

This paper describes objective technical results and analysis. Any subjective views or opinions that might be expressed in the paper do not necessarily represent the views of the U.S. Department of Energy or the United States Government.



SAND2020-12188C

Exceptional service in the national interest

EARLY EXPERIENCES WITH A64FX

PRESENTED BY

SIMON D. HAMMOND

SCALABLE COMPUTER ARCHITECTURES
SANDIA NATIONAL LABORATORIES

SDHAMMO@SANDIA.GOV

NOVEMBER 2020

AGENDA

- Overview of A64FX Testbeds at Sandia
- Initial Performance Results
- Update on ATSE Software Environment
- Other talks and resources at SC20



A64FX TESTBEDS AT SANDIA

Sandia has operated the ASC Advanced Architectures testbed collection for almost a decade

- Brings novel systems to the labs to enable testing for NNSA mini-apps, benchmarks and workloads
- Allows the NNSA to gain peeks into potential for future system designs and provides an environment to develop advanced software packages (e.g. Kokkos, Trilinos, etc)
- Names of machines follow themes of “American Firsts” – pioneers in science, engineering and American society

A64FX Testbed at Sandia is called *Inouye*

- Named after Senator Daniel Inouye (1924 – 2012), winner of the Medal of Freedom and the first Japanese American to serve in the US House of Representatives and in the US Senate
- 8 nodes of A64FX interconnected with InfiniBand
- Integrated by Penguin Computing





MEMORY BANDWIDTH TESTING

A64FX processor has 32GB of HBM

- One of the most attractive features of the processor is potential for very high memory bandwidth

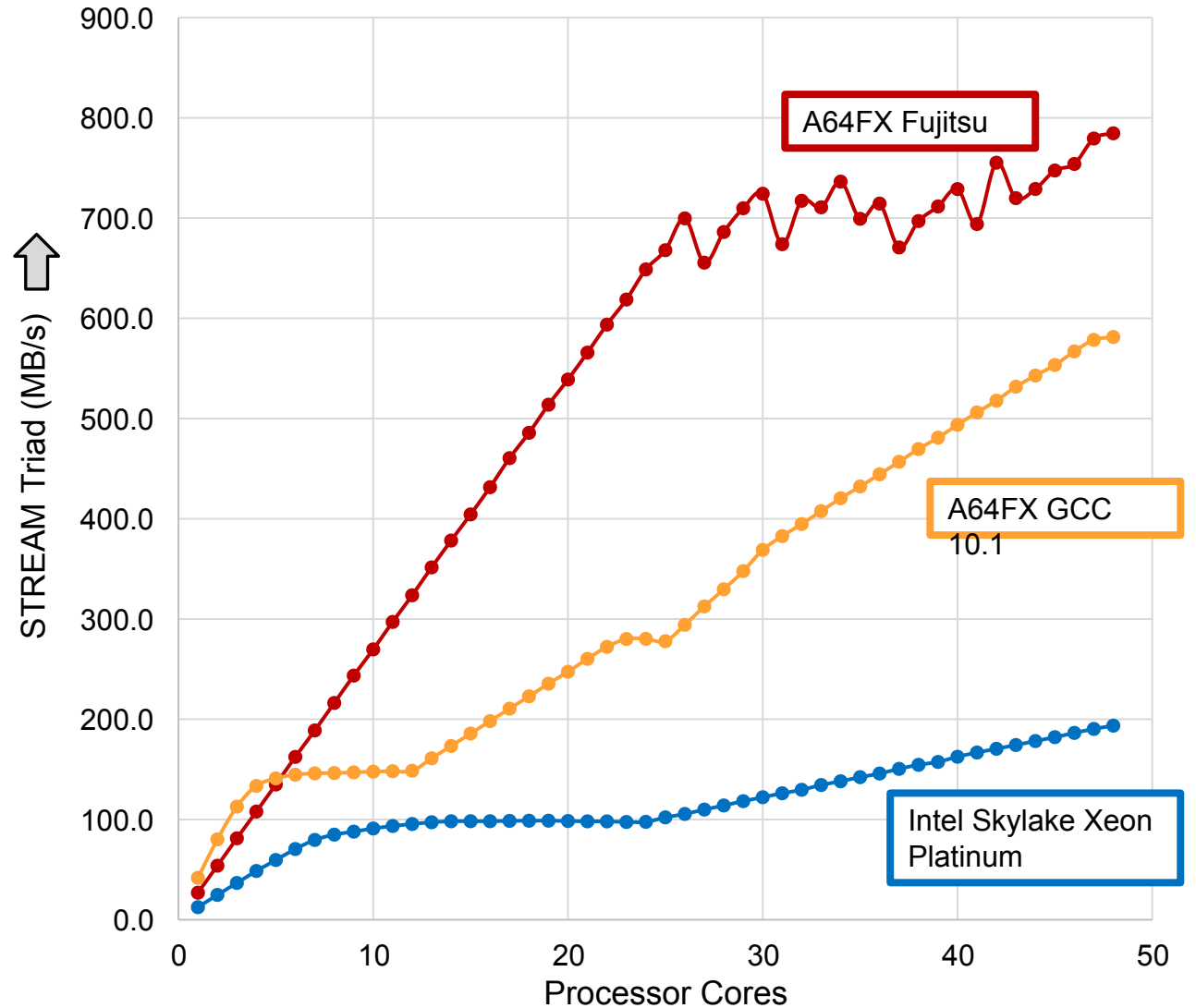
Achieves around 800GB/s in testing (have seen as high as 850GB/s)

