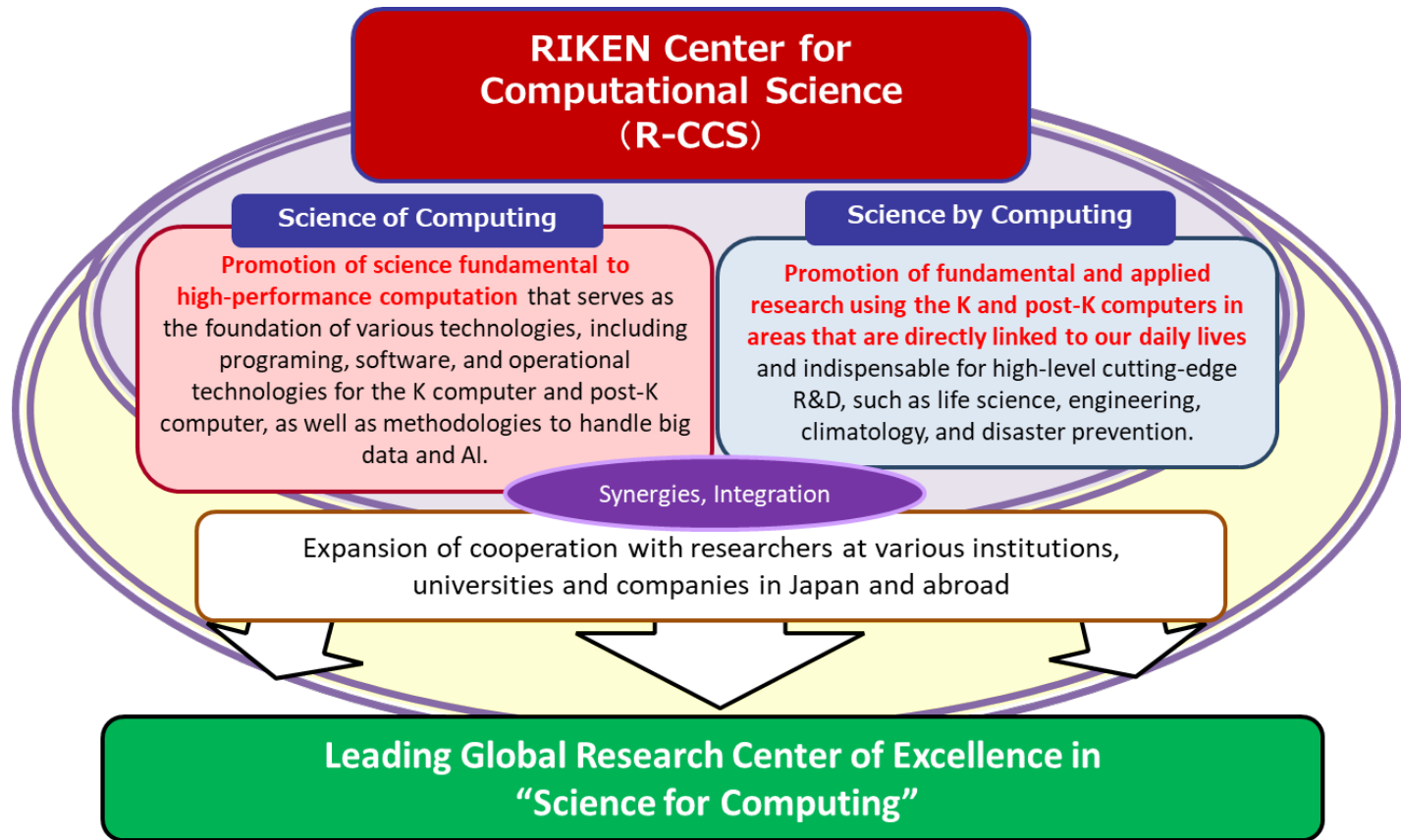# A64fx and Fugaku - A Game Changing, HPC / AI Optimized Arm CPU to enable Exascale Performance



- **Satoshi Matsuoka**
  - **Director, RIKEN Center for Computational Science & Professor, Tokyo Institute of Technology**
- **20190925 Linaro Connect @ San Diego, USA**

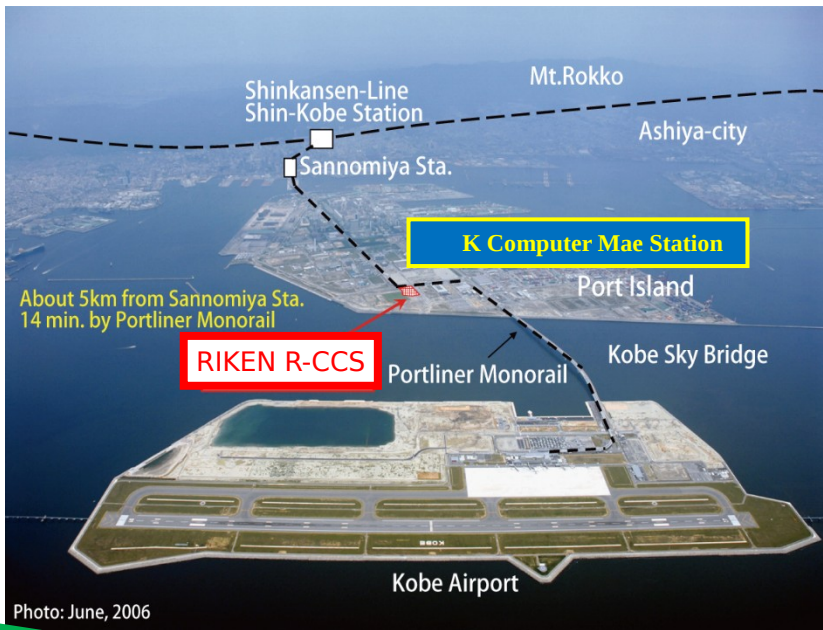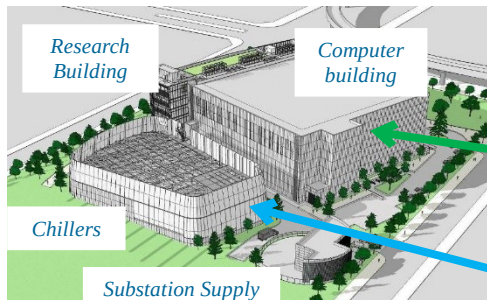# Riken R-CCS: Leadership HPC Research Center "Science of Computing by Computing for Computing"

# R-CCS with K Computer



Kobe

Tokyo

**423 km (263 miles)
west of Tokyo**

Research Building

Computer building

Chillers

Substation Supply

Mt.Rokko

Shinkansen-Line
Shin-Kobe Station

Sannomiya Sta.

