

Towards the Development of an Exascale Network Digital Twin

John K. Holmen, Oak Ridge National Laboratory (ORNL)

Md Nahid Newaz, Oakland University

Srikanth Yoginath, ORNL

Matthias Maiterth, ORNL

Amir Shehata, ORNL

Nick Hagerly, ORNL

Christopher Zimmer, ORNL

Wes Brewer, ORNL

**EXPERIENCE
ORNL**
MEET. EXPLORE. LEARN.

ORNL is managed by UT-Battelle, LLC for the US Department of Energy



U.S. DEPARTMENT OF
ENERGY

Motivation

- Resiliency a fundamental challenge for exascale systems
 - ~60M components in Frontier
- Component counts and complexity lead to more and new failures
- Important to ensure system functionality, performance, and usability
- Talk captures investigation of network digital twins



<https://www.flickr.com/photos/olcf/52117623843/in/album-72177720299483343/>

Frontier

- 9,408 HPE Cray EX235a nodes
- Theoretical peak of 2 Exaflop
 - Compute similar to 194,544 PS5s
- 74 cabinets weighing 8,000 pounds each
 - Total weight similar to a Boeing 747
- 90 miles of network cables
 - Perth to Wedge Island
- 700 PB of Lustre storage
 - 25 Mt. Everests of DVDs



<https://www.flickr.com/photos/olcf/52117623763/in/album-72177720299483343/>

Digital Twins

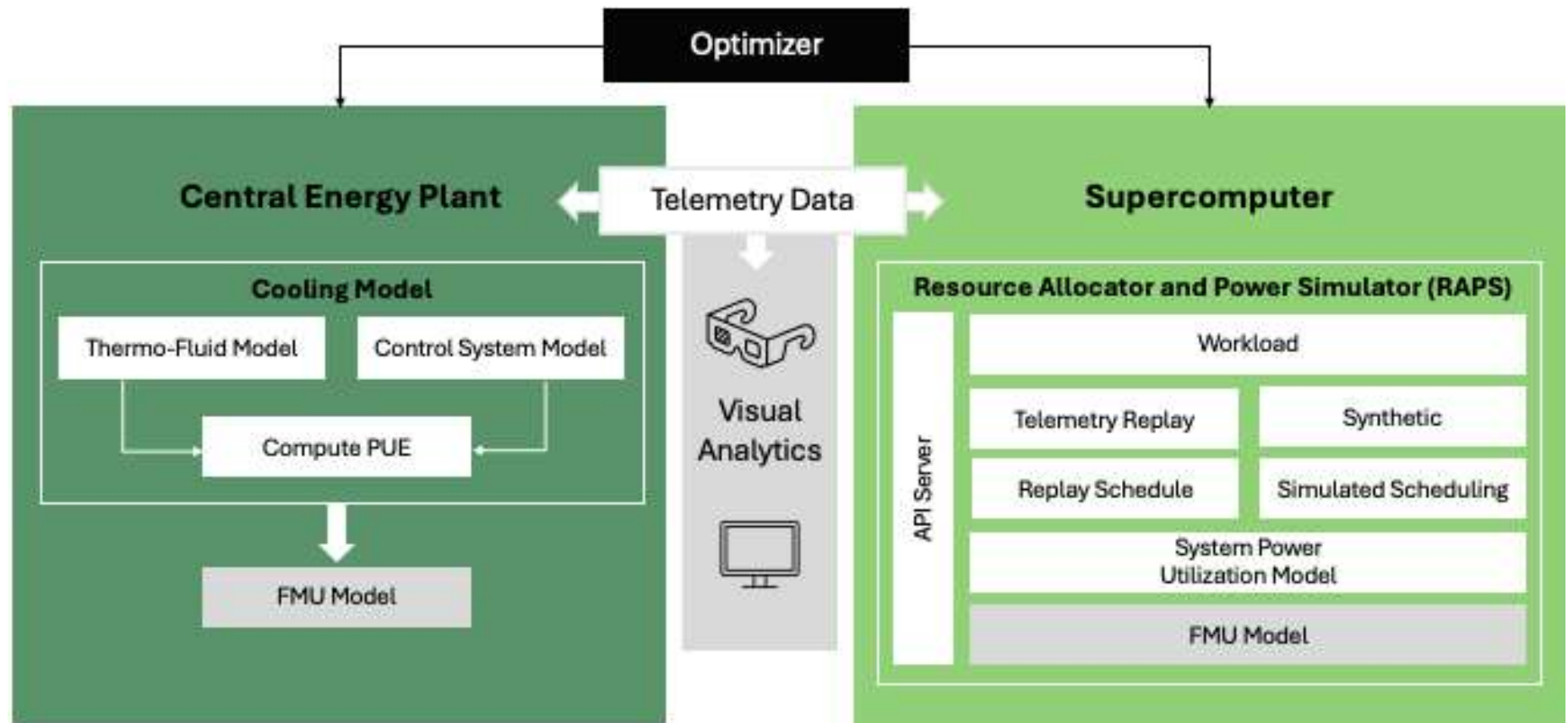
- National Academy of Sciences, Engineering, and Medicine (NASEM) define a digital twin as:
 - “A digital twin
 - is a set of virtual information constructs that mimics the structure, context, and behavior of a natural, engineered, or social system (or system-of-systems),
 - is dynamically updated with data from its physical twin,
 - has a predictive capability, and
 - informs decisions that realize value.
 - The bidirectional interaction between the virtual and the physical is central to the digital twin.”

