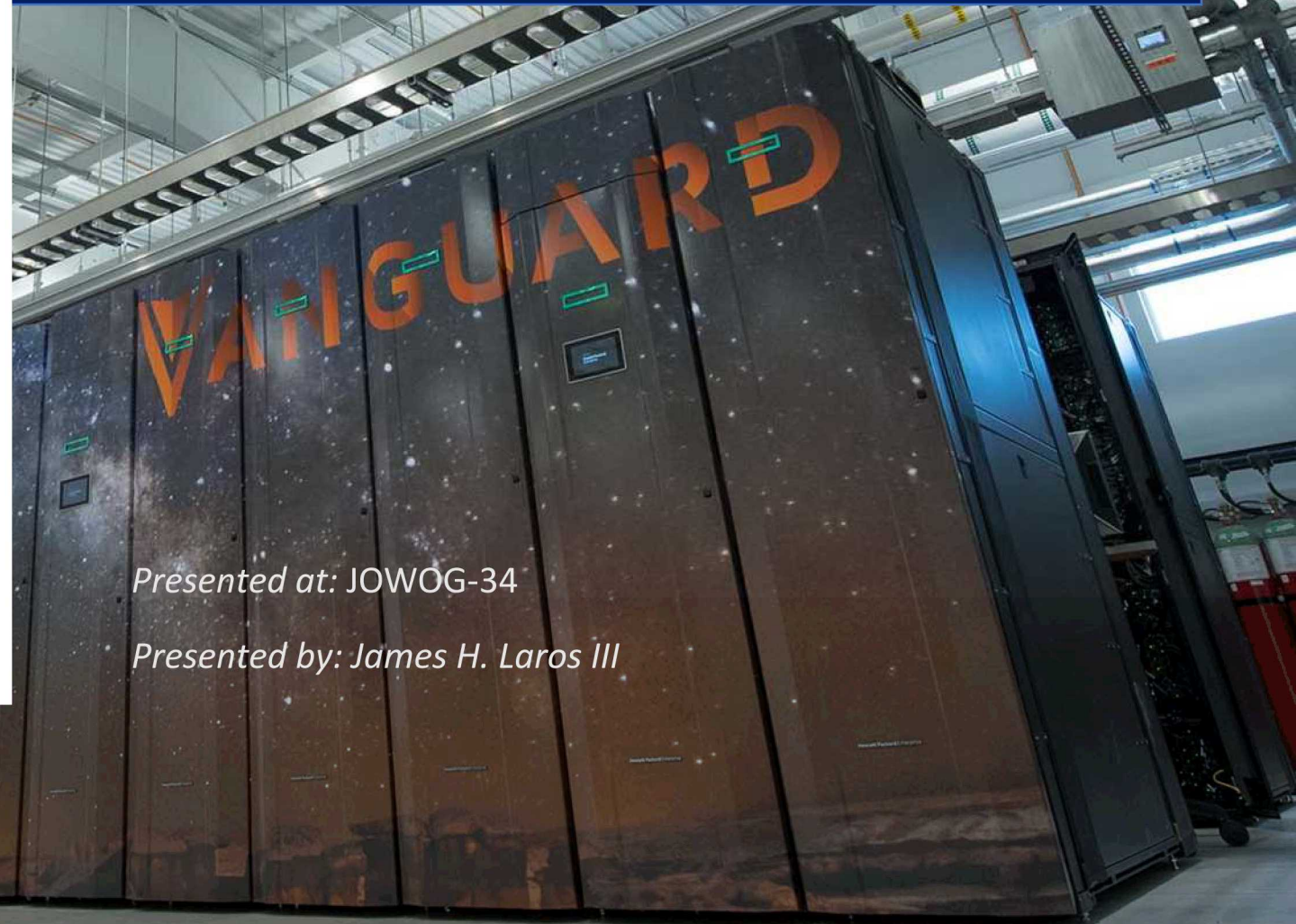


ARM SUPERCOMPUTER



Presented at: JOWOG-34

Presented by: James H. Laros III

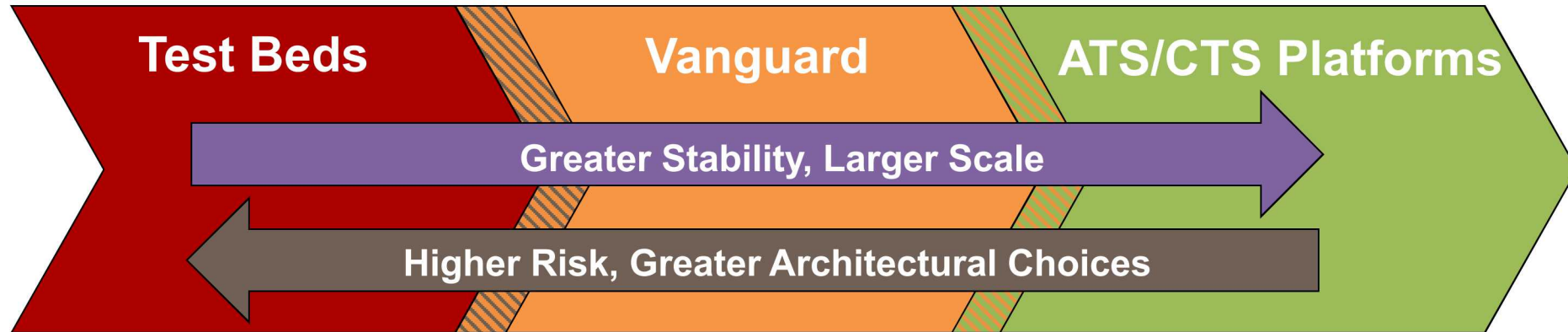
Vanguard Program: Advanced Technology Prototype Systems

- **Prove viability of advanced technologies for NNSA integrated codes, at scale**
- Expand the HPC-ecosystem by developing emerging yet-to-be proven technologies
 - Is technology viable for future ATS/CTS platforms supporting ASC mission?
 - Increase technology AND integrator choices
- Buy down risk and increase technology and vendor choices for future NNSA production platforms
 - Ability to accept higher risk allows for more/faster technology advancement
 - Lowers/eliminates mission risk and significantly reduces investment
- Jointly address hardware and software technologies
- First Prototype platform targeting Arm Architecture

Success achieved through Tri-Lab involvement and collaboration



Where Vanguard Fits



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
- Not Production
- **Tri-lab resource but not for ATCC runs**

