Software and Hardware for High Performance and Low Power Homogeneous and Heterogeneous Multicore Systems

Hironori Kasahara

Professor, Dept. of Computer Science & Engineering Director, Advanced Multicore Processor Research Institute Waseda University, Tokyo, Japan IEEE Computer Society President Elect 2017 URL: http://www.kasahara.cs.waseda.ac.jp/

Waseda Univ. GCSC

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 ∝ Frequency * Voltage² (Voltage ∝ Frequency) Power ∝ Frequency³

If Frequency is reduced to 1/4 (Ex. 4GHz→1GHz), Power is reduced to 1/64 and Performance falls down to 1/4 . <<u>Multicores</u>>

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



> 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

Trend of Peak Performances of Supercomputers





Compiler Co-designed Multicore RP2



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

Core#0	Core#1		Process Technology	90nm, 8-layer, triple- Vth, CMOS	
Core#2	Core#3		Chip Size	104.8mm ² (10.61mm x 9.88mm)	
Core#6	Core#7	VSWC	CPU Core Size	6.6mm ² (3.36mm x 1.96mm)	
Core#4	Core#5		Supply Voltage	1.0V–1.4V (internal), 1.8/3.3V (I/O)	
DDRPAD	GCPG		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"



Cancer Treatment Carbon Ion Radiotherapy

(Previous best was 2.5 times speedup on 16 processors with hand optimization)



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)

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.



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



Compi

(*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

Generation of Coarse Grain Tasks Macro-tasks (MTs)

- Block of Pseudo Assignments (BPA): Basic Block (BB)
- Repetition Block (RB) : natural loop

Subroutine Block (SB): subroutine



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



PRIORITY DETERMINATION IN DYNAMIC CP METHOD



Earliest Executable Conditions

EEC: Control dependence + Data Dependence Control dependences show executions of MTs are decided Data dependences show data accessed by MTs are ready MT2 may start execution after MT1 branches to MT2 and MT1 finish execution.

Macrotask No.	Earliest Executable Condition			
1				
2 MT3 may start exec	ution 1 2			
3 after MT1 branches	$\begin{array}{c} \text{to M13.} \\ \hline 2 \ 4 \ \text{OR} \ (1) \ 3 \end{array}$			
5 MT6 may start execu	tion (4) 5 AND [2 4 OR (1) 3]			
6after MT3 finish exect7MT2 branches to MT	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$			
8	(2) 4 OR (1) 3			
9	(8) 9			
10	(8) 10			
11	89 OR 810			
12	11 12 AND [9 OR (8) 10]			
13	11 13 OR 11 12			
14	(8) 9 OR (8) 10			
15	2 15			

V1.0 © 2014, IEEE All rights reserved

Automatic processor assignment in 103.su2cor

• Using 14 processors

Coarse grain parallelization within DO400



MTG of Su2cor-LOOPS-DO400 • Coarse grain parallelism PARA_ALD = 4.3



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



Inter-loop data dependence analysis in TLG

- Define exit-RB in TLG as Standard-Loop
- Find iterations on which a iteration of Standard-Loop is data dependent
 - e.g. K_{th} of RB3 is data-dep on K-1_{th},K_{th} of RB2, on K-1_{th},K_{th},K+1_{th} of RB1 indirectly.



Example of TLG

Decomposition of RBs in TLG

- Decompose GCIR into $DGCIR^p(1 \le p \le n)$
 - n: (multiple) num of PCs, DGCIR: Decomposed GCIR
- Generate CAR on which DGCIR^p&DGCIR^{p+1} are data-dep.
- Generate LR on which DGCIR^p is data-dep.





An Example of Data Localization for Spec95 Swim



(a) An example of target loop group for data localization

Data Layout for Removing Line Conflict Misses by Array Dimension Padding Declaration part of arrays in spec95 swim

before padding

after padding



Statement Level Near Fine Grain Task



Task Graph for FPPPP

Statement level near fine grain parallelism



Elimination of Redundant Synchronization for Shared Data on Centralized Shared Memory after Static Task Scheduling



Generated Parallel Machine Code for Near Fine Grain Parallel Processing



【W-CDMA Base Band Communication】 Near Fine Grain Parallel Processing of EAICH Detection Program on RP2 Multicore with 4 SH4A cores

- Hadamard transform often used in the signal processing
- Parallel Processing Method
 - Near fine grain parallel processing among statements
 - Static Scheduling



1.62 times speedup for 2cores, 3.45 times speedup for 4 cores for EAICH on RP2.



Code Generation Using OpenMP

- Compiler generates a parallelized program using **OpenMP API**
- One time single level thread generation
 - Threads are forked only once at the beginning of a program by OpenMP "PARALLEL SECTIONS" directive
 - Forked threads join only once at the end of program
- Compiler generates codes for each threads using static or dynamic scheduling schemes
- Extension of OpenMP for hierarchical processing is not required

Multicore Program Development Using OSCAR API V2.0



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.

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



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





V1.0 © 2014, IEEE All rights reserved

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





Model Base Designed Engine Control on V850 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.



OSCAR Compile Flow for Simulink Applications



Speedups of MATLAB/Simulink Image Processing on Various 4core Multicores

(Intel Xeon, ARM Cortex A15 and Renesas SH4A)



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

Buoy Detection : http://www.mathworks.co.jp/matlabcentral/fileexchange/44706-buoy-detection-using-simulink

Color Edge Detection : <u>http://www.mathworks.co.jp/matlabcentral/fileexchange/28114-fast-edges-of-a-color-image--actual-color--not-converting-to-grayscale-/</u>

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

Parallel Processing on Simulink Model

• The parallelized C code can be embedded to Simulink using C mex API for HILS and SILS implementation.



