

# BoF: The Virtual Institute for I/O and the IO-500

Julian M. Kunkel, Jay Lofstead, John Bent, George Markomanolis

German Climate Computing Center, Sandia National Lab, Seagate Government Solutions, KAUST

2017-11-15



# Outline

- 1 Overview
- 2 Comprehensive Data Center List (CDCL)
- 3 Roadmap
- 4 Summary

# Introduction

## Goals of the Virtual Institute for I/O

- Provide a platform for I/O enthusiasts for exchanging information
- Foster training and collaboration in the field of high-performance I/O
- Track and encourage the deployment of large storage systems by hosting information about high-performance storage systems

<https://www.vi4io.org>



# Introduction

## Philosophical cornerstones of the institute

- Treat every member and participant equally
- Allow free participation without any membership fee inclusive to all
- Be independent of vendors and research facilities

# Open Organization

- The organization uses a wiki as central hub
  - Everybody (registered users) can edit the content
  - Mayor changes should be discussed (see below)
  - The wiki uses tag clouds to link between similar entities
- Supported by mailing lists
  - Call-for-papers
  - Announce list for relevant information
  - Contribute list to discuss and steer organizational issues
- Mayor changes should be discussed on the contribute mailing list
- Members can vote for changes

***Everybody is welcome to participate***

# Wiki Content

- Groups involved in high-performance storage  
*Overview of research groups (evtl. companies involved in research)*
  - Product development the group is involved in
  - Research projects (with links to their source)
  - Tags for layers, products and knowledge
- Tools: *Overview of relevant tools with small descriptions*
  - Types of tools: analysis, benchmarking, I/O middleware
  - Tags for layers and features
- Data Comprehensive Center List (CDCL) / High-Performance storage list  
*Characteristics of data center systems*
  - Editable and owned by the community
- Internal section  
*Provides templates and describes rules for editing the page*

# Comprehensive Data Center List (CDCL)

The CDCL contains system characteristics for sites, supercomputer and storage

## System Model

- The system model has been refined since ISC
  - Now based on an extensible JSON schema, optimized editor
  - Supports now (all) logical components and subcomponents
- Characteristics and peak values
- Measured values \*-500

## Components with characteristics

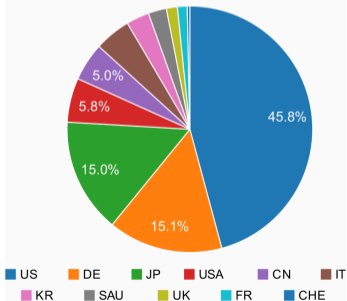
- Site, supercomputer, online storage, tape archives
- Compute nodes, storage nodes, local storage, accelerators, ...
- Supporting: e.g., CPU type, memory available, ...

# CDCL Storage View 2017

## Features

- Table view with selectable columns
- Flexible metrics selection/aggregation
- Multi-year analysis will be supported

