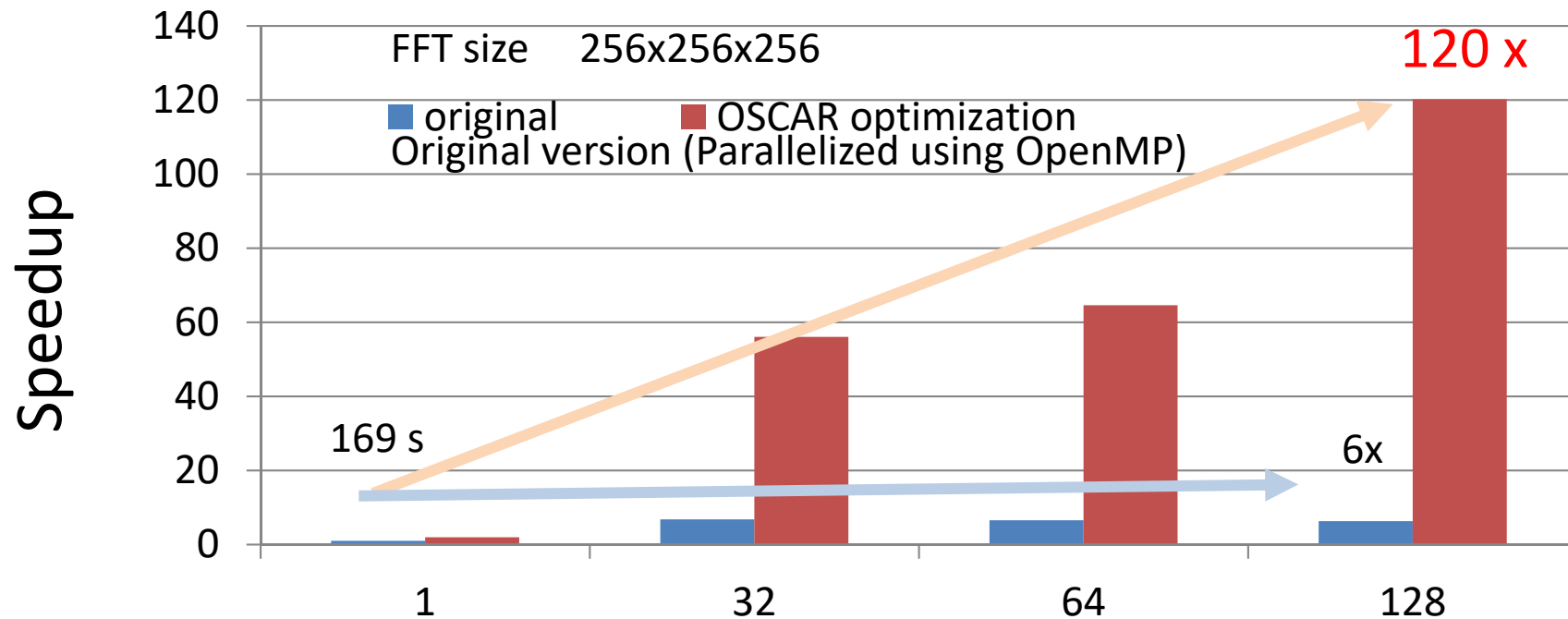


放射線医学研究所
施設の費用: 120億円



- Original **sequential execution time 2948 sec (50 minutes)** using GCC was reduced to **9 sec with 144 cores** (327.6 times speedup)
- Reduction of treatment cost and reservation waiting period is expected

Parallelization of 3D-FFT for New Magnetic Material Computation on Hitachi SR16000 Power7 CC-Numa Server



OSCAR optimization

- reducing number of data transpose with interchange, code motion and loop fusion

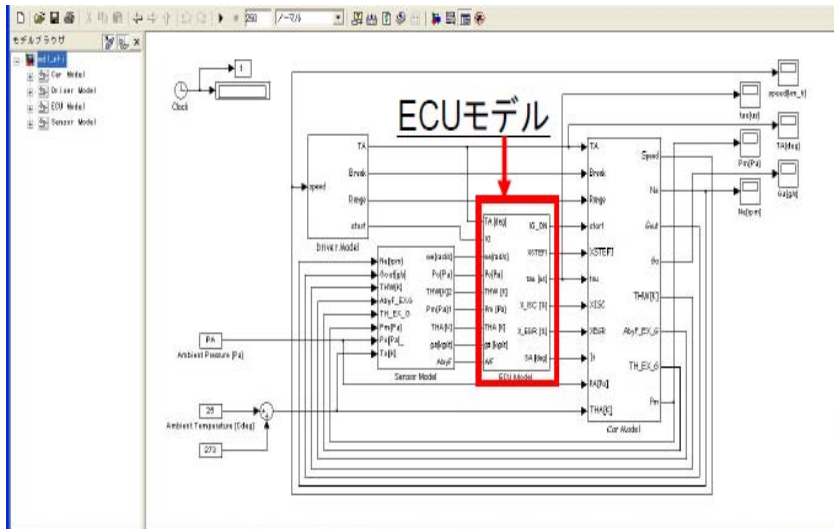
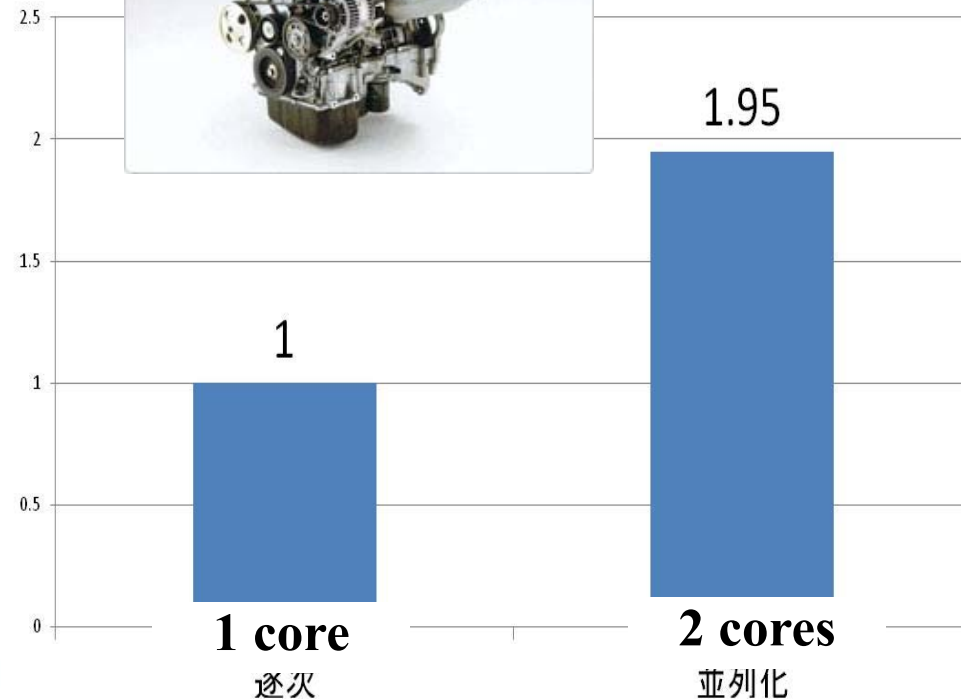