ATS/CTS Platforms

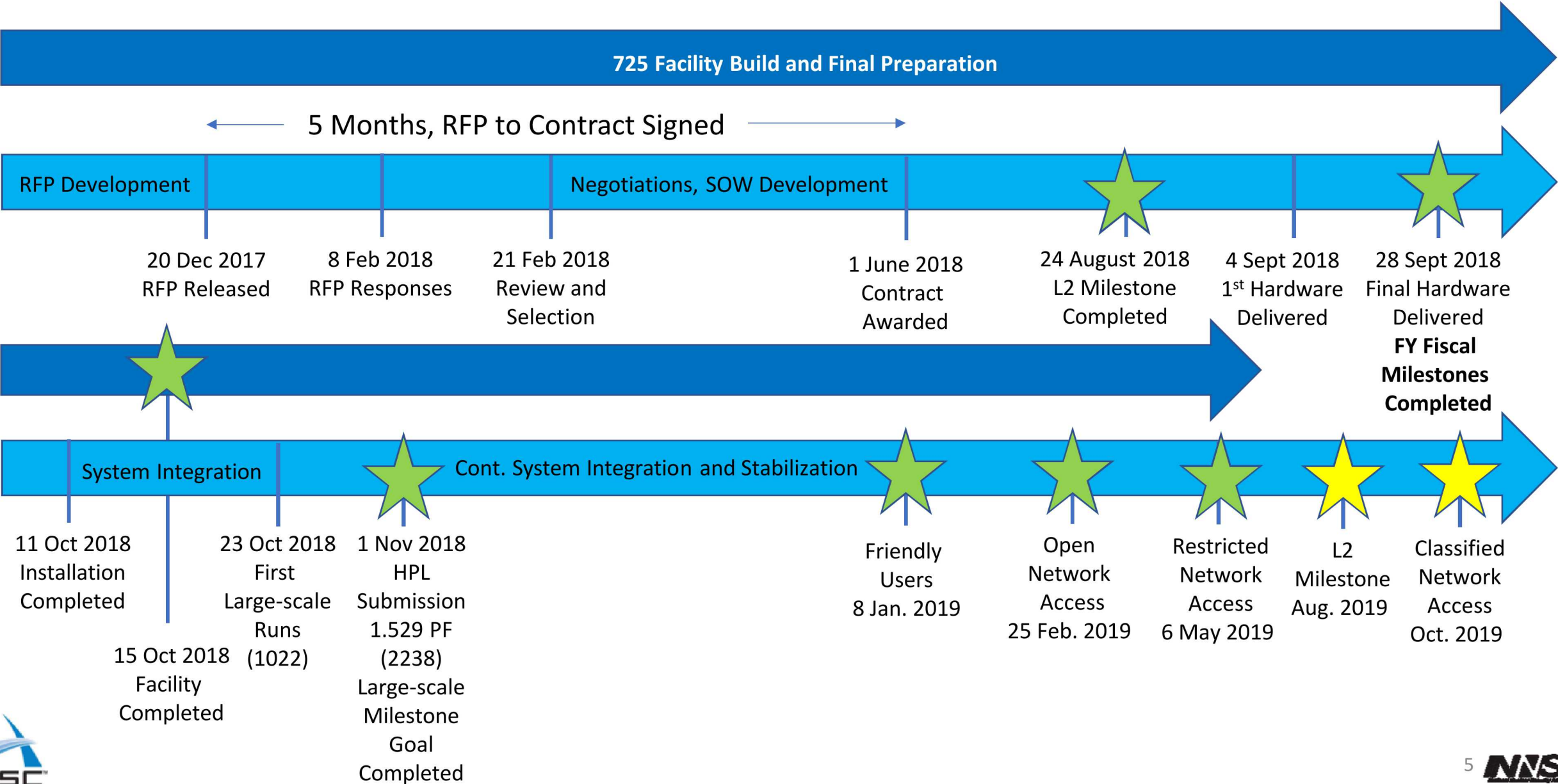
- Leadership-class systems (Petascale, Exascale, ...)
- Advanced technologies, sometimes first-of-kind
- Broad user-base
- **Production use**

Vanguard Astra: At a Glance

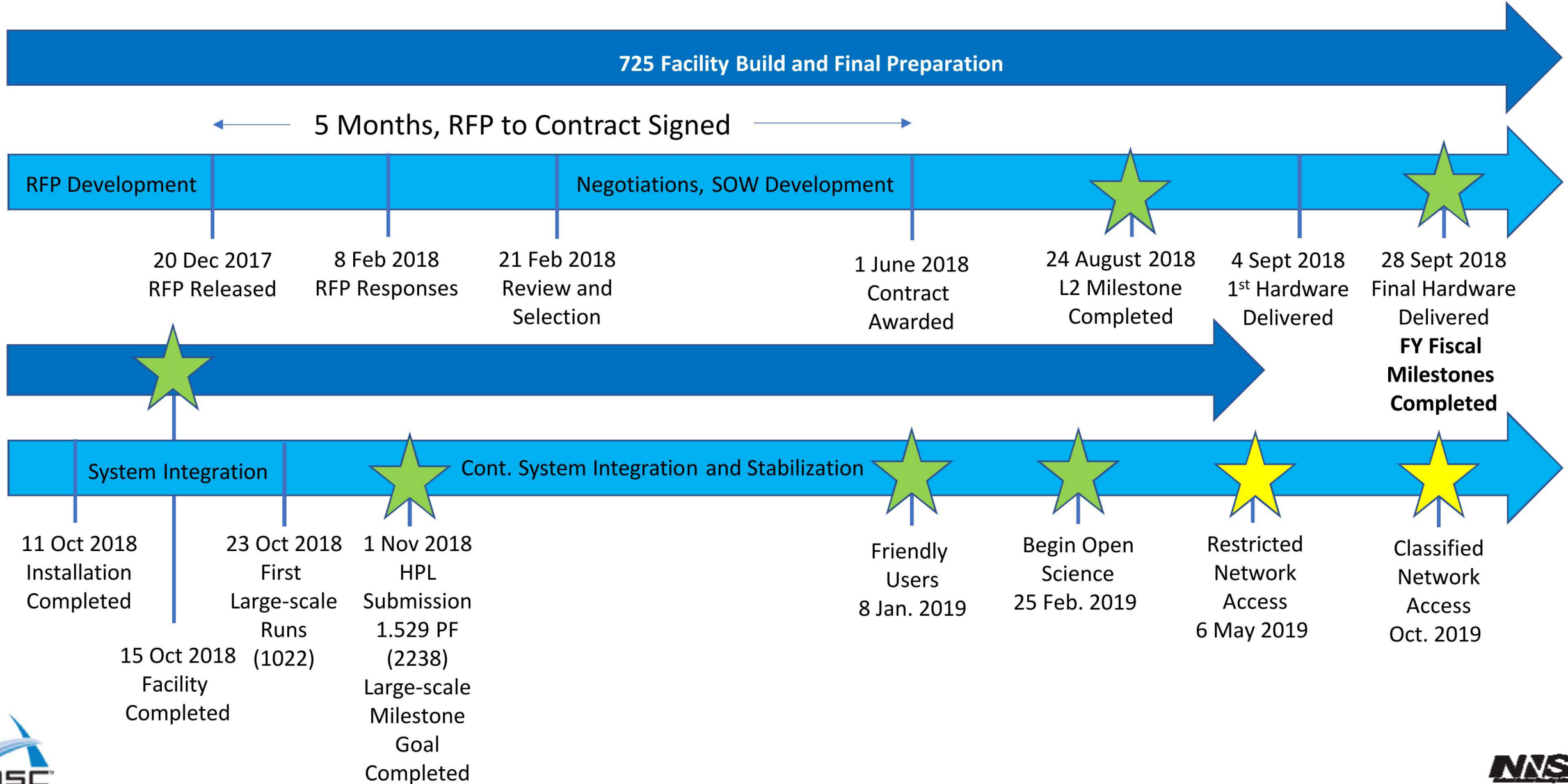
- 2,592 HPE Apollo 70 compute nodes
 - 5,184 CPUs, 145,152 cores, 2.3 PFLOPs (peak)
- Cavium Thunder-X2 ARM SoC, 28 core, 2.0 GHz
- Memory per node: 128 GB
 - 16 x 8 GB DDR DIMMs
 - Aggregate capacity: 332 TB, 885 TB/s (peak)
 - 247 GB/s per node STREAM
- Mellanox IB EDR, ConnectX-5
 - 112 36-port leaf, 3 648-port spine switches
- ATSE software stack
 - TOSS Base Operating system
- HPE Apollo 4520 All-flash Lustre storage
 - Storage Capacity: 403 TB (usable)
 - **Upgrade to 3X memory in preparation for move to classified**
 - Storage Bandwidth: 250 GB/s
 - 400 GB/s stunt mode, 432 GB/s peak



