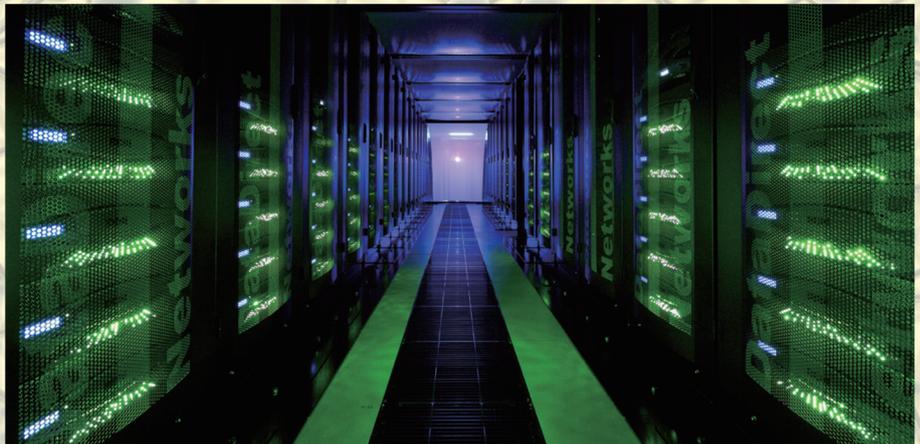




THE UNIVERSITY OF TOKYO

# SCD/ITC

**Supercomputing Division,  
Information Technology Center, The University of Tokyo**





# Scientific Computing & Numerical Algorithms/Libraries

## HACApK : Distributed Memory $\mathcal{H}$ -matrices Library

Low-rank structured matrices represented by hierarchical matrices ( $\mathcal{H}$ -matrices) have recently received attention as fast computational techniques for dense matrices arising from scientific simulations. For matrix size  $N$ , the memory complexity of low-rank structured matrices is at worst  $O(N \log N)$ , which is much lower than that of dense matrices  $O(N^2)$ .

We have been developing an open-source  $\mathcal{H}$ -matrices library named HACApK. Development of the library is started in 2012 as a part of the ppOpen-HPC project. To exploit recent supercomputer systems, HACApK is designed for a distributed memory cluster of SMP. From 2017, our proposals to enhance functions of the library are accepted as an international joint research project of JHPCN and JSPS KAKENHI projects. In the projects, we address wide spectrum of issues, such

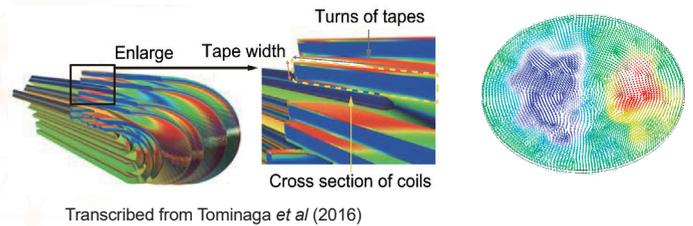
as porting the library to GPU and FPGA, increasing available matrix arithmetic functions and improvement of  $\mathcal{H}$ -matrices for massively parallel processing. The latest version of the HACApK library is available at the web-page of the ppOpen-HPC project.

HACApK library is employed in practical simulations, such as electric field, earthquake cycle, superconductor and micromagnetics. The use of the library enables us to conduct large scale simulations. Moreover, we have been challenging to open a new frontier to apply the library. For example, we address the application of HACApK to the elastodynamic boundary integral equation method to investigate the earthquake rupture dynamics, as a general JHPCN project from 2017.

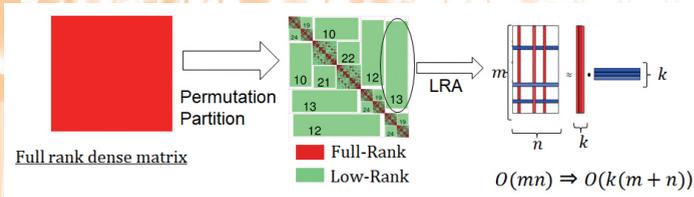
### Example analyses

• Superconductor

• Spin Torque Oscillator



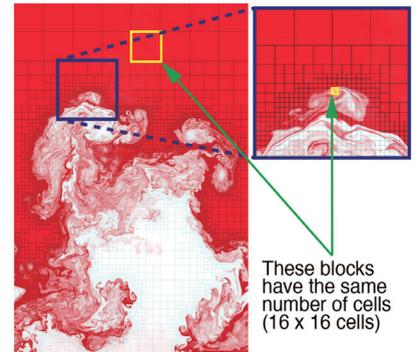
### Overview of Low-rank structured matrices



