

Shared Data Facility (SDF) Overview and Update

Yemi Adesanya

May 8th 2020, Future Computing Frontiers at SLAC



- Deliver **common shared computing infrastructure** to tackle massive throughput data analytics at SLAC
- Enable critical, data-heavy computing workflows in several key mission areas:
 - **SLAC users facilities (LCLS, UED, CryoEM, SSRL), Machine Learning, HEP, FES**
- The SDF infrastructure would offer:
 - **High-throughput** and **high capacity storage**
 - Comprehensive set of frameworks, tools and services
 - **Baseline capabilities for all SLAC users**
 - A cost model for stakeholders with demands that exceed the baseline

- The benefits of a centrally integrated hardware architecture:
 - **Increased operational efficiencies** (lower administration overhead)
 - Coordinated procurements for Economies of Scale
 - Increased utilization by leveraging 'idle/free' compute cycles
- Promote a model for **sustainable scientific computing services**
 - Drive lifecycle and continued support for modern, capable solutions to deliver the science
- **Strong Alignment to Science Goals and Priorities**
 - Partner with science via SDF steering and advisory committees

The Challenge: Raise the Bar for “Baseline” Scientific Computing

- What do we consider to be “Baseline”?
 - A Service that addresses a **common computing requirement** (not unique to any specific project or application)
 - A Service typically **managed and supported through a central organization** (OCIO SCS team)
 - A level of **service the user community expects from the lab as a “birthright” entitlement (for free)**
- Why is Baseline Scientific Computing important?
 - **Baseline services are at the core of many (critical) Scientific applications**
 - **Baseline capabilities help seed new science initiatives before any project-specific grants are awarded**
 - **Baseline solutions foster labwide collaboration and partnership**

The Challenge: Raise the Bar for “Baseline” Scientific Computing

- What is the risk posed by lack of support for Baseline?
 - **No ongoing strategy to address the current and future core computing needs of the lab**
 - **No sustainable lifecycle or modernization**
 - **Decentralization leads to inefficiencies, lack of governance, policy, etc**
 - **Science and collaboration suffers**

Making the Case for Shared Integrated Infrastructure

- **LCLS-II and CryoEM** applications/workflows demand similar high-throughput solutions
- **LCLS-II infrastructure could potentially contribute to the Baseline Capability**
 - 70% of LCLS-II compute time could run other science without impacting LCLS-II operations
- **SLAC Machine Learning initiative** also requires optimal bandwidth between compute (GPU) and storage
- Integrate compute and storage hardware projections from these facilities/projects
- Architect a common infrastructure and consider scalability and total operating cost

	Stage 1 (2019-2024)	Stage 2 (2025-2028)	Main Driver
CPU Compute	1 PFLOPS		LCLS-II
GPU	1 to 10 PFLOPS	> 10 PFLOPS	Cryo-EM + LCLS-II + ML
Disk Storage	10 to 30 PB	50 to 100 PB	Cryo-EM + LCLS-II
Tape Archive	10 to 100 PB	100 to 500 PB	HEP + LCLS-II
Border network	200 Gb/s	1 Tb/s	LCLS-II

SDF is NOT about deploying Siloed Solutions

- Our **existing siloed solutions**:
 - Are **Inefficient** in terms of scalability, utilization and support
 - **Hinder sustainability** of compute, network and storage resources
 - Prevent implementation of **baseline services** to provide meaningful resources for all users; complicates use
 - Impact **long-term planning**

Silos limit our ability to collaborate and align on Computing Strategy!

So what exactly is SDF?

SDF is more than a “facility”, it’s an overarching Computing Strategy

- **An integrated hardware design that includes Storage, Compute, GPU and Fast Networking?**
 - Yes, all of the above. The focus is on fast access to storage from the compute servers
- **A funding model for all of this hardware?**
 - Yes, SDF will standardize on limited number of hardware configurations and coordinate combined purchases with stakeholders / business managers
- **A Datacenter Strategy?**
 - Yes (See Christian Pama). We need to carefully plan for the future infrastructure as it scales over time
- **An Organization?**
 - We’ll develop a matrixed organization of talent distributed across the lab. It will take an entire village to pull this off!
 - SDF will be overseen by a steering committee comprised of key science representatives to ensure alignment with Mission requirements and priorities
- **A set of policies and best practices?**
 - SDF must ensure resources are managed effectively through policies and controls
 - (examples: storage quotas, hardware lifecycle refresh, data retention periods)
- **Raise the bar for Baseline Scientific Computing**
 - Seek lab funding to sustain the baseline
 - Share project resources (LCLS-II) when feasible
 - Continual engagement with the Science Community to stay aligned with evolving requirements

SDF Phase 1 Deployment

Yemi Adesanya

May 8th 2020, Future Computing Frontiers at SLAC

SDF Phase 1 Storage

We're putting the finishing touches on the new storage:

