# Vanguard Astra: Maturing the ARM Software Ecosystem for U.S. DOE/ASC Supercomputing

ExaComm'18
June 28, 2018

Kevin Pedretti, Jim H. Laros III, Si Hammond
ktpedre@sandia.gov

SAND2018-7066 C

*Exceptional service in the national interest*
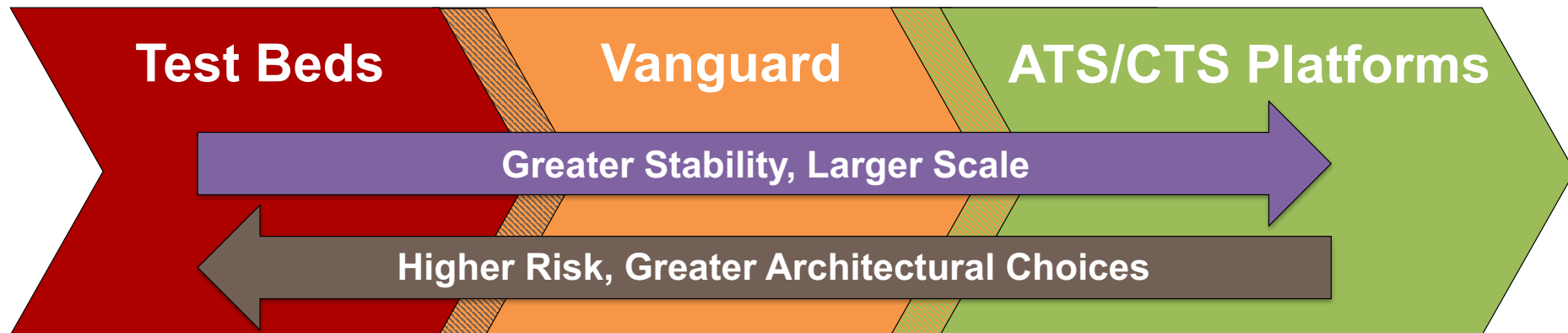
**Sandia National Laboratories**

# Outline

- **Vanguard prototype systems**
- Vanguard Astra ARM-based supercomputer
- Advanced Tri-lab Software Environment (ATSE)
- R&D directions
- Conclusion

# Vanguard:
# Large-scale Prototype Systems

- Expand the HPC ecosystem by developing emerging, yet-to-be-proven, technologies
  - Is technology viable for future production platforms supporting ASC integrated codes?
  - Increase technology choices
- Address hardware and software technologies together
  - If hardware technology is new, gaps in software stack are certain
- Buy down risk before commitment on capability/capacity class investment

# Where Vanguard Fits



**Test Beds** → **Vanguard** → **ATS/CTS Platforms**

Greater Stability, Larger Scale →

← Higher Risk, Greater Architectural Choices

**Test Beds**
- Small testbeds (~10-100 nodes)
- Breadth of architectures Key
- **Brave users**

**Vanguard**
- Larger-scale experimental systems
- Focused efforts to mature new technologies
- Broader user-base
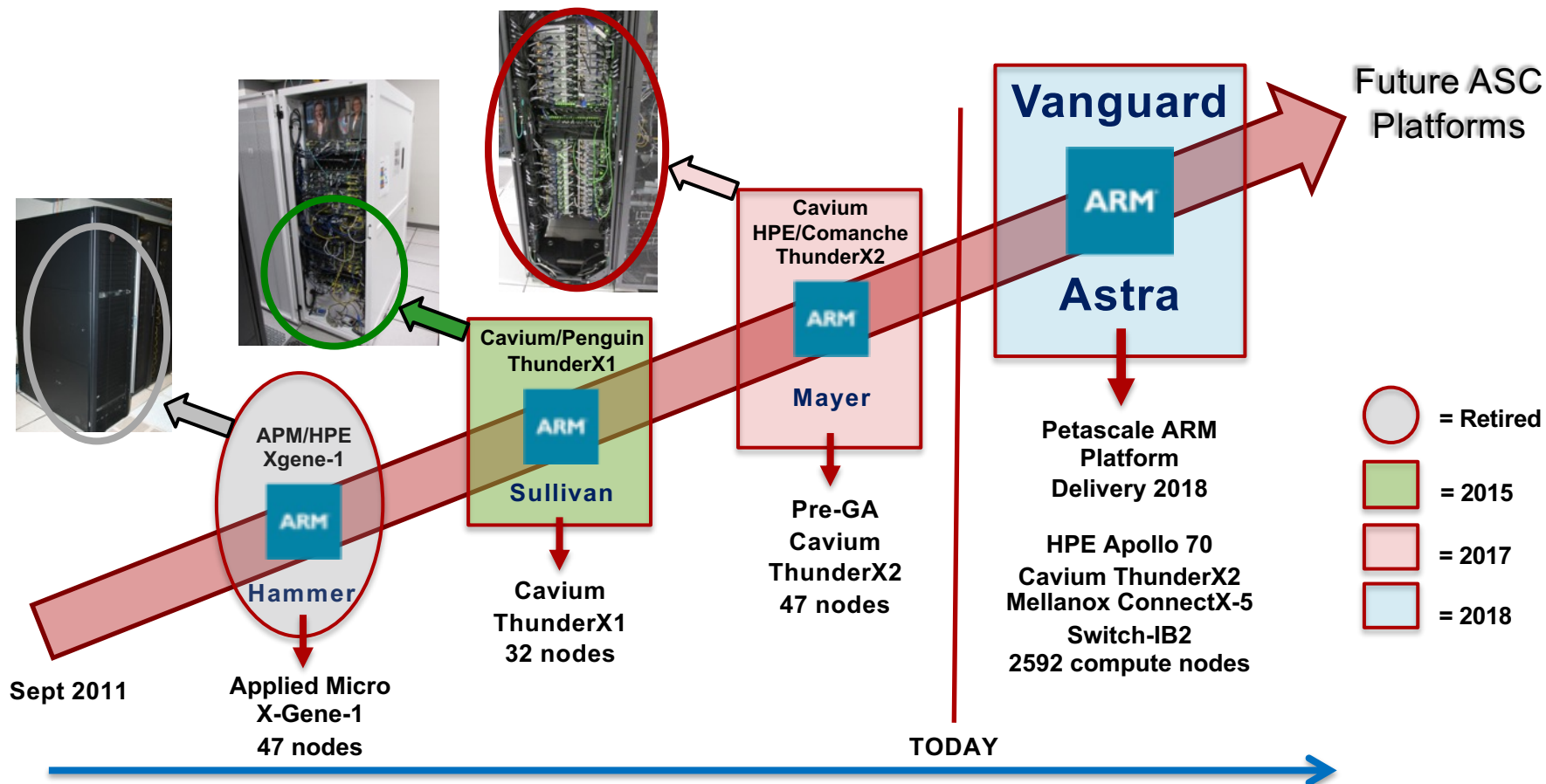- **Demonstrate viability for production use**
- NNSA Tri-lab resource

