# The DataLake model for the scientific computing of the next decade



As discussed at the DOMA kickoff meeting at CERN with the XDC perspective

> Daniele Cesini – INFN-CNAF Patrick Fuhrmann - DESY



eXtreme DataCloud is co-funded by the Horizon2020 Framework Program – Grant Agreement 777367 Copyright © Members of the XDC Collaboration, 2017-2020



Data Management for extreme scale computing

### Mainly a report form the Data Organization Management Access KickOff





WLCG DOMA kick-off: https://indico.cern.ch/event/729930/

2018-06-14

#### 2018-06-14

D.Cesini / P. Fuhrmann - The DataLake model - INFN CCR WS 2018 - Rimini

Reduce Cost: local + Stobal Scale out: Shared Infrastucture (Science ada coperation) Cost: Local site adm QOO: DBS team larg Site mostly UR Rugs Resource Usage Optimisation RAIN-X -> (LAD) SiteA N-Roy Q.S Q.S: Q.S. O.S. Import/W

### DOMA



WLCG DOMA kick-off

Photo: Courtesy Xavier

### **The DataLake**





# Data Lakes in Wikipedia



X A data lake is a method of storing data within a system or repository

••• in its natural format

---+ facilitates the collocation of data in various structural forms

.... usually object blobs or files.

### X The idea of data lake is to have a single store of all data

··· → raw data

••• transformed data for Reporting, Visualization, Analytics and machine learning

### ✗ The data lake includes

- → structured data, i.e DB
- semi-structured data, i.e. CSV, logs, XML, JSON)
- unstructured data, i.e. emails, documents, PDFs

→ binary data, i.e. images, audio, video



2018-06-14

D.Cesini / P. Fuhrmann - The DataLake mode

# Data Lakes in Wikipedia





# **Data Lakes in WLCG**



X Storage Consolidation

- Storage services are the most complicated to operate in WLCG
  - ---- Investigating the possibility to reducing the number of storage endpoints
- ---- Consolidate, especially at the national level, different endpoints into a single distributed instance
- X Data lakes are an extension of storage consolidation
- See Geographically distributed storage centers, potentially deploying different storage technologies, are operated and accessed as a single entity
- X Storage centers in a lake should be connected through high bandwidth links (> 10Gb/s in 2017) and relatively low latency (< 50 ms), consequently it is expected that a single data lake would not span more than one continent.

2018-06-14

# On the definition of the Data Lake



Very diverse understanding of the expression data lake.
 Attempts to define it by 'name space' or 'region' or 'country' all failed.

- X Good suggestion by Maria Girone: We should avoid using the expression
  Data lakes 2
- X Looks like one thing, but is composed of many



Or...



### X Looks like one thing and is composed by one thing...



X CachingX Network Performance

2018-06-14

# **Data Lakes in WLCG**



Experiments already handle very complex data distribution workflows and rely on a certain level of non co-location between data and processing units

X They all implemented techniques and workflows for remote reading of data

X The point here is to move this "intelligence" and "orchestration" of geographically distributed data **at the infrastructure level** 

# **QoS: What is that ?**



- X Ongoing work started with INDIGO. Great progress based on work from KIT and CNAF.
- Instead of caring about disk and tape on the experiment framework level, it would be less complex and more future proof define the quality of local storage, concerning persistency, access latency and price
- X These 'storage classes' need to be well and commonly defined

### X Example:

Storage centers define two (or more) classes with names and properties:
'archive': Probability of data loss < (1 in 1.000.000) and max access latency 2 hours</p>
'archive-online' same data loss probability but max access latency 1 ms.

X QoS classes need to be discoverable

# **QoS Advantages**



- X We don't care how storage provides fulfill their QoS SLA: e.g. Tape, CEPH with 20 copies or engraved in stone.
- X We are prepared for storage technology changes, w/o changing our data persistency model.