## AMR Framework with multiple GPUs to Realize Effective High-Resolution Simulations

Recently grid-based physical simulations with GPU require effective methods to adapt grid resolution to certain sensitive regions of simulations. An adaptive mesh refinement (AMR) method is one of the effective methods to compute certain local regions that demand higher accuracy with higher resolution. To develop the applications adopting AMR effectively with maintaining high performance on multiple GPUs, we are developing a block-based AMR framework for stencil applications. Programmers just write C++11 lambda expressions called the stencil functions that update a grid point on Cartesian grid and the framework executes them over a tree-based AMR data structure effectively. An entire computational domain is divided into a large number of the small uniform grid blocks with the same size recursively. The computation for all grid blocks can be solved with a single execution of a conventional stencil calculation for Cartesian grid regardless of their resolutions. The framework provides the halo exchange between GPUs based on the temporal blocking method, which contributes to performance improvement. It also provides mesh refinement mechanism and data migration that are required for AMR applications with a dynamic load balance technique.

The framework-based application for compressible flow has achieved to reduce the computational time to less than 15% with 10% of memory footprint in the best case compared to the equivalent computation running on the fine uniform grid. It also has demonstrated good weak scalability with 84% of the parallel efficiency on multiple NVIDIA P100 GPUs on the TSUBAME3.0 supercomputer.



**A snapshot of the Rayleigh-Taylor instability simulation obtained by compressible flow computation with AMR**

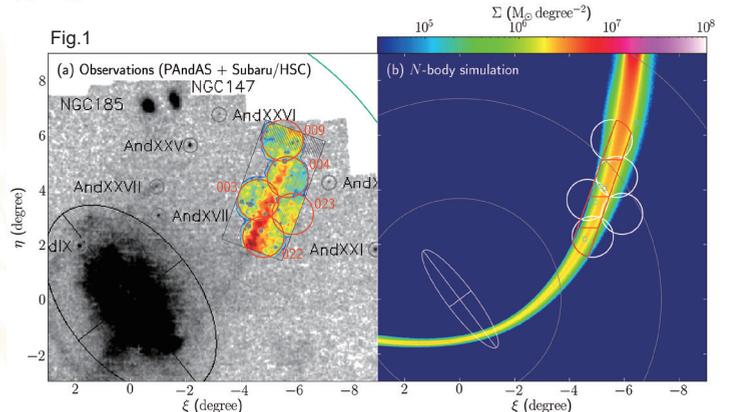
## Software development for numerical astrophysics

Collisionless N-body simulations are frequently employed to explore the formation and evolution of galaxies. We have developed an initial-condition generator and a gravitational N-body code and have performed galactic merger simulations by using them (Fig.1).

MAGI (MAny-component Galaxy\_INITIALIZER) is an initial-condition generator that supports various types of density models, their superposition, and the presence of multiple disks. We tested the dynamical stability of systems generated by MAGI representing elliptical and disk galaxies and confirmed that the model galaxies maintained their initial distributions for over a billion years. MAGI is publicly and freely available at <https://bitbucket.org/ymiki/magi>.

GOTHIC (Gravitational Oct-Tree code accelerated by Hierarchical time step Controlling) is a gravitational N-body code, which includes both the tree method and the hierarchical time step. The code runs entirely on GPU and is optimized for from the Fermi to the Volta GPU architectures. The performance measurements on Tesla V100, the current flagship GPU by NVIDIA, revealed that the N-body simulations of the Andromeda galaxy model with 8388608 particles took 33 ms per step. It corresponds to 3.1 TFlop/s, which is 20% of the single-precision theoretical peak performance. Tesla V100 achieves a 1.4 to 2.2-fold acceleration in comparison with Tesla P100, the flagship GPU in the

previous generation. The observed speed-up of 2.2 is greater than 1.5, which is the ratio of the theoretical peak performance of the two GPUs. The independence of the units for integer operations from those for floating-point number operations enables the overlapped execution of both operations leading to the speed-up rate above the theoretical peak performance ratio.



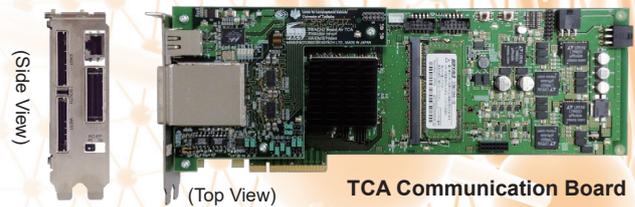
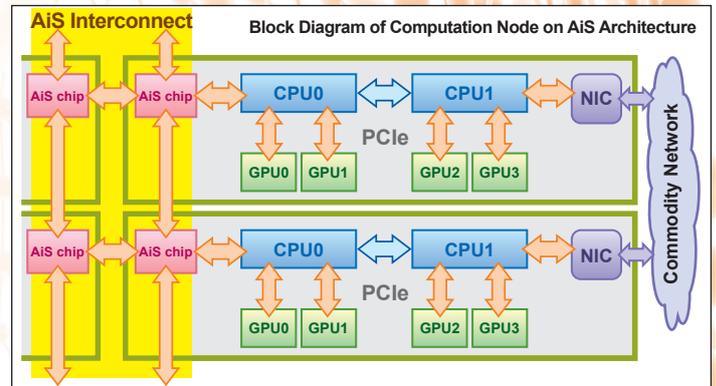
# System, tools & hardware