Vanguard-Astra: Timeline

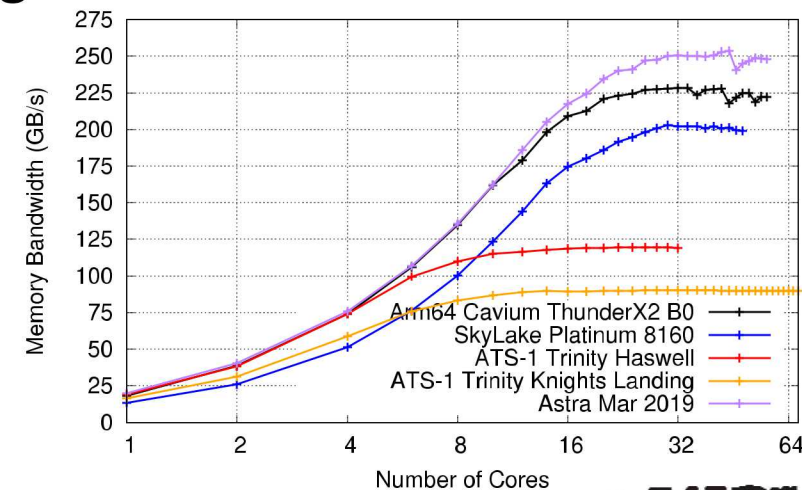


Vanguard-Astra: Timeline



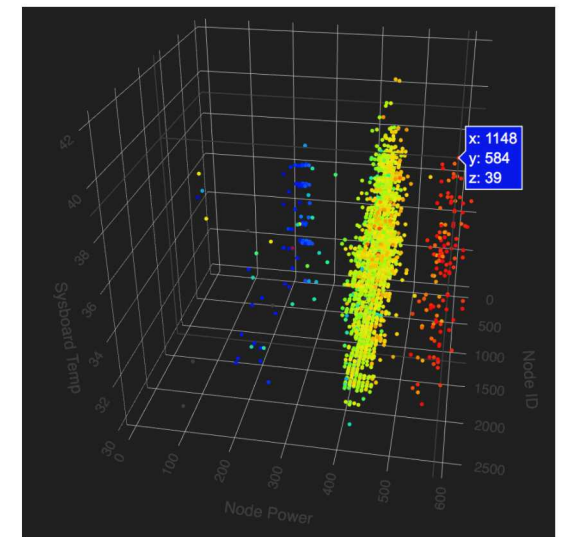
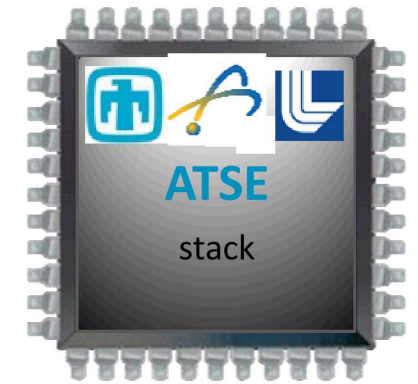
Vanguard Astra: Lessons Learned or Reasons to Prototype new Technologies

- “Right sizing” technical review team worked well - Smaller can be better
 - reinforced at Crossroads TAT
- “Right sizing” system to meet goals
 - Depends on the technology targeted
 - Depends on budget
 - At scale reasonably shows will run on both CTS and ATS
 - Test entire software environment, systems management included, at scale
 - Vendor interest (depends on vendor)
- Unless you are into hard hats and safety vests consider finishing the building before system delivery
- Problems are not observed in isolation
- Initial implementation of Thermal solution required tweaking
 - Met our target of 21-22 C water, 24-25 inlet air temp
 - Lowered water temp to 20 C - led to increased stability
- Delivered processor part originally did not meet SOW memory bandwidth requirements – Does now



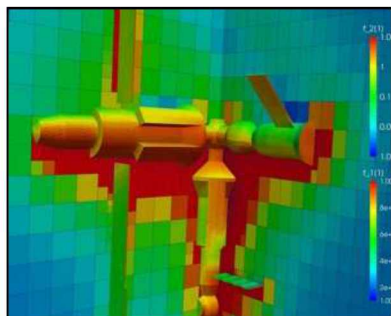
Vanguard Astra: Lessons Learned *or* Reasons to Prototype new Technologies

- Similar to CTS, NALU and other applications were forcing out of spec voltage swings
 - In this case memory bus
- Pesky fabric instability
 - Lots of hands on nodes probable cause
- Pioneering new systems management solution (HPCM) with vendor
 - Combined with new software stack (ATSE)
- At scale testing reveals previously unseen issues
 - kworker bug not seen on Comanche or x86 platforms
- Systems monitoring is CRITICAL (debug/analyze many of above)
- Early hardware requires frequent and quick iterations of software stack
 - Tension with accelerated move to classified where this is a challenge
- Keeping system in sync (updating software) a challenge -> future work with containers



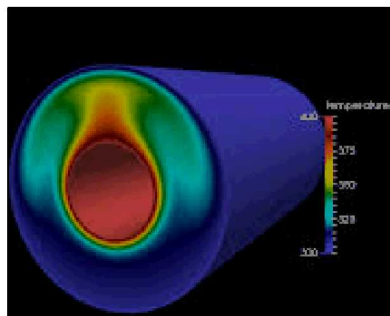
Vanguard Astra: Lessons Learned *or* Reasons to Prototype new Technologies

- ThunderX2 is less reliant on vectorization to utilize available memory bandwidth.
 - Cores can consume available memory bandwidth without vectorized code.
 - Downside: vector units are small so compute-dense code may run slower, extra cores help offset this when comparing node-to-node
- Most of our complex solver libraries and applications compile with GCC or Arm compilers without significant issues.
 - Functional portability for broad code portfolio without significant code rework (NALU, SPARC, CTH, etc.)
 - Acid test is getting the performance out of generated code
- Cache performance will likely impact some of our codes that have reasonable locality
 - Suspect that caches simply perform slower on TX2 versus Xeon
 - Lack of support for gather operations
- Most packages ported and running on the platform, ATSE environment has worked out well



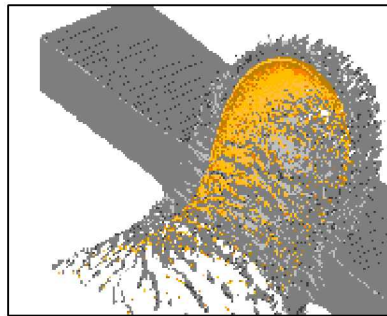
Monte Carlo

1.60X



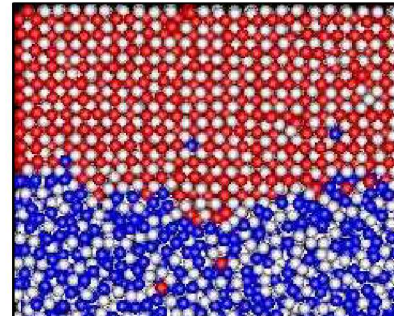
CFD Models

1.45X



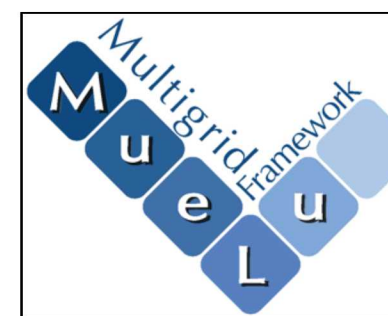
Hydrodynamics

1.30X



Molecular Dynamics

1.42X



Linear Solvers

1.87X

Vanguard Astra

- Arm processor – first time used for HPC platform at scale
- All-flash Lustre storage – not a first but a first with Arm clients
- HPE - New integrator, lots of experience but little specific to DOE and these types of platforms
- Vanguard program will certainly have lowered risk for all for follow on platforms
- New technology and integrator possibilities for NNSA
 - Well this was true...
 - **HPE acquires supercomputer giant Cray Inc for \$1.3bn**

The logo for Arm, consisting of the lowercase letters 'arm' in a bold, blue, sans-serif font.The logo for Hewlett Packard Enterprise, featuring a green rectangular outline above the text 'Hewlett Packard Enterprise' in a bold, black, sans-serif font.

