



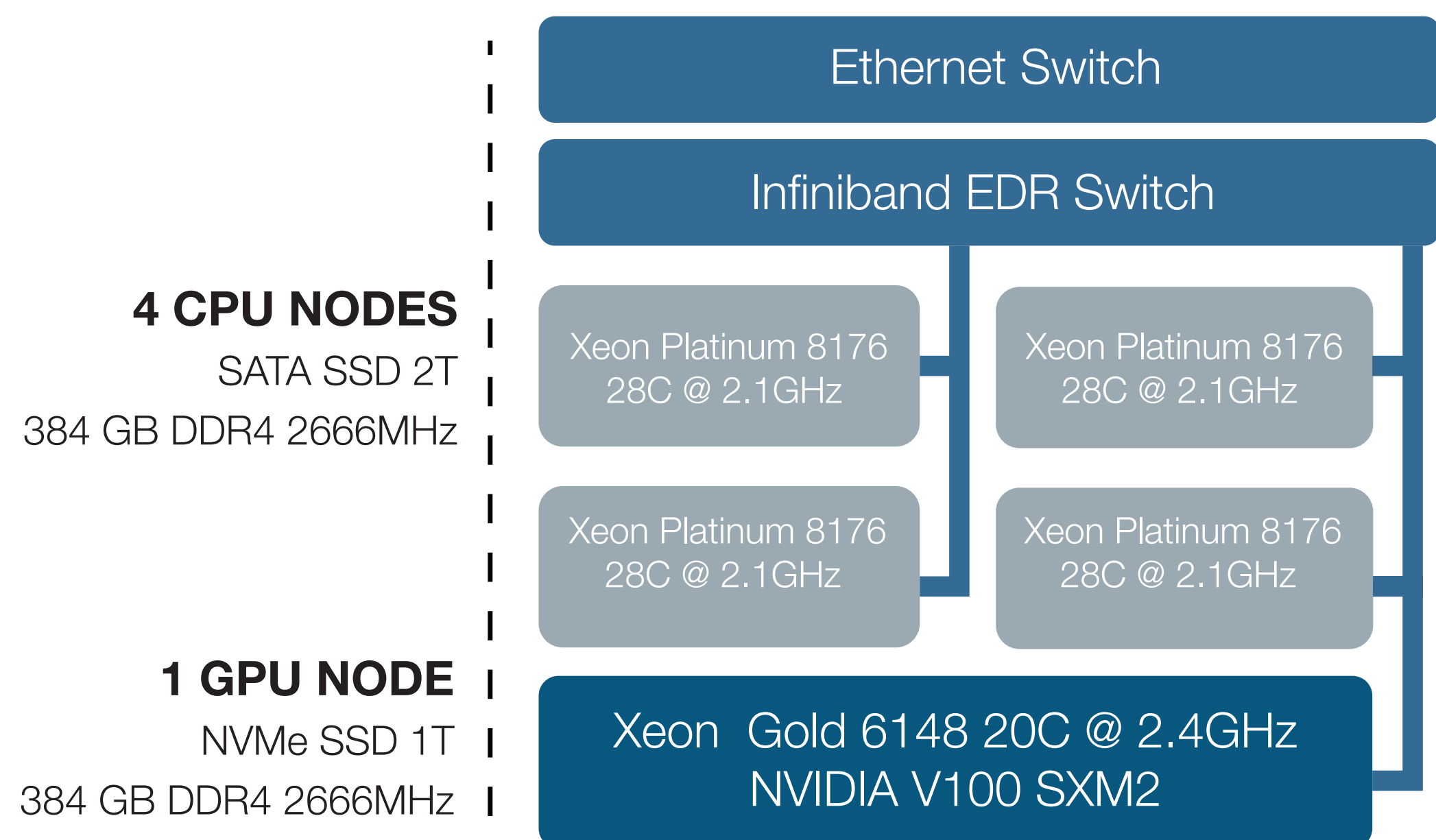
# National Tsing Hua University

## SC18 Student Cluster Competition

YuHsuan Cheng, KengJui Hsu, ChiChen Yang, HungHsin Chen, ShaoFu Lin and YuanChing Lin