## Accelerators in Switch (AiS) Concept: Coordination among CPU, GPU, and FPGA

GPU computing is now widely performed for accelerating computational science and engineering applications to improve performance significantly with less power consumption. In order to reduce the latency of inter-node GPU communication, TCA (Tightly Coupled Accelerators) had been developed to enable direct communication using PCI Express. PEACH2/3 (PCI Express Adaptive Communication Hub ver. 2/3) chip was implemented in cooperation with University of Tsukuba and Keio University. Flexible control is realized by the using FPGA (Field Programmable Gate Array). We achieved the direct communication among inter-node GPUs based on PCI Express protocol with up to 8x improvement for latency.

In the TCA, FPGA was used only for the communication device. Thus, many unused logics in FPGA are helpful as a small accelerator for streaming computation through CPU-GPU or GPU-GPU with reconfiguration. Accelerators in Switch (AiS) is the concept to build such a cluster using accelerators like an enhancement of TCA architecture. AiS approach would become one solution towards post-Moore era. This research includes the following topics:

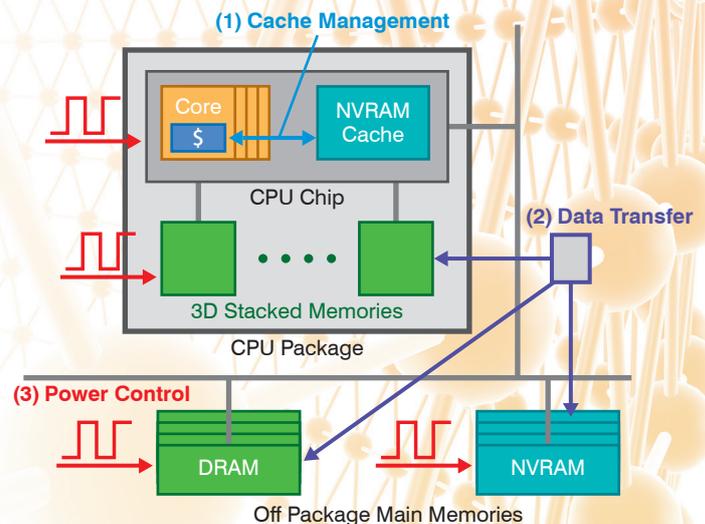
- Feasibility study on porting practical applications to FPGA using OpenCL and other programming methods as AiS programming environment
- Investigation for a new platform and interface technologies



## Improving Energy Efficiency using Emerging Memory Technologies

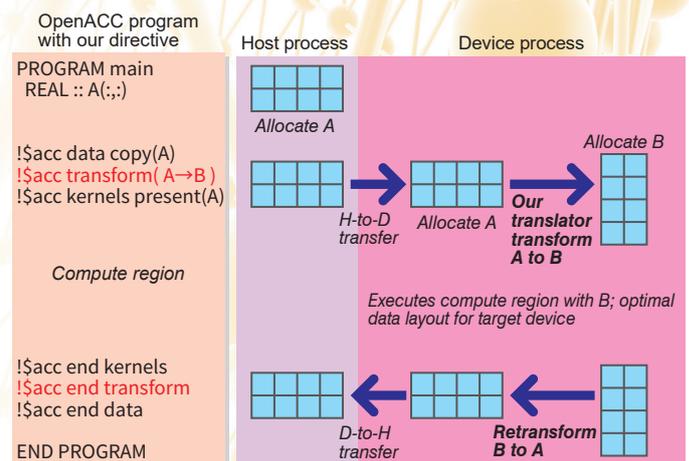
It is critical to improve the energy efficiency (performance per watts) of supercomputing nodes to build future exa-scale systems. This is because recent top-class supercomputers have already reached the power supply limit of a few tens of megawatts; therefore, we cannot merely scale the number of nodes to gain performance. In addition, VLSI technology scaling, which is the major contributor to energy-efficiency improvements, is coming to an end; therefore, alternative approaches must be developed.

To this end, we are focusing on the node architecture which comprises emerging memory technologies (e.g., 3D stacking and NVRAM), and developing software/hardware techniques to exploit the energy efficiency of the nodes. These memory technologies are helpful to scale bandwidth or capacity with smaller power budgets, and are indispensable to improve the performance of various memory intensive applications. Particularly, we are developing (1) data management on NVRAM-based cache hierarchies, (2) data transfer optimizations on hybrid main memories which comprise multiple different memory technologies, and (3) power control techniques for such systems. By combining these techniques, the energy efficiency can be considerably improved.