storage capacity



#	Site		Supercomputer		Storage		
	Name	nationality	Name	memory_capacity	Name	capacity	
			in PFLOPs	in TiB		in PiB	
1	LANL	US	Trinity	11.08	1,919.03	Lustre	72.83
2	DKRZ	DE	Mistral	3.12	204.00	Lustre02 Lustre01 HPSS	52.00
3	LLNL	US	Sequoia	20.10	1,364.24	Grovo	48.85
4	RKEN	JP	K Computer	10.62	1,136.87	Lustre FEFS	39.77
5	NCAR	USA	Cheyenne	5.33	184.40	HPSS GPFS	37.00
6	NERSC	US	Cori Phase I	4.90	204.00	Lustre	30.00
7	ORNL	US	Titan	27.10	645.74	Spider 2	28.00
8	NCSA	US	Blue Waters	13.40	1,500.00	HPSS Lustre	26.40
9	JCAHPC	JP	Oakforest-PACS	24.91	836.09	Lustre Burst Buffer	24.10
10	CINECA	IT	Marconi A2 Fermi	12.93	413.97	GPFS GPFS	23.71
11	ANL	US	Mira	10.00	698.49	GPFS	21.32
12	JSC	DE	Juqueen	5.90	407.45	HPSS JUST	20.30
13	JAMSTEC	JP	Earth Simulator	1.31	291.04	Home Data Work Archive	19.62
14	KMA	KR	Miri	2.90	0.00	Lustre	19.27
15	NSCC	CN	TaihuLight	125.00	1,191.44	Sunway	17.76
16	AFRL	US	Thunder	5.61	406.54	Lustre	15.54
17	KAUST	SAU	Shaheen II	7.20	718.50	Lustre HPSS	15.28
18	LRZ	DE	SuperMUC Phase 2	3.58	176.44	GPFS	15.00
19	NASA	US	Pleiades	4.97	603.90	Lustre	14.21
20	NSCG	CN	Tianhe-2 Tianhe-1A	59.60	1,169.61	Tianhe-2 H2FS Tianhe-2 Lustre Lustre	14.18
21	TACC	US	Stampede	9.60	245.56	Lustre	12.43
22	ERDC DSRC	US	Topaz	4.57	401.63	Lustre	10.66
23	HLRS	DE	Hazel Hen	7.40	876.75	HPSS Lustre	8.88
24	TEP	FR	Pangea	6.71	49.11	Lustre	8.17
25	GSIC	JP	Tsubame 2.5	5.76	67.67	Lustre	6.93
26	ENI	IT	HPC2	4.60	0.00	GPFS	6.66
27	PGS	US	Abel	5.37	531.14	Lustre	5.33
28	Nagoya University	JP	PRIMEHPC	3.20	83.67	Lustre	5.33
29	ECMWF	UK	Cray XC40	4.25	0.00	HPSS Lustre	5.33
30	ARL	US	Excalibur	3.70	385.63	Lustre	4.09
31	EPCC	UK	Archer	2.55	0.00	Lustre	3.91
32	PNL	US	Cascade	3.40	167.35	Lustre	2.40
33	CSCS	CHE	Piz Daint	7.79	153.70	Lustre	2.22



# Roadmap for 2018

## Web page development

- Provide a Javascript for embedding into any data center web page
  - Allowing the site to describe and visualize their system
  - Hosted by the site directly
  - Allowing a simple export into VI4IO data center list
  - ⇒ Towards a **standardized presentation** of systems !
- Polish presentation of site's information
- Benchmarking section
  - Enabling upload of any benchmark's result and linking them with systems
  - Supporting views to benchmarks
  - For example, useful for IO-500

# Roadmap for 2018

## Supported community activities

- Roadmaps for community benchmarks (ior, mdtest, ...)
- Standardized presentation of systems
- Standardization of lossy compression specifications
- Stabilization of IO-500 + presentation of its results
- New training page linking resources for learning high-performance storage

# Summary

- The Virtual Institute for I/O is a community hub
  - Open to everybody and free to join
- It contains information about
  - Tools, benchmarks
  - Research groups
  - Standardization efforts
- It hosts the Comprehensive Data Center List (CDCL)
  - Covers many metrics and allows flexible visualization
  - Will track metrics across years
  - Can be updated by members
- Contact me if you are interested in **standardized** system presentation
- ***We need you to participate!***

# Appendix



# Collected Information

## Peak Performance

- Theoretical value based on hardware limits
  - e.g. network (server) throughput, SATA limits
- Best performance of one server x number of servers.
- Describe in the text how the peak is computed

## Sustained Performance

- Actually observed performance with an application or benchmark
- You can use any benchmark and measurement protocol
- Just make sure you are not measuring cache effects
- Describe in the text how the value has been measured

## Tags

- Describe hardware and software features individually
- Include coarse grained and fine grained information
  - Lustre, Lustre 2.7, DNE Phase 1
  - Infiniband, FDR-14, fat-tree, blocking 2:2:1
- A taxonomy is needed – but overkill so far
  - Approach: check existing tags and manually fix tag incompatibility

# Tracking Data Across Multiple Years

## Strategy

- Every begin of a year, systems from the last list are copied over
- Decomission: 5 years after installation, systems are removed from the list

## Dealing with hardware upgrades

- Procurement in phases: a small system is delivered first, later a big one
  - If both systems work as one big system, you can first add “NAME phase 1”, then later add the system “NAME”
    - Combine the characteristics
  - If not, then you can keep “NAME phase 1” and “NAME phase 2” systems
- Minor upgrades: e.g., more storage, more compute nodes
  - Just update the system characteristics of this year’s supercomputer
  - Keep the older lists as they are

# Some More Analysis: Relationship Storage/Memory Capacity

- 33 sites are in the list
- Correlation storage cap. vs.
  - memory capacity = 0.64
  - compute peak = 0.057
- Mean(storage/mem capacity) = 59

