


An Overview of Post-K Development

Yutaka Ishikawa
RIKEN AICS

9:00– 9:50 31st of January, 2018

28th-31st, January, 2018
HPC Asia 2018, Tokyo, Japan



Outline



- FLAGSHIP2020 Project
- Co-Design and R&D organization
- CPU Architecture and Post-K System Software
- Light-Weight OS Kernel: IHK/McKernel

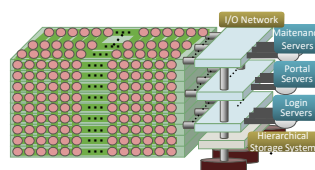


FLAGSHIP2020 Project



□ Missions

- Building the Japanese national flagship supercomputer, post K, and
- Developing wide range of HPC applications, running on post K, in order to solve social and science issues in Japan



□ Project organization

- Post K Computer development
- RIKEN AICS is in charge of development
- Fujitsu is vendor partner.
- International collaborations: DOE, JLESC, CEA
(NCSA, ANL, UTK, JSC, BSC, INRIA, RIKEN)

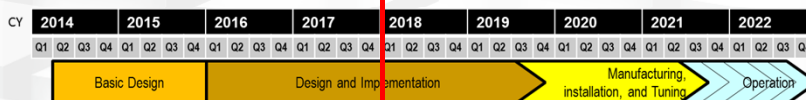
□ Current Status

- In the middle of “Design and Implementation” phase.
- CPU has been designed and is being implemented
 - now evaluated by simulators and compilers
- System software stack has been designed and is being implemented

• Applications

- The government selected
 - 9 social & scientific priority issues
 - 4 exploratory issues
- and their R&D organizations.

NOW



20018/1/31

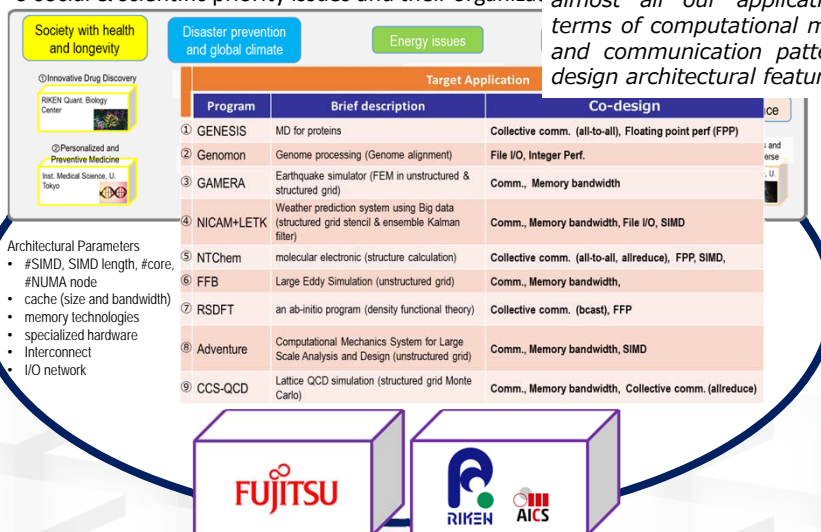
3

Co-design: *How we are running under the concept*



9 social & scientific priority issues and their organization

Those are representatives of almost all our applications in terms of computational methods and communication patterns to design architectural features.

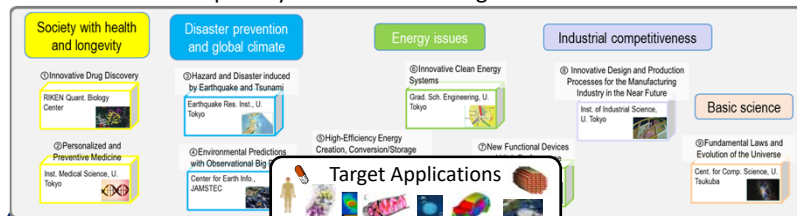


4

Co-design: *How we are running under the concept*



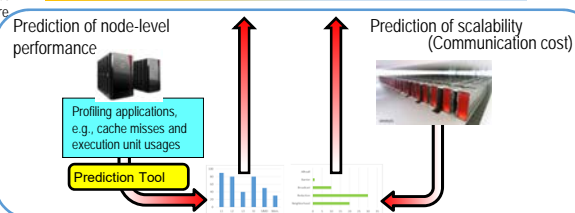
9 social & scientific priority issues and their organizations