Engine Control by multicore with Denso

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 using local memories
- Millions of lines C codes consisting conditional branches and basic blocks

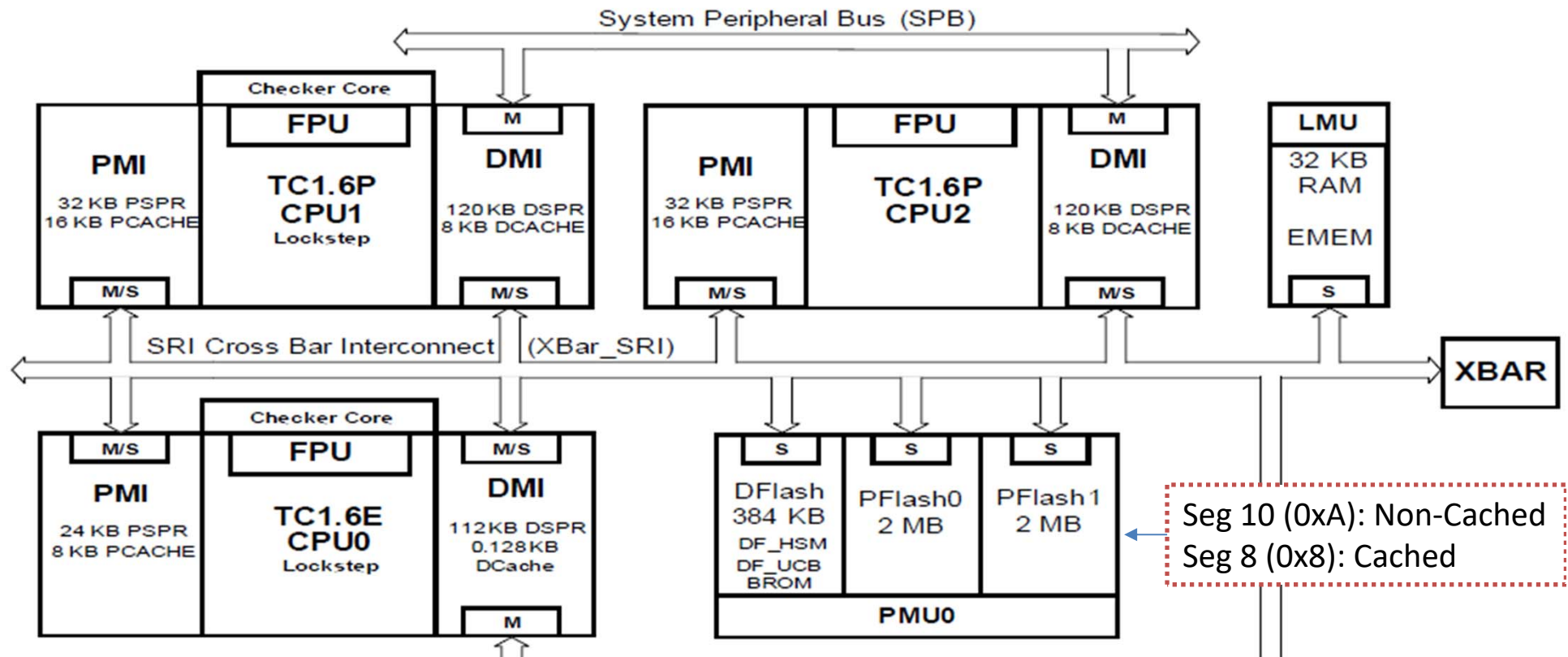


Infineon AURIX TC277



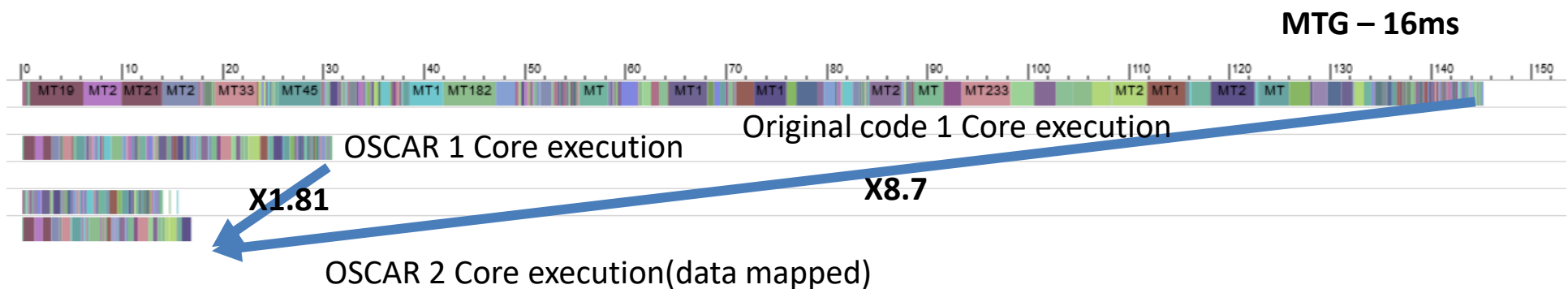
Abbreviations :