**ATS/CTS Platforms**
- Leadership-class systems (Petascale, Exascale, ...)
- Advanced technologies, sometimes first-of-kind
- Broad user-base
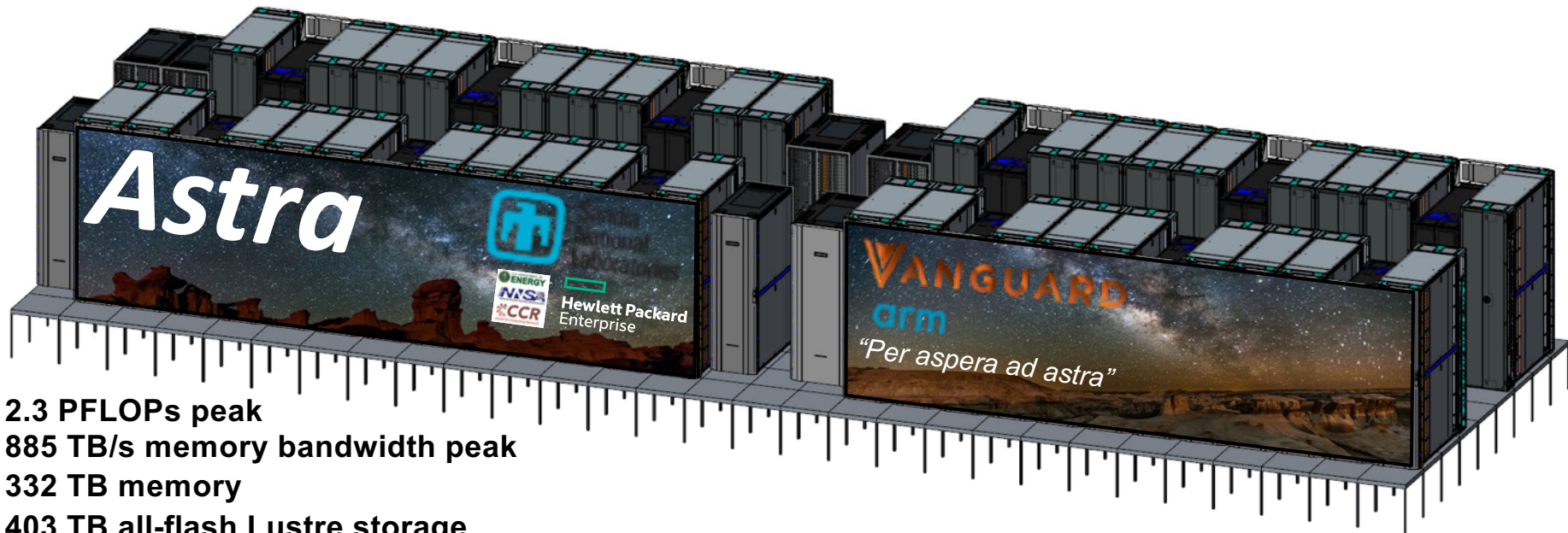- **Production use**

# Outline

- Vanguard prototype systems
- Vanguard Astra ARM-based supercomputer
- Advanced Tri-lab Software Environment (ATSE)
- R&D directions
- Conclusion