## OpenACC Extension for Performance Portability

OpenACC is gaining momentum as an implicit and portable interface in porting legacy CPU-based applications to heterogeneous, highly parallel computational environment involving many-core accelerators such as GPUs and Intel Xeon Phi. OpenACC provides a set of loop directives similar to OpenMP for the parallelization and also to manage data movement, attaining functional portability across different heterogeneous devices; however, the performance portability of OpenACC is said to be insufficient due to the characteristics of different target devices, especially those regarding memory layouts, as automated attempts by the compilers to adapt is currently difficult. We are currently working to propose a set of directives to allow compilers to have better semantic information for adaptation; here, we particularly focus on data layout such as Structure of Arrays, advantageous data structure for GPUs, as opposed to Array of Structures, which exhibits good performance on CPUs. We propose a directive extension to OpenACC that allows the users to flexibility specify optimal layouts.

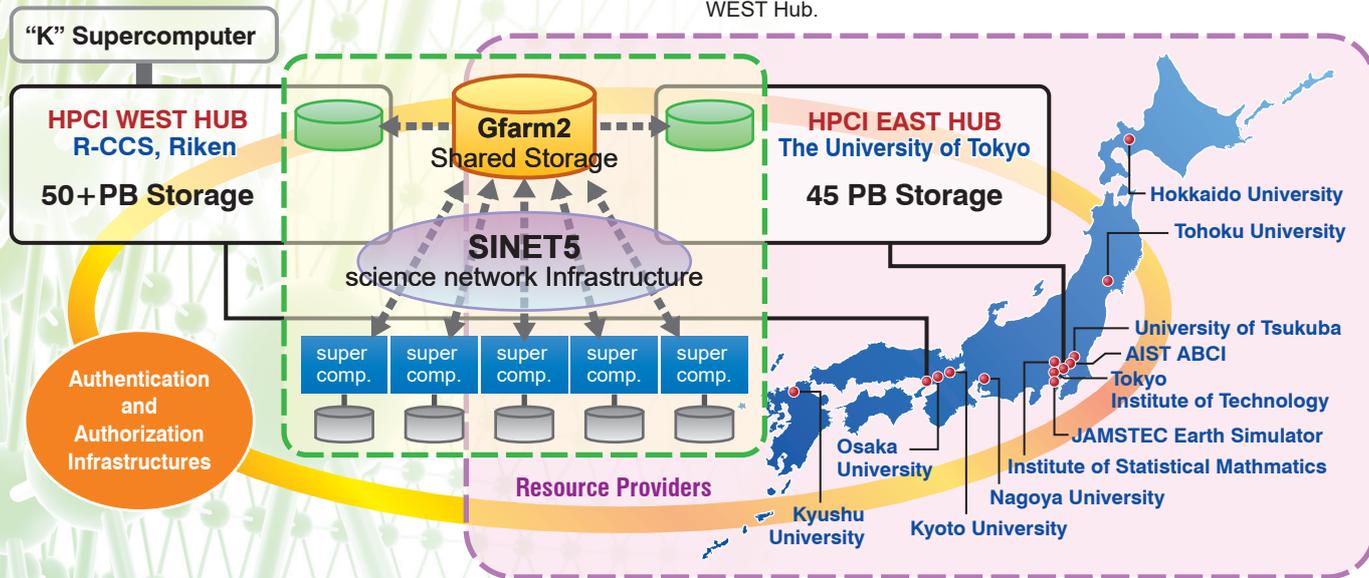


# Supercomputers in SCD/ITC

## HPCI: High Performance Computing Infrastructure

High performance computing infrastructure (HPCI) is an environment that enables an easy usage of flagship “K” and “Fugaku” supercomputers and other computation resources (tier-2) in Japan. In addition, HPCI is expected to match a user’s needs and computational resources to accelerate exploratory research, large-scale research, and industrial use of HPC. HPCI comprises 13 computational resource providers; nine

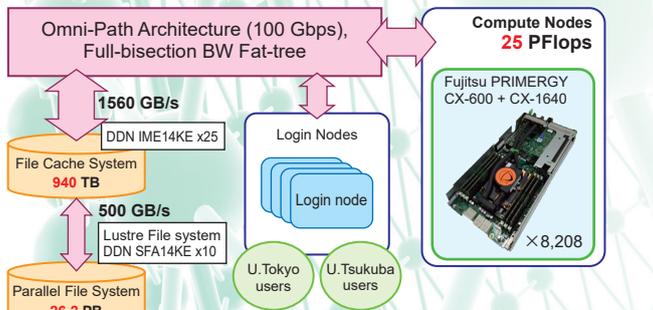
of these providers are supercomputing centers at national universities and four are governmental research institutes. These resource suppliers are connected via SINET5, which is a high-speed academic backbone network with 100 Gbps. SCD/ITC participates in this project as a hub resource provider in the Kanto region (the HPCI EAST Hub). The HPCI EAST Hub provides a 45-PB storage system in combination with the WEST Hub.