Ashiya-city

**K Computer Mae Station**

About 5km from Sannomiya Sta.
14 min. by Portliner Monorail

RIKEN R-CCS

Portliner Monorail

Port Island

Kobe Sky Bridge

Photo: June, 2006

Kobe Airport

**Computer room  50 m x 60 m = 3,000 m²**
**Electric power up to 15 MW**
**Water cooling system**

**Gas-turbine co-generation 5 MW x 2**

# K computer (decommissioned Aug. 16, 2019)

## Specifications
- Massively parallel, general purpose supercomputer
- No. of nodes:  88,128
- Peak speed:     11.28 Petaflops
- Memory:        1.27 PB
- Network: 6-dim mesh-torus (Tofu)

## Top 500 ranking
LINPACK measures the speed and efficiency of linear equation calculations.
Real applications require more complex computations.
- No.1 in Jun. & Nov. 2011
- No.20 in Jun 2019

*First supercomputer in the world to retire as #1 in major rankings (Graph 500)*

## Graph 500 ranking
"Big Data" supercomputer ranking
Measures the ability of data-intensive
- No. 1 for 9 consecutive editions since 2015

## HPCG ranking
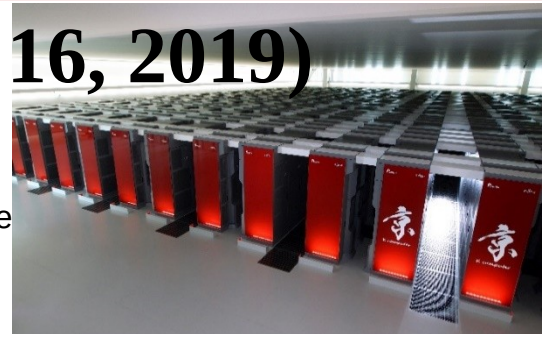Measures the speed and efficiency of solving linear equation using HPCG
Better correlate to actual applications
- No. 1 in Nov. 2017, No. 3 since Jun 2018

## ACM Gordon Bell Prize
**"Best HPC Application of the Year"**
- Winner 2011 & 2012. several finalists

# K-Computer Shutdown Ceremony 30 Aug 2019

# The Nex-Gen "Fugaku"
# 富岳 Supercomptuer

*Mt. Fuji representing the ideal of supercomputing*

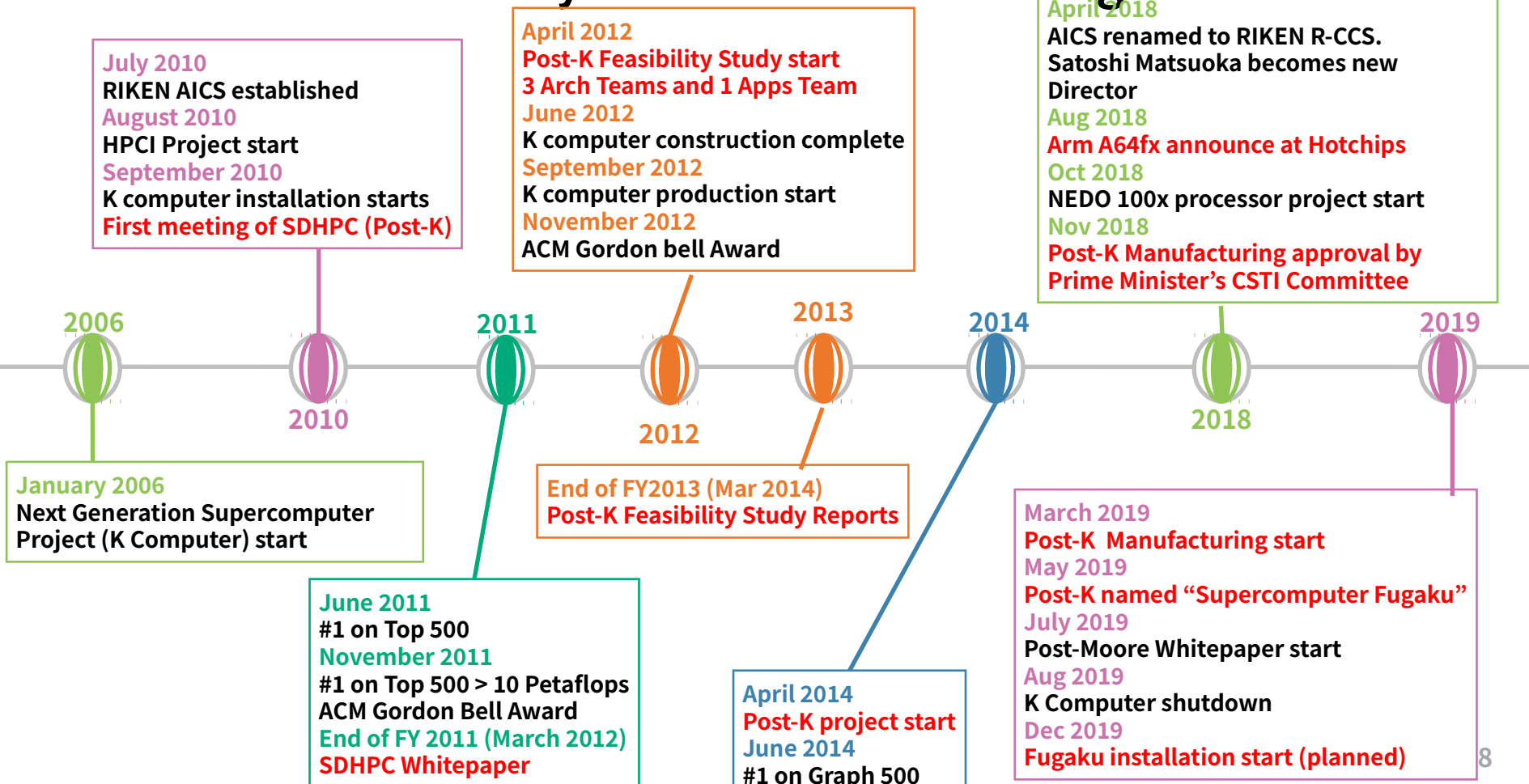High-Peak --- Acceleration of Large Scale Application (Capability)

**Broad Base --- Applicability & Capacity**
**Broad Applications: Simulation, Data Science, AI, …**
**Broad User Bae: Academia, Industry, Cloud Startups, …**

# Arm64fx & Fugaku（富岳）are:

- **Fujitsu-Riken design A64fx ARM v8.2 (SVE), 48/52 core CPU**
  - *HPC Optimized:* Extremely high package high memory BW (1TByte/s), on-die Tofu-D network BW (˜400Gbps), high SVE FLOPS (˜3Teraflops), various AI support (FP16, INT8, etc.)
  - Gen purpose CPU – Linux, Windows (Word), other SCs/Clouds
  - Extremely power efficient – > **10x power/perf efficiency for CFD benchmark** over current mainstream x86 CPU
- **Largest and fastest supercomputer to be ever built circa 2020**
  - > 150,000 nodes, superseding LLNL Sequoia
  - > 150 PetaByte/s memory BW
  - Tofu-D 6D Torus NW, 60 Petabps injection BW (10x global IDC traffic)
  - 25˜30PB NVMe L1 storage
  - ˜10,000 endpoint 100Gbps I/O network into Lustre
  - The first 'exascale' machine (not exa64bitflops but in apps perf.)

# Brief History of R-CCS towards Fugaku

**July 2010**
RIKEN AICS established
**August 2010**
HPCI Project start
**September 2010**
K computer installation starts
**First meeting of SDHPC (Post-K)**

**April 2012**
**Post-K Feasibility Study start**
**3 Arch Teams and 1 Apps Team**
**June 2012**
K computer construction complete
**September 2012**
K computer production start
**November 2012**
ACM Gordon bell Award

**April 2018**
AICS renamed to RIKEN R-CCS.
Satoshi Matsuoka becomes new
Director
**Aug 2018**
**Arm A64fx announce at Hotchips**
**Oct 2018**
NEDO 100x processor project start
**Nov 2018**
**Post-K Manufacturing approval by
Prime Minister's CSTI Committee**

2006   2011   2013   2014   2019

2010   2012   2018

**January 2006**
**Next Generation Supercomputer
Project (K Computer) start**

**End of FY2013 (Mar 2014)**
**Post-K Feasibility Study Reports**

**June 2011**
**#1 on Top 500**
**November 2011**
**#1 on Top 500 > 10 Petaflops
ACM Gordon Bell Award**
**End of FY 2011 (March 2012)**
**SDHPC Whitepaper**

**April 2014**
**Post-K project start**
**June 2014**
**#1 on Graph 500**

**March 2019**
**Post-K  Manufacturing start**
**May 2019**
**Post-K named "Supercomputer Fugaku"**
**July 2019**
Post-Moore Whitepaper start
**Aug 2019**
K Computer shutdown
**Dec 2019**
**Fugaku installation start (planned)**

