Automated Network Services for Exascale Data Movement

Frank Würthwein, Jonathan Guiang, Aashay Arora, Diego Davila, John Graham, Dima Mishin, Thomas Hutton, Igor Sfiligoi, Harvey Newman, **Justas Balcas**, Preeti Bhat, Tom Lehman, Xi Yang, Chin Guok, Oliver Gutsche, Phil Demar



May 8, 2023





But let's keep in mind that this work is extensible to any collaboration that uses Rucio (adaptable to other data management systems too) to move data across sites

Frank Würthwein, Jonathan Guiang, Aashay Arora, Diego Davila, John Graham, Dima Mishin, Thomas Hutton, Igor Sfiligoi, Harvey Newman, **Justas Balcas**, Preeti Bhat, Tom Lehman, Xi Yang, Chin Guok, Oliver Gutsche, Phil Demar



May 8, 2023



Motivation = HL-LHC

Caltech

CMS expects more than half **an exabyte of new data for each year of LHC** operations during the High-Luminosity LHC era from about 2029-2040

- One annual processing workflow of a few hundred PBs
- Every 3 years, exabyte scale re-processing workflow

Total aggregate data flows are expected to be **dominated by the largest flows**.

- Given that:
 - WAN Network resources are not unlimited!
 - We don't have infinite \$\$; We cannot just throw hardware to the problem

We need to be smart and **use our resources** efficiently







Make Rucio capable to **schedule transfers on the network**. **Improve accountability.**

Predetermined transfer speed and quality of service (time to completion).

Fine-grain managed transfers can be also fine-grain monitored since they travel alone within a well-identified network channel.

Comparing Achieved V.S. Allocated bandwidth will make network & endpoint issues evident.

Quick review: How transfers work nowadays? Caltech





Caltech

Let's STOP using the network as a black box and instrument it!

There can be multiple paths between 2 SEs, but tools don't get to pick which one to use. The routing algorithms know nothing about our priorities



*TPC - Third Party Copy, *AAA - Any Data, Any Time, Anywhere

prioritization



SENSE Architecture



SENSE Architecture



Configure **VPNs** between Routers so we can enforce a given path to be used for specific set of transfers



Implement **QoS** so that we can prioritize certain transfers at the DTN level/Router Level

10Gbps	Everything else	
40Gbps	Medium-important Dataset	XRootD
50Gbps	VERY-important dataset	

Caltech

Rucio/DMM/SENSE Workflow



Proof of Concept Testing



Currently working toward ~400 Gbps site-to-site. Only a few hosts needed for these rates. Working through some technical issues in the areas of End System QoS, FTS/XRootD configurations.



Higher Speed transfers using FTS/XRootD



- Once FTS Transfers are submitted, FTS Slowly increase number of active transfers (see red line).
- Due to this, XRootD endpoints do not get enough streams to reach
 >200gbps.
- Working to increase transfer rates
- Need a dynamic way to control submission rate (to FTS and FTS to XRootD)



Source	Destination	i≣ vo	Submitted	Active	Staging	S.Active	Archiving	Finished	Failed	Cancel	1h)	Thr.		
+ davs://xrootd- sense-ucsd- redirector.sdsc	davs://sense- redir- 01.ultralight.or	CMS	1284	190	-	-	-	14117	-	-	100.00 %	5223.57 MiB/s	al	۲
			1284	190	0	0	0	14117	0	0	100.00 %			

Key takeaways

Caltech

- Today, science workflows view the network as an opaque infrastructure inject data and hope for an acceptable Quality of Experience
- We should allow workflow agents to interact with the network ask questions, see what is possible, get flow specific data and resources
- Science workflow planning should be able to include the networks as a firstclass resource (alongside compute, storage, instruments)
- This requires collaborative cross-discipline teams for workflow co-design
- The same mechanisms that allow the above can also be used by individual networks to distribute traffic more efficiently across entire infrastructure

Thank You!

Want to try out? Join SENSE Testbed? Ask questions? Drop an email to SENSE Group:



sense-info@es.net









This ongoing work is partially supported by the US National Science Foundation (NSF) Grants OAC-2030508, OAC-1841530, OAC-1836650, MPS-1148698, and PHY-1624356. In addition, the development of SENSE is supported by the US Department of Energy (DOE) Grants DE-SC0015527, DE-SC0015528, DE-SC0016585, and FP-00002494. Finally, this work would not be possible without the significant contributions of collaborators at ESNet, Caltech, and SDSC.

Backup slides





SENSE/Rucio/FTS Workflow





What's next?

Caltech

- Development Goals:
 - DMM Development and policies. Allow it be adaptable and define importance of data transfer.
 - Add more sites US (Fermilab (T1), Nebraska (T2), Vanderbilt (T2)), Brazil Sprace (T2), CERN (T2).
 Looking for more European site(s).
 - More NOS (Network Operating Systems) support in SiteRM (Dell OS 10, FreeRTR, Juniper)
 - Quality of Service (Hard QoS, Soft QoS) What to do once underutilized/oversubscribed?
 - Link weights on WAN:
 - Caltech-LasVegas-CERN (130ms, 10gbit max); Caltech-SFO-CERN (163ms, 20gbit max)
 - Policy for fair-share between experiments. Who gets how much and what?
 - Automated End-to-End troubleshooting, monitoring, alarming. (pin-point exact hop failing, alerting)
 - Other experiment use cases and support in SENSE.
- Participate in the WLCG Data Challenge 2024

FTS Transfers via SENSE Path logged in MONIT (using CERN FTS3@Pilot Instance)





UCSD to Caltech testing at higher speeds



- Using FDT (Not FTS/XRootD)
- Green background traffic, Yellow Priority path requested via SENSE
- Total Capacity between UCSD-Caltech (300gbps). Background 200G, Priority 100G.
- Working through some issues with Linux TC, Kubernetes/Multus Private NS Issues. Also evaluating use of BPF and Smart NICs for end-system options, P4TC.

TC – Traffic Control, NS – Namespaces, BPF – Berkeley Packet Filter, P4TC – Protocol-Independent Packet Processor Traffic Control

J.Balcas "Automated Network Services for Exascale Data Movement", CHEP 2023

Caltech