

1 Outline

1.1 Supercomputer System

Before SY 2015, the ISSP supercomputer center provided users with three supercomputer systems: NEC-SX9 (System A), SGI Altix ICE 8400EX (System B), and FUJITSU PRIMEHPC FX10 (System C). At the beginning of SY 2015, Systems A and B were replaced by SGI ICE XA/UV hybrid system, and this system will be called System B for the next five years. The (new) System B and System C are installed in the main building of ISSP.

System B - SGI ICE XA/UV hybrid system is a massively-parallel supercomputer with three types of compute nodes: 19 “Fat” nodes, 1584 “CPU” nodes, and 288 “ACC” nodes. “Fat” nodes are each comprised of four Intel Xeon E5-4627v3 CPUs (10 cores/CPU) and 1 TB of memory. “CPU” nodes have two Intel Xeon E5-2680v3 CPUs (12 cores/CPU) and 128 GB of memory. “ACC” nodes have two nVIDIA Tesla K40 GPUs in addition to two Xeon E5-2680v3 CPUs and 128 GB of memory. System B achieves 2.6 PFLOPS in theoretical peak performance with high power efficiency. The subsystem comprised of only CPU nodes ranks 61st on the November 2015 Top 500 List, which is a ranking based on total performance measured by the HPL benchmark. The subsystem of ACC nodes ranks 104th on the Top 500 List, and it also ranks 23rd on the Green 500 List, which is a ranking based on performance per watt of electrical power consumption. The compute nodes communicate to each other through FDR Infiniband. The Fat nodes are interconnected in fat tree topology, while the CPU and ACC nodes are connected in enhanced hypercube topology. System B entered official operation on Aug. 21, 2015.

System C - FUJITSU PRIMEHPC FX10 has been in service since April, 2013. It is highly compatible with K computer, the largest supercomputer in Japan. System C consists of 384 nodes, and each node has 1 SPARC64TM IXfx CPU (16 cores) and 32 GB of memory. The total system achieves 90.8 TFlops theoretical peak performance.

For further details, please contact ISSP Supercomputer Center (SCC-ISSP).

[Correspondence: center@issp.u-tokyo.ac.jp]

1.2 Project Proposals

The ISSP supercomputer system provides computation resources for scientists working on condensed matter sciences in Japan. All scientific staff members (including post-docs) at universities or public research institutes in Japan can submit proposals for projects related to research activities on materials and condensed matter sciences. These proposals are peer-reviewed by the Advisory Committee members (see Sec. 1.3), and then the computation resources are allocated based on the review reports. The leader of an approved project can set up user accounts for collaborators. Other types of scientists, including graduate students, may also



Figure 1: Supercomputer System at the SCC-ISSP

be added. Proposal submissions, peer-review processes, and user registration are all managed via a web system.

The computation resources are distributed in a unit called “point”, determined as a function of available CPU utilization time and consumed disk resources. There were calls for six classes of research projects in SY 2015. The number of projects and the total number of points that were applied for and approved in this school year are listed in Table 1.

In addition, from SY 2010, ISSP Supercomputer is providing 20% of its computational resources for Computational Materials Science Initiative (CMSI), which aims at advancing parallel computations in condensed matter, molecular, and materials sciences on the 10-PFlops K Computer. The points for projects run by CMSI are distributed in accord with this policy. Computer time has also been allotted to Computational Materials Design (CMD) workshops run by CMSI, as well as for Science Camps held in ISSP for undergraduate students.

1.3 Committees

In order to fairly manage the projects and to smoothly determine the system operation policies, the Materials Design and Characterization Laboratory (MDCL) of the ISSP has organized the Steering Committee of the MDCL and the Steering

Table 1: Classes of research projects in SY 2015

Class	Maximum Points		Application	# of Proj.	Total points			
	Sys-B	Sys-C			Applied		Approved	
					Sys-B	Sys-C	Sys-B	Sys-C
A	100	100	any time	13	1.3k	1.3k	1.3k	1.3k
B	1k	500	twice a year	57	43.8k	8.6k	26k	7.8k
C	10k	2.5k	twice a year	162	1394.8k	181.3k	475k	147.8k
D	10k	2.5k	any time	2	18k	4k	14k	2.5k
E	30k	2.5k	twice a year	5	150k	7.5k	79k	6.2k
S	–	–	twice a year	0	0	0	0	0
CMSI				33	125k	139k	125k	139k