# Co-Design Activities in Fugaku

Multiple Activities since 2011

## Science by Co-Authoring Computing
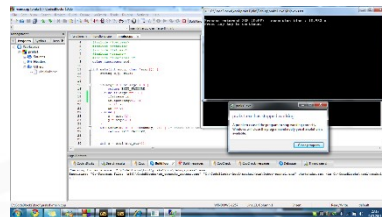
· 9 Priority App Areas High Concern to General Public: Medical/Pharma, Environment/Disaster, Energy, Manufacturing, …
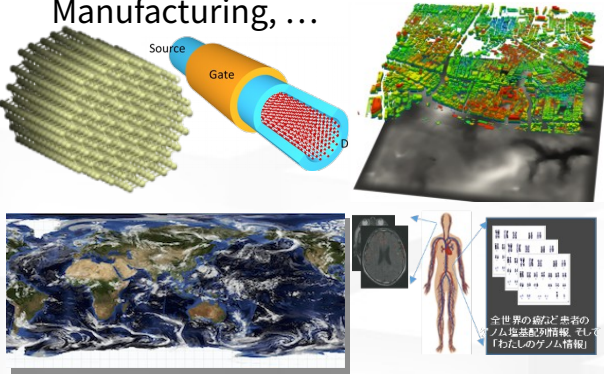
## Science of Computing

**A 6 4 f x**

For the Post-K supercomputer

Select representatives from 100s of applications signifying various computational characteristics

Design systems with parameters that consider various application characteristics

- Extremely tight collaborations between the Co-Design apps centers, Riken, and Fujitsu, etc.
- Chose 9 representative apps as "target application" scenario
- Achieve up to x100 speedup c.f. K-Computer
- Also ease-of-programming, broad SW ecosystem, very low power, …

# Research Subjects of the Post-K Computer

**The post K computer will expand the fields pioneered by the K computer, and also challenge new areas.**
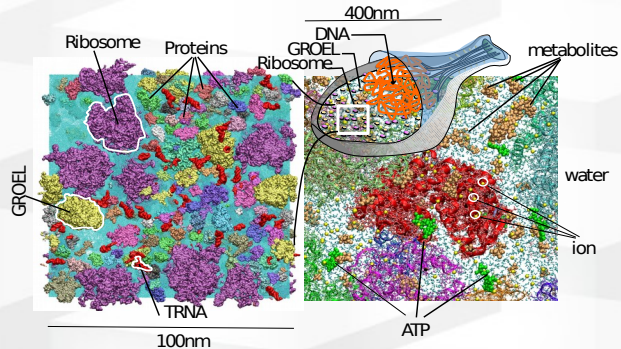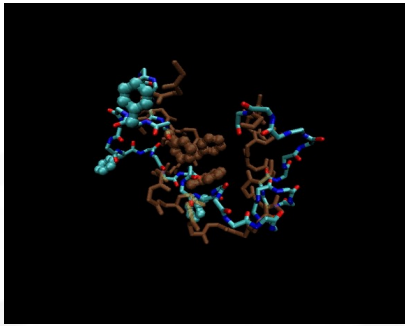
# Genesis MD: proteins in a cell environment
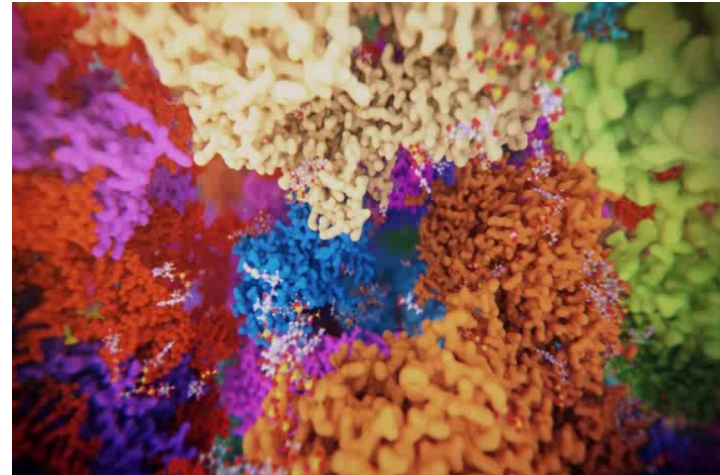
## Protein simulation before K

- Simulation of a protein in isolation

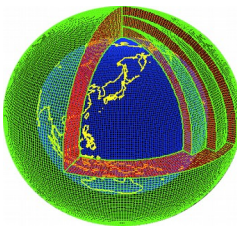Folding simulation of Villin, a small protein with 36 amino acids





## Protein simulation with K

- all atom simulation of a cell interior
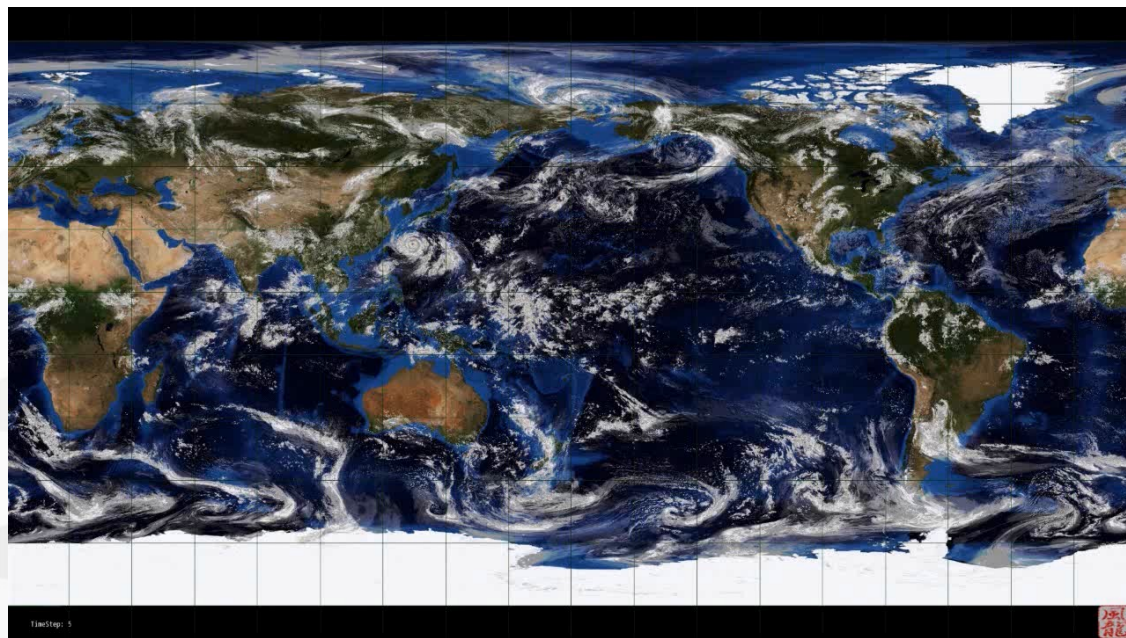- cytoplasm of Mycoplasma genitalium