# Sandia's NNSA/ASC ARM Platforms



**Future ASC Platforms**

**Vanguard**

ARM

**Astra**

Petascale ARM Platform Delivery 2018

HPE Apollo 70
Cavium ThunderX2
Mellanox ConnectX-5
Switch-IB2
2592 compute nodes

**Cavium HPE/Comanche ThunderX2**

ARM

**Mayer**

Pre-GA
Cavium
ThunderX2
47 nodes

**Cavium/Penguin ThunderX1**

ARM

**Sullivan**

Cavium
ThunderX1
32 nodes

**APM/HPE Xgene-1**

ARM

**Hammer**

Applied Micro
X-Gene-1
47 nodes

**Sept 2011**

**TODAY**

| | |
|---|---|
| ⬭ | = Retired |
| 🟩 | = 2015 |
| 🟥 | = 2017 |
| 🟦 | = 2018 |

# *per aspera ad astra*

## through difficulties to the stars

**Astra**

**Vanguard**
arm
*"Per aspera ad astra"*

**2.3 PFLOPs peak**
**885 TB/s memory bandwidth peak**
**332 TB memory**
**403 TB all-flash Lustre storage**
**1.2 MW**

## Demonstrate viability of ARM for U.S. DOE NNSA Supercomputing
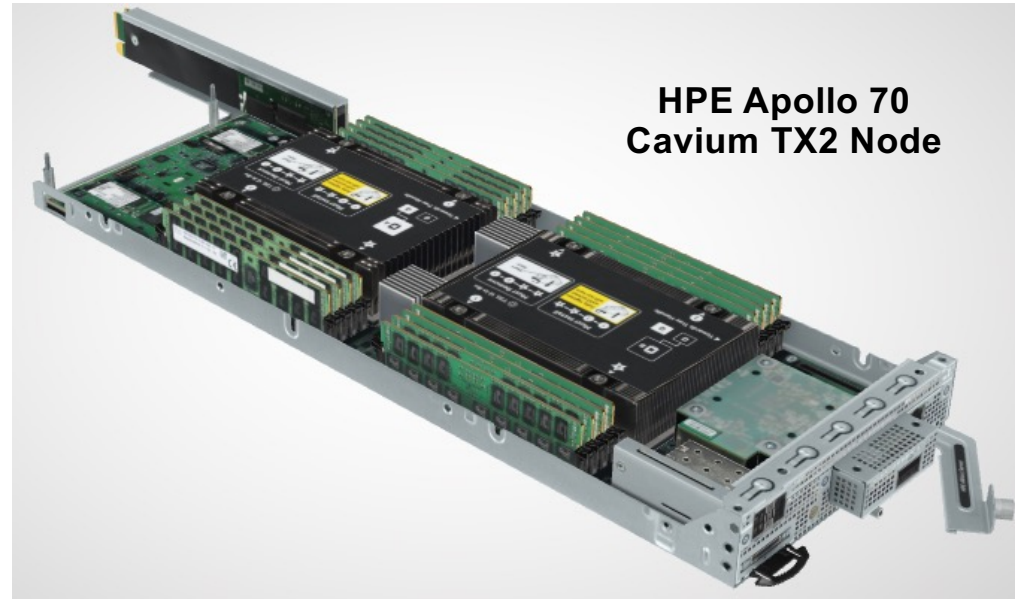
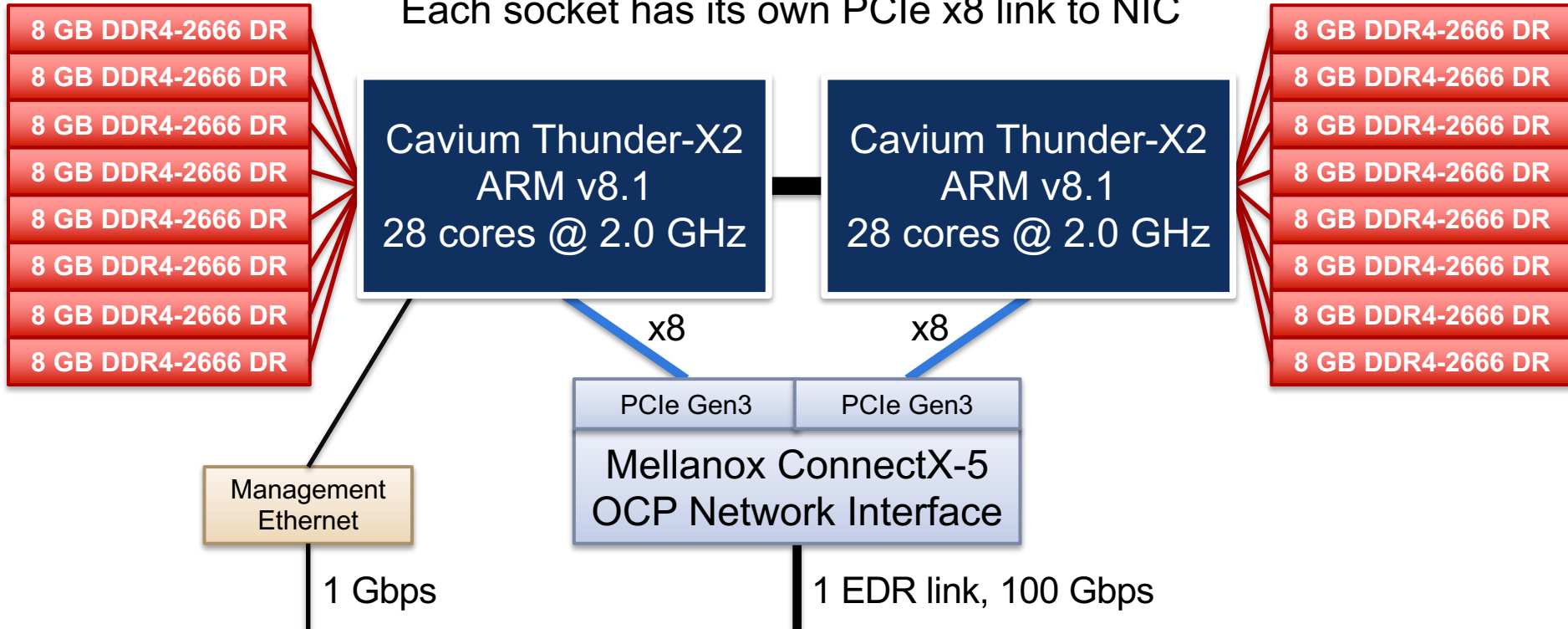# Vanguard-Astra Compute Node Building Block