Backup

Challenges

- Initial implementation of thermal design required modification
 - Basic design, 1 cooling unit per 3 compute racks (4 rack set)
 - Initially, 4 rack sets were not isolated from each other
 - Thought was that this configuration would provide some redundancy in event of failure
 - Resulted in turbulence and low pressure regions in input (cold air) side
 - Remedied by isolating each 4 rack set and adding sealing material to maintain pressure differential
 - Prevents hot air re-circulation
 - We are currently seeing a thermally stable system, operating within design specifications
 - 21-22 degree C input water temperature to cooling unit, 24-25 degree C input air
- Memory bandwidth below performance targets
 - Delivered SKU did not provide 247GB/sec STREAM performance, per node
 - After problem was analyzed, determined that majority of existing processors could support specified memory bandwidth with modification
 - Processors that did not pass screening are being replaced
- Various (small number) of other failure modes being root-caused and addressed
 - All inclusive less than 200 nodes are currently being tracked



Artist Rendering



Construction Completed < 1 Year

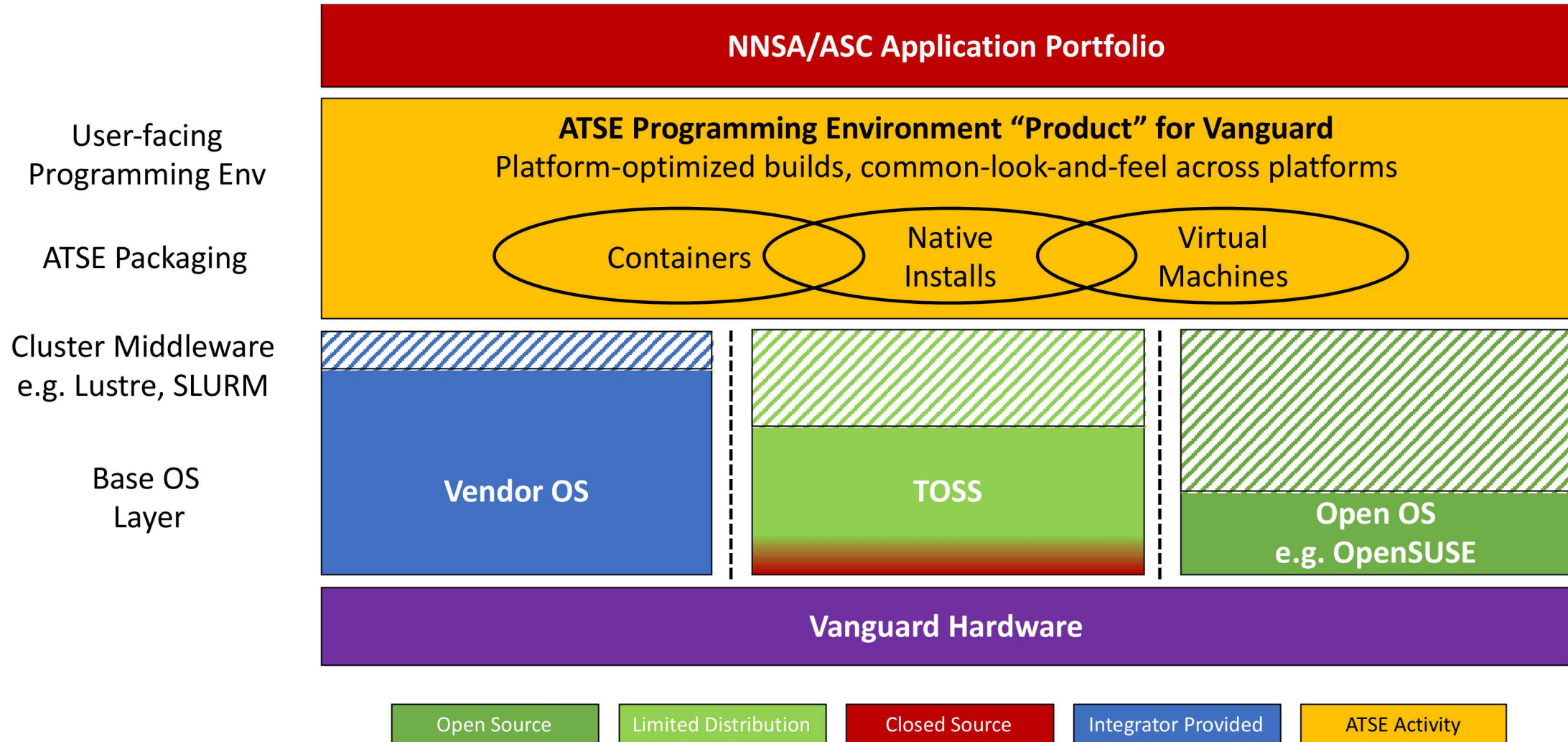


Celebrity Groundbreaking



Delivery and Integration Overlapped Construction

ATSE: Integration with Multiple Base Operating Systems



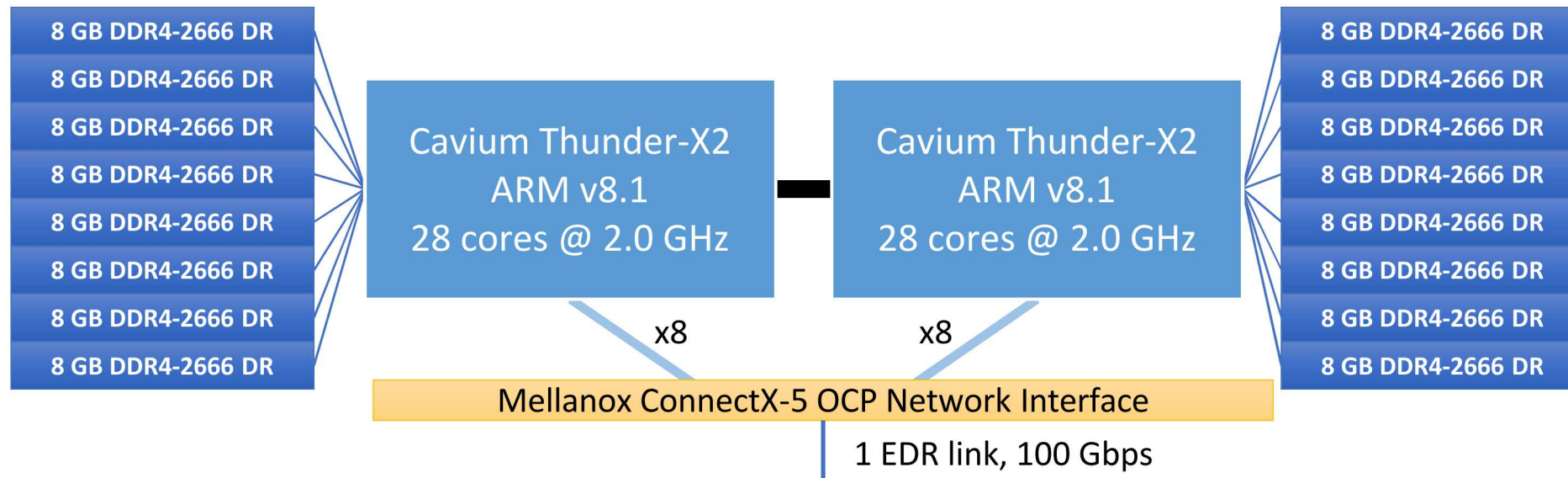
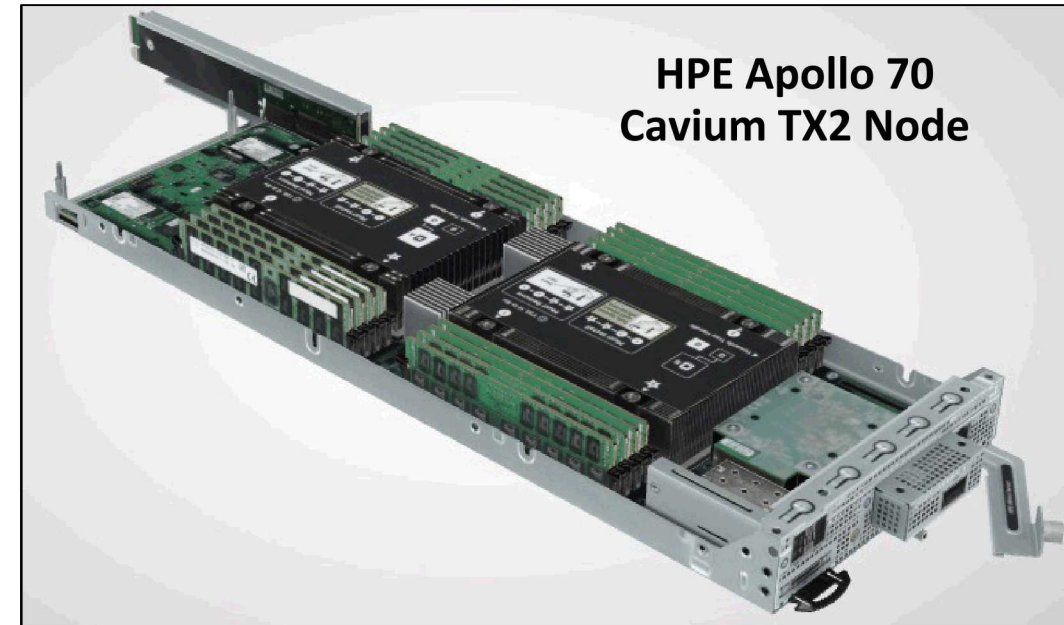
Vanguard Program: Tri-Lab Software Effort (ATSE)

- Advanced Tri-lab Software Environment
- 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 (HPC-optimized Linux, network stack, filesystems, containers/VMs, ...)
 - Job scheduling and management (WLM, scalable app launcher, user tools, ...)
 - System management (OS image management, boot, system monitoring, ...)
- Leverage prototype aspect of system for scalable system software R&D

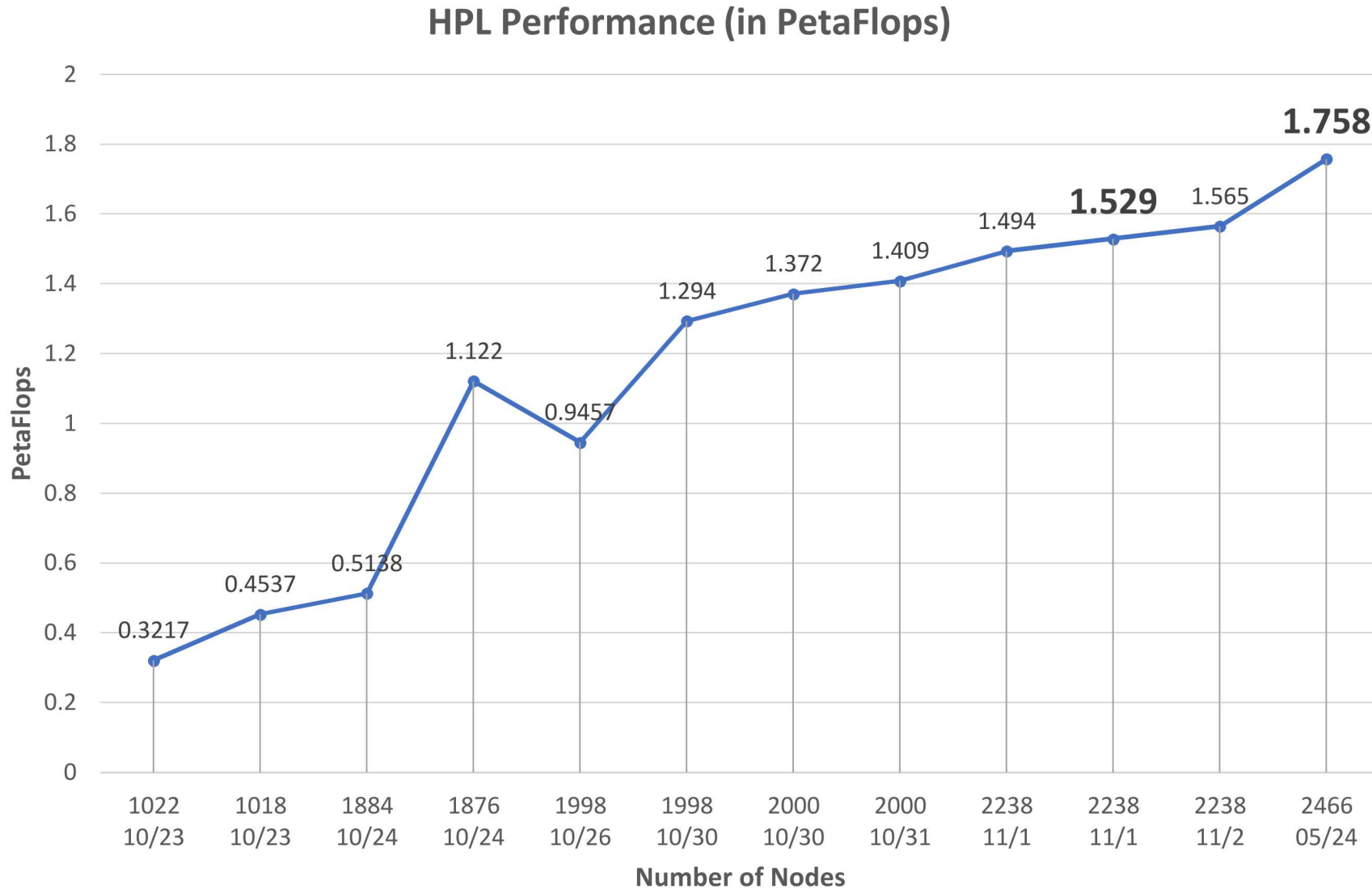
Improve 0 to 60 time... Vanguard-Astra arrival to useful work done