ExaDigiT

- Multidisciplinary international collaboration
 - Academia, HPC centers, industry
- Community-driven effort to design open-source framework for digital twins of liquid-cooled supercomputers
- Effort guided by 8 working groups
- Differing progress across groups, e.g.,:
 - Established cooling and power models
 - Early investigation for network model



Architecture Overview



Resource Allocator and Power Simulator (RAPS)

- Enables replay and simulation of jobs
- Shows jobs running through queue
- Captures cooling, power, etc.
- Calculates efficiency, carbon emissions, etc.
- Replayed 183 days of Frontier data

Pressure and Flow Rates									
Output					Average Value				
Work Done by CDUs (kW)					2.5				
Facility Supply Pressure (psig)					63.6				
Facility Return Pressure (psig)					38.0				
Facility Flowrate (gpm)					127.6				
Rack Flowrate (gpm)					273.9				
Rack Supply Pressure (psig)					51.9				
Rack Return Pressure (psig)					24.3				

Power and Temperature									
CDU	Rack 1 (kW)	Rack 2 (kW)	Rack 3 (kW)	Sum (kW)	Loss (kW)	Facility Supply Temperat_ (°C)	Facility Return Temperat_ (°C)	Rack Supply Temperat_ (°C)	Rack Return Temperat_ (°C)
1	235	235	227	700	40	22.0	20.1	27.8	31.0
2	133	151	145	429	31	22.0	20.0	28.7	29.1
3	148	131	129	408	30	22.0	20.7	28.5	28.8
4	227	129	128	484	36	22.0	20.4	28.0	29.6
5	115	217	118	450	33	22.0	20.6	28.6	29.0
6	120	115	118	353	27	22.0	20.0	27.4	29.1
7	125	124	125	374	29	22.0	20.4	28.1	28.4
8	130	129	128	379	29	22.0	20.4	28.5	28.4
9	134	131	124	389	30	22.0	20.5	28.4	28.6
10	167	164	205	536	38	22.0	20.7	27.0	30.0
11	286	158	248	692	41	22.0	20.3	27.9	30.7
12	137	135	135	413	30	22.0	20.0	28.1	29.0
13	138	138	141	417	30	22.0	20.5	28.1	29.0
14	138	138	0	276	21	22.0	20.6	28.0	27.5
15	138	142	142	422	30	22.0	20.1	27.0	29.3
16	147	148	135	422	30	22.0	20.4	27.8	29.6
17	255	140	140	535	37	22.0	20.0	27.3	30.2
18	151	152	162	465	34	22.0	20.2	28.4	29.5
19	150	158	151	459	35	22.0	20.2	28.9	29.4
20	150	155	150	455	33	22.0	20.2	28.6	29.4
21	155	167	167	489	35	22.0	20.4	28.0	29.7
22	175	186	178	539	39	22.0	20.0	27.1	30.1