- Dual socket
  Cavium Thunder-X2 CN99xx
  28 cores @ 2.0 GHz

- 8 DDR4 controllers per socket

- One 8 GB DDR4-2666 dual-rank DIMM per controller

- Mellanox EDR InfiniBand ConnectX-5 VPI OCP

- Tri-Lab Operating System Stack based on RedHat 7.5+

**HPE Apollo 70 Cavium TX2 Node**

# Vanguard-Astra Compute Node

8 DDR4 channels/socket, 1 DIMM/channel
Each socket has its own PCIe x8 link to NIC

# Vanguard-Astra System Packaging

**HPE Apollo 70 Chassis: 4 nodes**



**36 compute racks**
**(9 scalable units, each 4 racks)**

**2592 compute nodes**
**(5184 TX2 processors)**

**3 IB spine switches**
**(each 540-port)**

**HPE Apollo 70 Rack**



**18 chassis/rack**

**72 nodes/rack**

**3 IB switches/rack**
**(one 36-port switch**
**per 6 chassis)**
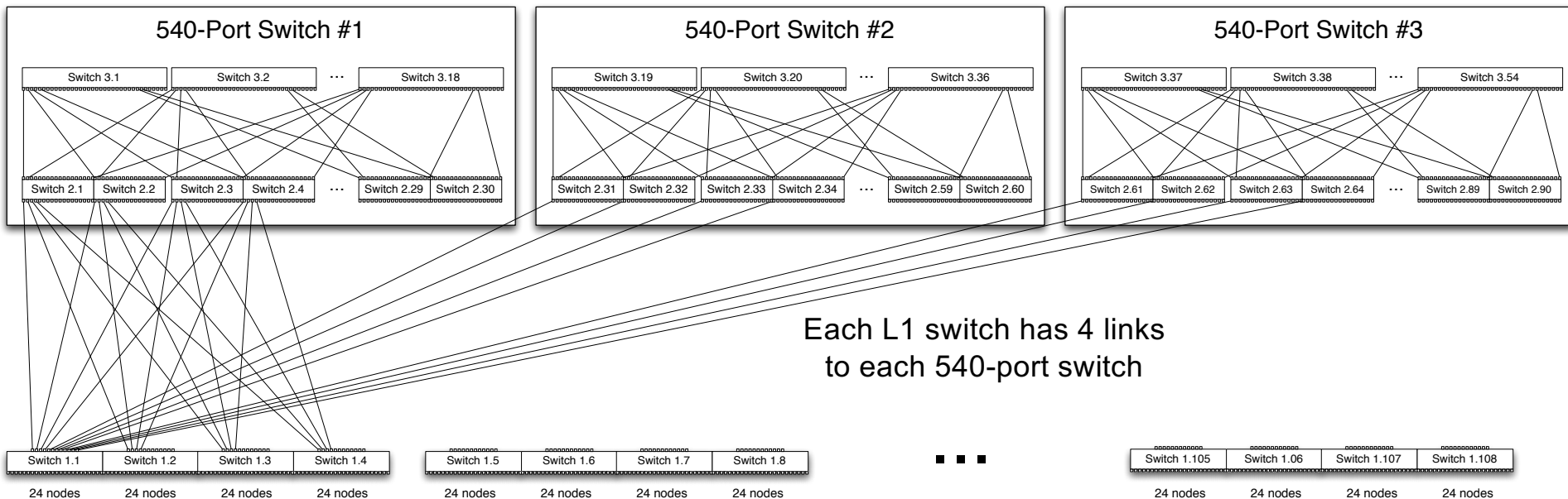


**VANGUARD**
*Astra*

# Network Topology Visualization

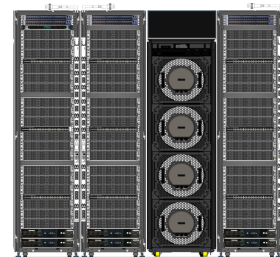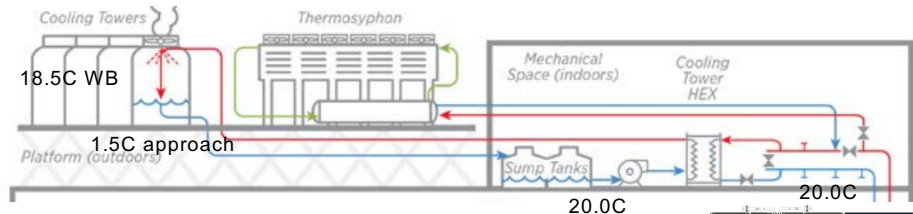Mellanox Switch-IB2 EDR, Radix 36 switches, 3 level fat tree, 2:1 taper at L1, SHARP

| 540-Port Switch #1 | 540-Port Switch #2 | 540-Port Switch #3 |
|---|---|---|

Switch 3.1 | Switch 3.2 | ··· | Switch 3.18

Switch 3.19 | Switch 3.20 | ··· | Switch 3.36

Switch 3.37 | Switch 3.38 | ··· | Switch 3.54

Switch 2.1 | Switch 2.2 | Switch 2.3 | Switch 2.4 | ··· | Switch 2.29 | Switch 2.30

Switch 2.31 | Switch 2.32 | Switch 2.33 | Switch 2.34 | ··· | Switch 2.59 | Switch 2.60

Switch 2.61 | Switch 2.62 | Switch 2.63 | Switch 2.64 | ··· | Switch 2.89 | Switch 2.90

Each L1 switch has 4 links
to each 540-port switch

Switch 1.1 | Switch 1.2 | Switch 1.3 | Switch 1.4

Switch 1.5 | Switch 1.6 | Switch 1.7 | Switch 1.8

■ ■ ■

Switch 1.105 | Switch 1.06 | Switch 1.107 | Switch 1.108

24 nodes | 24 nodes | 24 nodes | 24 nodes

24 nodes | 24 nodes | 24 nodes | 24 nodes

24 nodes | 24 nodes | 24 nodes | 24 nodes

108 L1 switches * 24 nodes/switch = 2592 compute nodes

# Vanguard-Astra Advanced Power & Cooling

**Power and Water Efficient:**

- Total 1.2 MW in the 36 compute racks are cooled by only 12 fan coils
- These coils are cooled without compressors year round. No evaporative water at all almost 6000 hours a year
- 99% of the compute racks heat never leaves the cabinet, yet the system doesn't require the internal plumbing of liquid disconnects and cold plates running across all CPUs and DIMMs
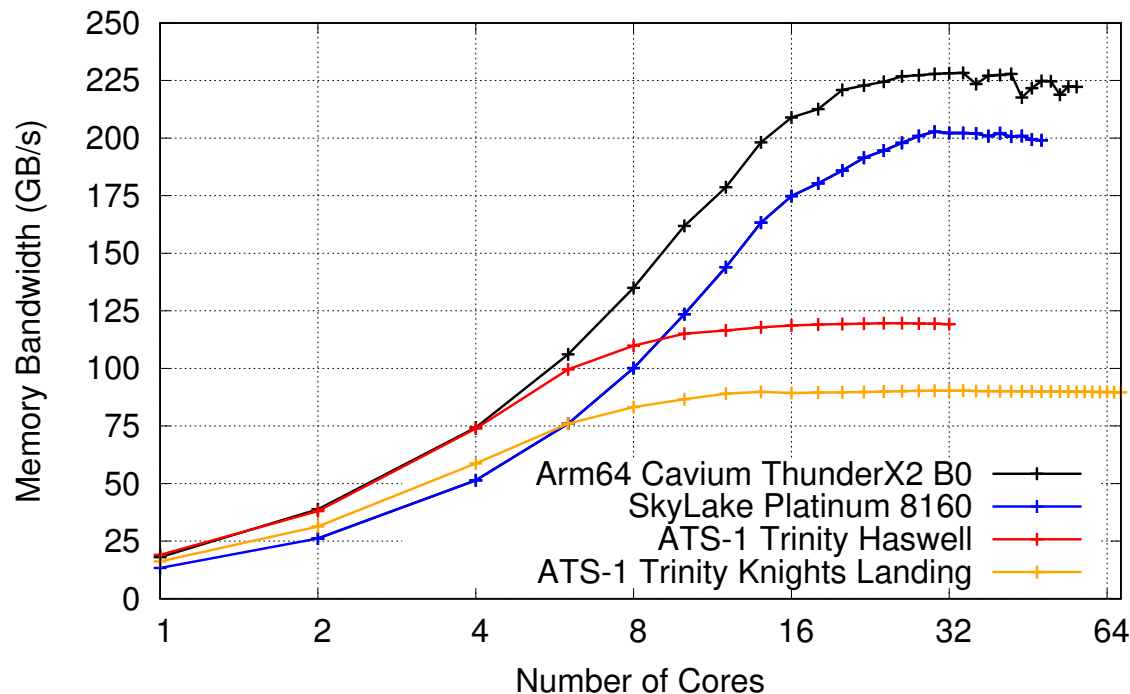
- Builds on joint work by NREL and Sandia: https://www.nrel.gov/esif/partnerships-jc.html