# NICAM: Global Climate Simulation

- Global cloud resolving model **with 0.87 km-mesh** which allows resolution of cumulus clouds
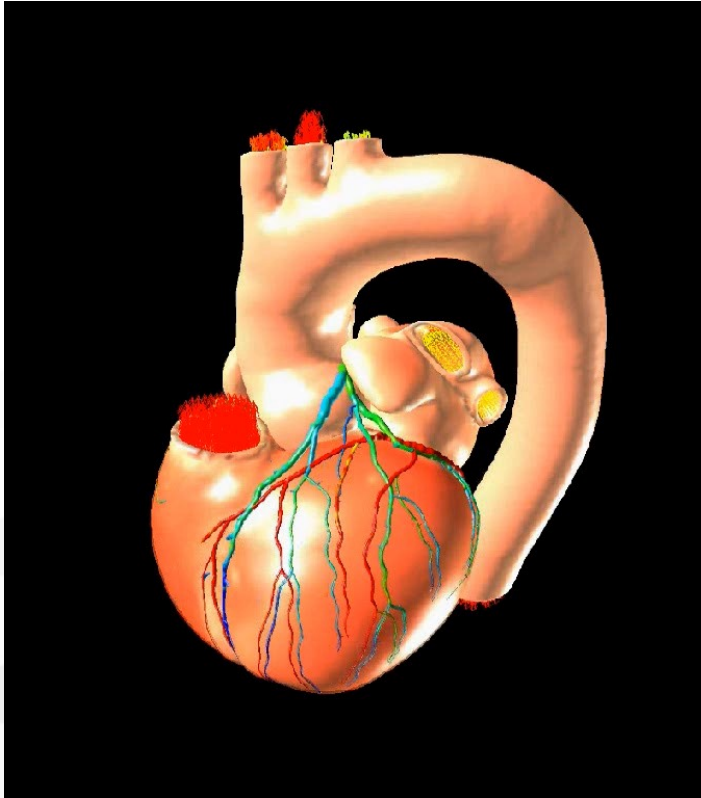- Month-long forecasts of Madden-Julian oscillations in the tropics is realized.
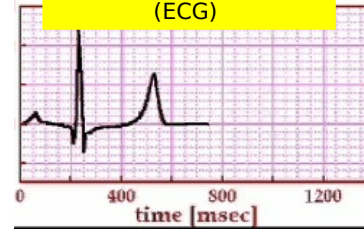
Global cloud resolving model



Miyamoto et al (2013) , Geophys. Res. Lett., 40, 4922–4926, doi:10.1002/grl.50944.
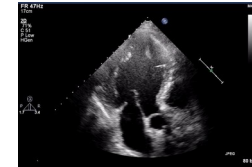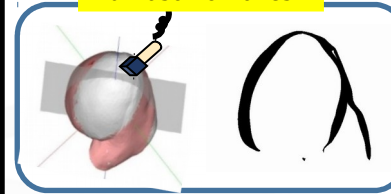
# Heart Simulator



electrocardiogram (ECG)



ultrasonic waves



- Multi-scale simulator of heart starting from molecules and building up cells, tissues, and heart

- Heartbeat, blood ejection, coronary circulation are simulated consistently.

- Applications explored
  - congenital heart diseases
  - Screening for drug-induced irregular heartbeat risk
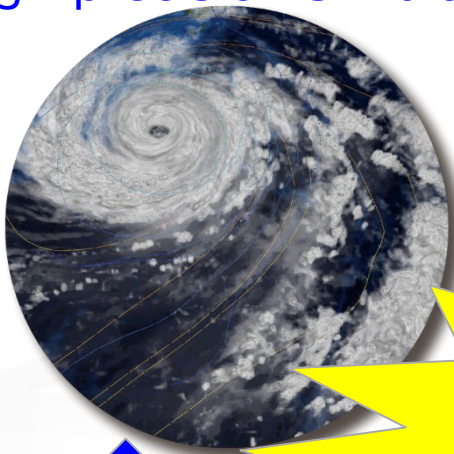
UT-Heart, Inc., Fujitsu Limited

# Fugaku: The Game Changer

**1. Heritage of the K-Computer, HP in simulation via extensive Co**

- High performance: up to x100 performance of K in real applications
- Retain BYTES/FLOP of K (0.4~0.5) for real application performance
- Simultaneous high performance and ease-of-programming

**2. New Technology Innovations of Fugaku**

- **High Performance, esp. via high memory BW**

  Performance boost by "factors" c.f. mainstream CPUs in many HPC & Society5.0 apps via <u>BW & Vector acceleration</u>

- **Very Green e.g. extreme power efficiency**

  Ultra Power efficient design & various power control knobs

- **Arm Global Ecosystem & SVE contribution**

  Top CPU in ARM Ecosystem of 21 billion chips/year, SVE co-design and world's first implementation by Fujitsu

- **High Perf. on Society5.0 apps incl. AI**

  Architectural features for high perf on Society 5.0 apps based on Big Data, AI/ML, CAE/EDA, Blockchain security, etc.

**Global leadership not just in the machine & apps, but as cutting edge IT**

FUJITSU A64FX™

ARM: Massive ecosystem from embedded to HPC

*...ogy not just limited to Fugaku, but into societal IT infrastructures e.g. C*

# A64FX Leading-edge Si-technology



- **TSMC 7nm FinFET & CoWoS**
  - Broadcom SerDes, HBM I/O, and S RAMs
  - 8.786 billion transistors
  - 594 signal pins

# A64FX technologies: Scalable architecture

- Core Memory Group (CMG)

  - 12 compute cores for computing and
     an assistant core for OS daemon, I/O, etc.

  - Shared L2 cache

  - Memory controller

- Four CMGs maintain cache coherence w/ on-chip directory

  - Threads binding within a CMG allows linear speed up of cores' performance

**CMG configuration**

| core | core | core | core | core | core | core |
| core | core | core | core | core | core | |

**X-Bar connection**

**L2 cache 8MiB 16-way**

**Memory controller**

**HBM2**  **Network on chip**

**A64FX chip configuration**

| Tofu controller | PCIe controller |

HBM2 ↔ CMG | Network on chip | CMG ↔ HBM2

HBM2 ↔ CMG | | CMG ↔ HBM2

# Fugaku's FUjitsu A64fx Processor is···

- **an Many-Core ARM CPU···**
  - 48 compute cores + 2 or 4 assistant (OS) cores
  - Brand new core design
  - Near Xeon-Class Integer performance core
  - ARM V8 --- 64bit ARM ecosystem
  - Tofu-D + PCIe 3 external connection



- **···but also an accelerated GPU-like processor**
  - SVE 512 bit x 2 vector extensions (ARM & Fujitsu)
    - Integer (1, 2, 4, 8 bytes) + Float (16, 32, 64 bytes)
  - Cache + scratchpad-like local memory (sector cache)
  - HBM2 on package memory – Massive Mem BW (Bytes/DPF ˜0.4)
    - Streaming memory access, strided access, scatter/gather etc.
  - Intra-chip barrier synch. and other memory enhancing features

- **World's first implementation of SVE, high performance, low power**

18

# Himeno Benchmark (Fortran90)

- Stencil calculation to solve Poisson's equation by Jacobi method



† "Performance evaluation of a vector supercomputer SX-aurora TSUBASA",
SC18,  https://dl.acm.org/citation.cfm?id=3291728

# "Fugaku" CPU Performance Evaluation (3/3)

**FUJITSU**

- **WRF: Weather Research and Forecasting model**
  - Vectorizing loops including IF-constructs is key optimization
  - Source code tuning using directives promotes compiler optimizations



WRF v3.8.1 (48-hour,12km, CONUS) on 48 cores

×1.32

×1.56

1

Intel Xeon Platinum 8168 2 CPUs | Fugaku A64FX 1 CPU (asis) | Fugaku A64FX 1 CPU (w/ src tuning)

1 / (Iteration time) Normalized by Xeon

W 800 mm
D1400 mm
H2000 mm
384 nodes

230 mm

CMU

280 mm

60 mm

60 mm

**CPU Package**

FUJITSU

**CPU**

A64fx

**A0 Chip Booted in June Undergoing Tests**

# A64FX: Tofu interconnect D

■ Integrated w/ rich resources

- ■ Increased TNIs achieves higher injection BW & flexible comm. patterns
- ■ Increased barrier resources allow flexible collective comm. algorithms

■ Memory bypassing achieves low latency

- ■ Direct descriptor & cache injection

| | TofuD spec |
|---|---|
| Port bandwidth | 6.8 GB/s |
| Injection bandwidth | 40.8 GB/s |
| **Measured** | |
| Port throughput | 6.35 GB/s |
| Ping pong latency | 0.49~0.54 μs |