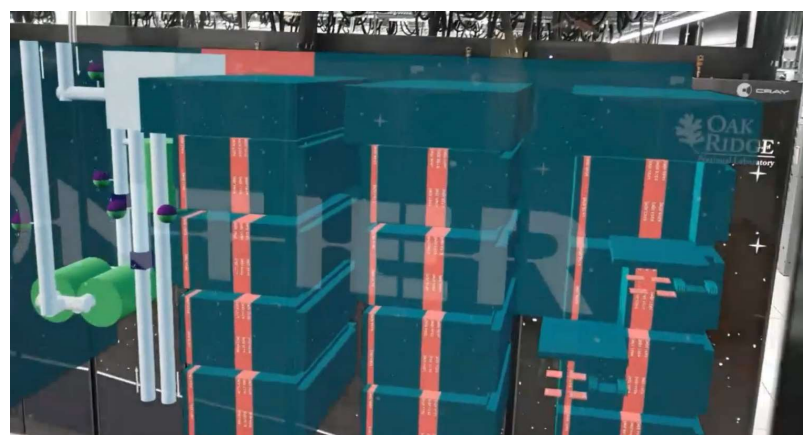
Total Power									
11,910 MW (Loss 0.843 MW - 7.08%) (PUE 1.1)									

Job Queue						
JOBID	WALL TIME	NAME	ST	NODES	NODE SEGMENTS	TIME
1577563	8:42		R	200	35	6:44
1576316	9:40		R	184	64	6:44
1578713	4:00		R	180	42	3:57
1576222	5:56		R	160	41	3:36
1577696	6:00		R	120	72	3:19
1578714	4:00		R	180	67	3:19
1577703	6:00		R	144	54	3:03
1578913	4:01		R	160	47	2:00
1578941	2:00		R	192	79	1:42
1578949	2:01		R	192	78	1:32
1578952	2:00		R	70	14	1:29
1578958	2:00		R	2	1	1:26
1578961	2:05		R	192	89	1:21
1578973	2:00		R	20	5	1:16
1578972	2:00		R	20	5	1:16
1578974	2:00		R	20	12	1:16
1578975	2:00		R	20	10	1:16
1578977	2:00		R	20	9	1:16
1578976	2:00		R	20	11	1:16
1578978	2:00		R	20	13	1:16
1578986	2:00		R	20	9	1:13
1578985	2:00		R	20	12	1:13
1578984	2:00		R	20	3	1:13
1578983	2:00		R	20	6	1:13
1578982	2:00		R	20	4	1:13
1578981	2:00		R	20	3	1:13
1578980	2:00		R	20	3	1:13
1578979	2:00		R	20	4	1:13
1578994	2:00		R	48	11	1:07
1578990	2:00		R	20	9	1:07
1578999	2:00		R	2	1	1:07
1578997	2:00		R	2	2	1:07

Status Update					
Time	Num Jobs Running	Num Jobs Queued	Active Nodes	Free Nodes	Down Nodes
6:45	62	0	5537	3935	0

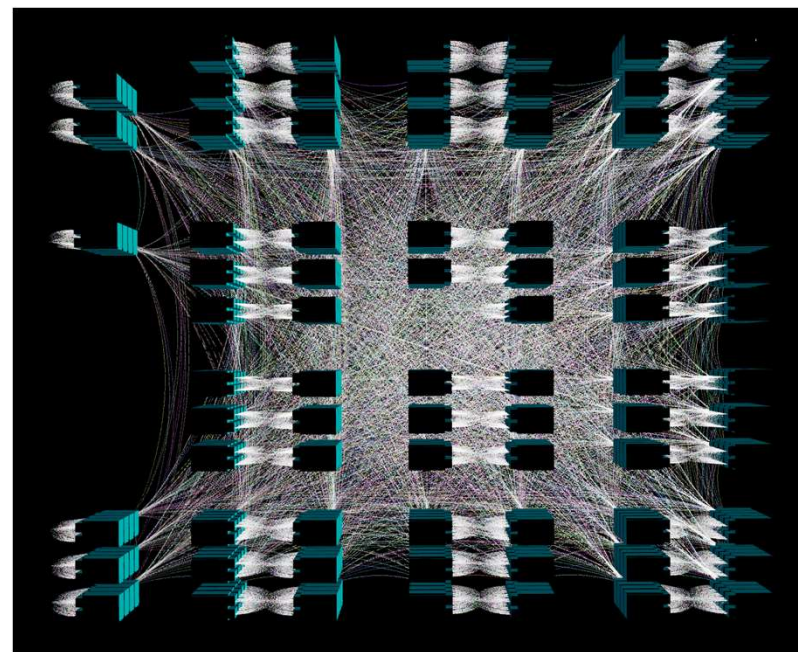
Frontier Augmented Reality (AR) Environment