- Compares to around 200GB/s for DDR4-based dual-socket Intel Xeon systems (although these have substantially higher capacity)

Fujitsu compilers delivers significant gains over GCC

- Up to 200GB/s (~30%) extra performance

STREAM Bandwidth





LINEAR SOLVERS

Linear solvers traditionally are heavy users of memory bandwidth

- HPCG, MiniFE, AMG/HYPRE etc
- Experience shows HBM can provide significant workload improvements

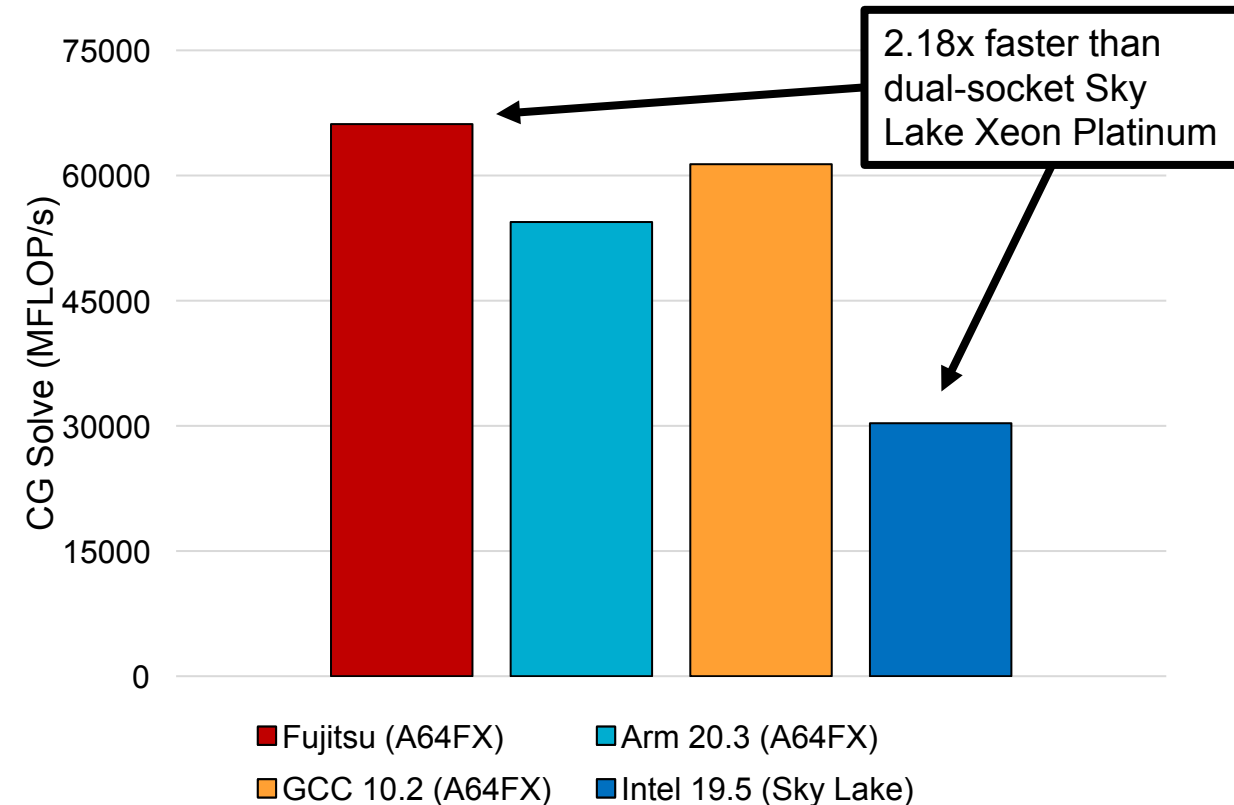
A64FX with HBM is over 2x faster than a dual-socket Sky Lake Xeon with DDR4