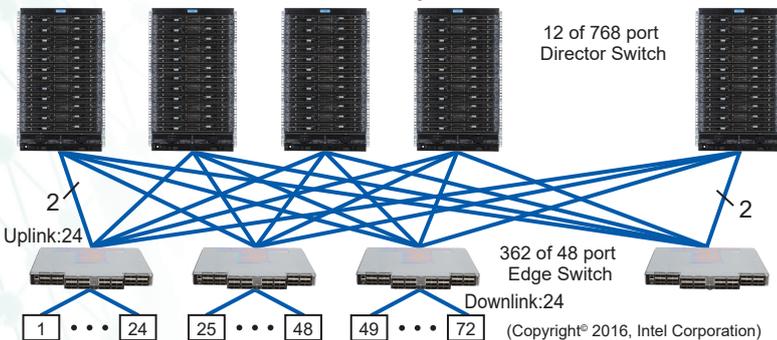
## Oakforest-PACS (FUJITSU PRIMERGY)

Oakforest-PACS is the first supercomputer introduced by JCAHPC (Joint Center for Advanced HPC) which is established by SCD/ITC and Center for Computational Sciences, U. Tsukuba (CCS). The system consists of 8,208 nodes of Intel Xeon Phi (Knights Landing) as a host processor, and Omni-Path Architecture provides 100 Gbps interconnection. In addition, the system employs the parallel file system with 26 PB, and file cache system of 940 TB with BW of over 1.5 TB/sec. This system is located at Kashiwa campus and operated by Fujitsu since December 2016. Oakforest-PACS has been offering computing resource to researchers in Japan and their international collaborators through various types of programs, such as by HPCI, by MEXT’s Joint Usage/Research Centers, and by each of CCS and ITC. It is expected to contribute to drastic developments of new frontiers of various field of studies, including computational science and engineering (CSE). This system is utilized for education and training of students and young researchers in both CSE and high-performance computing (HPC) as well. Both of CCS and ITC continue to make further social contributions through operations of the Oakforest-PACS.

Entire	Theoretical peak	25 PFLOPS
	Main memory	128TB (High BW)+ 770 TB (Low BW)
	Number of nodes	8,208
Compute node	FUJITSU PRIMERGY CX600 M1 + CX1640 M1	
Processor	Intel Xeon Phi (Knights Landing) 7250	
CPU (Core)	68 core, 1.4 GHz, 2 x AVX512	
Theoretical peak	3.05 TFLOPS	
Main memory	16 GB (High BW) + 96 GB (Low BW)	
Memory bandwidth	490 GB/sec (High BW, effective) + 115 GB/sec (Low BW)	
Cache memory	L2: 1 MB / tile (2 cores)	
Interconnect	Intel OmniPath Architecture (100 Gbps)	
Interconnect Topology	Full-bisection BW Fat Tree	
Parallel file system	Lustre Filesystem(DDN SFA14KE x10) 26PB, 500GB/sec	
File cache system	Burst buffer(DDN IME14K x25) 940TB, 1560GB/sec	



Full bisection bandwidth Fat-tree by Intel® Omni-Path Architecture



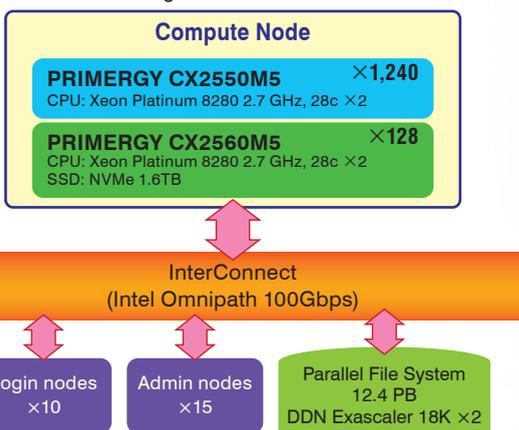
Power consumption	4.2 MW (including cooling)	
Number of racks	102	
Cooling system	Compute node	Warm-water cooling
		Direct cooling (CPU)
		Rear door cooling (except CPU)
Others	Facility	Cooling tower & Chiller
	Type	Air cooling
	Facility	PAC

Ranking		
Top500 (Nov. 2016)	#6 (#1 in Japan)	13.55 PFLOPS
Green500 (Nov. 2016)	#6 (#2 in Japan)	4985.7 MFLOPS/W
HPCG (Nov. 2016)	#3 (#2 in Japan)	385.5 TFLOPS
IO500 (Nov. 2018)	#1	137.78(BW 560.10 GiB/s, MD 33.89 kiOPS)

# Supercomputers in SCD/ITC

## Oakbridge-CX (FUJITSU PRIMERGY)

Oakbridge-CX (OBCX) is the Massively Parallel Supercomputer System using Intel Xeon CascadeLake CPUs with the total performance of 6.61 PFLOPS. We have started the operation since July 2019, and plan to offer the computing service from Oct. 2019. The 128 nodes of compute node employ a NVMe SSD in each node for supporting staging, checkpointing, and data-intensive applications. Moreover, SSDs on designated nodes can be dynamically converged as a single shared file system using BeeGFS on Demand (BeeOND). OBCX will become a prototype for the next-generation applications based on the fusion of "Simulation+Data+Learning."