- Interactive scene visualizing Frontier
 - Implemented using Unreal Engine 5
 - Interaction using Microsoft HoloLens
- Color-coded visualization of various data gathered on Frontier
- Filters allow viewing of:
 - Cabinet interior (i.e., nodes),
 - Cooling infrastructure,
 - Network infrastructure,
 - Power infrastructure, etc.



AR Network Infrastructure

- Recently added network visualization capabilities
- Connections can be color-coded
- Currently displaying all connections
 - Filter logic needed
- Network-related data not yet incorporated
 - Working to identify and gather such data as a part of this work



Network Digital Twins

- Mimics network infrastructure and transmission of data
 - Dynamically updated with data from its physical twin
- Example Use Cases:
 - Understanding and mitigating network congestion
 - e.g., congestion studies and routing optimization
 - Application fingerprinting
 - e.g., characterize and model workloads
 - Improve parallel discrete event simulator models (e.g., those in SST)
 - e.g., validate network models at first-of-kind scales

Target Use Case

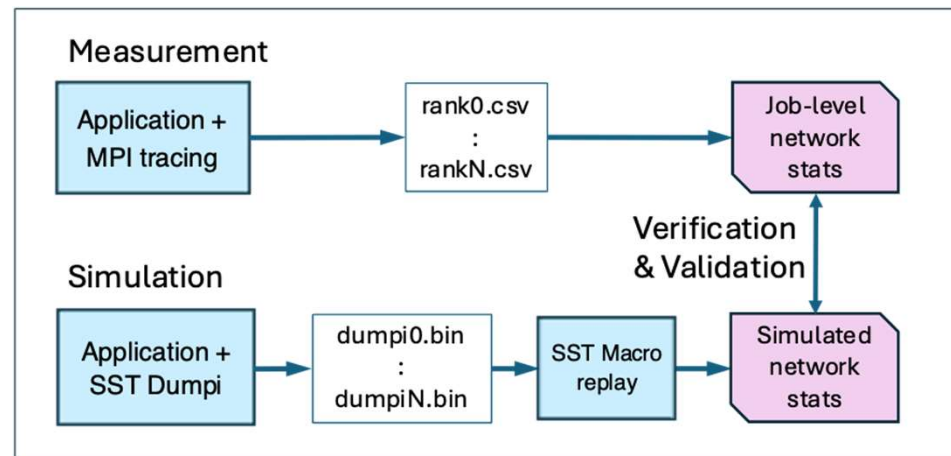
- Short-Term: Extend RAPS functionality
 - RAPS uses CPU and GPU utilization to model system power
 - Improve workload modeling by adding network utilization
- Long-Term: Application fingerprinting
 - Aim to collect:
 - Message sizes and counts for transfers across network
 - Network-related hardware power consumption
 - Combine the two for application fingerprinting
 - i.e., identify types of workloads stressing the network

Tools

- SST/macro (Structural Simulation Toolkit)
 - Platform to simulate full-scale machines and evaluate changes
 - Models to estimate processing and network component performance
- SST DUMPI Trace Library
 - Trace collection and replay tools for MPI applications
 - SST/macro uses trace data to simulate machine variations
- Custom MPI Tracers
 - fi_hook utility and PMPI used to trace calls
- Also explored CrayPat and Darshan+autoperf
 - Too coarse-grained (e.g., summarized statistics after run)