- But the problem size that can fit in a node is quite a lot smaller (typically 128 – 256GB on Xeon vs. 32GB on A64FX HBM)

Fujitsu compiler provides best optimization for HBM

- Several compiler options tell the compiler to tune its output for HBM and memory B/W
- Changes cost model for how operands are loaded and which instructions used for writes to memory (see: zfill instructions)

MiniFE CG Solver Mini-App Comparison



Translate into around 1.6X improvement in HPCG performance compared to Sky Lake Xeon nodes



LULESH HYDRODYNAMICS

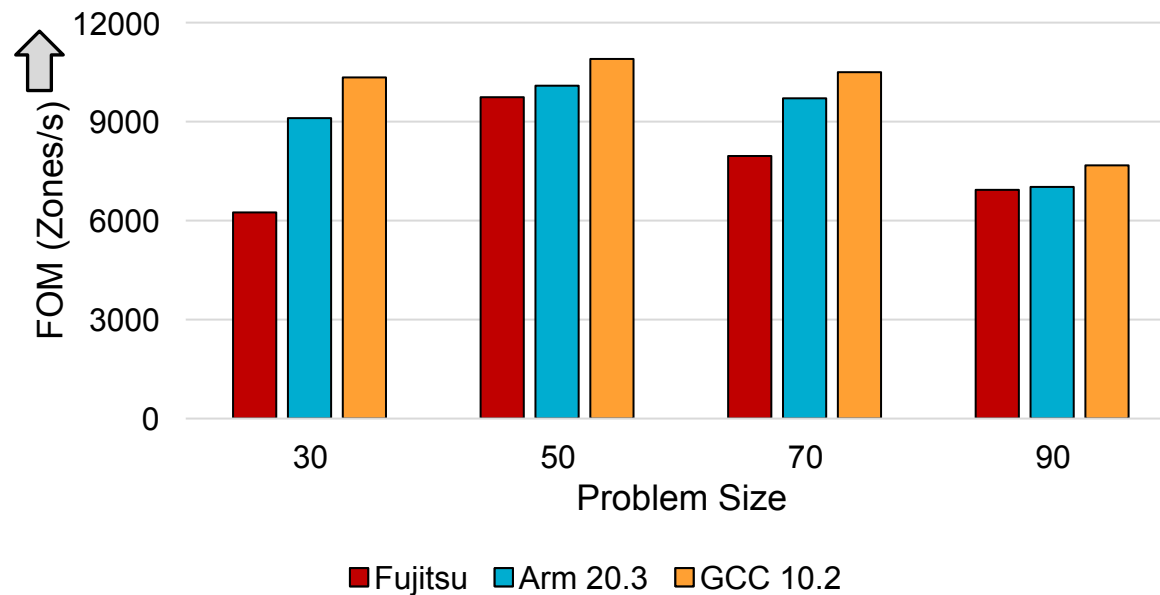
LULESH is an unstructured hydrodynamics benchmark mini-app from LLNL