- Class A is for trial use by new users; proposals for Class A projects are accepted throughout the year.
- Proposals for projects in Classes B (small), C (mid-size), E (large-scale), and S (exceptional) can be submitted twice a year. Approved projects in Classes A, B, C, and E continue to the end of the school year.
- In Class D, projects can be proposed on rapidly-developing studies that need to perform urgent and relatively large calculations. An approved project continues for 6 months from its approval.
- Class S is for projects that are considered extremely important for the field of condensed matter physics and requires extremely large-scale computation. The project may be carried out either by one research group or cooperatively by several investigators at different institutions. A project of this class should be applied with at least 10,000 points; there is no maximum. We require group leaders applying for Class S to give a presentation on the proposal to the Steering Committee of the SCC-ISSP. Class S projects are carried out within one year from its approval.
- Project leaders can apply for points so that the points for each system do not exceed the maximum point shown in this table.

Committee of the SCC-ISSP, under which the Supercomputer Project Advisory Committee (SPAC) is formed to review proposals. The members of the committees in SY 2015 were as follows:

Steering Committee of the MDCL

HIROI, Zenji	ISSP (Chair person)
KATO, Takeo	ISSP
KAWASHIMA, Naoki	ISSP
MORI, Hatsumi	ISSP
NAKATSUJI, Satoru	ISSP
NOGUCHI, Hiroshi	ISSP
SUGINO, Osamu	ISSP
SUEMOTO, Toru	ISSP
KIMURA, Kaoru	Univ. of Tokyo
YOSHIMOTO, Yoshihide	Univ. of Tokyo
HASEGAWA, Tadashi	Nagoya Univ.
MIYASAKA, Hitoshi	Tohoku Univ.
MORIKAWA, Yoshitada	Osaka Univ.
OKAMOTO, Yuko	Nagoya Univ.
OTSUKI, Tomi	Sophia Univ.
TAKEDA Mahoto	Yokohama Natl. Univ.

Steering Committee of the SCC-ISSP

NOGUCHI, Hiroshi	ISSP (Chair person)
HARADA, Yoshihisa	ISSP
KAWASHIMA, Naoki	ISSP
SUGINO, Osamu	ISSP
TAKADA, Yasutami	ISSP
TSUNETSUGU, Hirokazu	ISSP
KASAMATSU, Shusuke	ISSP
MORITA, Satoshi	ISSP
SHIBA, Hayato	ISSP
WATANABE, Hiroshi	ISSP
HATANO, Naomichi	Univ. of Tokyo
IMADA, Masatoshi	Univ. of Tokyo
NAKAJIMA, Kengo	Univ. of Tokyo
TSUNEYUKI, Shinji	Univ. of Tokyo
YOSHIMOTO, Yoshihide	Univ. of Tokyo
MOHRI, Tetsuo	Tohoku Univ.
MORIKAWA, Yoshitada	Osaka Univ.
ODA, Tatsuki	Kanazawa Univ.
OKAMOTO, Yuko	Nagoya Univ.
OTSUKI, Tomi	Sophia Univ.

SUZUKI, Takafumi	Univ. of Hyogo
YATA, Hiroyuki	ISSP
FUKUDA, Takaki	ISSP

Supercomputer Project Advisory Committee

NOGUCHI, Hiroshi	ISSP (Chair person)
HARADA, Yoshihisa	ISSP
KATO, Takeo	ISSP
KAWASHIMA, Naoki	ISSP
OZAKI, Taisuke	ISSP
SUGINO, Osamu	ISSP
TAKADA, Yasutami	ISSP
TSUNETSU, Hirokazu	ISSP
KASAMATSU, Shusuke	ISSP
MORITA, Satoshi	ISSP
SHIBA, Hayato	ISSP
WATANABE, Hiroshi	ISSP
AOKI, Hideo	Univ. of Tokyo
ARITA, Ryotaro	Univ. of Tokyo
HATANO, Naomichi	Univ. of Tokyo
HUKUSHIMA, Koji	Univ. of Tokyo
IKUHARA, Yuichi	Univ. of Tokyo
IMADA, Masatoshi	Univ. of Tokyo
IWATA, Jun-Ichi	Univ. of Tokyo
MIYASHITA, Seiji	Univ. of Tokyo
MOTOME, Yukitoshi	Univ. of Tokyo
NAKAJIMA, Kengo	Univ. of Tokyo
OGATA, Masao	Univ. of Tokyo
OSHIYAMA, Atsushi	Univ. of Tokyo
TODD, Synge	Univ. of Tokyo
TSUNEYUKI, Shinji	Univ. of Tokyo
WATANABE, Satoshi	Univ. of Tokyo
YOSHIMOTO, Yoshihide	Univ. of Tokyo
NEMOTO, Koji	Hokkaido Univ.
AKAGI, Kazuto	Tohoku Univ.
KAWAKATSU, Toshihiro	Tohoku Univ.
KURAMOTO, Yoshio	Tohoku Univ.
MOHRI, Tetsuo	Tohoku Univ.
SHIBATA, Naokazu	Tohoku Univ.
KIM, Kang	Niigata Univ.
ISHIBASHI, Shoji	AIST
MIYAMOTO, Yoshiyuki	AIST
OTANI, Minoru	AIST
KOBAYASHI, Kazuaki	NIMS