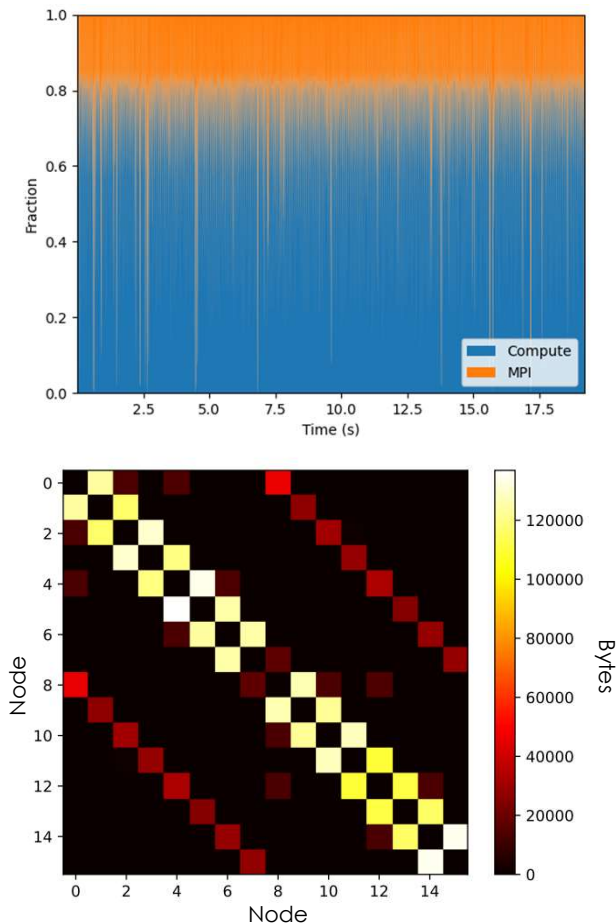
Goal

- Trace applications with SST DUMPI
- Replay traces with SST/macro
- Trace applications with MPI tracer(s)
- Validate simulations



Progress

- Generated fixed-time quanta charts with SST (top)
 - Time-dependent histogram showing split between compute and MPI
 - 1764 rank NAS Parallel Benchmark run
 - Block tri-diagonal solve
- Generated spyplots with SST (bottom)
 - Visualizes message counts or bytes between network endpoints
 - 16 node miniVite run
 - Graph analytics benchmark tool



Progress (cont.)

- Explored CrayPat and Darshan
 - Moved to tracing for more fine-grained data
- Evaluating two approaches
 - Collecting traces using libfabric's hook fabric provider utility, fi_hook
 - Collecting traces using PMPI
 - <https://github.com/hagertrn/mpi-trace>

- Example output:

[Rank 63] MPI_lrecv started 1713984327.205686331, ended 1713984327.205687046 (elapsed 0.000000715), moved 1572864 bytes from source 59
[Rank 63] MPI_lrecv started 1713984327.205688477, ended 1713984327.205688715 (elapsed 0.000000238), moved 1572864 bytes from source 62