- Represents unstructured accesses to a mesh in memory
- Forces use of gathers and scatters in vectorized code because elements are not necessarily adjacent in memory

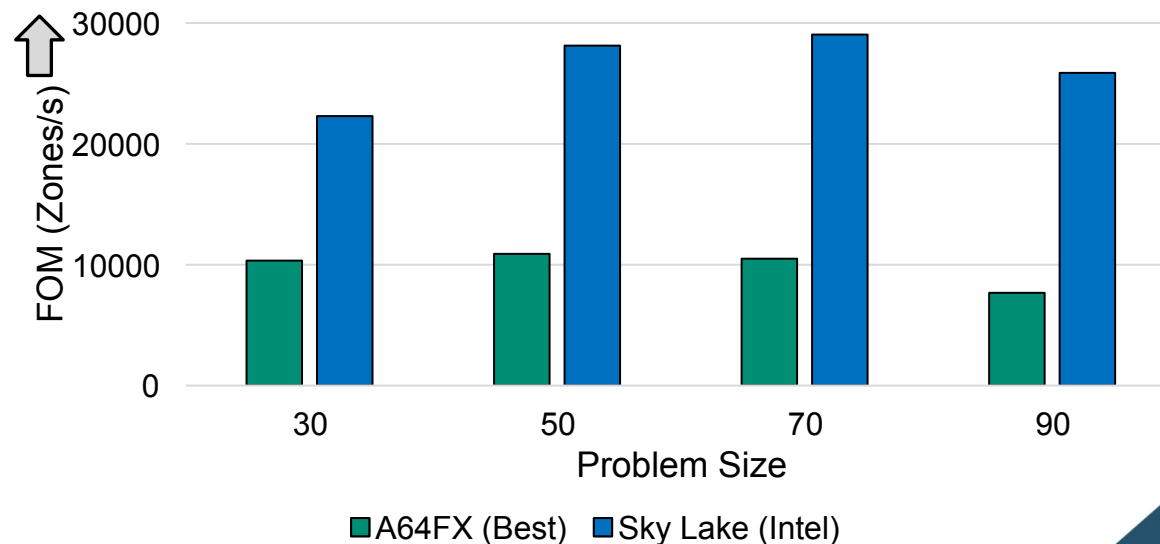
Mixed results for compilers in the platform

- See top graph (GCC 10.2 outperforms)
- Intel Sky Lake Xeon provides strong performance because there is much larger amounts of L2 cache (and L3) on the machine
- Represents results we see on many highly unstructured problems

Comparison of Compilers on A64FX



A64FX vs. Sky Lake Xeon





UPDATE ON ATSE (ADVANCED TRILAB SOFTWARE ENVIRONMENT)

Collaboration with HPE, OpenHPC, Arm and NNSA Labs

- Develop a robust software environment for Arm (and other) platforms
- Well tested combination of software products which are supported and optimized for NNSA production code teams