Architectural Parameters

- #SIMD, SIMD length, #core
- cache (size and bandwidth)
- memory technologies
- specialized hardware
- Interconnect
- I/O network

- ❑ Mutual understanding both computer architecture/system software and applications
- ❑ Looking at performance predicted by the following method

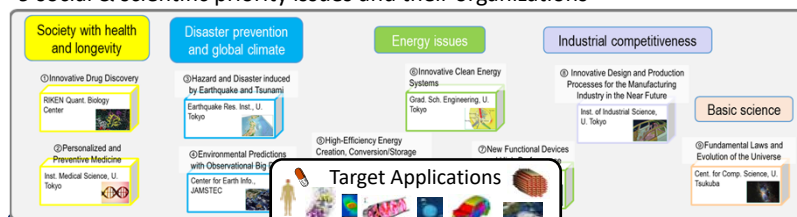


5

Co-design: *How we are running under the concept*



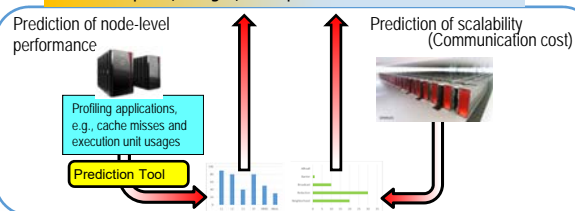
9 social & scientific priority issues and their organizations



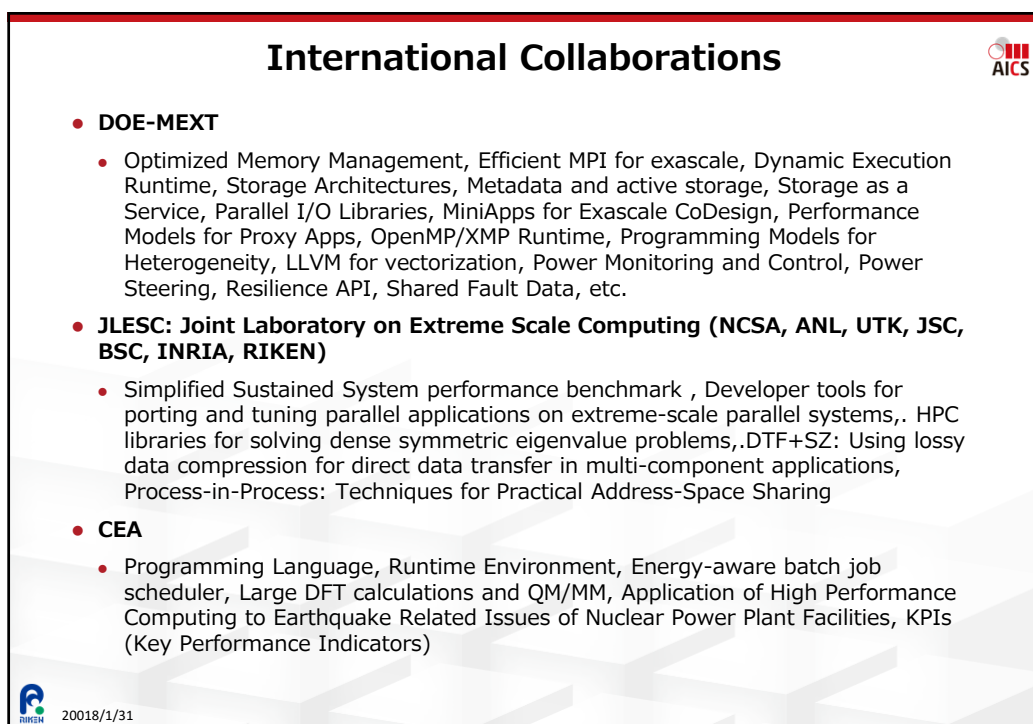
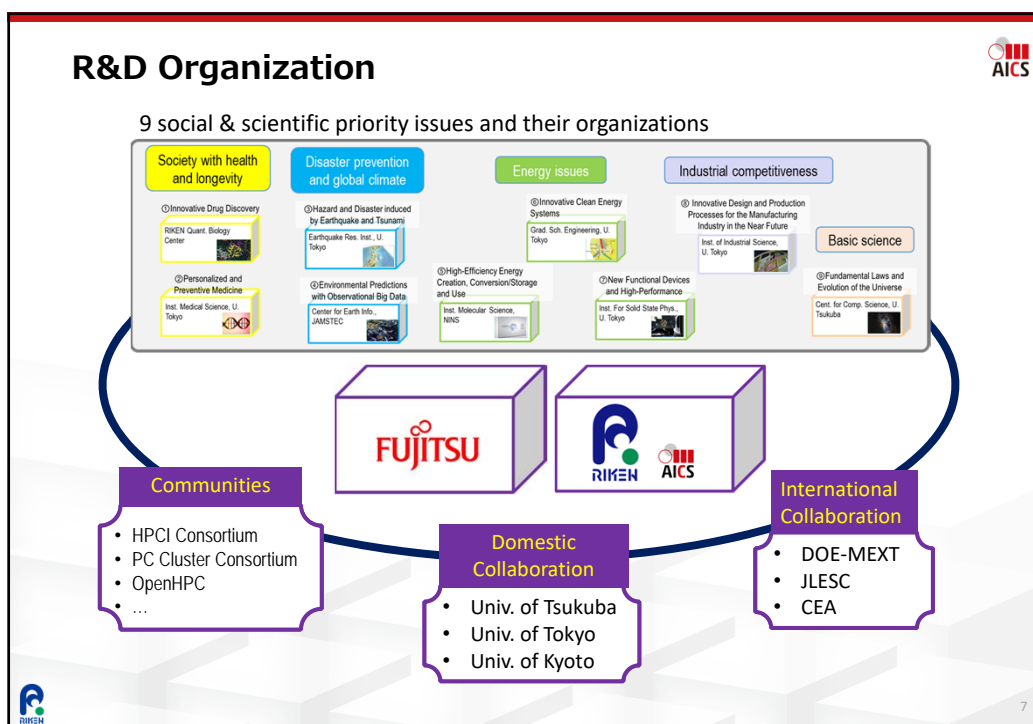
Architectural Parameters