- Validation of SST a work in progress

Challenges

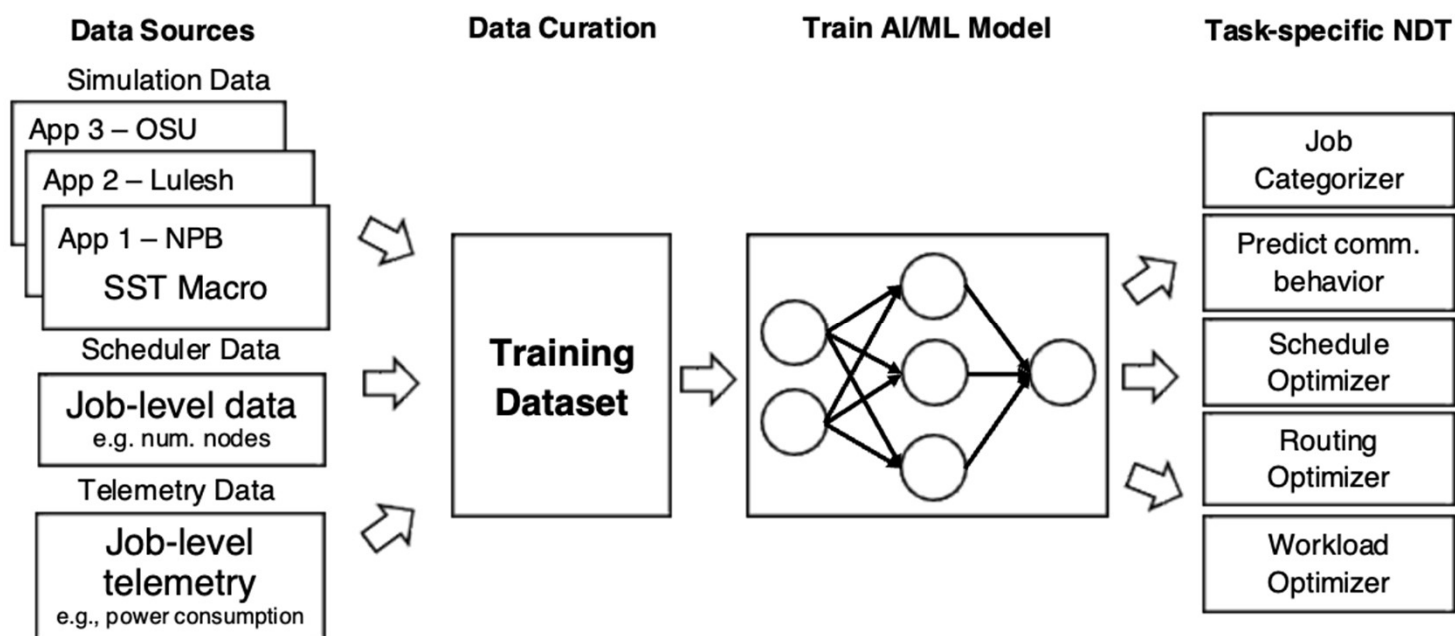
- Data Availability
 - Cooling and power data readily available
 - Network data hasn't been as readily available
 - Manual data collection is time intensive
 - Complicated by uncertainty in which tools are useful
- Data Visualization
 - Node-level details easier to visualize
 - e.g., GPU interconnectedness not needed when visualizing power
 - Network visualization complex due to interconnectedness
 - Difficult to gather meaningful insight from full view

Challenges (cont.)

- Differing Scales
 - Digital twin operates on 15-second intervals
 - Network time scale much faster
 - e.g., 100-350 ns for switching technology
 - Unclear how to incorporate network data
 - Aggregate data to align with 15-second interval?
 - Operate network data on separate interval?
- Simulation Time
 - Explored use of simulation to inform network digital twin
 - Tens of hours to simulate runs with SST/macro
 - Turn-around time not feasible

ExaDigit Integration Goals

- Explore AI/ML as simulation alternative
 - Train on validated simulation data



ExaDigit Integration Goals (cont.)

- Extend Resource Allocator and Power Simulator capabilities
 - Capture network-related data
 - Study relationship between network and power
- Extend augmented reality scene
 - Visualize messages sent/received
 - Filter network components shown to meaningful subsets
- Explore ways to passively gather network data
 - Eliminate need to manually gather data
 - Capture network-related data across all jobs

Conclusions

- Network digital twin important to overall digital twin
 - Frontier's scale helpful for validation
- Unclear how to best integrate a network digital twin
 - Separate operating intervals?
- Challenging to find tools aligned with goals
 - Suggestions?
- Tracers found most helpful
 - Anticipate manual data collection and processing

ExaDigiT Collaborations

- Monthly meetings and standalone working group meetings:
 - Visual Analytics
 - Application Fingerprinting
 - Power & Cooling
 - Use Cases & Architectures
 - Documentation
 - Networking
 - AI/ML/RL
 - Verification, Validation, and Uncertainty Quantification (VVUQ)
- If interested, contact Wes Brewer: brewerwh@ornl.gov

Questions?

This research used resources of the Oak Ridge Leadership Computing Facility, which is a DOE Office of Science User Facility supported under Contract DE-AC05-00OR22725.

holmenjk@ornl.gov

ORNL is managed by UT-Battelle, LLC for the US Department of Energy



U.S. DEPARTMENT OF
ENERGY

**EXPERIENCE
ORNL**
MEET. EXPLORE. LEARN.