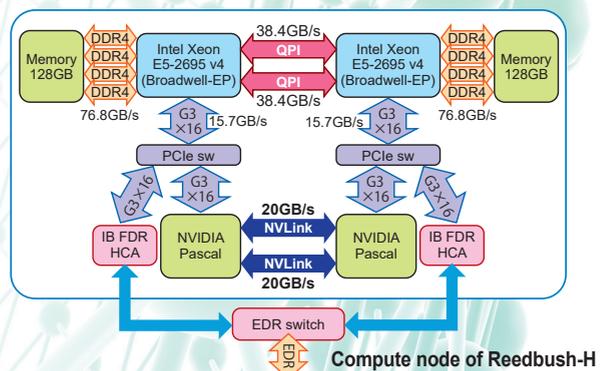
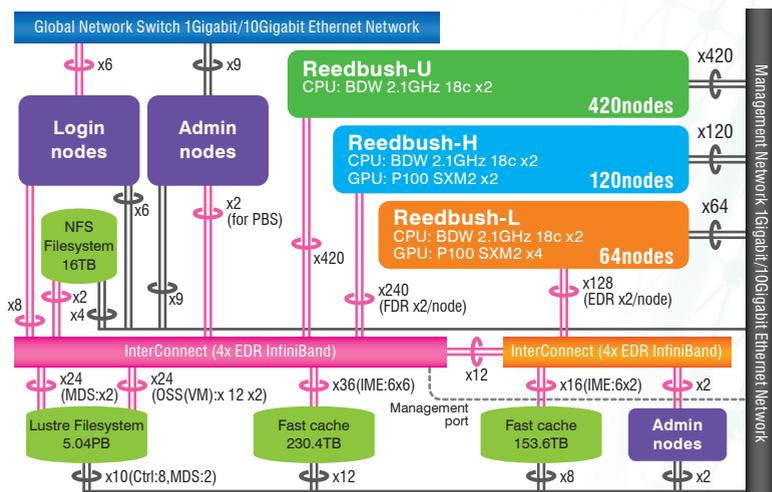
Peak performance	6.61 PFlops	
Total memory size	256.5 TByte	
Number of nodes	1240	128
Compute node	Fujitsu PRIMERGY CX2550 M5	Fujitsu PRIMERGY CX2560 M5
CPU	Intel Xeon Platinum 8280 (CascadeLake, 28 cores, 2.7 GHz) 4.83 TFLOPS	
Memory	192 GB (DDR4)	
Interconnect	Intel Omni-Path (100 Gbps)	
Interconnect topology	Full Bisection BW Fat Tree	
SSD	—	1.6 TB(NVMe, Read: 3.20 GB/s, Write: 1.32 GB/s)
Parallel file system	Lustre Filesystem (DDN SFA18KE x2) 12.4 PB, 98 GB/s	

## Reedbush (SGI Rakable system)

Reedbush is the first supercomputer system using accelerators in SCD/ITC, and its total peak performance is up to 3.3 PFLOPS. We started computing service with Reedbush-U (CPUs only) in July 2016 and continue to operate this part until June 2020. In terms of the GPU part, the Reedbush-H (with 2 GPUs per node), and the Reedbush-L (with 4 GPUs per node) subsystem have been available since March 2017 and July 2017, respectively. This system is located at the Asano campus in Tokyo and is operated by HPE (ex-SGI). This system has the following missions.

	Reedbush-U	Reedbush-H	Reedbush-L
Peak performance	509 TFlops	1417 TFlops	1433 TFlops
Number of nodes	420	120	64
Total memory size	105 TByte	30 TByte + 3.75 TByte	16 TByte + 4 TByte
Compute node	SGI Rackable C2112-4GP3		SGI Rackable C1102-GP8
CPU	Intel Xeon E5-2695v4 (Broadwell-EP, 18 core, 2.1 GHz) x 2 socket 1209.6 GFlops		
Memory	256 GB (DDR4-2400 x 4ch x 2), 153.6 GB/sec		
GPU	None	NVIDIA Tesla P100 (Pascal, 5.3 TFlops, 16 GB, 720 GB/sec) x 2	NVIDIA Tesla P100 (Pascal, 4.8-5.3 TFlops, 16 GB, 720 GB/sec) x 4
Interconnect	InfiniBand EDR 4x (100 Gbps)	InfiniBand FDR 4x 2 link (56 Gbps x2)	InfiniBand EDR 4x 2 link (100 Gbps x2)
Interconnect topology	Full-bisection BW Fat Tree		Full-bisection BW Fat Tree
Parallel file system	Lustre Filesystem (DDN SFA14KE x3) 5.04 PB, 145.2 GB/sec		
File cache system	Burst buffer (DDN IME14K x6) 230.4 TB, 385.2 GB/sec		
			Burst buffer (DDN IME240 x8) 153.6 TB, 166.4 GB/sec