Astra Architecture

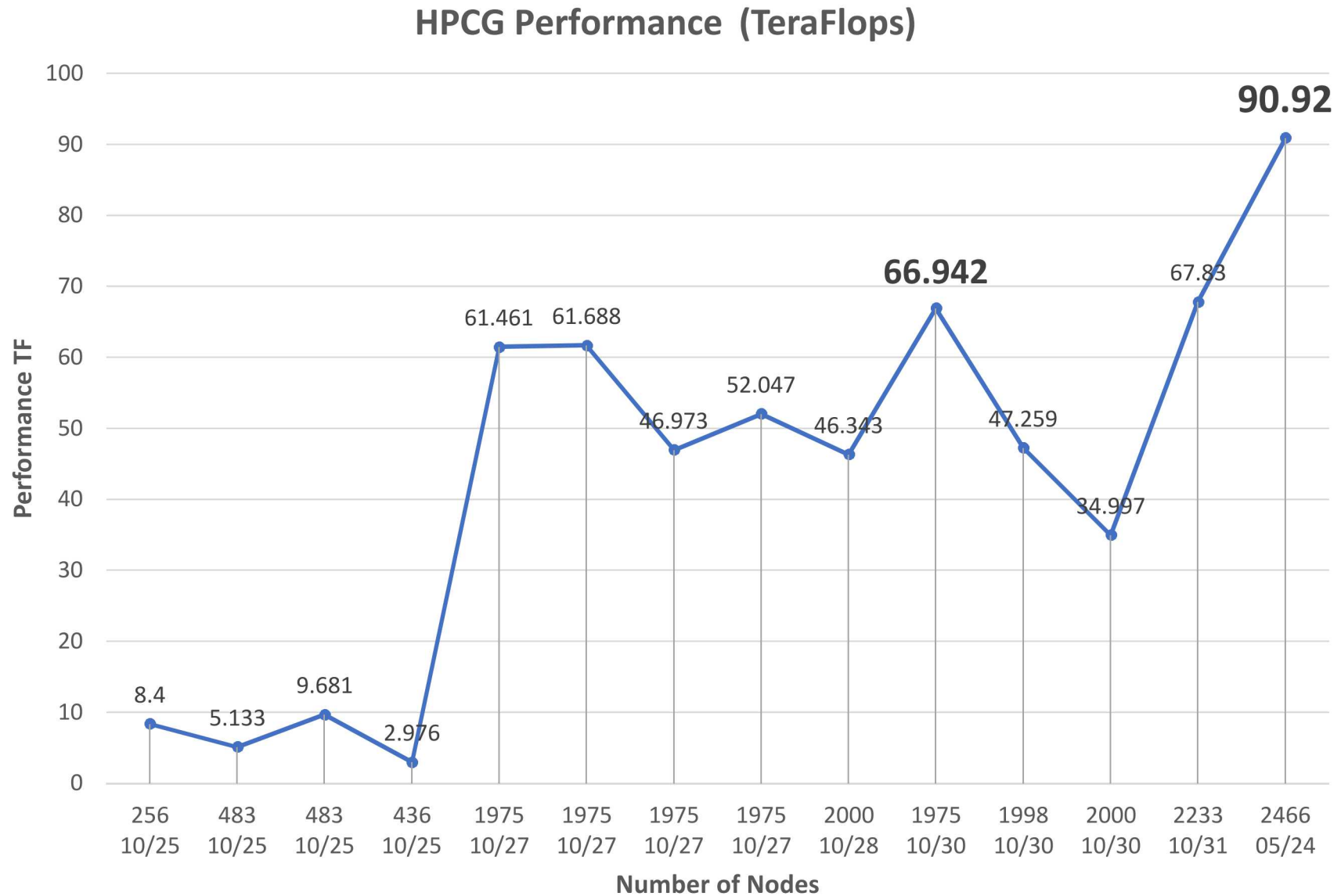
- **2,592** HPE Apollo 70 compute nodes
 - Cavium Thunder-X2 **Arm** SoC, 28 core, 2.0 GHz
 - 5,184 CPUs, 145,152 cores, 2.3 PFLOPs system peak
 - 128GB DDR Memory per node (**8 memory channels per socket**)
 - Aggregate capacity: 332 TB, Aggregate Bandwidth: 885 TB/s
- Mellanox IB EDR, ConnectX-5
- HPE Apollo 4520 All-flash storage, Lustre parallel file-system
 - Capacity: 403 TB (usable)
 - Bandwidth 244 GB/s



Initial Large Scale Testing and Benchmarks (HPL)

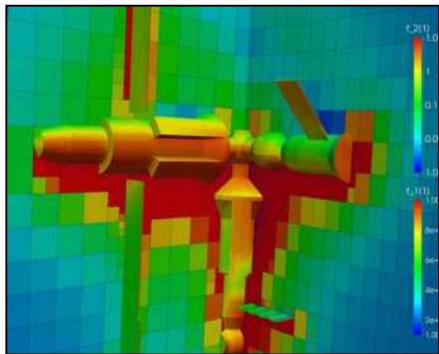


Initial Large Scale Testing and Benchmarks (HPCG)