TATEYAMA, Yoshitaka	NIMS
HATSUGAI, Yasuhiro	Univ. of Tsukuba
KOBAYASHI, Nobuhiko	Univ. of Tsukuba
OKADA, Susumu	Univ. of Tsukuba
YABANA, Kazuhiro	Univ. of Tsukuba
ODA, Tatsuki	Kanazawa Univ.
SAITO, Mineo	Kanazawa Univ.
HIDA, Kazuo	Saitama Univ.
NAKAYAMA, Takashi	Chiba Univ.
FURUKAWA, Nobuo	Aoyama Gakuin Univ.
MATSUKAWA, Hiroshi	Aoyama Gakuin Univ.
TAKANO, Hiroshi	Keio Univ.
YAMAUCHI, Jun	Keio Univ.
YASUOKA, Kenji	Keio Univ.
TOMITA, Yusuke	Shibaura Inst. Tech.
OTSUKI, Tomi	Sophia Univ.
OBATA, Shuji	Tokyo Denki Univ.
ANDO, Tsuneya	Tokyo Tech.
TADA, Tomofumi	Tokyo Tech.
HOTTA, Takashi	Tokyo Metropolitan Univ.
TOHYAMA, Takami	Tokyo Univ. of Sci.
WATANABE, Kazuyuki	Tokyo Univ. of Sci.
HAGITA, Katsumi	National Defense Academy
KONTANI, Hiroshi	Nagoya Univ.
OKAMOTO, Yuko	Nagoya Univ.
SHIRAISHI, Kenji	Nagoya Univ.
TANAKA, Yukio	Nagoya Univ.
KAWAKAMI, Norio	Kyoto Univ.
MASUBUCHI, Yuichi	Kyoto Univ.
YAMAMOTO, Ryoichi	Kyoto Univ.
YANASE, Yoichi	Kyoto Univ.
KAWAMURA, Hikaru	Osaka Univ.
KUROKI, Kazuhiko	Osaka Univ.
KUSAKABE, Koichi	Osaka Univ.
MORIKAWA, Yoshitada	Osaka Univ.
OGUCHI, Tamio	Osaka Univ.
SHIRAI, Koun	Osaka Univ.
YOSHIDA, Hiroshi	Osaka Univ.
YUKAWA, Satoshi	Osaka Univ.
HARIMA, Hisatomo	Kobe Univ.
SAKAI, Toru	Univ. of Hyogo
SUGA, Seiichiro	Univ. of Hyogo
SUZUKI, Takafumi	Univ. of Hyogo
TATENO, Masaru	Univ. of Hyogo
HOSHI, Takeo	Tottori Univ.

YASUDA, Chitoshi

Univ. of the Ryukyus

1.4 Staff

The following staff members of the SCC-ISSP usually administrate the ISSP Supercomputer.

NOGUCHI, Hiroshi	Associate Professor (Chair person)
KAWASHIMA, Naoki	Professor
SUGINO, Osamu	Associate Professor
WATANABE, Hiroshi	Research Associate
KASAMATSU, Shusuke	Research Associate
NOGUCHI, Yoshifumi	Research Associate
SHIBA, Hayato	Research Associate
MORITA, Satoshi	Research Associate
YATA, Hiroyuki	Technical Associate
FUKUDA, Takaki	Technical Associate
ARAKI, Shigeyuki	Technical Associate

2 Statistics (School Year 2015)

2.1 System and User Statistics

In the following, we present statistics for operation time taken in the period from April 2015 to March 2016 (SY 2015). In Table 2, we show general statistics of the supercomputer system in SY 2015. The total number of compute nodes in System B, and C is 1891 and 384 respectively. Consumed disk points amount to about 1% and 7% of the total consumed points in System B and C respectively. Roughly 20% of the total points in System B 60% of that in System C were consumed by CMSI projects. This means that about 20% of the total computational resources in this school year were actually used by CMSI projects.