- #SIMD, SIMD length, #core
- cache (size and bandwidth)
- memory technologies
- specialized hardware
- Interconnect
- I/O network

- ❑ Mutual understanding both computer architecture/system software and applications
- ❑ Looking at performance predicted by the following method
- ❑ Finding out the best solution with constraints, e.g., power consumption, budget, and space



6



An Overview of Post K



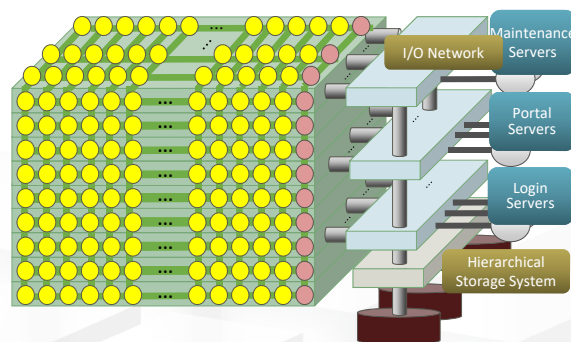
- **Compute Node, Compute + I/O Node**

- Armv8-A SVE + Fujitsu Extension

- **6D mesh/torus Interconnect**

- **3-level hierarchical storage system**

- Silicon Disk
- Magnetic Disk
- Storage for archive



20018/1/31

9

An Overview of System Software Stack



Easy of use is one of our KPIs (Key Performance Indicators)

Providing wide range of applications/tools/libraries/compilers

Linux Distribution

Fortran, C/C++, OpenMP, Java, ...

Math libraries

Tuning and Debugging Tools

Parallel Programming Environments
XMP, FDPS, ...