| Projected power of the system by component | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| per constituent rack type (W) | | | | | total (kW) | | | | |
| | wall | peak | nominal (linpack) | idle | | racks | wall | peak | nominal (linpack) | idle |
| Node racks | 39888 | 35993 | 33805 | 6761 | | 36 | 1436.0 | 1295.8 | 1217.0 | 243.4 |
| MCS300 | 10500 | 7400 | 7400 | 170 | | 12 | 126.0 | 88.8 | 88.8 | 2.0 |
| Network | 12624 | 10023 | 9021 | 9021 | | 3 | 37.9 | 30.1 | 27.1 | 27.1 |
| Storage | 11520 | 10000 | 10000 | 1000 | | 2 | 23.0 | 20.0 | 20.0 | 2.0 |
| utility | 8640 | 5625 | 4500 | 450 | | 1 | 8.6 | 5.6 | 4.5 | 0.5 |
| | | | | | | | 1631.5 | 1440.3 | 1357.3 | 274.9 |

# Cavium Arm64 Providing Best-of-Class Memory Bandwidth



STREAM TRIAD

TX2 DDR4-2400
SkyLake 8160

Trinity Haswell
Trinity KNL DDR
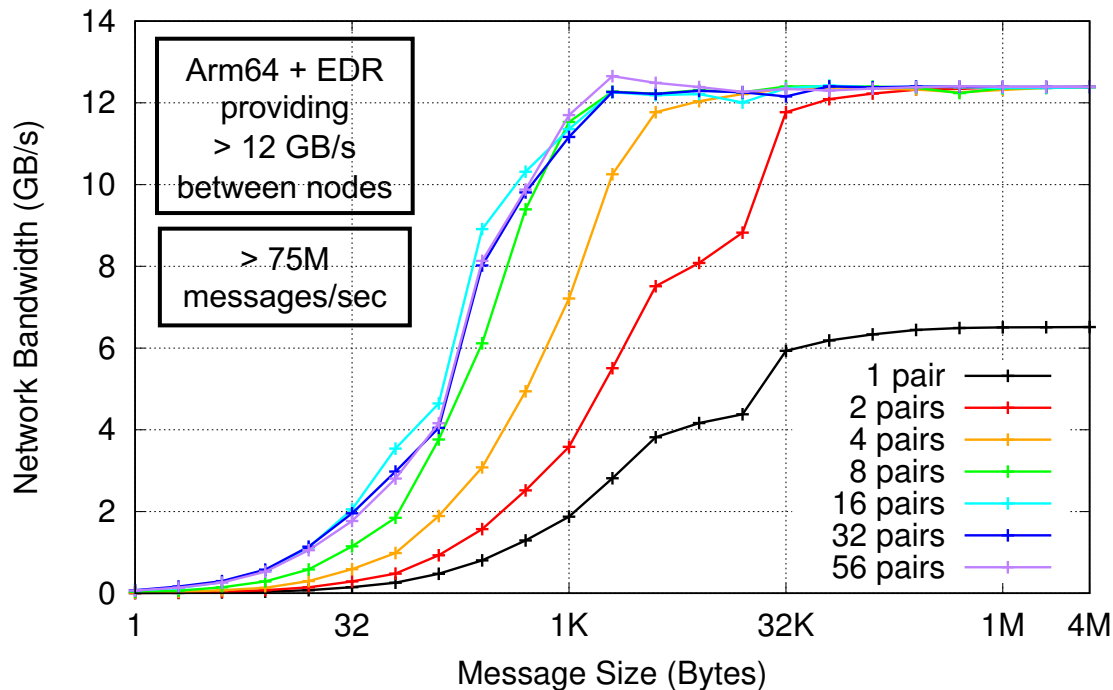
# Network Bandwidth on ThunderX2 + Mellanox MLX5 EDR with Socket Direct

## Socket Direct – Each socket has dedicated path to the NIC



## OSU MPI Multi-Network Bandwidth



Arm64 + EDR providing > 12 GB/s between nodes

> 75M messages/sec
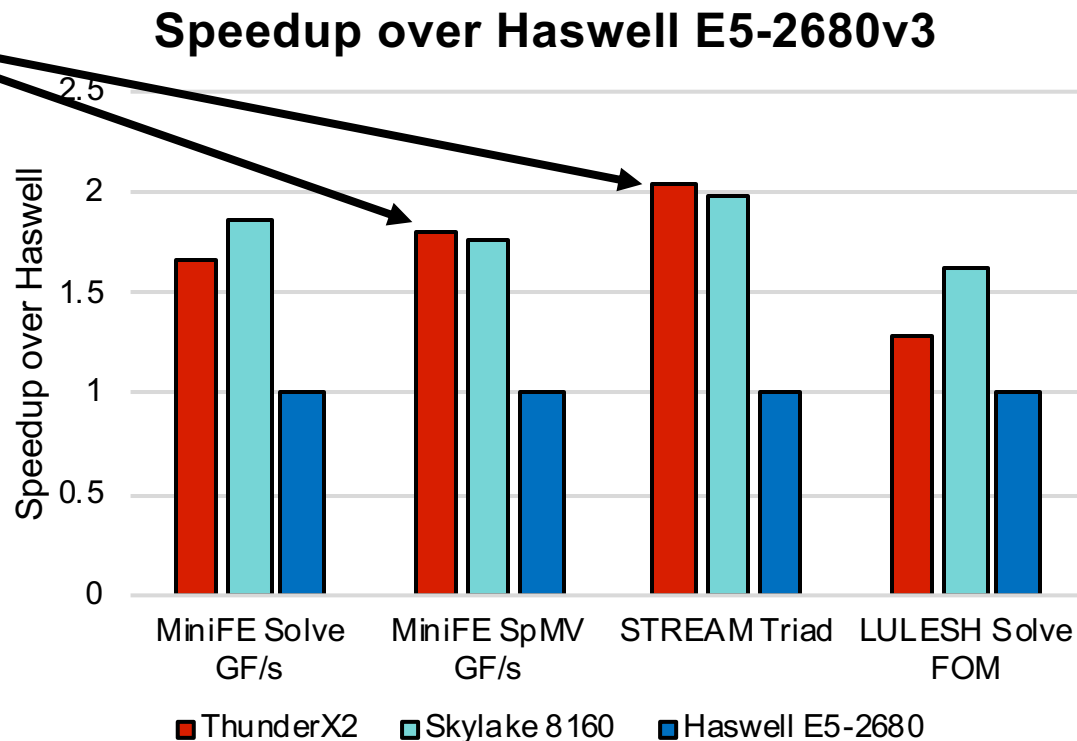
# Mini-App Performance on Cavium ThunderX2