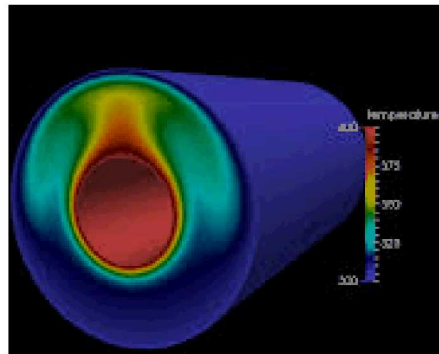
Early Results from Astra

Baseline: Trinity ASC Platform (Current Production), dual-socket Haswell



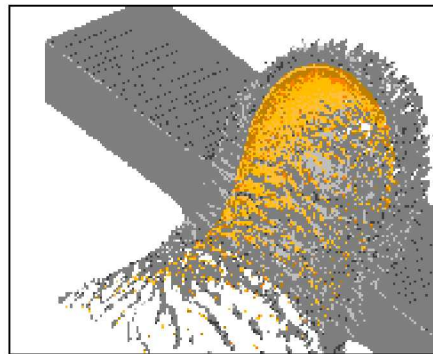
Monte Carlo

1.60X



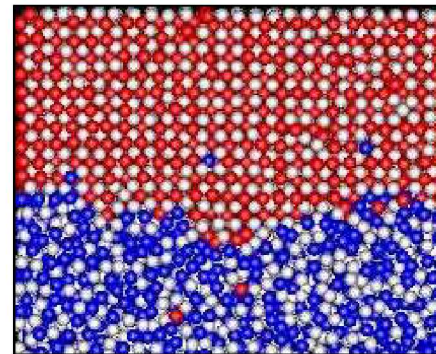
CFD Models

1.45X



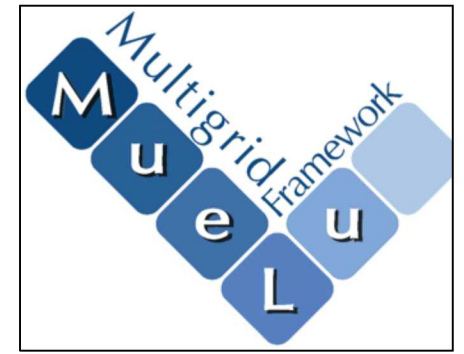
Hydrodynamics

1.30X



Molecular Dynamics

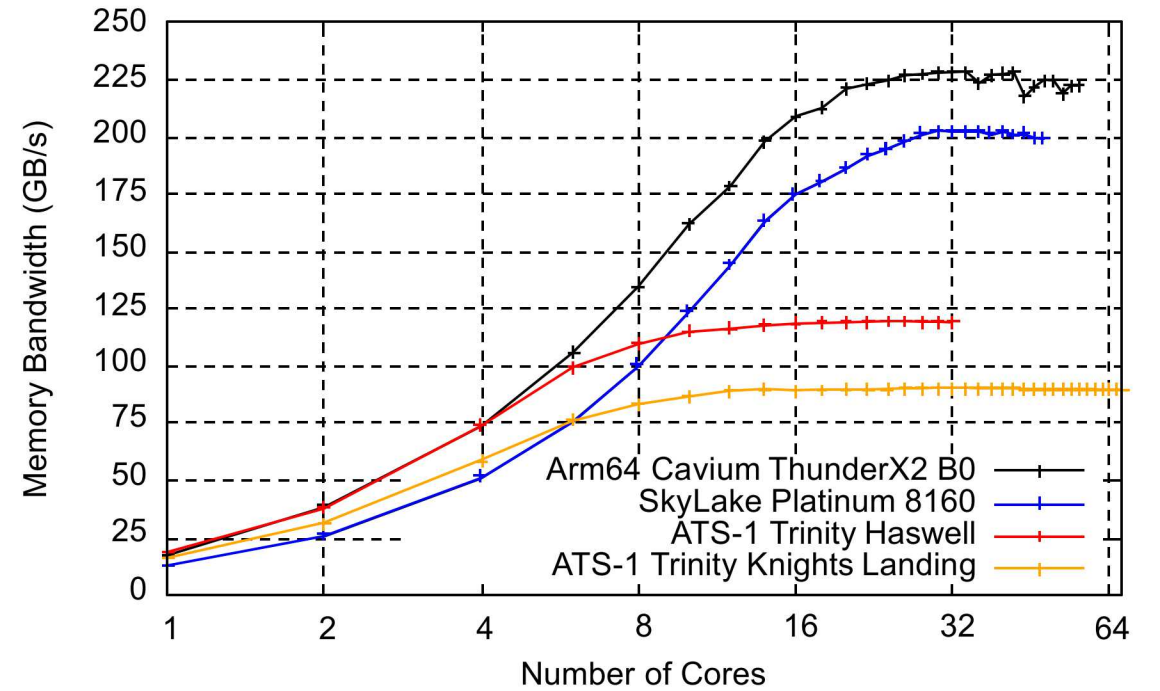
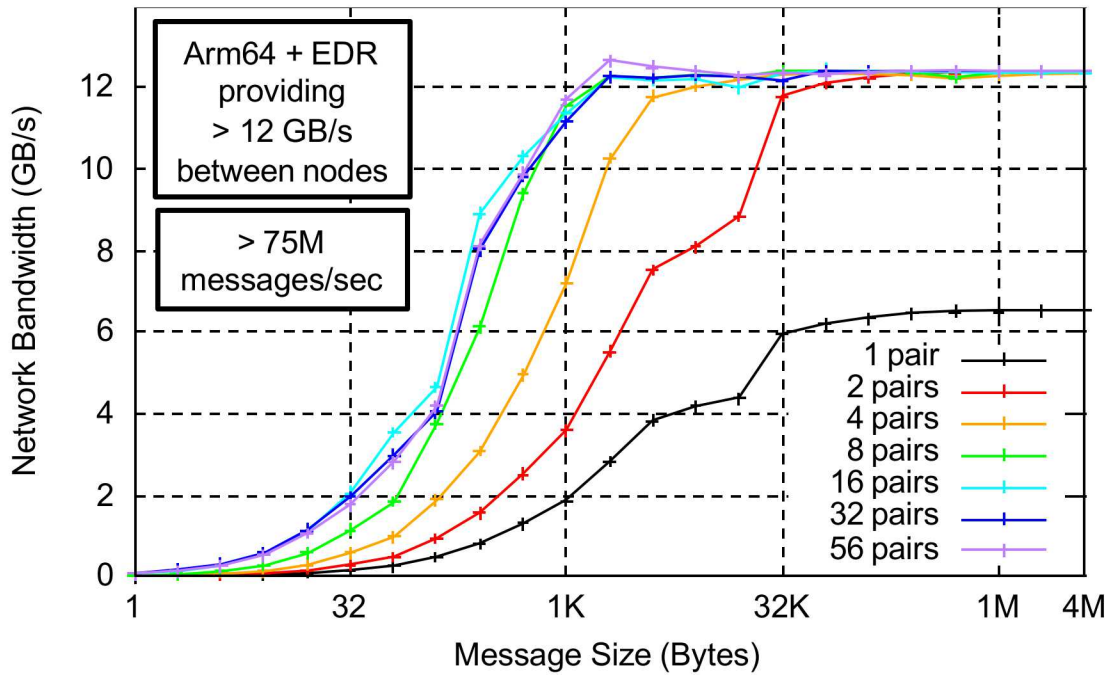
1.42X



Linear Solvers

1.87X

Astra Early Results

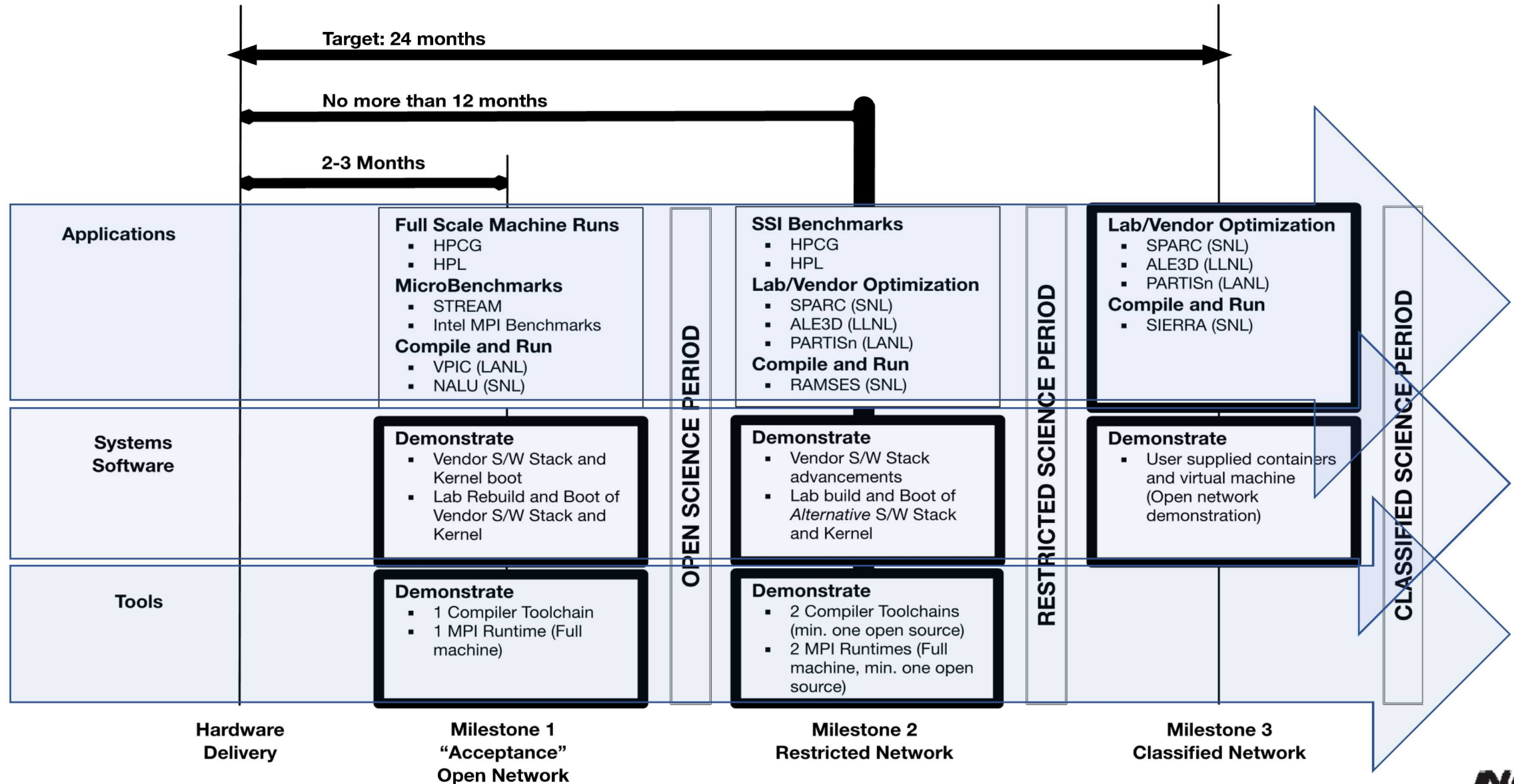


What is ATSE?