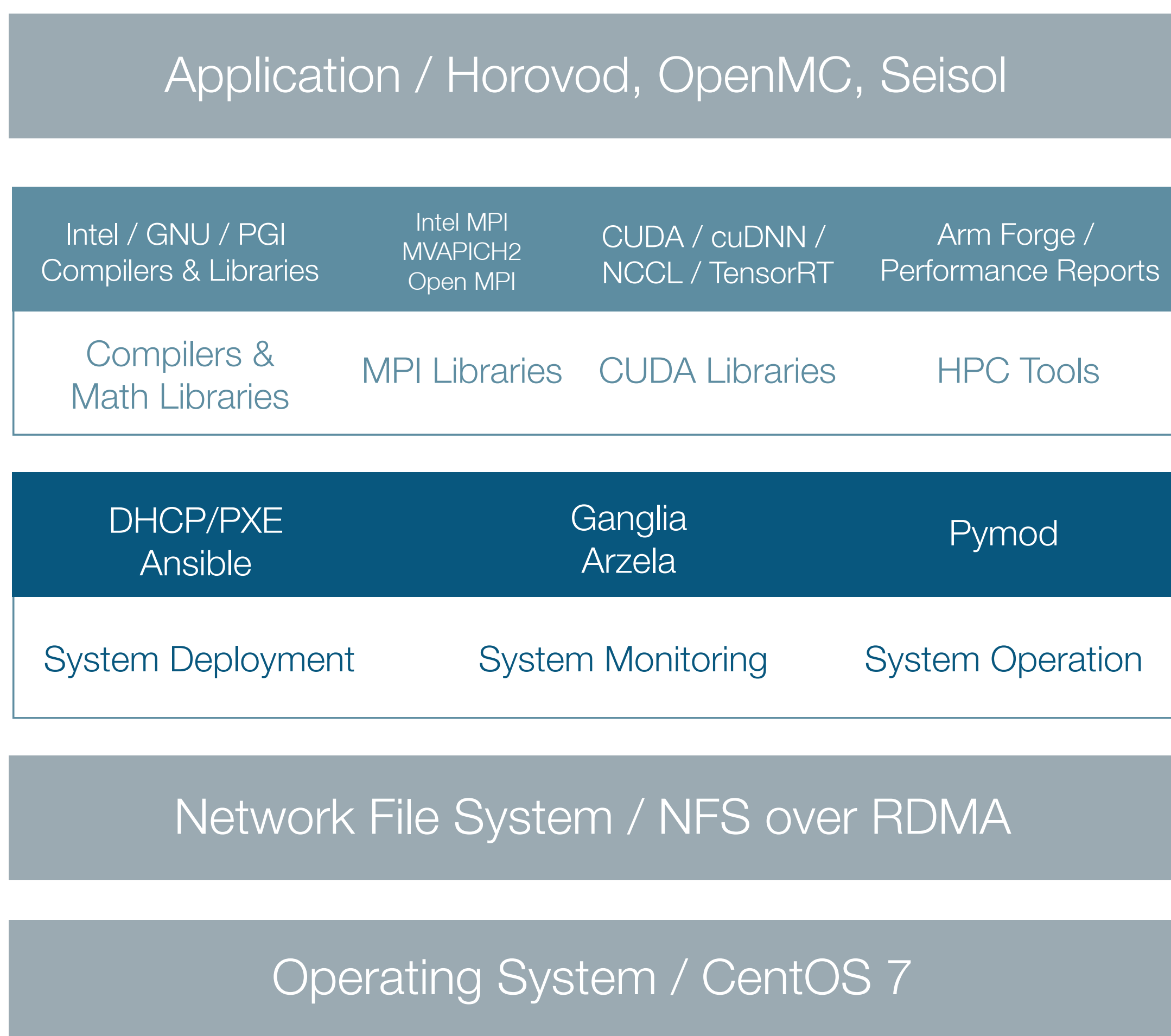
### Architecture and System

Both Intel Xeon Platinum 8176 and Intel Xeon Gold 6148 have many cores and balanced multi-core to power ratio. We choose NVIDIA Tesla V100 with NVLINK to maximize the double precision floating point operation with 46% speedup on the bandwidth. Mellanox EDR enables a 100Gbps network throughput and a low latency interconnection.