# CMU: CPU Memory Unit

- A64FX CPU x2 (Two independent nodes)
- QSFP28 x3 for Active Optical Cables
- Single-side blind mate connectors of signals &
- ~100% direct water coo

A64FX™

Water

Water

AOC    QSFP28 (Z)

AOC    QSFP28 (Y)

AOC    QSFP28 (X)

Electrical signals

# Fugaku system configuration

■ Scalable design



CPU    CMU    BoB    Shelf    Rack    System

| Unit | # of nodes | Description |
| --- | --- | --- |
| CPU | 1 | Single socket node with HBM2 & Tofu interconnect D |
| CMU | 2 | CPU Memory Unit: 2x CPU |
| BoB | 16 | Bunch of Blades: 8x CMU |
| Shelf | 48 | 3x BoB |
| Rack | 384 | 8x Shelf |
| System | 150k+ | As a Fugaku system |

# "Fugaku" Chronology

*(Disclaimer: below includes speculative schedules and subject to change)*

- May 2018 A0 Chip came out, almost bug free
- 1Q2019 B0 Chip on hand, bug free, exceeded perf. target
- Mar 2019 "Fugaku" manufacturing budget approval by the Diet, actual manufacturing contract signed **(now w/Society 5.0 AI mission also)**
- Aug 2019 End of K-Computer operations
- 4Q2019 "Fugaku" installation starts
- 1H2020 "Fugaku" preproduction operation starts
- 1˜2Q2021 "Fugaku" production operation starts (hopefully)
- And of course we move on…

# Overview of Fugaku System & Storage

- **3-level hierarchical storage**
  - 1st Layer: GFS Cache + Temp FS **(25~30 PB NVMe)**
  - 2nd Layer: Lustre-based GFS (a few hundred PB HDD)
  - 3rd Layer: Off-site Cloud Storage
- **Full Machine Spec**
  - **>150,000 nodes ~8 million High Perf. Arm v8.2 Cores**
  - **> 150PB/s memory BW**
  - **Tofu-D 10x Global IDC traffic @ 60Pbps**
  - **~10,000 I/O fabric endpoints**
  - **> 400 racks**
  - **~40 MegaWatts Machine+IDC PUE ~ 1.1 High Pressure DLC**
  - **NRE pays off: ~= 15~30 million state-of-the art competing CPU Cores for HPC workloads (both dense and sparse problems)**

# Prepping the 40+MW Facility (actual photo)

# Fugaku Performance Estimate on 9 Co-Design Target Apps

☐ **Performance target goal**

✓ 100 times faster than K for some applications (tuning included)
✓ 30 to 40 MW power consumption

☐ **Peak performance to be achieved**

|  | PostK | K |
|---|---|---|
| Peak DP (double precision) | >400+ Pflops (34x +) | 11.3 Pflops |
| Peak SP (single precision) | >800+ Pflops (70x +) | 11.3 Pflops |
| Peak HP (half precision) | >1600+ Pflops (141x +) | -- |
| Total memory bandwidth | >150+ PB/sec (29x +) | 5,184TB/sec |

☐ **Geometric Mean of Performance Speedup of the 9 Target Applications over the K-Computer**

> 37x+

| Category | Priority Issue Area | Performance Speedup over K | Application | Brief description |
|---|---|---|---|---|
| Health and longevity | 1. Innovative computing infrastructure for drug discovery | 125x + | GENESIS | MD for proteins |
| Health and longevity | 2. Personalized and preventive medicine using big data | 8x + | Genomon | Genome processing (Genome alignment) |
| Disaster prevention and Environment | 3. Integrated simulation systems induced by earthquake and tsunami | 45x + | GAMERA | Earthquake simulator (FEM in unstructured & structured grid) |
| Disaster prevention and Environment | 4. Meteorological and global environmental prediction using big data | 120x + | NICAM+LETKF | Weather prediction system using Big data (structured grid stencil & ensemble Kalman filter) |
| Energy issue | 5. New technologies for energy creation, conversion / storage, and use | 40x + | NTChem | Molecular electronic simulation (structure calculation) |
| Energy issue | 6. Accelerated development of innovative clean energy systems | 35x + | Adventure | Computational Mechanics System for Large Scale Analysis and Design (unstructured grid) |
| Industrial competitiveness enhancement | 7. Creation of new functional devices and high-performance materials | 30x + | RSDFT | Ab-initio simulation (density functional theory) |
| Industrial competitiveness enhancement | 8. Development of innovative design and production processes | 25x + | FFB | Large Eddy Simulation (unstructured grid) |
| Basic science | 9. Elucidation of the fundamental laws and evolution of the universe | 25x + | LQCD | Lattice QCD simulation (structured grid Monte Carlo) |

*As of 2019/05/14*

# Fugaku Programming Environment

- **Programing Languages and Compilers provided by Fujitsu**
  - Fortran2008 & Fortran2018 subset
  - C11 & GNU and Clang extensions
  - C++14 & C++17 subset and GNU and Clang extensions
  - OpenMP 4.5 & OpenMP 5.0 subset
  - Java
    - GCC and LLVM will be also available
- **Parallel Programming Language & Domain Specific Library provided by RIKEN**
  - XcalableMP
  - FDPS (Framework for Developing Particle Simulator)
- **Process/Thread Library provided by RIKEN**
  - PiP (Process in Process)

- **Script Languages provided by Linux distributor**
  - E.g., Python+NumPy, SciPy
- **Communication Libraries**
  - MPI 3.1 & MPI4.0 subset
    - Open MPI base (Fujitsu), MPICH (RIKEN）
  - Low-level Communication Libraries
    - uTofu (Fujitsu), LLC(RIKEN）
- **File I/O Libraries provided by RIKEN**
  - Lustre
  - pnetCDF, DTF, FTAR
- **Math Libraries**
  - BLAS, LAPACK, ScaLAPACK, SSL II（Fujitsu）
  - EigenEXA, Batched BLAS（RIKEN）
- **Programming Tools provided by Fujitsu**
  - Profiler, Debugger, GUI
- **NEW: Containers (Singularity) and other Cloud APIs**
- **NEW: AI software stacks (w/ARM)**
- **NEW: DoE Spack Package Manager**

# OSS Application Porting @ Arm HPC Users Group

(http://arm-hpc.gitlab.io/)

| Application | Lang. | GCC | LLVM | Arm | Fujitsu |
|---|---|---|---|---|---|
| LAMMPS | C++ | Modified | Modified | Modified | Modified |
| GROMACS | C | Modified | Modified | Modified | Modified |
| GAMESS* | Fortran | Modified | Modified | Modified | Modified |
| OpenFOAM | C++ | Modified | Modified | Modified | Modified |
| NAMD | C++ | Modified | Modified | Modified | Modified |
| WRF | Fortran | Modified | Modified | Modified | Modified |
| Quantum ESPRESSO | Fortran | Ok in as is | Ok in as is | Ok in as is | Modified |
| NWChem | Fortran | Ok in as is | Modified | Modified | Modified |
| ABINIT | Fortran | Modified | Modified | Modified | Modified |
| CP2K | Fortran | Ok in as is | Issues found | Issues found | Modified |
| NEST* | C++ | Ok in as is | Modified | Modified | Modified |
| BLAST* | C++ | Ok in as is | Modified | Modified | Modified |

# Fugaku Cloud Strategy

- **Industry use of Fugaku via intermediary cloud SaaS vendors. Fugaku as IaaS**

- **A64fx and other Fugaku Technology being incorporated into the Cloud**



Industry User 1

Industry User 2

Industry User 3

HPC SaaS Provider 1

HPC SaaS Provider 2

HPC SaaS Provider 3

Other IaaS Commercial Cloud

Extreme Performance Advantage

富岳

Various Cloud Service API for HPC

KVM/ Singularity, Kubernetes,

A64FX

Cloud Vendor 1

Cloud Vendor 2

Cloud Vendor 3

Cloud Workload Becoming HPC (including AI)
↓
Significant Performance Advantage
↓
Millions of Units shipped to Cloud