In the left column of Fig. 2, availabilities, utilization rates, and consumed points in each system are plotted for each month. Throughout the school year, the utilization rates were very high. Especially in System B, they were exceeding 90% throughout most of the year.

The user statistics are shown in the right column of Fig. 2. The horizontal axis shows the rank of the user/group arranged in the descending order of the execution time (hour \times nodes). The execution time of the user/group of the first rank is the longest. The vertical axis shows the sum of the execution time up to the rank. From the saturation points of the graphs, the number of “active” users of each system is around 250, and 70 for System B and C respectively. The maximum ranks in the graphs correspond to the number of the users/groups that submitted at least one job.

Table 2: Overall statistics of SY 2015

	System-B	System-C
total service time ($\times 10^3$ node·hours)	9728.4	3363.8
number of executed jobs	183721	22278
total consumed points ($\times 10^3$ point)	337.1	113.84
CPU points ($\times 10^3$ point)	334.6	106.0
disk points ($\times 10^3$ point)	2.5	7.9
points consumed by CMSI ($\times 10^3$ point)	63.3	59.5
total exec. time ($\times 10^3$ node·hours)	8803.1	2464.4
availability	96.0%	96.3%
utilization rate	90.5%	76.1%

2.2 Queue and Job Statistics

Queue structures of System B and C in SY 2015 are shown in Table 3. In System B, users can choose from three types of compute nodes; jobs submitted to queues with “cpu”, “acc”, and “fat” at the end of their queue names are submitted to CPU, ACC, and Fat nodes, respectively. See Sec. 1.1 for a description of each type of compute node. The user then has to choose the queue according to the number of nodes to use and the duration of their calculation jobs. Queue names starting with “F” are for jobs taking 24 hours or less, while those starting with “L” can run much longer up to 120 hours. More nodes are allotted to “F” queues in order to maximize the turnaround time of user jobs. The queue names starting with “i” are used for interactive debugging of user programs and the elapsed time limit is 30 minutes. The number following “F”, “L”, or “i” correspond to the number of nodes that can be used by one user job.

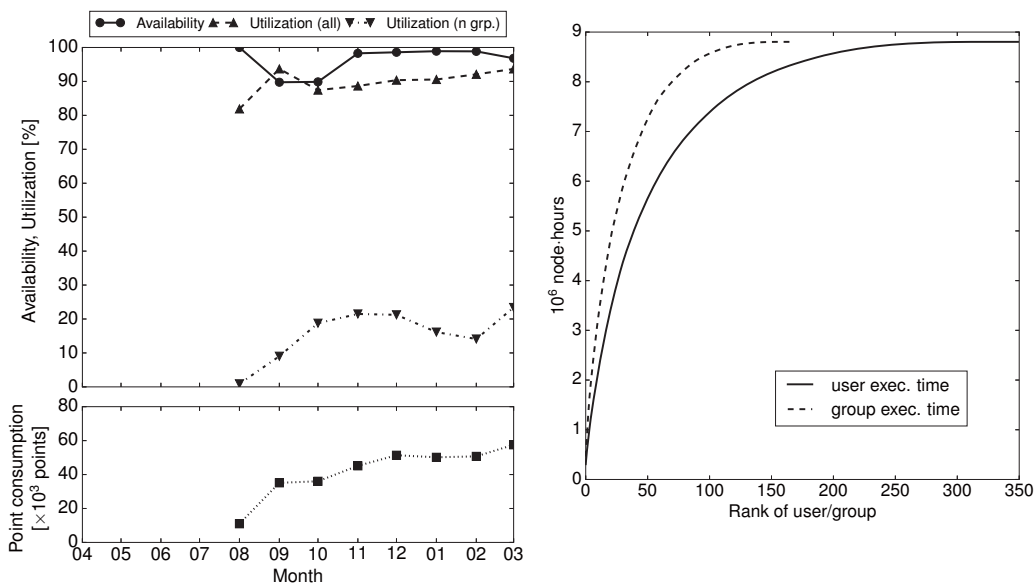
In System C, the “F” and “L” queues are set up similarly to System B. In addition, a debug queue is set up for short debugging jobs utilizing 1 to 4 CPUs, and an interactive queue that can use 1 to 4 CPUs is also available.

To prevent overuse of the storage, points are charged also for usage of disk quota in the three systems, as shown in Table 4. Disk points are revised often for optimal usage of the resources by examining usage tendencies each year.