Call parallelized C code from the S-Function block

OSCAR API Ver. 2.0 for Homogeneous/Heterogeneous Multicores and Manycores

List of Directives (22 directives)

- Parallel Execution API
 - parallel sections (*)
 - flush (*)
 - critical (*)
 - execution
- Memoay Mapping API
 - threadprivate (*)
 - distributedshared
 - onchipshared
- Synchronization API
 - groupbarrier
- Data Transfer API
 - dma_transfer
 - dma_contiguous_parameter
 - dma_stride_parameter
 - dma_flag_check
 - dma_flag_send

(* from OpenMP)

- Power Control API
 - fvcontrol
 - get_fvstatus
- Timer API
 - get_current_time
- Accelerator
 - accelerator_task_entry
- Cache Control
 - cache_writeback
 - cache_selfinvalidate
 - complete_memop
 - noncacheable
 - aligncache
 - 2 hint directives for OSCAR compiler
 - accelerator_task
 - oscar_comment

from V2.0

An Image of Static Schedule for Heterogeneous Multicore with Data Transfer Overlapping and Power Control





Power Reduction by Power Supply, Clock Frequency and Voltage

Control by OSCAR Compiler Frequency and Voltage (DVFS), Clock and Power gating of each cores are scheduled considering the task schedule since the dynamic power proportional to the cube of F (F^3) and the leakage power (the static power) can be reduced by the power gating (power off).



An Example of Machine Parameters for the Power Saving Scheme

- Frequency of each proc. is changed to several levels
- Voltage is changed together with frequency
- Each proc. can be powered on/off

state	FULL	MID	LOW	OFF
frequency	1	1/2	1/4	0
voltage	1	0.87	0.71	0
dynamic energy	1	3/4	1 / 2	0
static power	1	1	1	0

State transition overhead

state	FULL	MID	LOW	OFF	state	FULL	MID	LOW	OFF
FULL	0	40k	40k	80k	FULL	0	20	20	40
MID	40k	0	40k	80k	MID	20	0	20	40
LOW	40k	40k	0	80k	LOW	20	20	0	40
OFF	80k	80k	80k	0	OFF	40	40	40	0
delay time [u.t.]				energy overhead [µJ]					

Power Reduction Scheduling



46

Low-Power Optimization with OSCAR API



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



Automatic Power Reduction for MPEG2 Decode on Android Multicore ODROID X2 ARM Cortex-A94 cores http://www.youtube.com/channel/UCS43INYEIkC8i_KIgFZYQBQ



• 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.



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).

Automatic Parallelization of JPEG-XR for
Drinkable Inner Camera (Endo Capsule)
10 times more speedup needed after parallelization for 128 cores of
Power 7. Less than 35mW power consumption is required.





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



Solar Powered with compiler power reduction. Fully automatic parallelization and vectorization including local memory management

and data transfer.

Fujitsu VPP500/NWT: PE Unit



Performance of OSCAR Compiler Software Coherence Control

- Faster or Equal Processing Performance up to 4cores with hardware coherent mechanism on RP2.
- Software Coherence gives us correct execution without hardware coherence mechanism on 8 cores.



Automatic Local Memory Management 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

DLG0

CAR

DOI=34.35

DOI=34.34

LR

DO I=1.33

DO I=1.33

DO I=2.34

DO I=1.101 A(I)=2*I

ENDDO

DO I=1,100

ENDDO

DO I=2.100

ENDDO

C(I)=B(I)*B(I-1)

B(I)=B(I-1) +A(I)+A(I+1)





Adjustable Blocks

- Handling a suitable block size for each application
 - different from a fixed block size in cache



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, multidimensional 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']



8 Core RP2 Chip Block Diagram



Speedups by the Local Memory Management Compared with Utilizing Shared Memory on Benchmarks Application using RP2



20.12 times speedup for 8cores execution using local memory against sequential execution using off-chip shared memory of RP2 for the AACenc

Software Coherence Control Method on OSCAR Parallelizing Compiler

- Coarse grain task parallelization with earliest condition analysis (control and data dependency analysis to detect parallelism among coarse grain tasks).
- SCAR compiler automatically controls coherence using following simple program restructuring methods:
 - > To cope with stale data problems:

Data synchronization by compilers

- > To cope with false sharing problem:
 - Data Alignment

Array Padding

Non-cacheable Buffer



MTG generated by earliest executable condition analysis

Automatic Software Coherent Control for Manycores Performance of Software Coherence Control by

OSCAR Compiler on 8-core RP2



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





Future Multicore Products



Next Generation Automobiles

- Safer, more comfortable, energy efficient, environment friendly

- Cameras, radar, car2car communication, internet information integrated brake, steering, engine, moter control

Smart phones

endes es X)



-From everyday recharging to less than once a week

- Solar powered operation in emergency condition

- Keep health

Advanced medical systems

Personal / Regional Supercomputers



Cancer treatment, Drinkable inner camera

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



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

Summary

- To get speedup and power reduction on homogeneous and heterogeneous multicore systems, collaboration of architecture and compiler will be more important.
- 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", and "Wireless communication Base Band Processing" on various multicores.
 - For example, the automatic parallelization gave us 110 times speedup for "Earthquake Wave Propagation Simulation" on 128 cores of IBM Power 7 against 1 core, 327 times speedup for "Heavy Particle Radiotherapy Cancer Treatment" on 144cores Hitachi Blade Server using Intel Xeon E7-8890, 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.
- 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.
- For more speedup and power reduction, we have been developing a new architecture/compiler co-designed multicore with vector accelerator based on vector pipelining with vector registers, chaining, load-store pipeline, advanced DMA controller without need of modification of CPU instruction set.