- Development of new research field, and promotion for new users, such as big data and deep learning researchers
- Development of a pilot system of a next-generation supercomputer system for the integration and fusion of data analyses and scientific simulations



# Activities towards Society 5.0

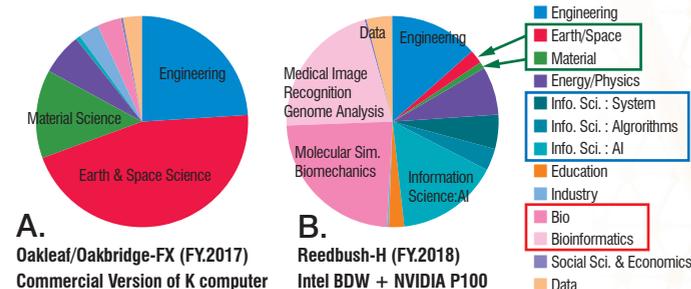
## What is Society 5.0 ?

[https://www8.cao.go.jp/cstp/english/society5\\_0/index.html](https://www8.cao.go.jp/cstp/english/society5_0/index.html)

Society 5.0 was proposed in the **5th Science and Technology Basic Plan by the Cabinet Office of Japan** as a future society that Japan should aspire to. It follows the hunting society (Society 1.0), agricultural society (2.0), industrial society (3.0), and information society (4.0). Society 5.0 is a human-centered society that balances economic advancement with the resolution of social problems by a system that highly integrates cyberspace and physical space, and will be achieved by Digital Innovation, such as IoT, AI, Big Data and etc.

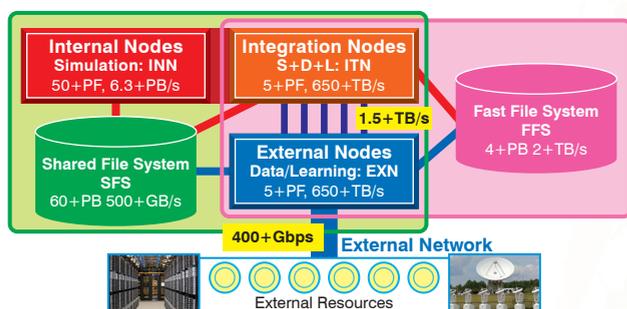
## New Directions in Supercomputing

Majority of SCD/ITC supercomputer system users belong to the fields of CSE, including engineering simulations (fluid dynamics, structural dynamics, and electromagnetics), earth sciences (atmosphere, ocean, solid earth, and earthquakes), and material sciences, as shown in the pie chart A, which shows usage rate of each research area on Oakleaf/Oakbridge-FX system (commercial version of the K computer) based on CPU hours in FY.2017. Recently, the number of users related to data science, machine learning, and artificial intelligence (AI) has been increasing, as shown in the pie chart B, which shows usage rate on Reedbush-H system with GPU's in FY.2018. Examples of new research topics are weather prediction by data assimilation, medical image recognition, and human genome analyses. Towards Society 5.0, a new type of method for solving scientific problems by integrations of "Simulation (S)", "Data (D)" and "Learning (L)" (S+D+L) is emerging.



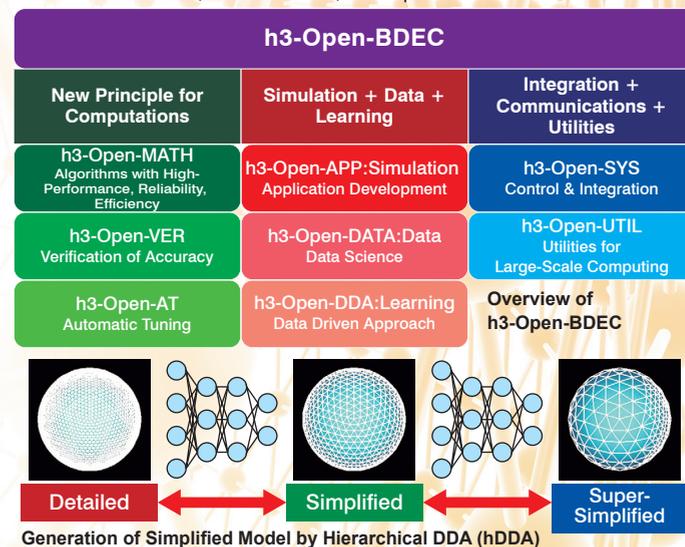
## BDEC System: Big Data & Extreme Computing