# A64fx in upcoming Stony Brook Cray System

## National Science Foundation
### WHERE DISCOVERIES BEGIN

SEARCH

| HOME | RESEARCH AREAS | FUNDING | AWARDS | DOCUMENT LIBRARY | NEWS | ABOUT NSF |

**Awards**

Search Awards
Recent Awards
Presidential and Honorary Awards
About Awards

**How to Manage Your Award**
Grant Policy Manual
Grant General Conditions
Cooperative Agreement Conditions
Special Conditions
Federal Demonstration Partnership
Policy Office Website

**Award Abstract #1927880**
**Category II : Ookami: A high-productivity path to frontiers of scientific discovery enabled by exascale system technologies**

| | |
|---|---|
| NSF Org: | OAC<br>Office of Advanced Cyberinfrastructure (OAC) |
| Initial Amendment Date: | July 11, 2019 |
| Latest Amendment Date: | August 29, 2019 |
| Award Number: | 1927880 |
| Award Instrument: | Cooperative Agreement |
| Program Manager: | Robert Chadduck<br>OAC Office of Advanced Cyberinfrastructure (OAC)<br>CSE Direct For Computer & Info Scie & Enginr |
| Start Date: | October 1, 2019 |
| End Date: | September 30, 2024 (Estimated) |
| Awarded Amount to Date: | $2,780,373.00 |
| Investigator(s): | Robert Harrison robert.harrison@stonybrook.edu (Principal Investigator)<br>Barbara Chapman (Co-Principal Investigator)<br>Matthew Jones (Co-Principal Investigator)<br>Alan Calder (Co-Principal Investigator) |
| Sponsor: | SUNY at Stony Brook<br>WEST 5510 FRK MEL LIB<br>Stony Brook, NY 11794-0001 (631)632-9949 |
| NSF Program(s): | Innovative HPC |
| Program Reference Code(s): | |
| Program Element Code(s): | 7619 |

**ABSTRACT**

The State University of New York proposes to procure and operate for at least four years the first computer outside of Japan with the A64fx processor developed by Fujitsu for the Japanese path to exascale computing (i.e., computers capable of 10^18 operations per second). The ARM-based, multi-core, 512-bit SIMD-vector processor with ultrahigh-bandwidth memory promises to retain familiar and successful programming models while achieving very high performance for a wide range of applications including simulation and big data. The testbed significantly extends current NSF-sponsored HPC technologies and will enable the community to evaluate and demonstrate the potential of this technology for deployment in multiple settings. Through integration with NSF's Extreme Science and Engineering Discovery Environment (XSEDE), the system will be widely accessible and fully leverages existing cyber infrastructure including the XDMoD monitoring system.

What does this mean for science? Compared with the best CPUs anticipated during the deployment period, A64fx offers 2-4x better performance on memory-intensive applications such as sparse-matrix solvers found in many engineering and physics codes.

## HPCwire
*Since 1987 - Covering the Fastest Computers in the World and the People Who Run Them*

Home
Technologies
Sectors
AI/ML/DL
Exascale
Specials
Resource Library
Podcast
Events
Job Bank

### Cray ARM-based 'Ookami' to Serve as Testbed for Computational Studies at Stony Brook
August 16, 2019

STONY BROOK, N.Y., August 16, 2019 – A $5 million grant from the National Science Foundation (NSF) to the Institute of Advanced Computational Science (IACS) at Stony Brook University will enable researchers nationwide to test future supercomputing technologies and advance computational and data-driven research on the world's most pressing challenges.

Serving as a testbed for advanced computer technologies, the Ookami system is expected to signal a new generation of high-speed U.S. supercomputers. Using a Cray ARM-based system, Ookami will deliver remarkably high performance for scientific applications, in part due to its blazing-fast memory. Robert J. Harrison, PhD, professor of applied mathematics and statistics and director of IACS, expects that these advanced technologies will enable researchers to more quickly and effectively conduct computational investigations. The project is led by IACS faculty in partnership with co-PI Matt Jones, PhD at the State University of New York at Buffalo, whose team will lead the capture of detailed operational metrics and provision of extensive

## Ookami

- Test bed for NSF researchers
  - First planned deployment of the Post-K processor outside of Japan
- Collaboration with Riken CCS
  - http://www.riken.jp/en/research/labs/r-ccs/
- Installation 3Q 2020
- $5M award NSF OAC 1942140 for purchase and operations

| Node | |
|---|---|
| Processor | A64FX |
| #Cores | 48+4 |
| Peak DP | 2.76 TOP/s |
| Peak INT8 | 22.08 TOP/s |
| Memory | 32GB@1TB/s |
| **System** | |
| #Nodes | 176 |
| Peak DP | 486 TOP/s |
| Peak INT8 | 3886 TOP/s |
| Memory | 5.6 TB |
| Disk | 0.5 PB |
| Comms | IB HDR-100 |

2

# Pursuing Convergence of HPC & AI (1)

- **Acceleration of Simulation (first principles methods) with AI (empirical method) :** *AI for HPC*
  - Interpolation & Extrapolation of long trajectory MD
  - Reducing parameter space on Paretho optimization of results
  - Adjusting convergence parameters for iterative methods etc.
  - AI *replacing* simulation when exact physical models are unclear, or excessively costly to compute
- **Acceleration of AI with HPC:** *HPC for AI*
  - HPC Processing of training data -data cleansing
  - Acceleration of (Parallel) Training: Deeper networks, bigger training sets, complicated networks, high dimensional data···
  - Acceleration of Inference: above + real time streaming data
  - Various modern training algorithms: Reinforcement learning, GAN, Dilated Convolution, etc.
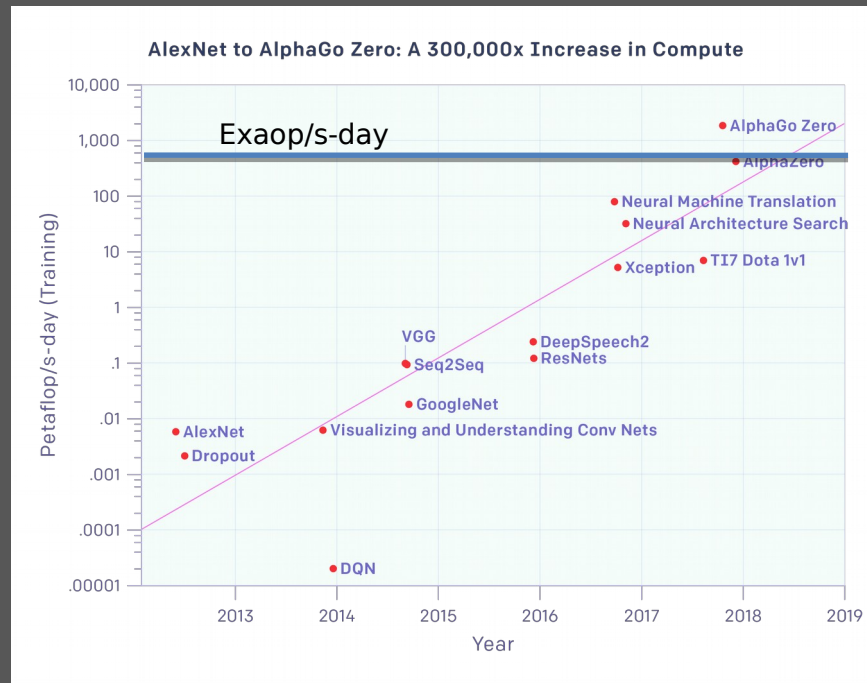
# Deep Learning Meets HPC
## 6 orders of magnitude compute increase in 5 years
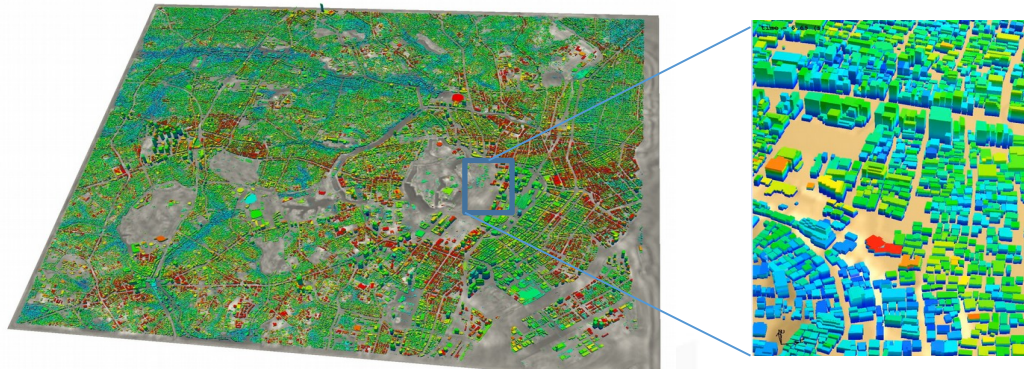[Slide Courtesy Rick Stevens @ ANL]