Communication
MPI

Application-oriented
File I/O

Process/Thread

Low Level Communication

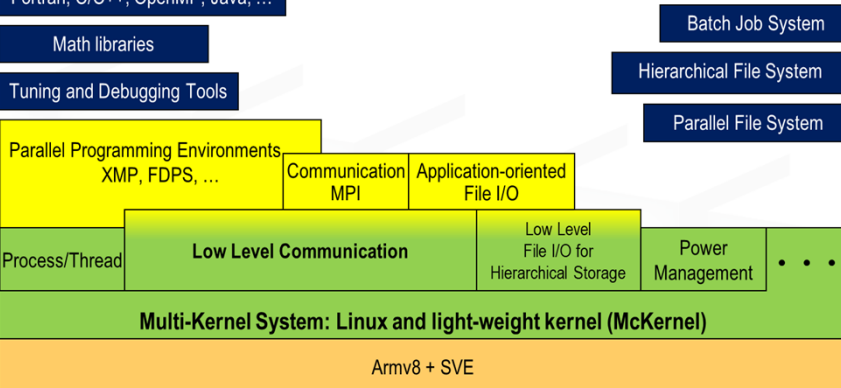
Low Level
File I/O for
Hierarchical Storage

Power
Management

...

Multi-Kernel System: Linux and light-weight kernel (McKernel)

Armv8 + SVE



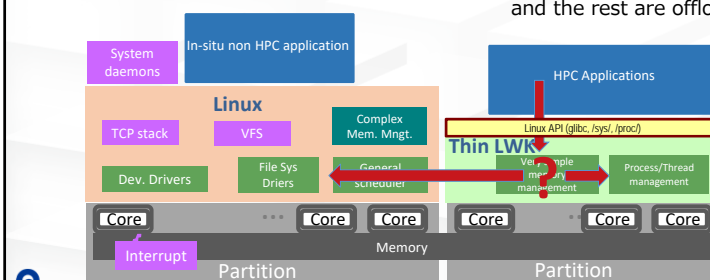
20018/1/31

10

IHK/McKernel developed at RIKEN



- **Partition resources (CPU cores, memory)**
- **Full Linux kernel on some cores**
 - System daemons and in-situ non HPC applications
 - Device drivers
- **Light-weight kernel(LWK), McKernel on other cores**
 - HPC applications
- **IHK: Linux kernel module**
Interface for Heterogeneous Kernels
 - Allows dynamically partitioning of node resources: CPU cores, physical memory, ...
 - Enables management of LWKs (assign resources, load, boot, destroy, etc..)
 - Provides inter-kernel communication, messaging and notification
- **McKernel: Light-weight kernel**
 - Is designed for HPC, noiseless, simple
 - Implements only performance sensitive system calls, e.g., process and memory management, and the rest are offloaded to Linux



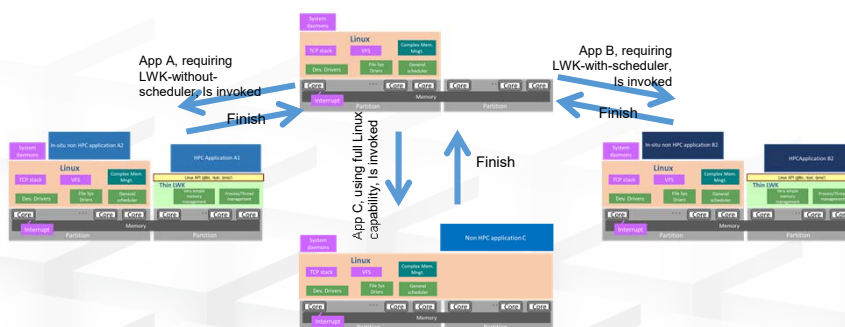
- Executes the same binary of Linux without any recompilation