PCACHE:	Program Cache
DCACHE:	Data Cache
DSPR:	Data Scratch-Pad RAM
PSPR:	Program Scratch-Pad RAM
BROM:	Boot ROM
PFlash:	Program Flash
DFlash:	Data Flash (EEPROM)
S :	SRI Slave Interface
M :	SRI Master Interface

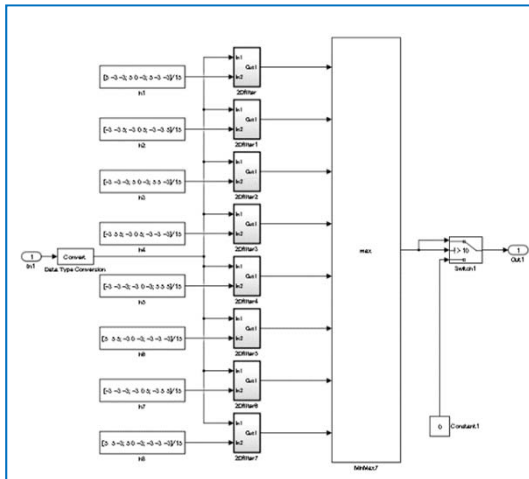


Automatic Parallelization of an Engine Control C Program with 400 thousands lines on AUTOSAR on 2 cores of Infineon AURIX TC277

- **Original sequential** execution time on 1 core: **145500** cycles
- **Sequential execution time by OSCAR** on 1 core: **29700** cycles
 - **4.9 times speedup on 1 core** against original execution by OSCAR Compilers automatic data allocation for local scratch pad memory, flush memory modules
- **2 core execution by OSCAR** Compiler: **16400** cycles
 - **1.81 times speedup with 2 core** against **1 core execution with OSCAR Compiler**
 - **8.7 times speedup against original sequential execution.**



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:
     * Import: '<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' */

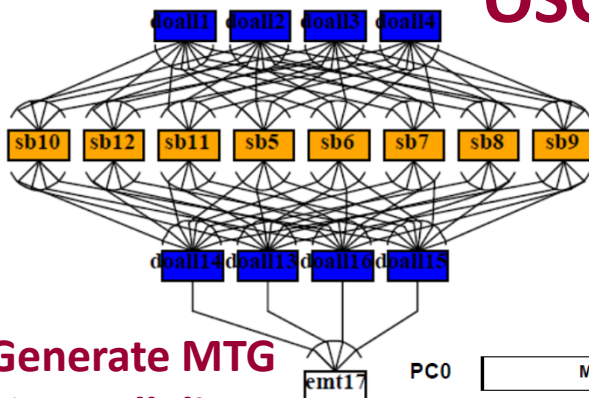
    /* Constant: '<S1>/h1' */
    VesselExtraction_Dfilter(VesselExtraction_B.DataTypeConversion,
        VesselExtraction_P.h1_Value, &VesselExtraction_B.Dfilter,
        (P_Filter_VesselExtraction_T *)&VesselExtraction_P.Dfilter);

    /* 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_Filter_VesselExtraction_T *)&VesselExtraction_P.Dfilter1);
}
    
```

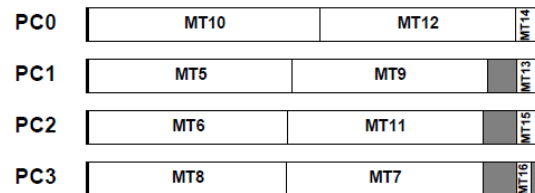
C code



OSCAR Compiler

(1) Generate MTG
→ Parallelism

(2) Generate gantt chart
→ Scheduling in a multicore



0.0E+00 4.0E-02
TIME [s]



```

void VesselExtraction_step ( )
{
    int thr1 ;
    int thr2 ;
    int thr3 ;

    void thread_function_001 ( void )
    {
        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 ) ;

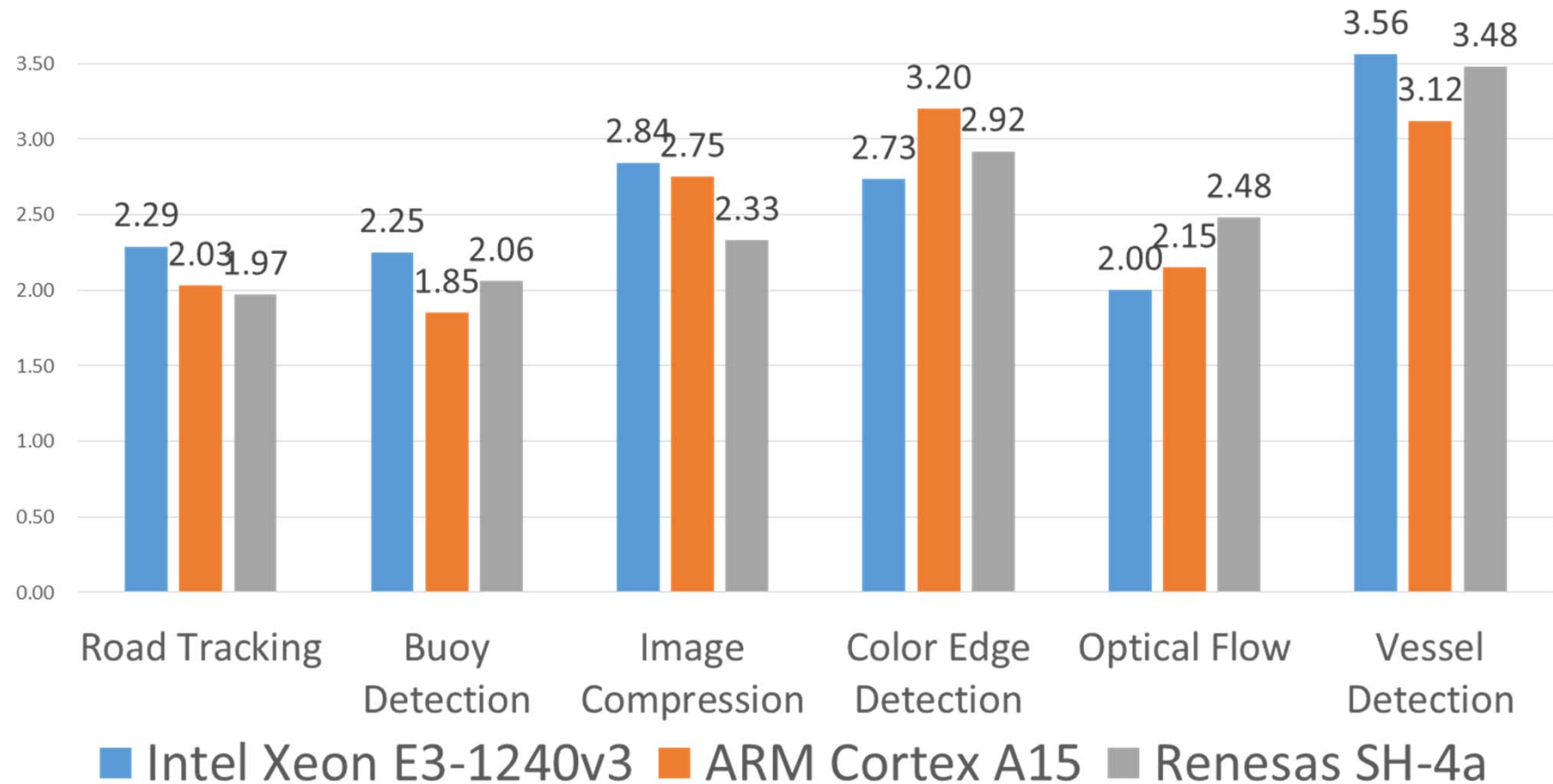
    VesselExtraction_step_PEO ( ) ;

