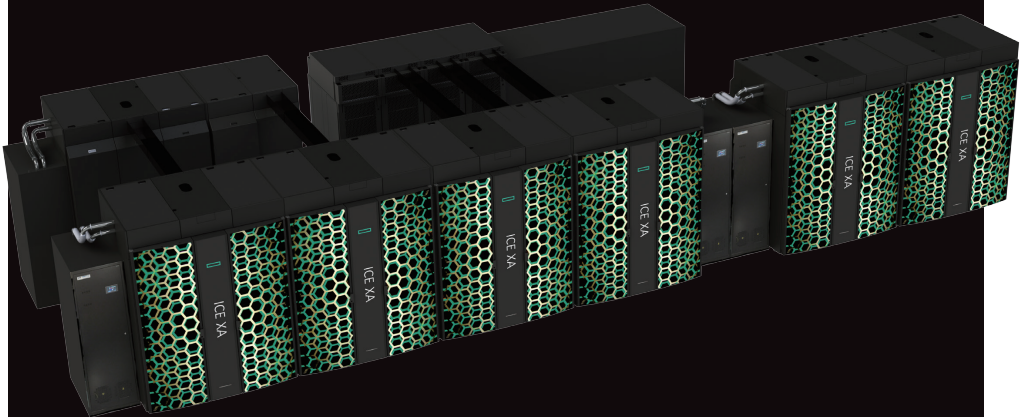


# HARDWARE SOFTWARE SPECIFICATIONS TSUBAME3.0

PREVIEW VERSION

- GPU acceleration for High Performance Computing and AI workloads
- High-Speed Network Interconnect
- High-Speed and Highly Reliable Storage Systems
- Low Power Consumption and Green Operation
- System and Application Software



# GPU-Equipped High-Performance Compute Nodes

TSUBAME3.0 system includes 540 compute nodes, which provides 12.15 PFlops performance in total. Each compute node is equipped with two CPUs and four GPUs in a compact design blade. In addition, high-speed network with four Omni-Path HFI and large capacity SSD accelerate applications' performance in Big Data and AI areas.

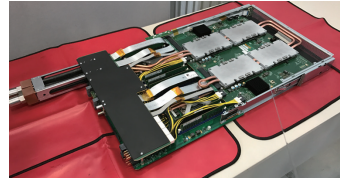
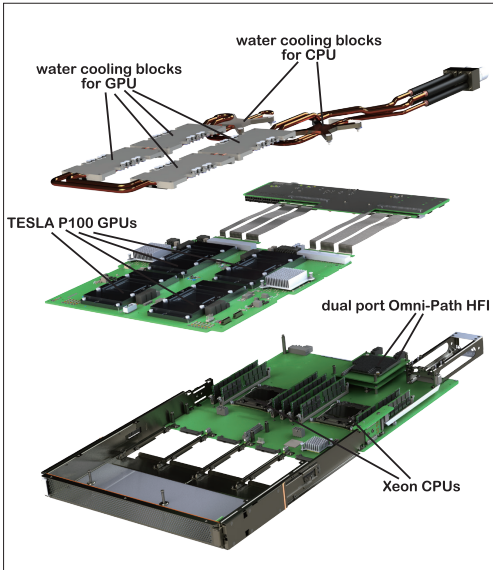
## HPE SGI ICE-XA (SGI 8600) IP139-SXM2 540 nodes

- CPU:** Intel Xeon E5-2680 V4 (Broadwell-EP, 2.4GHz) ×2 sockets  
14 cores per socket, total 28 cores per node.
- GPU:** NVIDIA TESLA P100 for NVLink-Optimized servers ×4.
- Memory:** 256GB (DDR4-2400 32GB module ×8)
- SSD:** Intel DC P3500 2TB (NVMe, PCI-E 3.0 ×4)
- Network:** Intel Omni-Path Architecture HFI (100Gbps) ×4



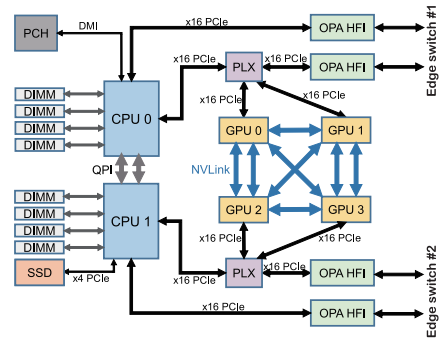
SGI ICE XA cabinet (front door opened)

Each E-Rack contains up to four chassis, and each chassis contains nine compute nodes.