- ThunderX2 providing high memory bandwidth
  - 6 channels (Skylake) vs. 8 in ThunderX2
  - See this in MiniFE SpMV and STREAM Triad
- Slower compute reflects less optimization in software stack
  - Examples – Non-SpMV kernels in MiniFE and LULESH
  - GCC and ARM versus Intel compiler

**Speedup over Haswell E5-2680v3**

Chart — Speedup over Haswell (y-axis), categories: MiniFE Solve GF/s, MiniFE SpMV GF/s, STREAM Triad, LULESH Solve FOM. Series: ThunderX2, Skylake 8160, Haswell E5-2680.

# Vanguard-Astra Acceptance Plan

## Milestone 1
Open Science
2-3 months

**Full Scale Machine Runs**
- HPCG
- HPL

**Micro-benchmarks**
- STREAM
- Intel MPI Benchmarks

**Compile and Run**
- **NALU (SNL)**
- **VPIC (LANL)**
- **PF3D (LLNL)**

## Milestone 2
Restricted Science
< 12 months

**SSI Benchmarks**
- HPCG
- HPL

**Lab/Vendor Optimization**
- **SPARC (SNL)**
- **PARTISn (LANL)**
- **ALE3D (LLNL)**

**Compile and Run**
- **RAMSES (SNL)**

## Milestone 3
Classified Science
Remainder of Life

**Lab/Vendor Optimization**
- **SPARC (SNL)**
- **PARTISn (LANL)**
- **ALE3D (LLNL)**
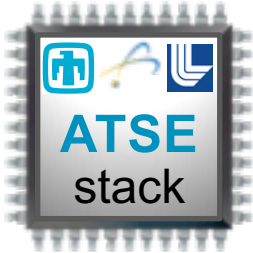
**Compile and Run**
- **SIERRA (SNL)**

**Demonstrate**
- **User-specified containers and virtual machines**

# Outline

- Vanguard prototype systems
- Vanguard Astra ARM-based supercomputer
- Advanced Tri-lab Software Environment (ATSE)
- R&D directions
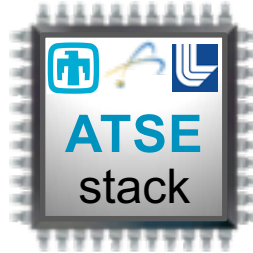- Conclusion

# Vanguard Tri-Lab Software Effort

- Accelerate maturity of ARM ecosystem for ASC computing
    - Prove viability for NNSA integrated codes running at scale
    - Harden compilers, math libraries, tools, communication libraries
        - Heavily templated C++, Fortran 2003/2008, Gigabyte+ binaries, long compiles
    - Optimize performance, verify expected results
- Build integrated software stack
    - Programming env (compilers, math libs, tools, MPI, OMP, SHMEM, I/O, ...)
    - Low-level OS (optimized Linux, network, filesystems, containers/VMs, ...)
    - Job scheduling and management (WLM, app launcher, user tools, ...)
    - System management (boot, system monitoring, image management, ...)

Improve 0 to 60 time... Astra arrival to useful work done
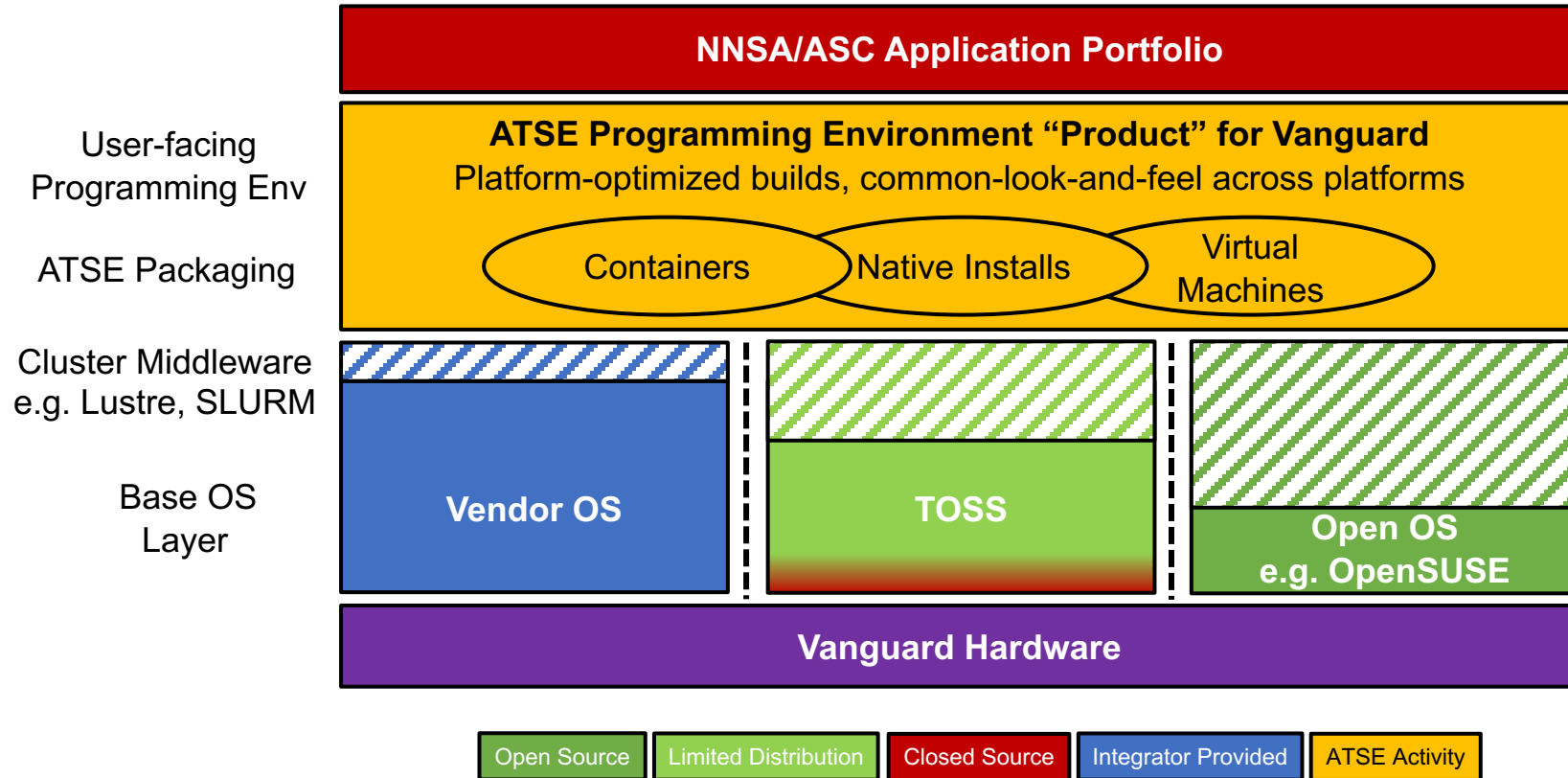
# Advanced Tri-lab Software Environment High-level Goals