    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)



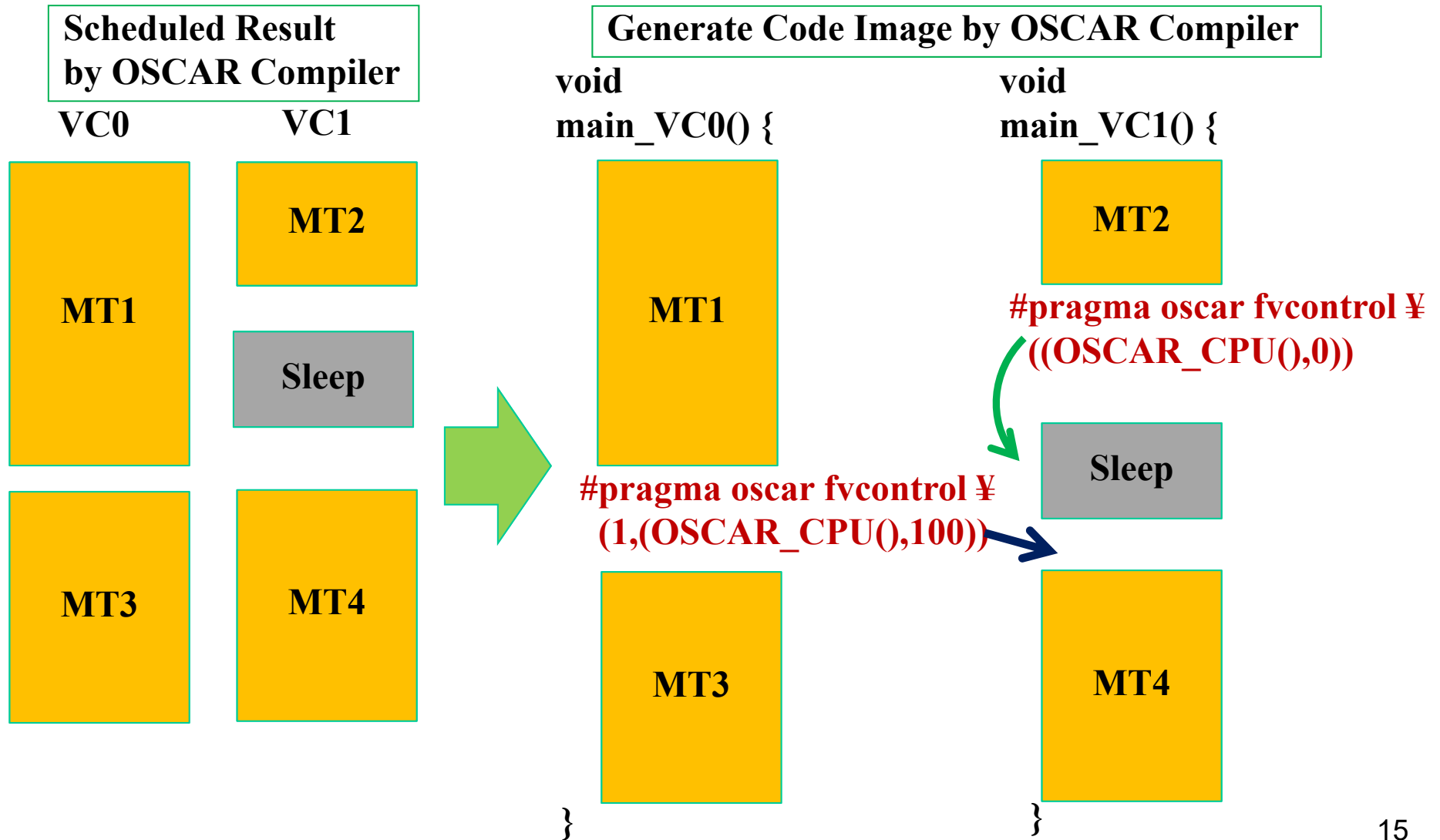
Road Tracking, Image Compression : <http://www.mathworks.co.jp/jp/help/vision/examples>