Early prototype of compute node

## Block Diagram



## Tesla P100 for NVLink-Optimized servers

### Peak performance :

- 5.3 TFLOPS (double precision)
- 10.6 TFLOPS (single precision)
- 21.2 TFLOPS (half precision)

Shader clock : 1328 MHz (1480MHz with boost)

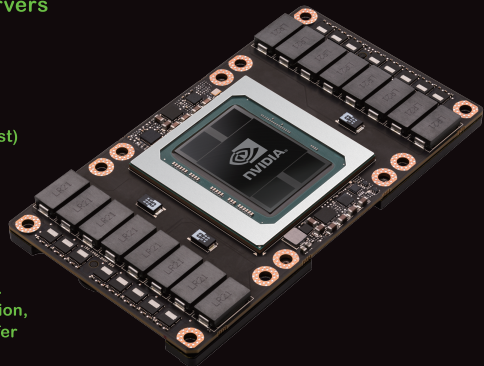
Number of CUDA cores (SP) : 3,584

Streaming Multiprocessors : 56

On-board memory: 16GB HBM2

Memory bandwidth : 720GB / sec

TDP: 300W

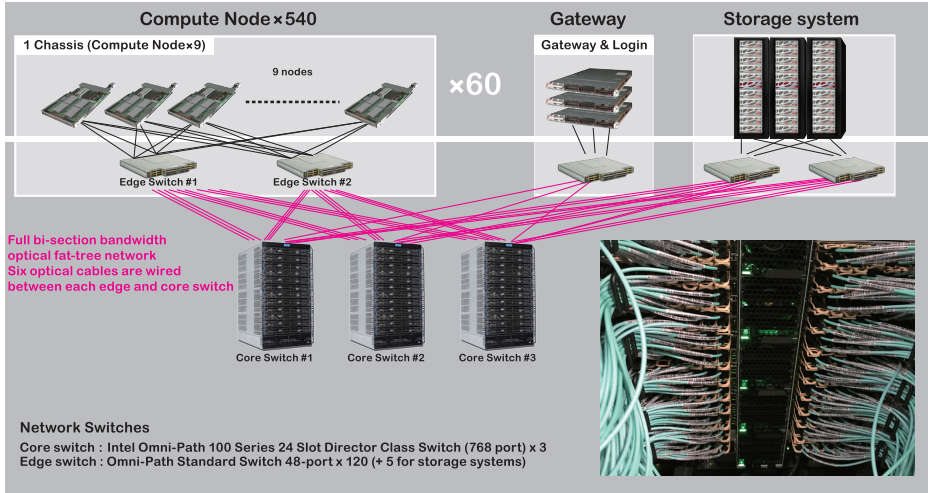


Pascal Architecture (NVIDIA GP100)

NVLink is direct interconnect between GPUs. Bandwidth of a link is 20GB/s for each direction, and four NVLinks enable 160GB/s data transfer in addition to the PCI-E bandwidth.

# High-Speed Interconnect

Compute nodes of TSUBAME3.0 interconnected with Omni-Path Architecture of full bi-section bandwidth fat-tree network achieving 432Tbps. End-to-End latency between the compute nodes is extremely low in microsecond-order time, resulting in high-speed performance and high-speed connection to highly reliable storages.