Although we do not mention here in detail, to promote utilization of the massively parallel supercomputer, background queues (“B4cpu”, “B36cpu”, “B144cpu”, “B18acc”, “B72acc”, and “B2fat”) which charge no points for the jobs have also been open in System B.

The number of jobs, average waiting time, and total execution time in each queue are shown in Table 5. In both System B and C, a large portion of jobs have been executed in “F” queues. The largest amount of the execution time has been consumed in the large-scale “F144cpu” queue, but substantial number of jobs were run in every queue, suggesting that a wide variety of user needs are met by this queuing scheme. In most of these queues, the queue settings meet the user’s tendencies in that the waiting times are on the order of the elapsed-time limit.

System B



System C

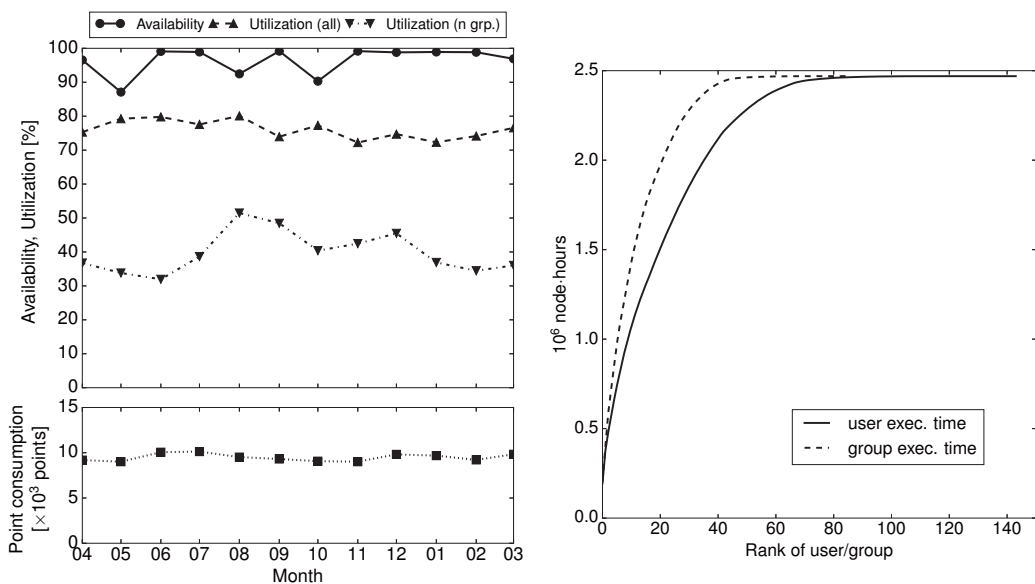


Figure 2: Left: Availabilities, utilization rates and point consumptions of each month during SY 2015. The utilization by CMSI projects (denoted by “n grp.”) is plotted in addition to the total utilization (denoted by “all”). Right: User statistics. The horizontal axis shows the rank of the user/group arranged in the descending order of the execution time (hour \times nodes). The vertical axis shows the sum of the execution time up to the rank.

Table 3: Queue structures in SY 2015

System-B					
queue name	Elapsed time limit (hr)	# of nodes /job	# of nodes /queue	Memory limit (GB)	job points /((node·day)
F4cpu	24	1-4	216	120/node	1
L4cpu	120	1-4	108	120/node	1
F36cpu	24	5-36	288	120/node	1
L36cpu	120	5-36	144	120/node	1
F144cpu	24	37-144	1008	120/node	1
L144cpu	120	37-144	144	120/node	1
i18cpu	0.5	1-18	72	120/node	1
F18acc	24	1-18	108	120/node	2
L18acc	120	1-18	54	120/node	2
F72acc	24	19-72	144	120/node	2
i9acc	0.5	1-9	36	120/node	2
F2fat	24	1-2	17	1000/node	4
L2 fat	120	1-2	6	1000/node	4
i1fat	0.5	1	2	1000/node	4

System-C				
queue name	Elapsed time limit (hr)	# of nodes /Job	# of nodes /queue	job points /((node·day)
debug	0.5	1-4	24	1
interactive	0.5	1-4	24	1
F12	24	2-12	60	1
F96	24	2-12	288	1
L12	120	24-96	24	1
L96	120	24-96	192	1

* In System C, the available memory size is limited to 28 GB per one CPU.

Table 4: Disk points of System B and C

		point/day
System B	/home	$0.001 \times \theta(q - 300)$
	/work	$0.0001 \times \theta(q - 3000)$
System C	/home	$0.05 \times \theta(q - 10)$
	/work	$0.005 \times \theta(q - 100)$

* q is denoted in unit of GB.