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/>

Low-Power Optimization with OSCAR API

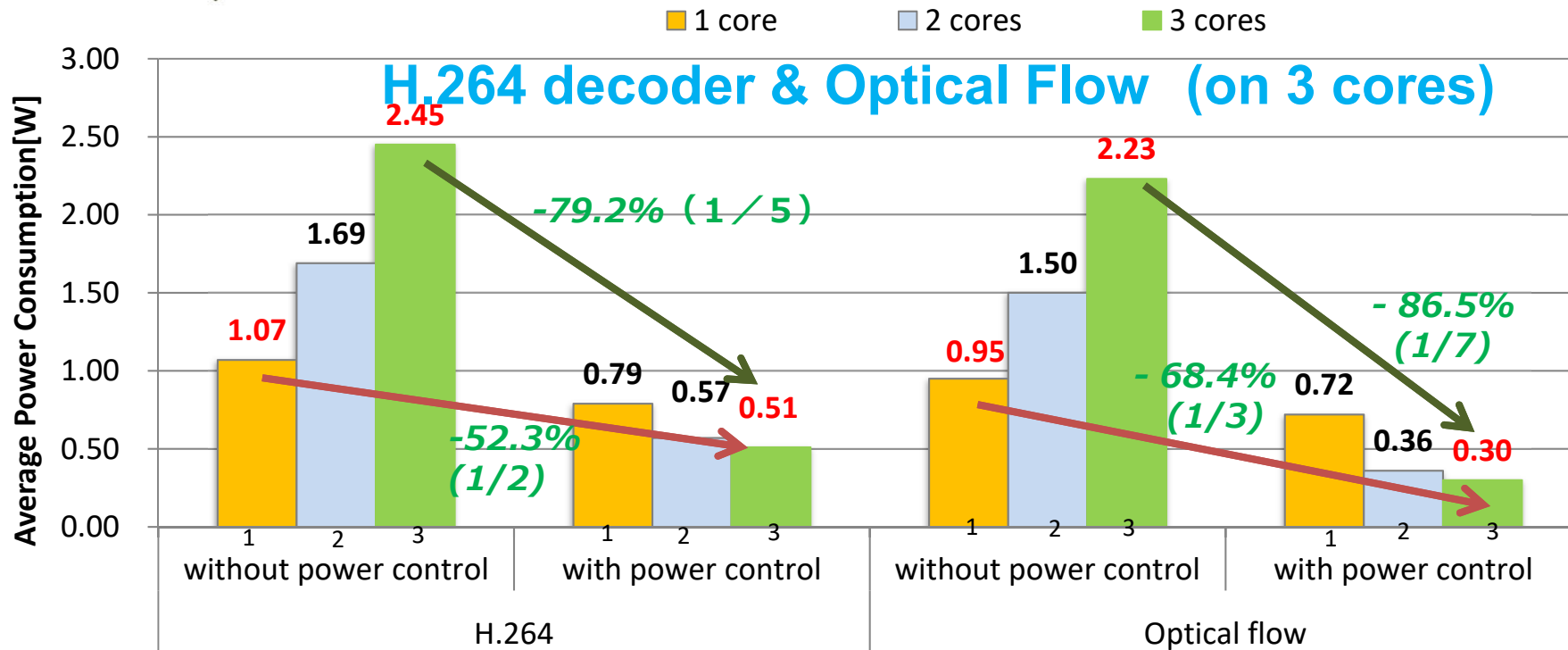


Automatic Power Reduction on ARM CortexA9 with Android

http://www.youtube.com/channel/UCS43INYEIkC8i_KIgfZYQBQ

ODROID X2

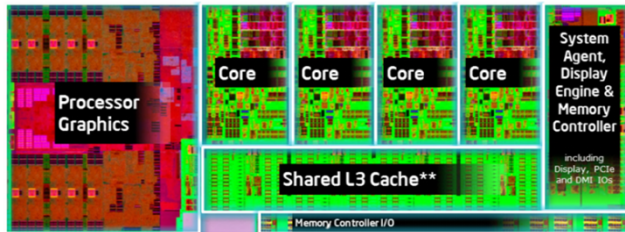
Samsung Exynos4412 Prime, ARM Cortex-A9 Quad core
1.7GHz~0.2GHz, used by Samsung's Galaxy S3



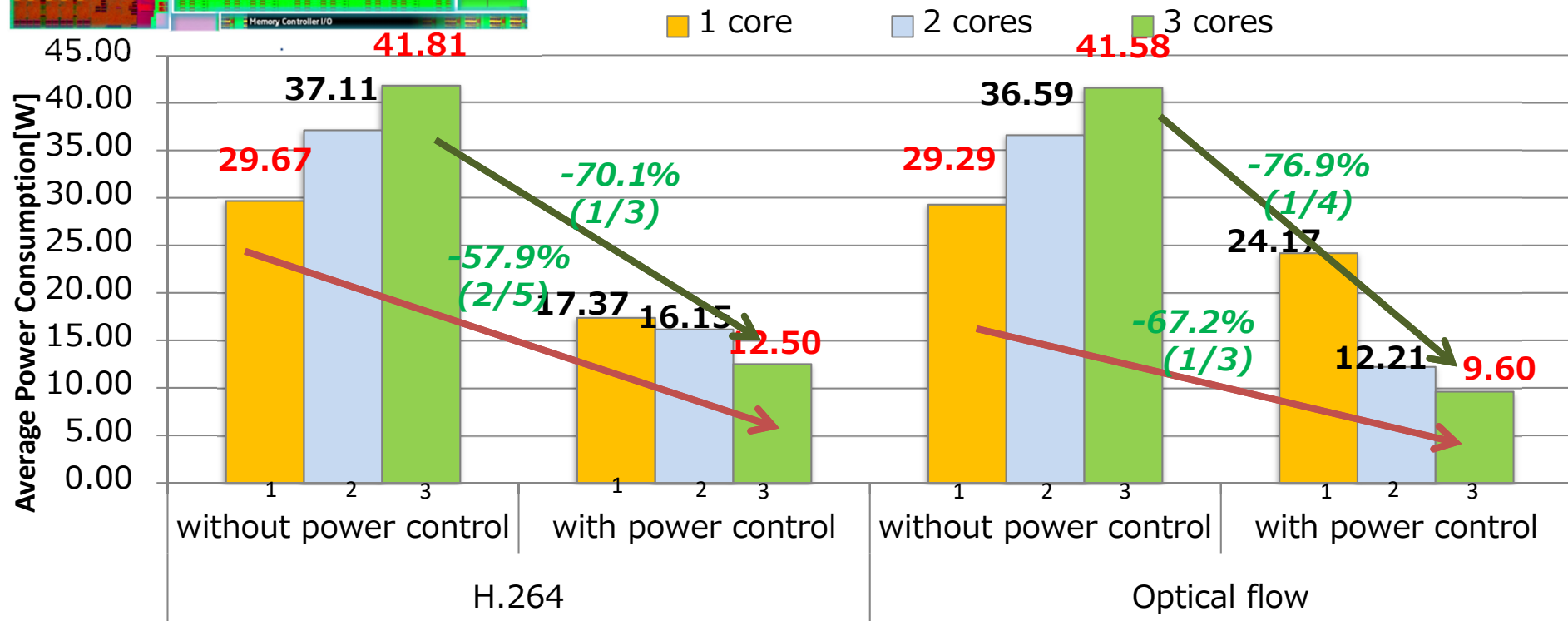
Power for 3cores was reduced to **1/5~1/7** against without software power control
Power for 3cores was reduced to **1/2~1/3** against ordinary 1core execution