### Software Stack

CentOS provides a massive number of stable packages in its official repositories. We develop our own system operation toolkits to control the power and monitor the temperature of our cluster. With the experiment of these applications, we decide to choose OpenMPI as our main MPI solution.



### Diversity

#### A. Gender diversity:

- ✓ Members of different genders
- ✓ Encourage sophomores or females to join our team

#### B. Diversity of experience:

- ✓ Encourage sophomores to join our team
- ✓ 3 senior 3 junior

#### C. Different majors:

- ✓ 1 member majors in Mathematics
- ✓ 5 members major in Computer Science
- ✓ Encourage students with interdisciplinary knowledge to join us

### 1 OpenMC

#### Strategies

##### A. Different execution configurations for different type of testcases

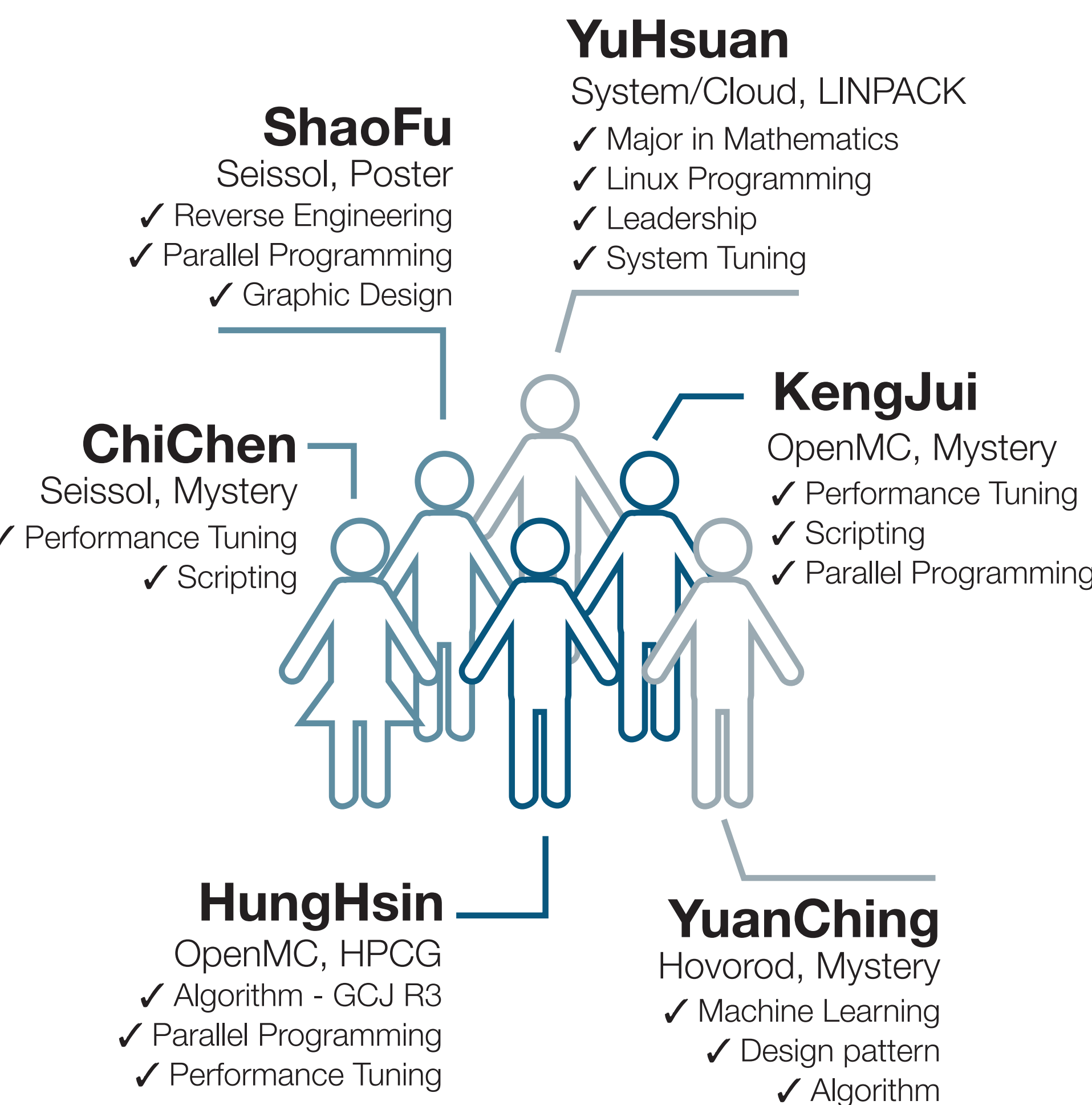
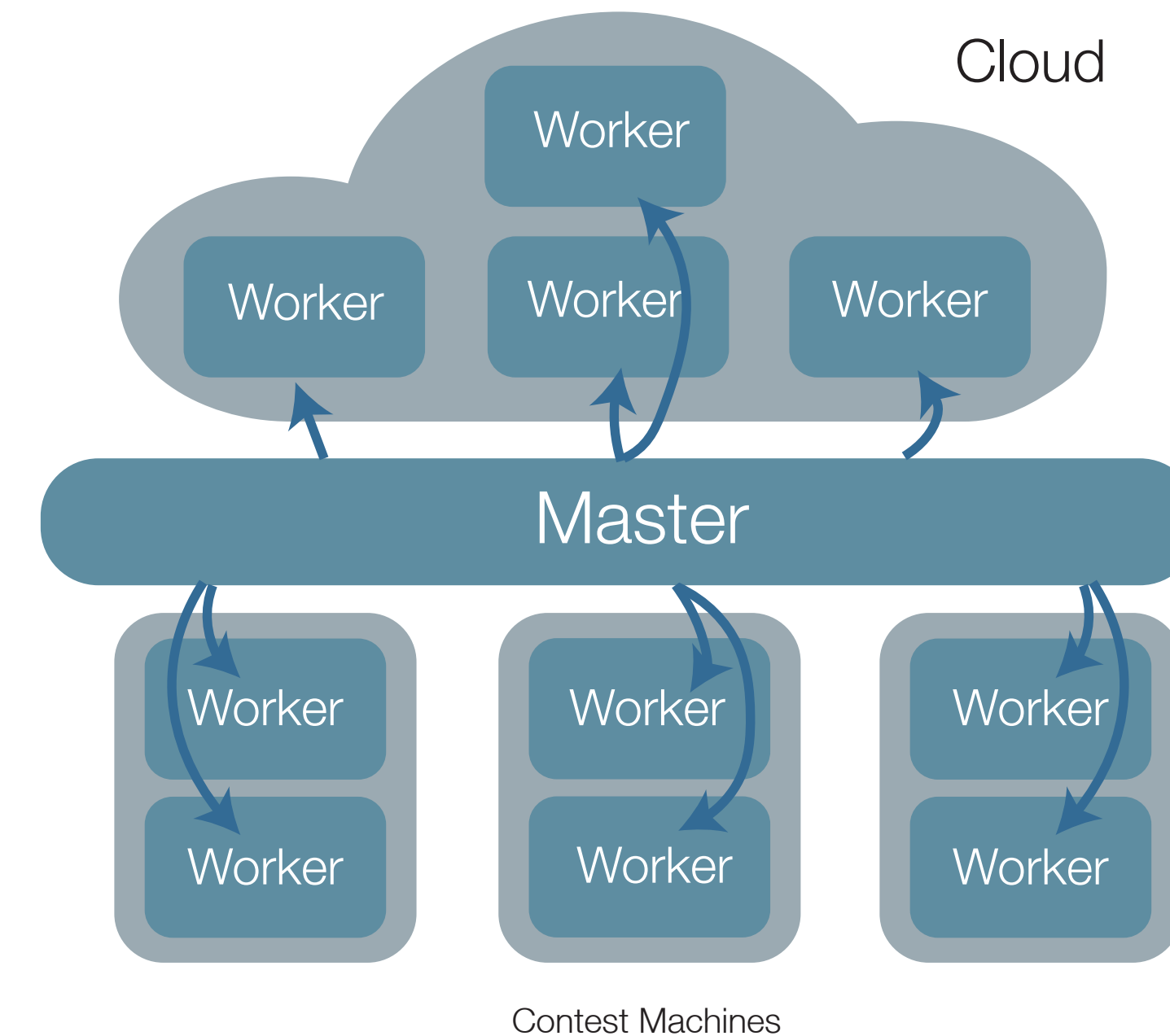
##### B. Well-designed job scheduler