Exascale Needs for Deep Learning

- Automated Model Discovery
- Hyper Parameter Optimization
- Uncertainty Quantification
- Flexible Ensembles
- Cross-Study Model Transfer
- Data Augmentation
- Synthetic Data Generation
- Reinforcement Learning



**AlexNet to AlphaGo Zero: A 300,000x Increase in Compute**

Exaop/s-day

Data points (top to bottom, by compute): AlphaGo Zero, AlphaZero, Neural Machine Translation, Neural Architecture Search, Xception, TI7 Dota 1v1, DeepSpeech2, ResNets, VGG, Seq2Seq, GoogleNet, Visualizing and Understanding Conv Nets, AlexNet, Dropout, DQN

Y-axis: Petaflop/s-day (Training), 10,000 to .00001
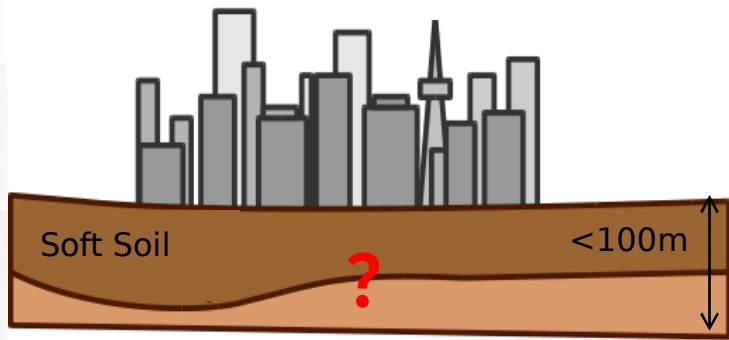X-axis: Year, 2013 to 2019
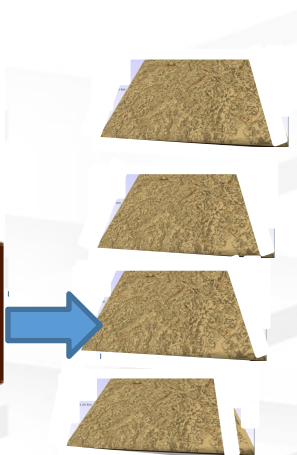
# Large Scale simulation and AI coming together
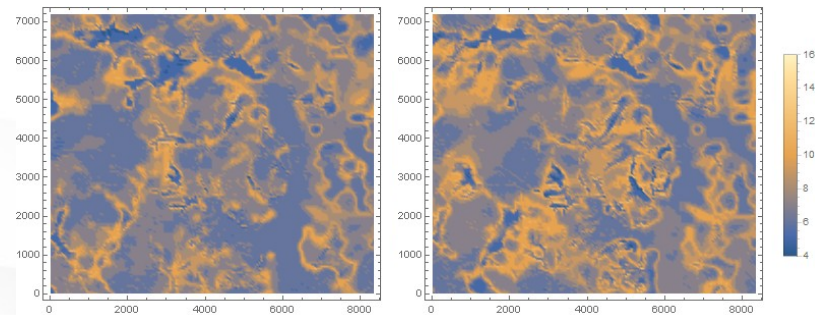[Ichimura et. al. Univ. of Tokyo, IEEE/ACM SC17 Best Poster 2018 Gordon Bell Finalist]



130 billion freedom earthquake of entire Tokyo on K-Computer (2018 ACM Gordon Bell Prize Finalist, SC16,17 Best Poster)



Soft Soil

<100m

?

Earthquake

Too Many Instances

Candidate Underground Structure 1

AI Trained by Simulation to generate candidate soft soil structure

Candidate Underground Structure 2

# 4 Layers of Parallelism in DNN Training

- Hyper Parameter Search
  - Searching optimal network configs & parameters
  - Parallel search, massive parallelism required

- Data Parallelism
  - Copy the network to compute nodes, feed different batch data, average => network reduction bound
  - TOFU: Extremely strong reduction, x6 EDR Infiniband

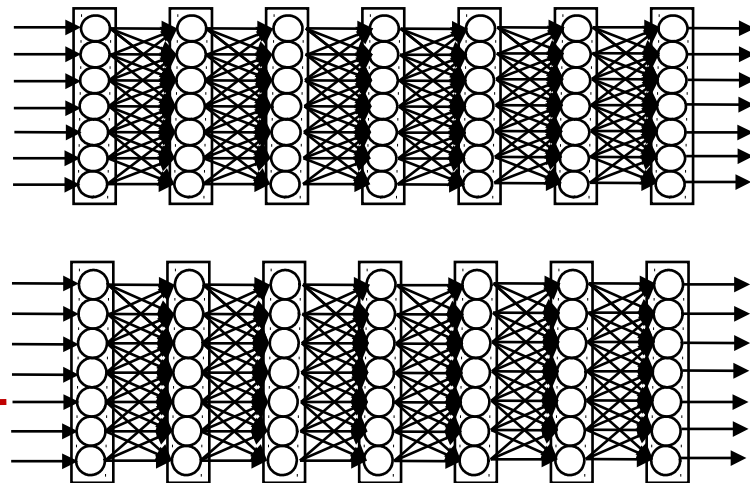**Inter-Node**

- Model Parallelism (domain decomposition)
  - Split and parallelize the layer calculations in propagation
  - Low latency required (bad for GPU) -> strong latency tolerant cores + low latency TOFU network

- Intra-Chip ILP, Vector and other low level Parallelism
  - Parallelize the convolution operations etc.
  - SVE FP16+INT8 vectorization support + extremely high memory bandwidth w/HBM2

**Intra-Node**

- Post-K could become world's biggest & fastest pl
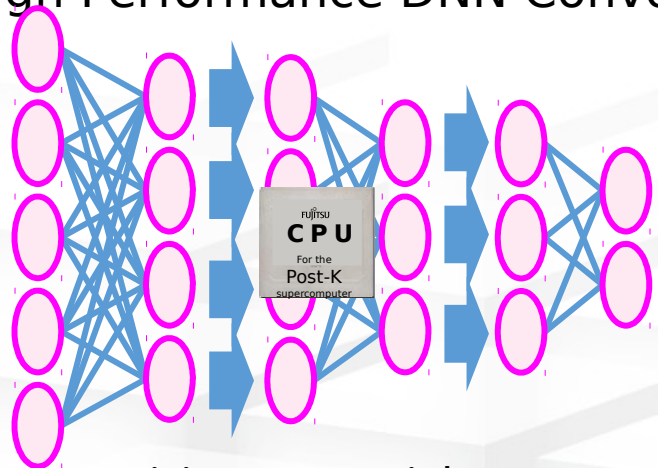
**Massive amount of total parallelism, only possible**