Provides broad range of software products

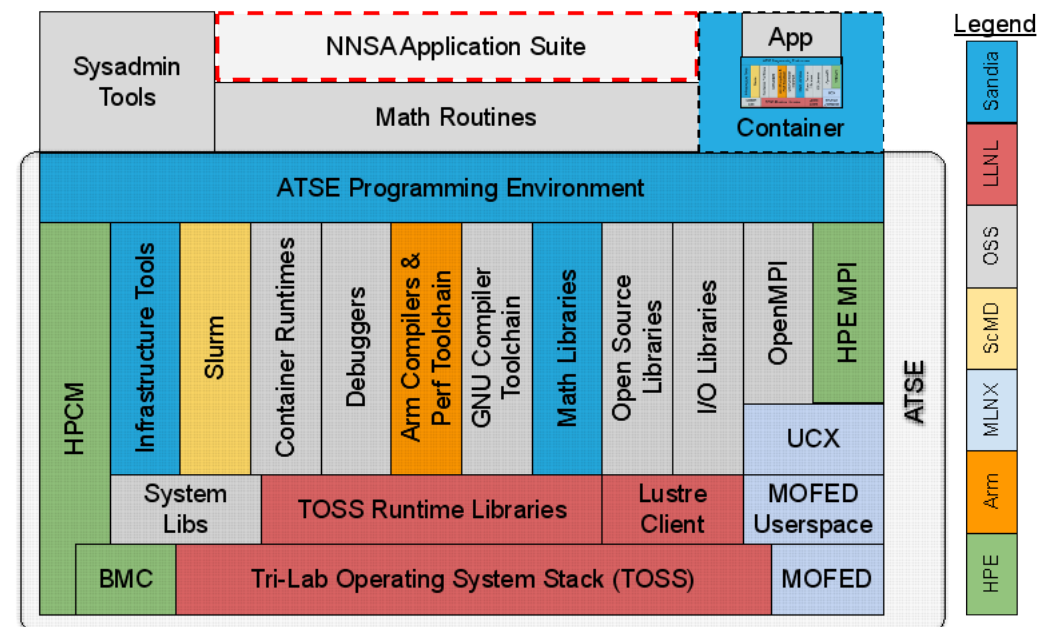
- Arm HPC Compiler, Math Libraries and Tools
- Several production tested GCC variants

Builds on the TOSS (Trilab Operating System) as a foundational Linux kernel

- Developed and supported by Lawrence Livermore Nat. Lab.



Hewlett Packard Enterprise



ATSE 1.2.5 Recipes Available @ <https://doi.org/10.5281/zenodo.4006668>



IN THE WORKS...

Kokkos C++ Performance Portable Programming Model

- Sandia is working on optimizations for Kokkos that specifically target A64FX, the Fujitsu compiler toolchains and SVE
- Include optimized math kernels via the Kokkos-Kernels math library

Initial versions of Trilinos Solvers already running on A64FX

- Not fully optimized and lots of areas for performance improvement

Open science applications building successfully on A64FX but still working on initial performance and benchmarking results

- Additional resources for our internal codes/systems



<https://kokkos.org/>



<https://trilinos.github.io>



MORE INFORMATION ON ASTRA AND ATSE AT SC20

SC20 Paper and Talk:

Chronicles of Astra: Challenges and Lessons from the First Petascale Arm Supercomputer