* $\theta(x)$ is equal to the Heaviside step function $H(x)$ multiplied by x , i.e., $xH(x)$.

Table 5: Number of jobs, average waiting time, total execution time, and average number of used nodes per job in each queue.

System-B				
queue	# of Jobs	Waiting Time (hour)	Exec. Time ($\times 10^3$ node-hour)	# of nodes
F4cpu	60680	20.14	667.69	1.92
L4cpu	5524	59.70	307.34	1.99
F36cpu	18095	31.29	955.91	10.24
L36cpu	2673	53.33	471.42	8.75
F144cpu	7965	15.28	3812.97	88.29
L144cpu	177	136.37	580.61	117.15
i18cpu	19750	0.54	33.35	6.66
F18acc	22329	14.14	288.57	2.09
L18acc	6097	39.16	111.13	1.41
F72acc	1617	5.69	188.79	43.99
i9acc	2436	0.58	2.68	4.14
F2fat	4243	8.53	24.22	1.11
L2fat	311	23.84	16.90	1.14
i1fat	1142	0.26	0.20	1.00

System-C				
queue	# of Jobs	Waiting Time (hour)	Exec. Time ($\times 10^3$ node-hour)	# of nodes
F12	8029	14.69	399.99	5.89
L12	64	85.11	10.72	5.59
F96	5180	43.50	2037.12	41.24
L96	6	555.42	15.15	96.00
debug	6966	0.16	1.87	2.61
interactive	1017	0.00	0.22	1.33

Table 6: List of supported software and project proposers for the GPGPU support service for SY 2015.

Software	Project Proposer	Affiliations
LargeScaleBoidsSimulator	Takashi Ikegami	The University of Tokyo
MDACP	Hiroshi Watanabe	The University of Tokyo
MODYLAS	Yoshimichi Andoh	Nagoya University
pTensor	Satoshi Morita	The University of Tokyo
RSCPMD	Yasuteru Shigeta	University of Tsukuba

The acc queues have relatively short waiting times, but we expect that to change as more users get accustomed to using GPGPUs.

2.3 GPGPU Support Service

As noted in Sec. 1.1, ACC nodes with graphics processing units (GPU) were introduced in System B in School Year 2015. Since GPUs were introduced in the ISSP Supercomputer center for the first time, many programs developed or utilized by users of this center have not been programmed for GPU computing. To help users take advantage of GPUs, the supercomputer center has started a service for porting users' materials science software to General Purpose GPUs (GPGPU)[1]. After a call for proposals (which will usually be in December), target programs for the next school year are selected by the Steering Committee of SCC. The porting service is carried out on each program for about two months; the coding is performed by engineers from the computer vender supplying the ISSP supercomputer system, and ISSP staff oversee the progress of the project and manage necessary communications with the proposer. Copyrights of the resulting software basically belong to the proposers, but the supported contents might be published under agreement with the proposer.

Since 2015 was the first year of the operation of System B with ACC nodes, the call for proposals was announced in May and target programs shown in Table 6 were selected in July. To determine the specific support contents and schedule, we first held kickoff meetings for each selected proposal from August to September.

In the actual service, analysis of the program was first performed to judge whether the target problems being solved can be accelerated using GPUs. After the analysis, the following support was carried out. For MDACP, MODY-LAS and RSCPMD, parallelization using OpenACC or CUDA was performed on hotspots. For pTensor, matrix-matrix multiplication, which originally used BLAS, was rewritten to use cuBLAS[2]. For LargeScaleBoidsSimulator, parallelization with multiple GPUs on a single process was extended to that with multiple GPUs on multiple processes. The contents and results of the porting service of MODY-LAS and pTensor were presented in ISSP joint-use supercomputer/CCMS symposium in 2016. Details on the porting service can be found on the ISSP supercomputer center web page[1].

References

- [1] <http://www.issp.u-tokyo.ac.jp/supercom/rsayh2/gpgpu>
- [2] cuBLAS, <http://docs.nvidia.com/cuda/cublas/>

Acknowledgments

The staffs would like to thank Prof. Takafumi Suzuki (now at University of Hyogo) for developing WWW-based system (SCM: SuperComputer Management System) for management of project proposals, peer-review reports by the SPAC committee, and user accounts. We also thank Ms. Reiko Yamaji for creating and maintaining a new WWW page of the ISSP Supercomputer Center.