- Two DDN Exascaler 18K controllers for LCLS-II, CryoEM, and Baseline (lab-funded)
 - 765 x 14TB NL SAS drives
 - 52 x 7.68TB SAS SSDs
 - 20 x 1.92TB SAS SSDs
 - NVMe Support
 - Up to 1800 HDDs per controller
 - Expand a single namespace across multiple controllers and storage pools
- ~7.5PB “/sdf” filesystem (Lustre version 2.12.x)
 - Data-On-Metadata ensures the first N bytes of every file are written to SSD
 - Home directories for the SDF compute nodes (no AFS)
 - Shared group space
 - Expect maximum aggregate throughput of ~60GB/sec
 - Final benchmarks once we have the new compute cluster online
 - We are not hitting the controller limit yet (so add more drives for more performance)
- Small shared scratch filesystem
 - For SDF compute nodes only
 - True ‘scratch’ = it will be purged



What about Storage-as-a-Service?

- SDF Baseline storage replaces StaaS
- SDF Baseline is lab-funded! But “free” only up to a certain point!
- We’ll develop some initial user and project quota limits and usage guidelines
- 25GB home directories for SDF Compute nodes
- Baseline will cover most StaaS migrations (30TB to 60TB)
- Baseline will cover *some* shared project space. How do we decide?
 - SDF Steering committee to review baseline storage requests
- Bigger demands (exceed baseline entitlement) will need project funding
 - Minimum buy-in is currently 45x14TB drive pool
- SDF Storage Cost Model coming soon
- We will limit access from legacy environments
- The intent is to build up SDF as we retire older clusters

Storage Cost Model (the tough part)

- How do we sustain storage in the long-term?
- What components should be treated baseline and funded with indirect?
- Purchase vs. Leasing?

Component	Cost
Storage Controller	\$129,864 each
Drive Enclosure	\$7,200 each
Enclosure Cables	\$151 each
14TB NL SAS drive	\$496 each
7.68TB SSD	\$2,704 each
1.92TB SSD	\$944 each
Vendor Support	Increases with the amount of storage

SDF Phase 1 CPU

- Dell PowerEdge C6525 Servers
- Node Specs:



- 2x 64-core *AMD Rome* EPYC 7702 CPUs @ 2.0GHz
 - AVX-256 SIMD
 - Up to 1 DP TFLOP
- 512GB RAM (4GB per core)
- Mellanox ConnectX-6 100Gb/s HDR100 InfiniBand Adapter
- 10GbE Base-T Ethernet
- 960GB SSD



“Rome” Cluster Stakeholders

- 11264 total cores or 176 TFLOPs

Project	# Cores
LCLS	2816
Fermi	2048
SUNCAT	5376
CryoEM	384
HPS	384
SuperCDMS	256

SDF Phase 1 GPU

Integrating 11 new Baseline funded GPU nodes as part of SLAC Machine Learning initiative.

Thank you, Daniel Ratner!

- Dual Intel Skylake 12-Core Processors
- 192GB RAM
- 10 x 2080Ti (11GB Mem)
- 6TB local SSD “scratch”

Existing CryoEM GPU servers will also be migrated to SDF

- New ethernet switches for SDF
 - “Rack-level” switches instead of fabric extenders to the central core
 - SDF VLANs with suitable access controls (security model)
- New 100Gb Infiniband Fabric
 - Optimal bandwidth between SDF Compute and DDN Storage
- Fiber trunk for ethernet and IB links between B050 1st floor <-> 2nd floor
- Switches and Fiber purchased with OCIO and LCLS funds

Datacenter Challenges

- We are deploying Phase 1 solutions in B050 Datacenter
 - Storage on the 2nd floor with generator-backed power
 - Compute on the 1st floor with house power
- Decommission EOL clusters and storage as we roll out SDF
- Repurposing racks and power infrastructure is a challenge requiring some trial and error
- Special thanks to Networking and Datacenter Team (Mark Foster, Christian Pama and Matt Wood)

SDF is a greenfield for modern solutions

SDF migration will not be seamless, but we need to modernize

- **No dependency on AFS**
 - Home directories on the new storage (“/sdf/home/...”)
- **Slurm is our batch scheduler for SDF**
 - Open Source
 - Comprehensive support for GPU scheduling (AKA fairshare)
 - Widely used within the research computing community: SRCC, NERSC, etc.
 - Slurm Support can be purchased from SchedMD
- **Active Directory authentication**
 - AD “Windows” accounts for logins
 - Integrate with open-source Identity Management framework
 - Avoid building dependencies on dated, homegrown infrastructure
 - More potential for streamlined and automated account provisioning
- **CentOS 7**
 - Run legacy RHEL6 applications in Singularity containers

Change is inevitable but it's an opportunity to deliver more for science

Stakeholder requirements will shape SDF capabilities

- The success of SDF will be measured on how we align with the science needs
- SDF stakeholders will provide requirements through a steering committee
- Hardware will be purchased/leased and lifecycled periodically based on current supported standards
- We anticipate a heterogenous (but controlled) hardware environment as technology and requirements evolve
- We must be flexible and responsive