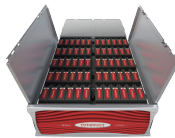
## Lustre



DDN EXAScaler x 3sets : Total 15.9PB, 150GB/s

Each EXAScaler consists of  
SFA14KXE + SS8462 x 10  
EF4024(MDS) x 2 + ED4024(MDT) x 2  
10TB 7.2Krpm NL-SAS HDD x 700 (incl. 20 spare)  
5.3PB effective capacity (20.4 physical)

## Home



DDN GridDirector : 45TB

SFA7700X  
300GB 2.5" 10Krpm SAS HDD x 226  
(Data: 200, Meta: 20, Spare: 6)  
NAS Gateway(NFS) x 4

## Campus storage



DDN GridDirector : 36TB

SFA7700X  
300GB 2.5" 10Krpm SAS HDD x 186  
(Data: 160, Meta: 20, Spare: 6)  
NAS Gateway(CIFS) x 4

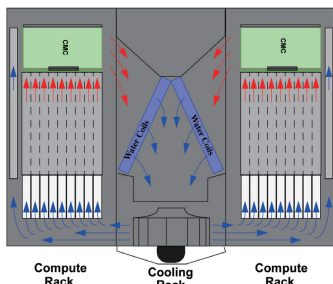
# High-Speed and Highly Reliable Storage Systems

TSUBAME3.0 provides massive storage volumes to serve various purposes, including 1.08 PB of SSDs embedded in compute nodes for scratch I/O, 15.9PB of parallel file systems such as Lustre for high speed parallel I/O, and 45 + 36 TB of home storage volumes for providing campus cloud storage services.

## Green Operation

Power performance in Linpack benchmark : TBA  
 Peak power consumption of system equipment : TBA  
 Average power consumption of system equipment: TBA  
 Idle power consumption for system equipment : TBA  
 Yearly average PUE : TBA

### Warm Water Cooling

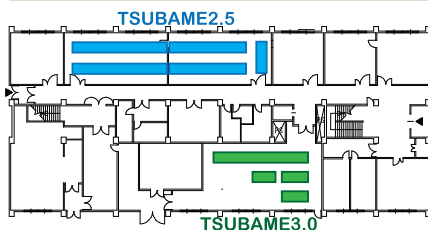


Cooling system requires additional electric power. In order to minimize it, we employ low-power evaporative cooling tower located on the rooftop of the building. It is estimated that the cooling tower provides less than 32 degrees water even in summer, and the compute nodes of TSUBAME3.0 are designed to work with the warm water.

CPUs and GPUs are direct water-cooled. The other components are air-cooled using rack-contained built-in heat exchanger (Water Coils in the figure) between water and air, allowing low-cost, high-maintainability, and high energy efficiency. The racks for storage and I/O systems employ rear door cooling technology.



### Small space installation



Highly efficient cooling of TSUBAME3.0 enables high density compute nodes. The space required for installation is less than half of TSUBAME2.5.

### System Software

### Docker support for efficient use of computing resources

The job management system and the cluster management system are working together to manage user environment as well as distributing computational resources to the insufficient part by taking from the node pool. Some jobs do not use all of computing resources (CPU cores, GPUs, memory, etc.) on allocated compute nodes. UNIVA Grid Engine's Docker support enables efficient node sharing between multiple jobs.

OS	SUSE Linux Enterprise Server 12 SP2
Batch System	UNIVA Grid Engine

### ISV (commercial) Software

(\* GPU full or partial support)

#### Compilers, Debuggers and Libraries

Intel Compiler (C/C++/Fortran)  
 PGI Compiler\*  
 (C/C++/Fortran, OpenACC, CUDA Fortran)  
 Allinea FORGE\*

#### Applications

ANSYS Workbench\*, Mechanical\* ABAQUS\*, ABAQUS CAE  
 ANSYS CFD, Fluent\*, HFSS\* MSC Nastran\*, Patran, Marc\*  
 COMSOL Multiphysics CST STUDIO SUITE\* (MW-Studio\*)  
 LS-DYNA AMBER\*  
 Gaussian\*, Gauss View Materials Studio, Discovery Studio  
 MATLAB\* Mathematica\*  
 AVS/Express, AVS/Express PCE Maple\*  
 Schrödinger Small-Molecule Drug Discovery Suite\*

Yellow: The license for all users White: The license for Tokyo Tech users Magenta: The license for industrial users

Published by Global Scientific Information and Computing Center, Tokyo Institute of Technology

2-12-1 Ookayama, Meguro-ku, Tokyo 152-8550, JAPAN

TEL : +81-3-5734-2087 FAX : +81-3-5734-3198 E-mail : tsubame@gsic.titech.ac.jp

<http://www.gsic.titech.ac.jp/>