ATSE
stack

- Build an open, modular, extensible, community-engaged, and vendor-adaptable ecosystem

- Prototype new technologies that may improve the DOE ASC computing environment (e.g., ML frameworks, containers, VMs, OS optimizations)

- Leverage existing efforts such as Tri-lab OS (TOSS), programing environments, and Exascale Computing Project software technologies

| Aug'17 Tri-lab Arm software team formed | Dec'17 ATSE Design Doc | Jul'18 Initial Release Target | Sep'18 First Use on Vanguard-Astra |
| --- | --- | --- | --- |

# Vanguard-Astra Software Stack


ATSE stack

**NNSA/ASC Application Portfolio**

User-facing Programming Env

**ATSE Programming Environment "Product" for Vanguard**
Platform-optimized builds, common-look-and-feel across platforms

ATSE Packaging

Containers   Native Installs   Virtual Machines

Cluster Middleware e.g. Lustre, SLURM

Base OS Layer

**Vendor OS**   **TOSS**   **Open OS e.g. OpenSUSE**

**Vanguard Hardware**

| Open Source | Limited Distribution | Closed Source | Integrator Provided | ATSE Activity |

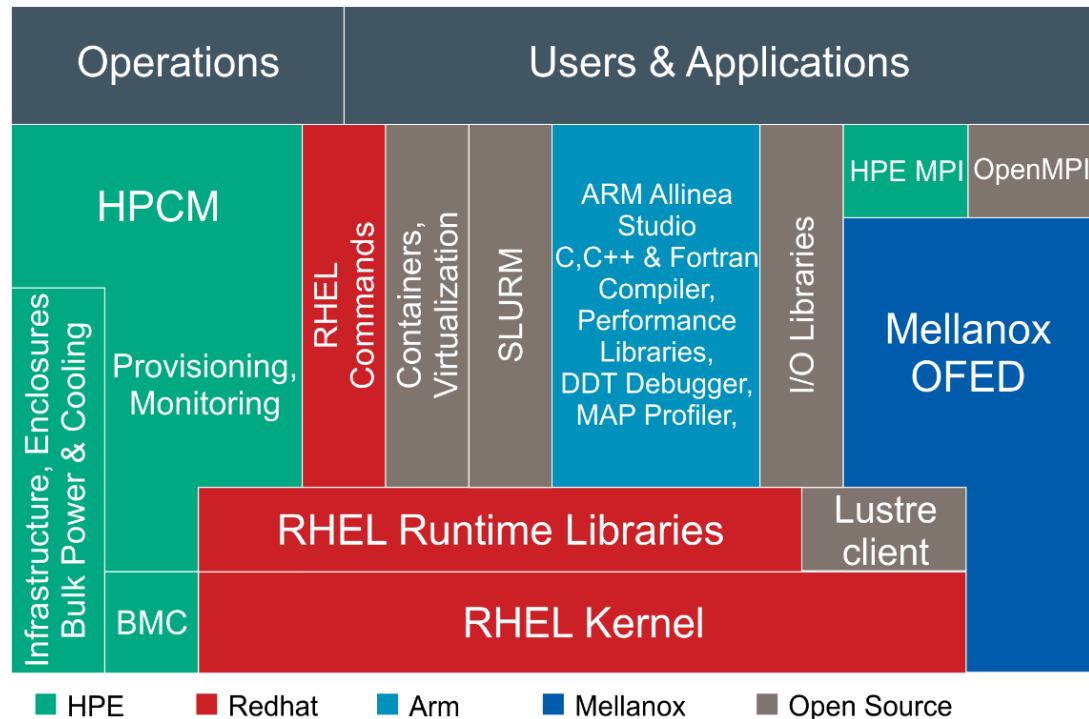# Integrate Components from Many Sources

# Close Collaboration with HPE Open Leadership Software Stack (OLSS) Effort

Hewlett Packard Enterprise

- HPE:
  - HPE MPI (+ XPMEM)
  - HPE Cluster Manager
- Arm:
  - Arm HPC Compilers
  - Arm Math Libraries
  - Allinea Tools
- Mellanox-OFED & HPC-X
- RedHat 7.x for aarch64



Operations | Users & Applications

HPCM

Infrastructure, Enclosures Bulk Power & Cooling

Provisioning, Monitoring

RHEL Commands

Containers, Virtualization

SLURM

ARM Allinea Studio C,C++ & Fortran Compiler, Performance Libraries, DDT Debugger, MAP Profiler,