Automatic Power Reuction on Intel Haswell

H.264 decoder & Optical Flow (3cores)



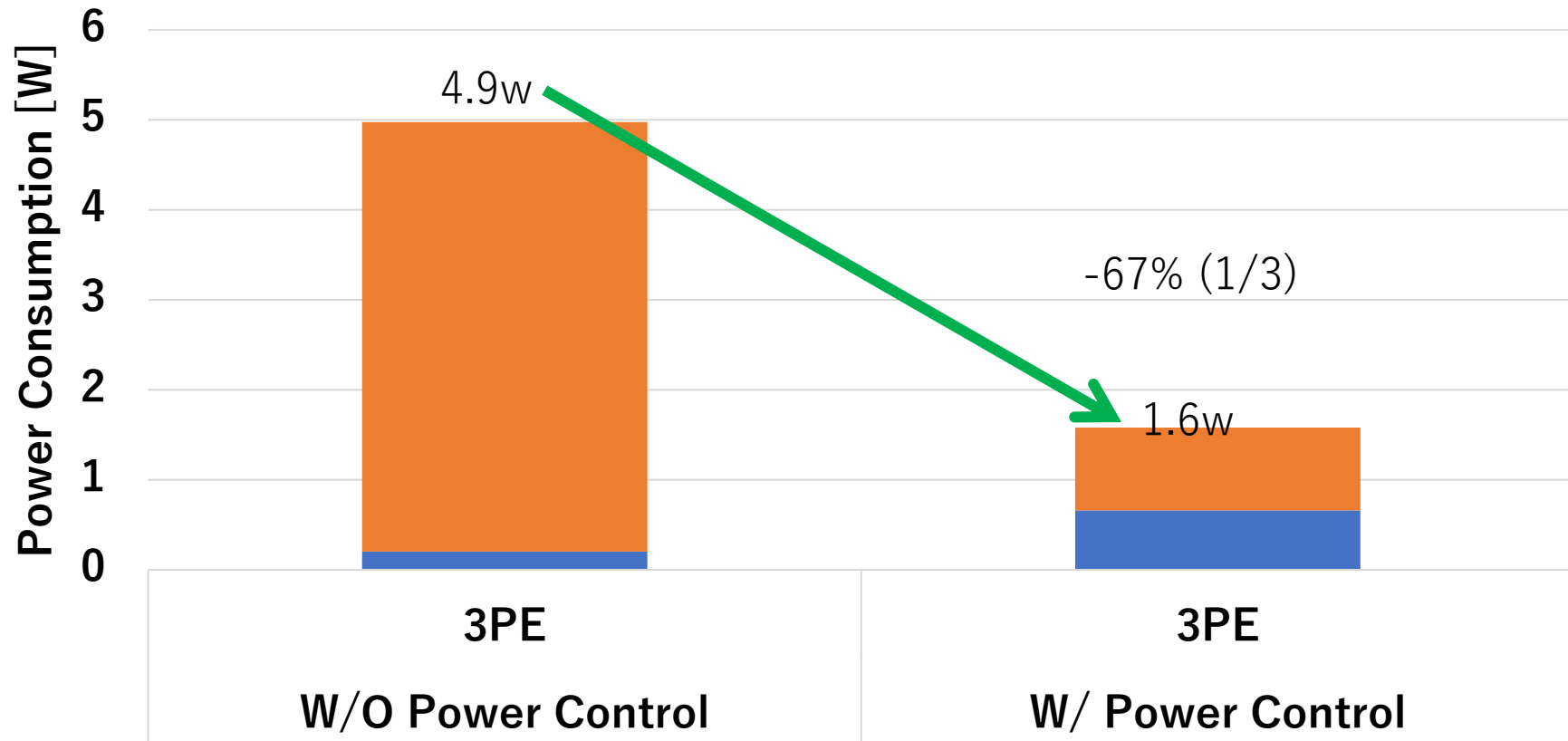
H81M-A, Intel Core i7 4770k
Quad core, 3.5GHz~0.8GHz



Power for 3cores was reduced to **1/3~1/4** against without software power control