# Massive Scale Deep Learning on Post-K

**Post-K Processor**
- ◆ High perf FP16&Int8
- ◆ **High mem BW for convolut**
- ◆ **Built-in scalable Tofu network**

High Performance DNN Convolution



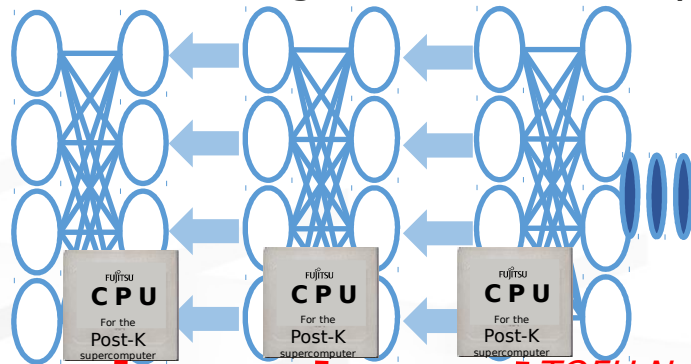Low Precision ALU + High Memory Bandwidth + Advanced Combining of Convolution Algorithms (FFT+Winograd+GEMM)

**Unprecedened DL scalability**

High Performance and Ultra-Scalable Network for massive scaling model & data parallelism



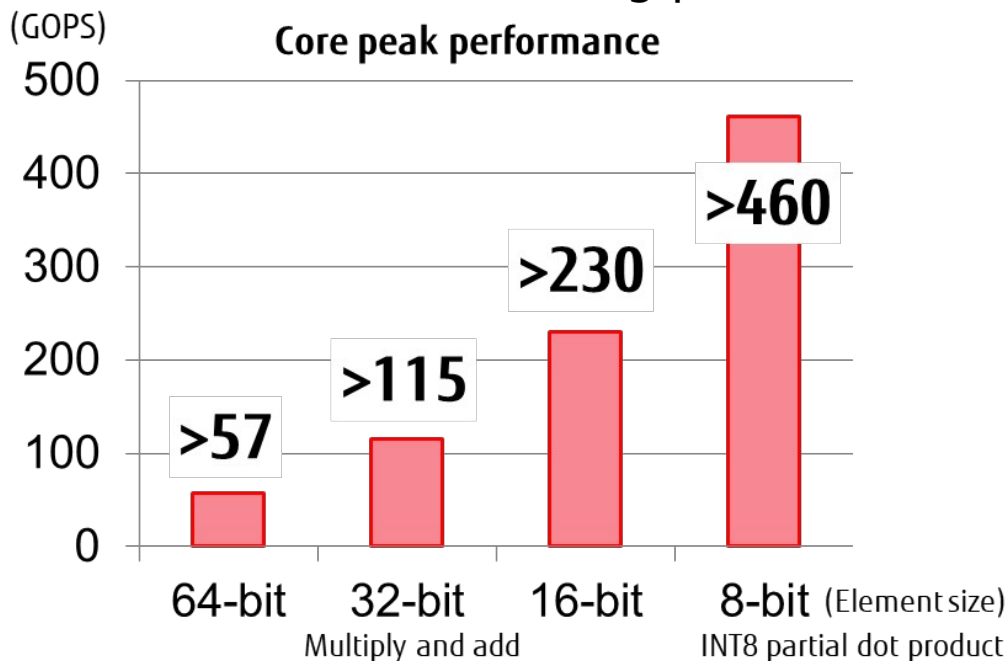*TOFU Network w/ high injection BW for fast reduction*

Unprecedented Scalability of Data/

# A64FX technologies: Core performance

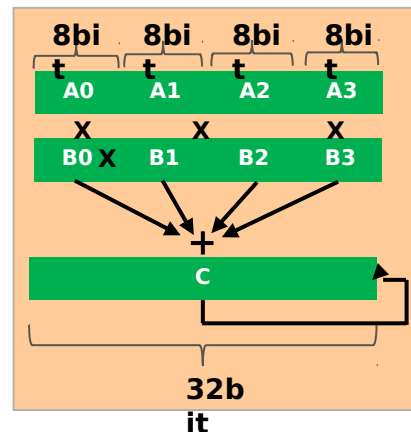High calc. throughput of Fujitsu's original CPU core w/ SVE

- 512-bit wide SIMD x 2 pipelines and new integer functions
- 8 bit inference => training possible?



(GOPS)

**Core peak performance**

| 64-bit | 32-bit | 16-bit | 8-bit | (Element size) |
| >57 | >115 | >230 | >460 | |
| Multiply and add | | INT8 partial dot product | | |

**INT8 partial dot product**

$$C = \Sigma \ (A_i \times B_i) + C$$

# "Isopower" Comparsion with the Best GPU



|  | NVIDIA Volta v100 | Fujitsu A64fx (2 A0 chip nodes) |
|---|---|---|
| **Power** | **400 W (incl. CPUs, HCAs DGX-1)** | **"similar"** |
| **Vectorized MACC Formats** | **FP 64/32/16, INT 32(?)** | **FP 64/32/16, INT 32/16/8 w/INT32 MACC** |
| **Multi-node Linpack** | **5.9 TF / chip (DGX-1)** | **> 5.3 TF / 2 chip blade** |
| **Flops/W Linpack** | **15.1 GFlops/W (DGX-2)** | **> 15 Glops/W** |
| **Stream Triad** | **855 GB/s** | **1.68 TB / s** |
| **Memory Capacity** | **16 / 32 GB** | **64 GB (32 x 2)** |
| **AI Performance** | **125 (peak) / ~95 (measured) Tflops FP16 Tensor Cores** | **~48 TOPS (INT8 MACC peak)** |
| **Price** | **~$11,000 (SXM2 32GB board only) ~$13,000 (DGX-1, per** | **Talk to Fujitsu ☺** |

# Large Scale Public AI Infrastructures in Japan

Inference
838.5PF
Training
86.9 PF

vs. Summit
Inf. 1/4
Train. 1/5

| | Deployed | Purpose | AI Processor | Inference Peak Perf. | Training Peak Perf. | Top500 Perf/ Rank | Green500 Perf/Rank |
|---|---|---|---|---|---|---|---|
| Tokyo Tech. TSUBAME 3 | July 2017 | HPC + AI Public | NVIDIA P100 x 2160 | 45.8 PF (FP16) | 22.9 PF / 45.8PF (FP32/FP16) | 8.125 PF #22 | 13.704 GF/W #5 |
| U-Tokyo Reedbush-H/L | Apr. 2018 (update) | HPC + AI Public | NVIDIA P100 x 496 | 10.71 PF (FP16) | 5.36 PF / 10.71PF (FP32/FP16) | (Unranked) | (Unranked) |
| U-Kyushu ITO-B | Oct. 2017 | HPC + AI Public | NVIDIA P100 x 512 | 11.1 PF (FP16) | 5.53 PF/11.1 PF (FP32/FP16) | (Unranked) | (Unranked) |
| AIST-AIRC AICC | Oct. 2017 | AI Lab Only | NVIDIA P100 x 400 | 8.64 PF (FP16) | 4.32 PF / 8.64PF (FP32/FP16) | 0.961 PF #446 | 12.681 GF/W #7 |
| Riken-AIP Raiden | Apr. 2018 (update) | AI Lab Only | NVIDIA V100 x 432 | 54.0 PF (FP16) | 6.40 PF/54.0 PF (FP32/FP16) | 1.213 PF #280 | 11.363 GF/W #10 |
| AIST-AIRC ABCI | Aug. 2018 | AI Public | NVIDIA V100 x 4352 | 544.0 PF (FP16) | 65.3 PF/544.0 PF (FP32/FP16) | 19.88 PF #7 | 14.423 GF/W #4 |