I/O Libraries

HPE MPI | OpenMPI

Mellanox OFED

RHEL Runtime Libraries

Lustre client

BMC

RHEL Kernel

HPE | Redhat | Arm | Mellanox | Open Source

# Early Application Porting

| Workload | GCC 7.2.0 | Arm HPC Compilers |
|----------|-----------|-------------------|
| LAMMPS | | |
| SPARTA | | |
| SPARC | | |
| NALU | | |
| CTH | | FORTRAN issue |
| Drekar | | |
| Xyce-UUR | | |
| VPIC | | |
| SNAP | | |

*Most codes build without trouble, optimization work remains*

**Placing collaborative vendor contracts to harden Arm64 compilers, math libraries, and tools – both for Astra and Arm ecosystem in general**

# Outline

- Vanguard prototype systems
- Vanguard Astra ARM-based supercomputer
- Advanced Tri-lab Software Environment (ATSE)
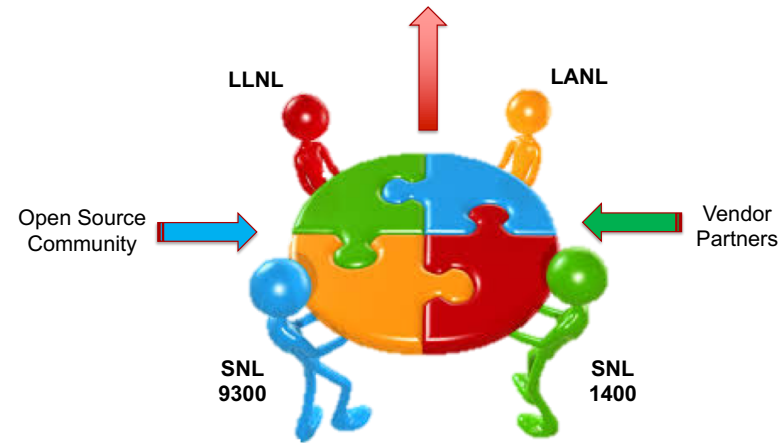- R&D directions
- Conclusion

# R&D Areas

- Workflows leveraging containers and virtual machines
  - Support for machine learning frameworks
  - ARMv8.1 includes new virtualization extensions, SR-IOV
- Evaluating parallel filesystems + I/O systems @ scale
  - GlusterFS, Ceph, BeeGFS, Sandia Data Warehouse, …
- Resilience studies over Astra lifetime
- Improved MPI thread support, matching acceleration
- OS optimizations for HPC @ scale
  - Exploring spectrum from stock distro Linux kernel to HPC-tuned Linux kernels to non-Linux lightweight kernels and multi-kernels
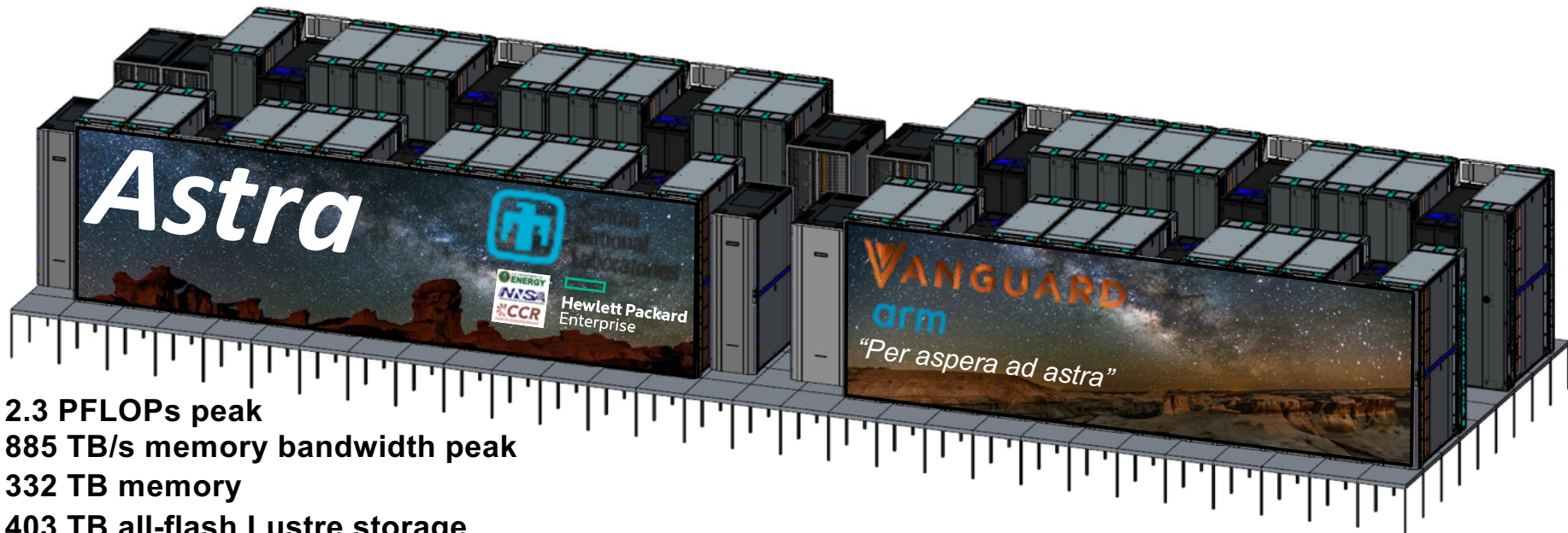  - Arm-specific optimizations

# Conclusion

- Vanguard expanding HPC ecosystem by developing emerging, yet-to-be-proven, technologies, taking appropriate risk
  - Mature new technologies for NNSA ASC integrated codes
- Vanguard-Astra will be one of the first Arm-based supercomputers
- NNSA Tri-lab team (Sandia, Los Alamos, Lawrence Livermore) is working in partnership with HPE, Arm, Cavium, RedHat, Mellanox, and others to develop the ATSE software stack for Astra

**Vanguard Collaboration**



LLNL

LANL

Open Source Community

Vendor Partners

SNL 9300

SNL 1400

# *per aspera ad astra*

## through difficulties to the stars



**Astra**

**Vanguard** arm

*"Per aspera ad astra"*

**2.3 PFLOPs peak**
**885 TB/s memory bandwidth peak**
**332 TB memory**
**403 TB all-flash Lustre storage**
**1.2 MW**

Demonstrate viability of ARM for U.S. DOE NNSA Supercomputing