Power for 3cores was reduced to **2/5~1/3** against ordinary 1core execution

Automatic Power Reduction of OpenCV Face Detection on big.LITTLE ARM Processor



- **ODROID-XU3**

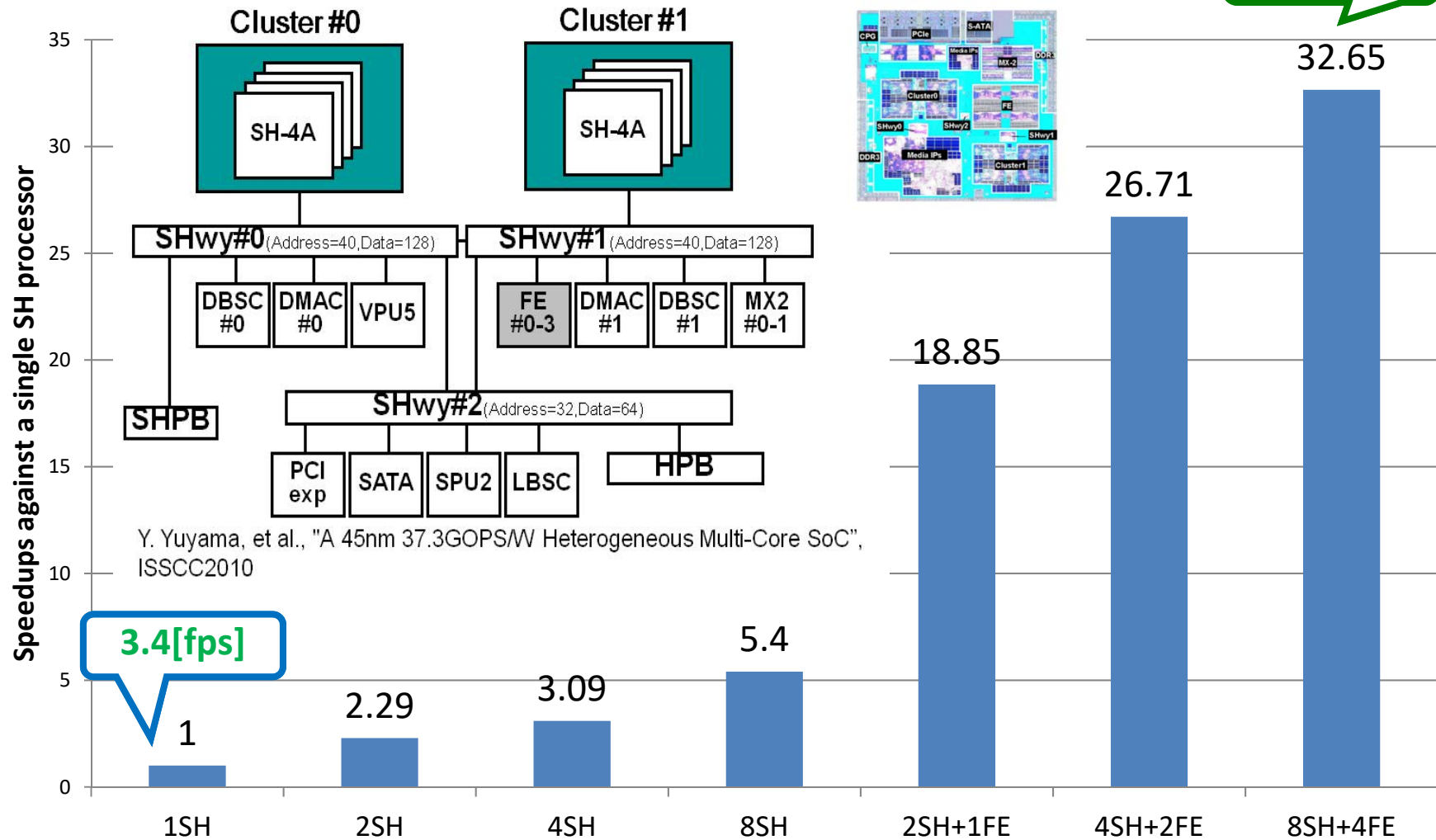
■ Cortex-A7 ■ Cortex-A15

- **Samsung Exynos 5422 Processor**

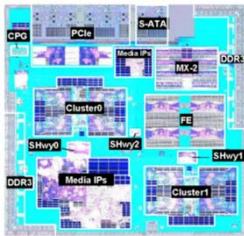
- 4x Cortex-A15 2.0GHz, 4x Cortex-A7 1.4GHz big.LITTLE Architecture
- 2GB LPDDR3 RAM
- Frequency can be changed by each cluster unit

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

111[fps]



3.4[fps]



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

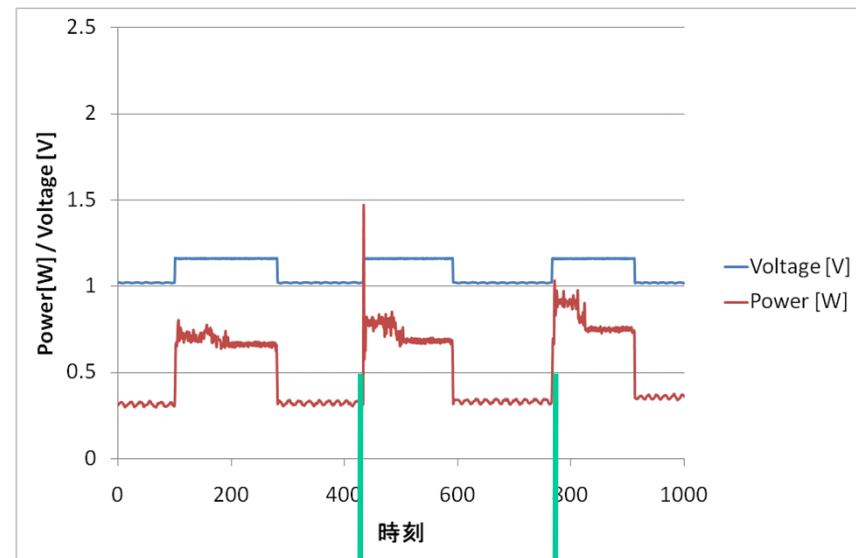
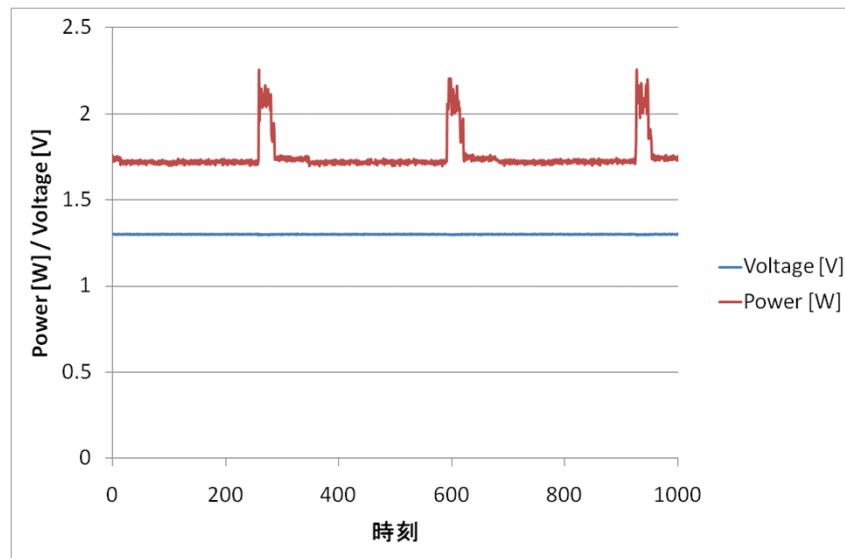
Without Power Reduction