- **Advanced Tri-lab Software Environment**
 - Sandia leading development with input from Tri-lab Arm team
 - Will be the user programming environment for Vanguard-Astra
 - Version 1.0 targets including the software components needed for Milestone 1
- **Lasting value**
 - Documented specification of:
 - Software components needed for ASC production applications
 - How they are configured (i.e., what features and capabilities are enabled) and interact
 - User interfaces and conventions
 - Reference implementation:
 - Deployable on multiple ASC systems with common look and feel
 - Tested against real ASC workloads (open networks and classified)
 - Something to point vendors at in procurements

ATSE is an integrated software environment for ASC workloads

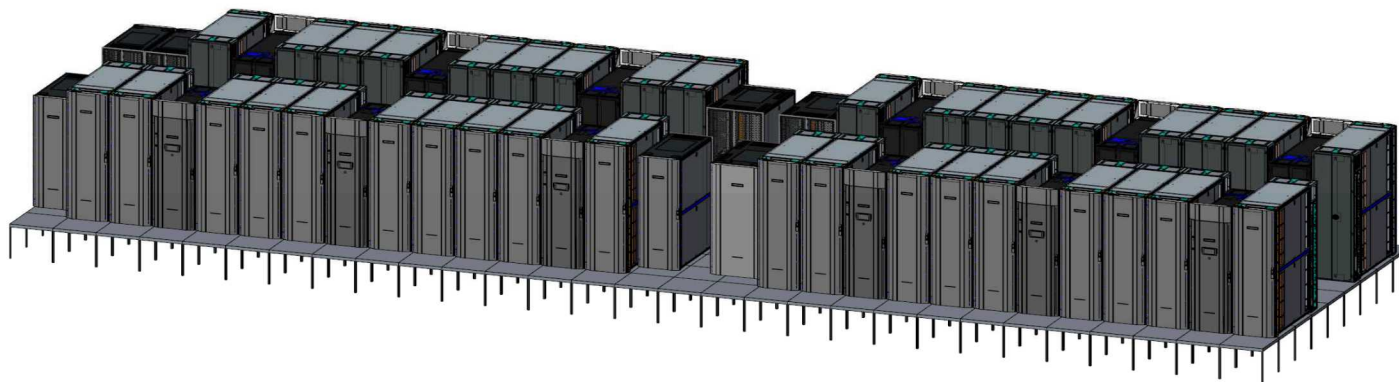
Vanguard-Astra: Project Milestones



Astra Advanced Power and Cooling

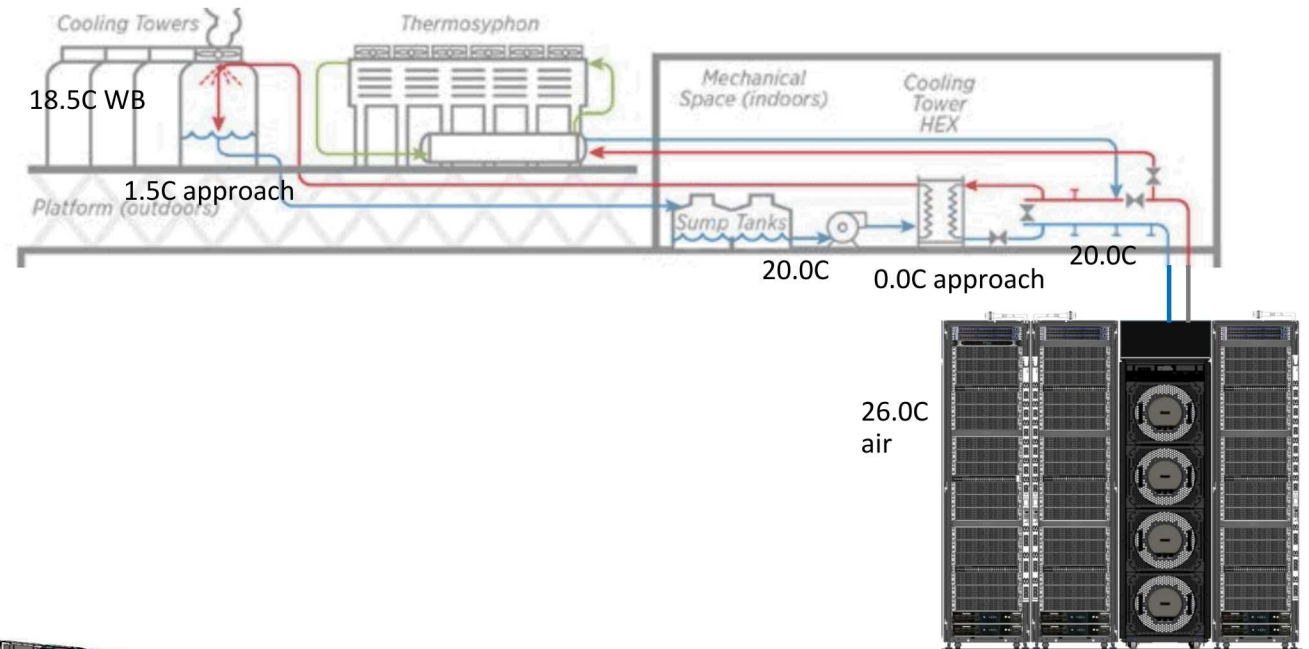
Extreme Efficiency:

- 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



Sandia Thermosyphon Cooler Hybrid System for Water Savings

Efficient tower and HEX can take hottest 36 hours of the year of 18.5C wetbulb to make 20C water to the fan coils



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

- Home
- Technologies
- Sectors
- AI/ML/DL
- Exascale
- Specials
- Resource Library
- Events
- Job Bank**
- About
- Solution Channels
- Subscribe



June 18, 2018

While the enterprise remains circumspect on prospects for Arm servers in the datacenter, the leadership HPC community is taking a bolder, brighter view of the x86 server CPU alternative. Amongst current and planned Arm HPC installations – i.e., the [innovative Mont-Blanc project](#), led by Bull/Atos, the ‘Isambard’ Cray XC50 going into the University of Bristol, and commitments from both Japan and France [among others](#) — HPE is announcing that it will supply the United States National Nuclear Security Administration (NNSA) with a 2.3 petaflops peak Arm-based system, named Astra. On track for deployment at Sandia National Labs later this year as part of the NNSA’s Arm-centric Vanguard project, Astra will be the world’s largest Arm system ever built. according to HPE.



TOP500 Lists

Up from 203

156	Sandia National Laboratories United States	Astra - Apollo 70, Cavium ThunderX2 ARM CN9975-2000 28C 2GHz, 4xEDR Infiniband HPE	138,096	1,758.0	2,209.5
-----	---	---	---------	---------	---------

Up from 36

29	156	Astra - Apollo 70, Cavium ThunderX2 ARM CN9975-2000 28C 2GHz, 4xEDR Infiniband , HPE Sandia National Laboratories United States	138,096	1,758.0	90.92
----	-----	---	---------	---------	-------

Sandia Labs News Releases

November 13, 2018

Astra supercomputer at Sandia Labs is fastest Arm-based machine on TOP500 list

Success suggests additional chip suppliers for supercomputing industry

ALBUQUERQUE, N.M. — Astra, the world's fastest [Arm](#)-based supercomputer according to the [TOP500](#) list, has achieved a speed of 1.529 petaflops, placing it 203rd on a ranking of top computers announced at The International Conference for High Performance Computing, Networking, Storage, and Analysis [SC18](#) conference in Dallas.



Exceptional Service in the National Interest