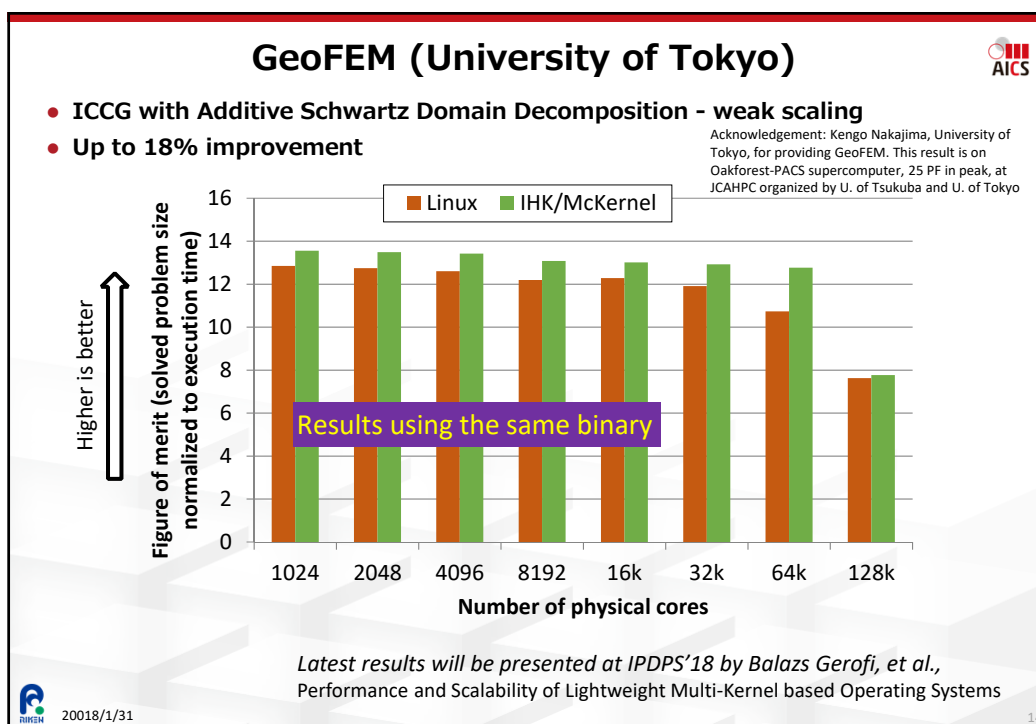
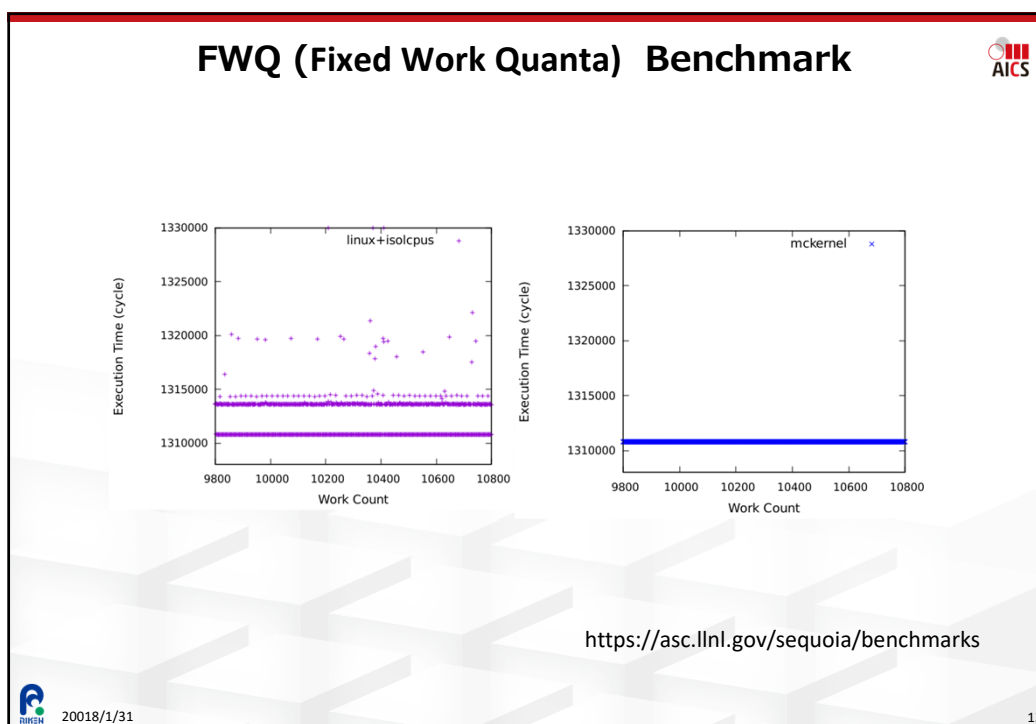
- IHK/McKernel runs on
 - Intel Xeon and Xeon phi
 - Fujitsu FX10 and FX100 (Experiments)