With Power Reduction by OSCAR Compiler
70% of power reduction

Average: 1.76[W]

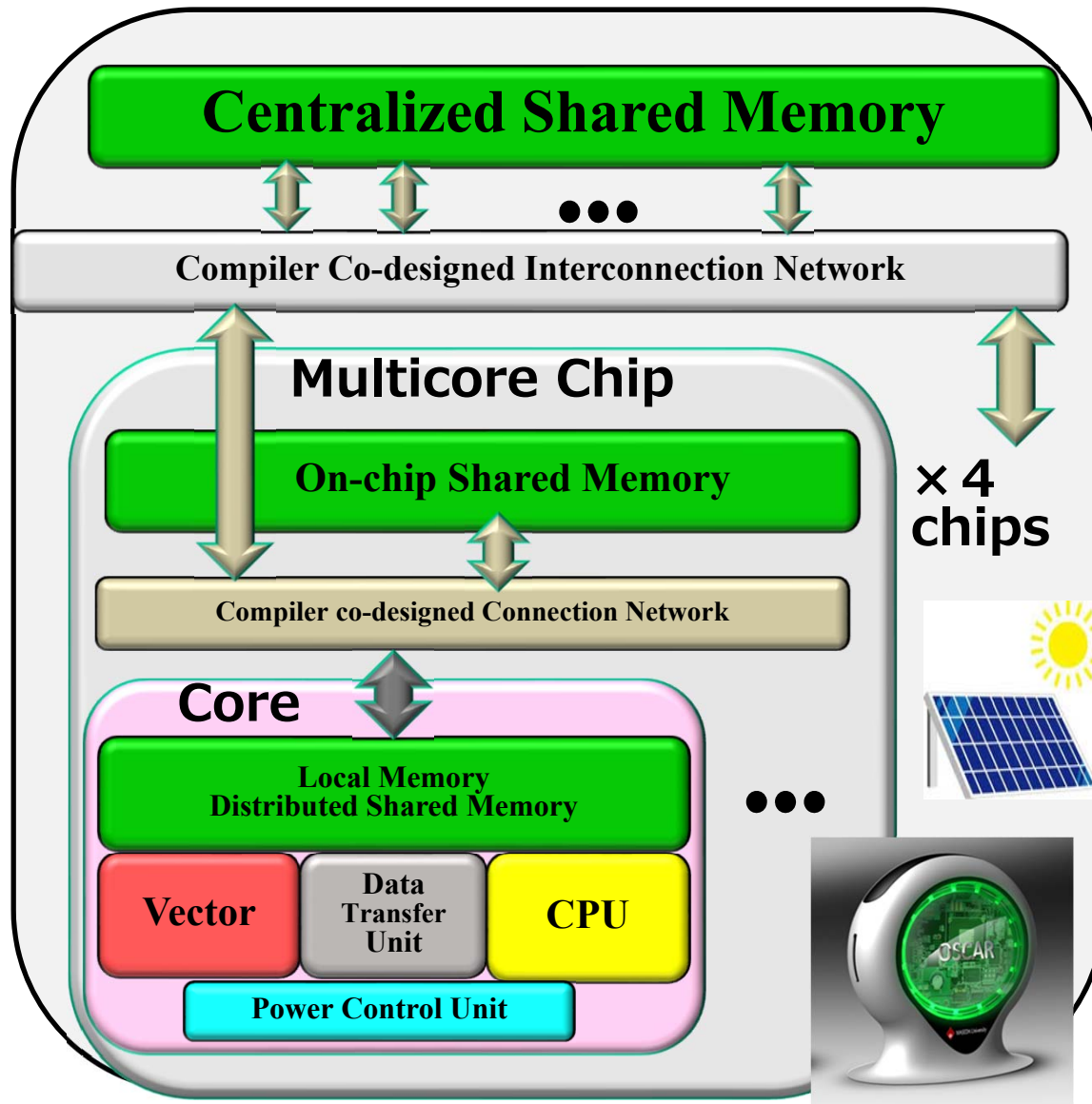


Average: 0.54[W]



**1cycle : 33[ms]
→30[fps]**

OSCAR Vector Multicore and Compiler for Embedded to Servers with OSCAR Technology



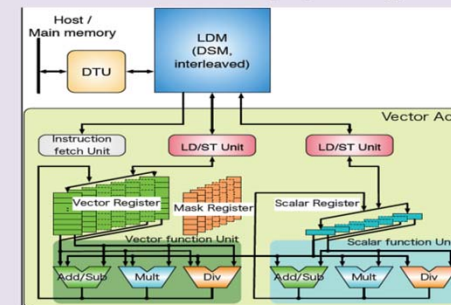
Target:

- Solar Powered
- 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



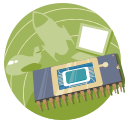
Function Units [tentative]

- 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





Future Multicore Products with Automatic Parallelizing Compiler



Next Generation Automobiles

- Safer, more comfortable, energy efficient, environment friendly
- Cameras, radar, car2car communication, internet information integrated brake, steering, engine, motor control

Smart phones



- From everyday recharging to less than once a week
- Solar powered operation in emergency condition
- Keep health

Advanced medical systems



Cancer treatment, Drinkable inner camera

- Emergency solar powered
- No cooling fan, No dust, clean usable inside OP room



Personal / Regional Supercomputers



Solar powered with more than 100 times power efficient : FLOPS/W

- Regional Disaster Simulators saving lives from tornadoes, localized heavy rain, fires with earth quakes