The BDEC system (Big Data & Extreme Computing), which is scheduled to be introduced to SCD/ITC in Summer 2021, is a Hierarchical, Hybrid, Heterogeneous (h3) system. The BDEC is the platform for integration of "Simulation, Data and Learning (S+D+L)", and consists of computing nodes for computational science, those for data science/machine learning, and those for integration. The aggregated peak performance of the BDEC system is expected to be 60+ PFLOPS, and it will comprise three types of compute nodes, "Internal Nodes (INN, peak performance: 50+PF, memory bandwidth: 6.30+PB/sec)" for traditional supercomputing applications, "External Nodes (EXN, 5+PF, 650+TB/sec)" for data and learning, and "Integration Nodes (ITN, 5+PF, 650+TB/sec)". Architecture of INN and ITN must be same, while that of EXN could be different. Each node of the EXN will be connected to external resources (e.g. data storage, servers, sensor networks, and etc.) directly through an external network (e.g., SINET, Japan). ITN and EXN will share the Fast File System (FFS, capacity: 4+PB, bandwidth: 2+TB/sec), while all nodes will can access the large-scale Shared File System (SFS, 60+PB, 500+GB/sec). The Reedbush-U/H/L systems and the Oakbridge-CX system are prototypes of the BDEC system.



## h3-Open-BDEC: Innovative Software Platform

We develop an innovative software platform "h3-Open-BDEC" for integration of (S+D+L), and evaluate the effects of integration of (S+D+L) on the BDEC system. The h3-Open-BDEC is designed for extracting the maximum performance of the supercomputers with minimum energy consumption focusing on (1) innovative method for numerical analysis with high-performance/high-reliability/power-saving based on the new principle of computing by adaptive precision, accuracy verification and automatic tuning, and (2) Hierarchical Data Driven Approach (hDDA) based on machine learning. This work will be supported by Japanese Government from FY.2019 to FY.2023 (JSPS Grant-in-Aid for Scientific Research (S), P.I.: Kengo Nakajima (ITC/U.Tokyo)). In Data Driven Approach (DDA), technique of machine learning is introduced for predicting the results of simulations with different parameters. DDA generally requires a lot of simulations for generation of teaching data. We propose the hDDA, where simplified models for generating teaching data are constructed automatically by machine learning with Feature Detection, MOR, UQ, Sparse Modeling and AMR. The h3-Open-BDEC is the first innovative software platform to realize integration of (S+D+L) on supercomputers in the Exascale Era, where computational scientists can achieve such integration without supports by other experts. Source codes and documents are open to public for various kinds of computational environments. This integration by h3-Open-BDEC enables significant reduction of computations and power consumptions, compared to those by conventional simulations. Idea of h3-Open-BDEC is that of "ppOpen-HPC (<https://github.com/Post-Peta-Crest/ppOpenHPC>)", which is part of a (five+three)-year project (FY.2011-2015, FY.2016-2018) supported by JST-CREST and DFG-SPPEXA in Germany. ppOpen-HPC is a framework for development of application with automatic tuning (AT). Possible applications on the BDEC system with h3-Open-BDEC are combined simulations/data assimilations for climate/weather simulations and earthquake simulations, and real-time disaster simulations, such as flood, earthquake and tsunami.



## Data Platform (DP)

The Data Platform (DP) is a project supported by Japanese Government, which develops a platform for utilization of various types of data sets towards Society 5.0. The eight national universities in JHPCN, NII (National Institute of Informatics) and AIST (National Institute of Advanced Industrial Science and Technology) are involved. NII operates the SINET (Science Information NETWORK). The Data Platform consists of computing nodes and storage. It is similar to the BDEC, but is focusing more on analyses and utilization of data. More flexible configuration and operation are possible compared to traditional supercomputer systems. Moreover, security of the system is a very critical issue for handling confidential data. Operation of the Data Platform will start in Spring 2021. The Data Platform and the BDEC will be installed in a same room of a new building in Kashiwa-II Campus of the University of Tokyo. The DP can access the Shared File System (SFS) of the BDEC.

# SCD/ITC

Supercomputing Division,  
Information Technology Center, The University of Tokyo



東京大学情報基盤センター

INFORMATION TECHNOLOGY CENTER, THE UNIVERSITY OF TOKYO

2-11-16 Yayoi, Bunkyo, Tokyo 113-8658, JAPAN TEL:03-5841-2710 FAX:03-5841-2708(G3) <https://www.itc.u-tokyo.ac.jp/en/>

2019.08