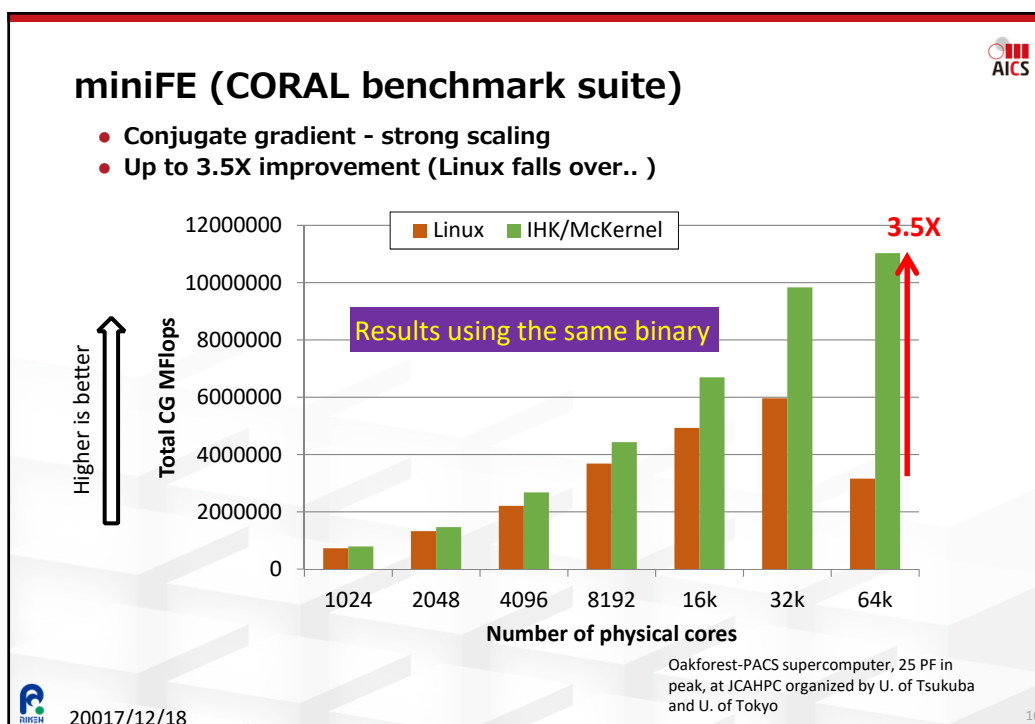
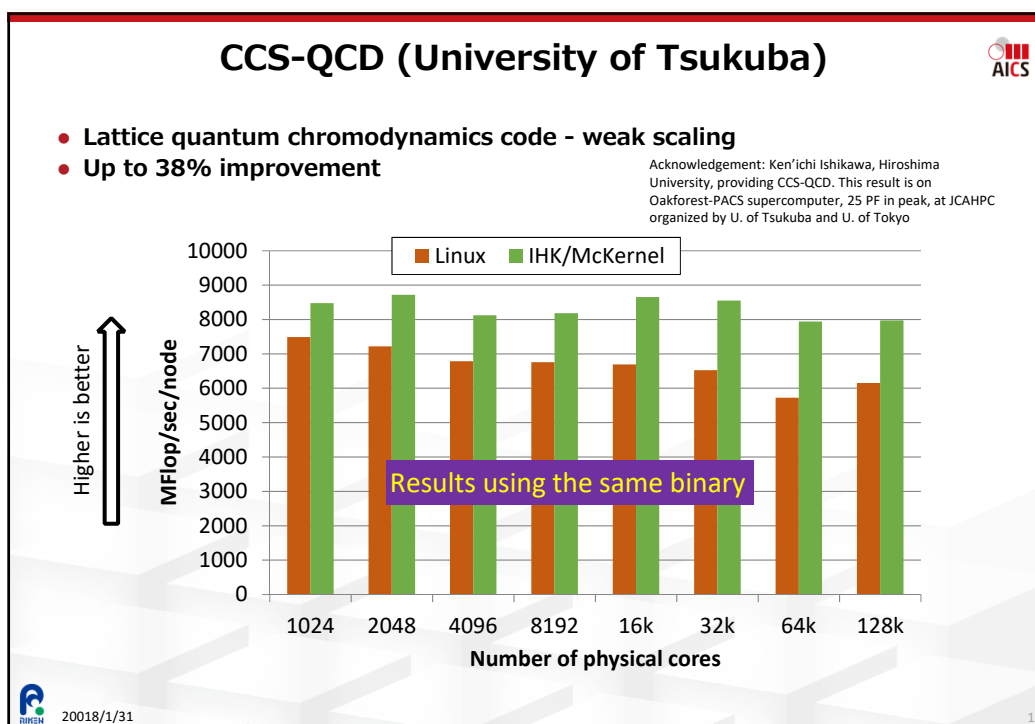
How to deploy IHK/McKernel