Presented by Andrew Younge and Kevin Pedretti from Sandia National Labs



Wednesday, 18 November 2020, 11:00 - 11:30 AM EST

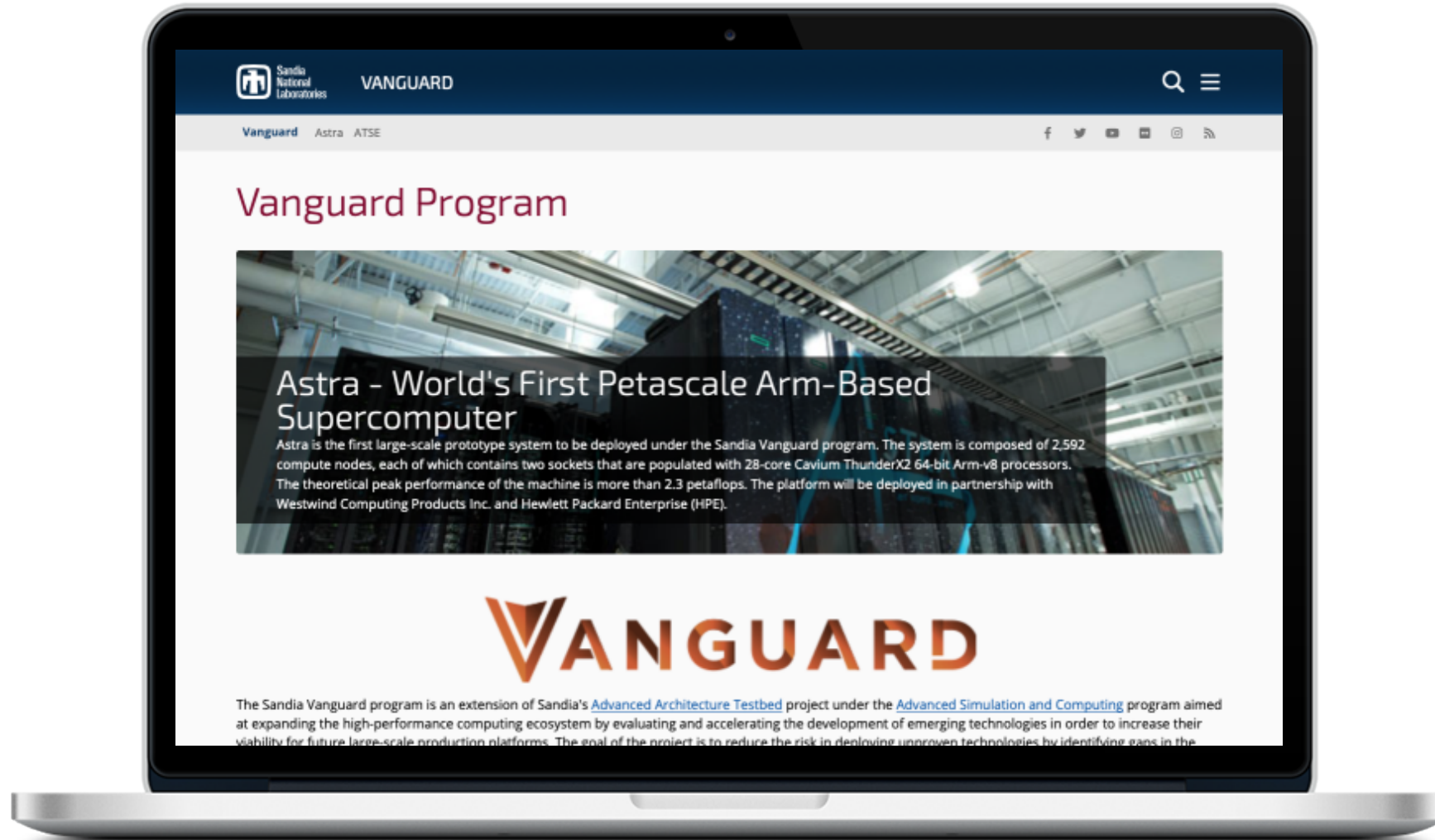


Architectures Technical Track





WANT TO KNOW MORE?



More information, links and papers at <https://vanguard.sandia.gov/>



ACKNOWLEDGEMENTS



Inouye builds on collaborations and work involving a wide collection of teams and laboratories, we are grateful to the following:

- Sandia ASC Advanced Architectures Testbed team and system admins
- Fujitsu and Penguin Computing for supporting early access to the A64FX processor
- Arm (especially the HPC Compiler team at Manchester UK and our performance experts Srinath and Olly)
- RIKEN Supercomputing Center (Japan)
- Our colleagues and team members at Los Alamos and Lawrence Livermore National Laboratories
- NNSA and ASC Headquarters for program support and for supporting our collaborations



Sandia
National
Laboratories

<http://www.sandia.gov>