- ✓ Automatic jobs dispatching
- ✓ Shortest job first
- ✓ 1 job per CPU socket
- ✓ Recover running state from statepoint
- ✓ Support for cloud job dispatching
- ✓ Add new jobs dynamically
- ✓ Centralized controller
  1. Manual job management
  2. Manual worker management
  3. Worker running configuration setting

##### C. Dynamic scheduling

##### D. Accurate time estimation

- ✓ Average runtime of several active batches
- ✓ Error within 1%



### System Strategies

#### Strategies

##### A. Power Management

- ✓ dynamically adjust CPU and GPU clock frequency
- ✓ controlling the cooling system
- ✓ power consumption estimation per application

##### B. Time Management

- ✓ automated job scheduler
- ✓ push jobs to cloud
- ✓ prepared working scripts
- ✓ collaboration between team members

##### C. Accurate time estimation

#### Cloud

Cyclecloud platform provides extra computing power without power constraints in the competition, and also avoids data loss in the power shut-off activity. OpenMC can take advantage of the Fv2 instance (Xeon Platinum 8168), which has a higher performance than our cluster. In addition, the RDMA feature of H16r instance enables multi-node speedup.

### 2 Horovod

Horovod is a distributed machine learning training framework. It aims to simplify the implementation of a distributed model with a higher scaling efficiency.

#### Strategies

##### A. Maximizing Multi-GPU Throughput

- ✓ Accelerated by 8x NVIDIA Tesla V100
- ✓ High speed interconnect with NVLINK

##### B. Power Consumption

- ✓ Balanced power performance ratio

##### C. Time Estimation

- ✓ Estimate convergence time for known models

Horovod implements parameters transferring upon the implementation of MPI and NCCL. We analyze the communication pattern by monitoring infiniband network bandwidth. Also, we tune the batch size and learning rate to fit the model well into our GPU memory.

### 3 SeisSol

SeisSol simulates an earthquake, and the SC version implemented some optimization on wave propagation, dynamic rupture and output.

#### Strategies

##### A. Prepare scripts:

- ✓ Easily run experiments and gather data from logs to plot figures.

##### B. Release machines:

- ✓ Start with the experiment that uses the most machines.

We reproduce figures with our own machines and then analyze it. In order to check our running method is correct, we had run Cori before the competition.

#### Why we will win?

- ✓ The passion toward parallel programming supports us during the long preparation.
- ✓ Our team leader Scott can use the knowledge in science domain to help us analyze problems and offer scientific insights.
- ✓ We not only work on performance tuning, but also study the background knowledge and the architecture of these applications.