- Linux Kernel with IHK kernel module is resident
 - daemons for job scheduler and etc. run on Linux
- McKernel is dynamically reloaded (rebooted) by IHK for each application
 - No hardware reboot



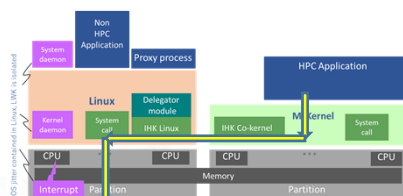




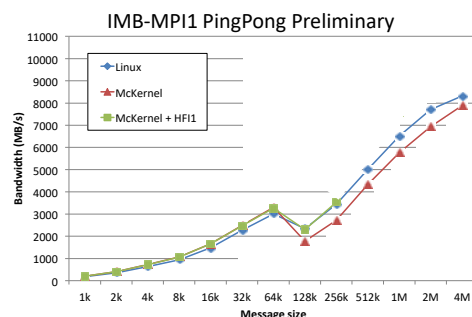
Ongoing R&D: Device driver integration



- In the current Omnipath, the `writew` systemcall is issued in large message sizes in order to utilize the RDMA capability.



This causes extra overhead, and thus communication latency and bandwidth is worse than Linux in the large message sizes.



- A partial Linux driver is ported to McKernel with less effort

- ✓ Making McKernel access the same kernel virtual address of Linux
- ✓ Providing the same Linux internal functions used by a device driver, such as `kmalloc` and synchronization primitives



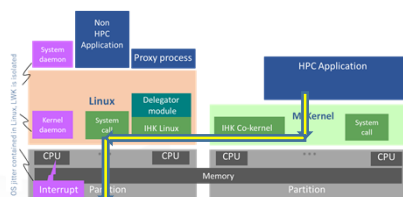
20018/1/31

17

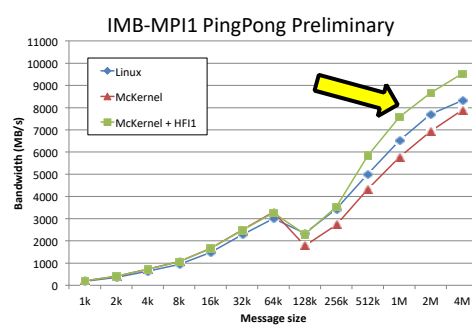
Ongoing R&D: Device driver integration



- In the current Omnipath, the `writew` systemcall is issued in large message sizes in order to utilize the RDMA capability.



This causes extra overhead, and thus communication latency and bandwidth is worse than Linux in the large message sizes.



- A partial Linux driver is ported to McKernel with less effort

- ✓ Making McKernel access the same kernel virtual address of Linux
- ✓ Providing the same Linux internal functions used by a device driver, such as `kmalloc` and synchronization primitives
- ✓ ...



20018/1/31

18

Concluding Remarks



- The system software stack for post-K is being designed and implemented with the leverage of international collaborations, CEA, DOE Labs, and JLESC (NCSA, INRIA, ANL, BSC, JSC, RIKEN)



- The software stack developed at RIKEN is open source
- It also runs on Intel Xeon and Xeon phi



20018/1/31

19