----> Eg. Tape could be replaced by disk arrays or cloud storage

- → HPC might need new storage qualities, like super low latency (BeeGFS)
- We could directly map to cloud providers
   Amazon: Glacier (cheap, safe, high latency), S3 (low latency, expensive)
- X QoS will become discoverable and orchestration middleware (Rucio) can use matching or AI algorithms to select the right combinations of 'classes' at the storage provider endpoints.
- X This kind of model is attractive to new (younger) communities. They essentially don't know anyway what a tape is.

# INDIGO-DataCloud/XDC CDMI QoS



### Quality of Service in storage – broker page

### Available Qualities of Storage

Name Access Latency [ms] Number of Copies Storage Lifetime Location Available Transitions Storage type 虛 disk 100 1 DE Processing tape, disk+tape DESY 18c disk+tape 100 2 DE Processing tape DiskAndTape 50 з DE SKIT TapeOnly 20 years Processing 2 DiskAndTape IT. 50 2 Processing **PSNC** SICIT DiskOnly 50 3 20 years DE Processing and . 1T DiskOnly 50 1 Processing profile1 10 3 20 years DE Processing profile2 10000 2 DE profile2 Archival profile1 -20 SSDDisk 1 IT Processing StandardDisk, Tape 10 Access Latency [ms] • 2 StandardDisk з 20 years 1000 Number of Copies  $\bullet$ **Storage Lifetime** • Location **Available Transitions** ۲ D.Cesini / P. Fuhrmann - The DataLake model - INFN CCF 13

KIT

**INFN** 

# **QoS in the Lake**



X We are looking for a solution in which the lake behaves as hierarchical system and optimizes the data organization based on policy (first) and usage (after) exploiting different QoS

Moreover, the solution should foresee the possibility to attach a volatile storage to the lake
 It should be used as tactical storage, hosting a redundant set of data, to optimize data access and the system should auto-recover in case the volatile storage disappears

# A QoS-based workflow in the Lake



- X One can imagine a data workflow where files that need several levels of processing start with a fast disk (hot) requirements and eventually transition to a nearline disk (warm) to finally end up archived on tape (cold) with a low likelihood they get recalled shortly after.
  - The three tiers mentioned before strongly relates to performance, reliability and resiliency which are metrics dominated by cost.
- X The transition between the storage tiers (or *hierarchy* change) can be either:
  - passive (ie. estimations on processing/access/replication time for a particular data/set)
  - ----- active by implementing triggers at metadata level once the needs for the data/set change.
    - This approach opens the possibility to steer data workflows in real-time through extended attributes "flags" on the namespace.

2018-06-14

### **The Data Lake**





# What are we told would justify changes ?

- X Reduce cost for storage
  - ···→ Global (WLCG level) and local (Site level)
  - → Hardware
  - Operations
  - ··· → Number of replicas is already below 2 on avg (Atlas reported 1.3)

### X Scale out

- ----> Does the current model really have a scale-out problem ?
- Which architecture would solve that ?

# X Or simply : Evolution forced by 'external' technologies or methodologies. (Best example is the 'cloud')

2018-06-14

# **On Cost**



- X Saving costs by reducing operational complexity. Or better: getting more storage for the same money
  - ----> Focus on larger sites, assuming small sites are ineffective
    - ----> There are no numbers available, how much that would save us
    - ----> There are ideas to run smaller sites 'operator-less' regarding storage
    - Probably it's better to avoid drying-out universities
- Operational cost savings are not automatically transferred to more storage space. That funding could simply disappear
- X Saving cost by providing fewer copies or smarter caching of data
  - ---- Improved high level orchestration. (data placement)
  - Well defined storage retention (QoS, see later)
  - ----> Balancing network against storage (depends on network costs)
  - → With or w/o smarter caching.

2018-06-14



## Some scenarios (Taken from XDC)

2018-06-14

# Data Lake, simple cache site In production for years e.g. NDGF, Michigan Data Management Center Data Pool Controller Satellite Site **Remote Data Pool**

2018-06-14

# **Complex multi provider scenario 1**





### **Advantages**

no additional software stack needed at sites.

### Still disadvantage

Local data not accessible in case data link is down or central service not available

2018-06-14



2018-06-14



# **Requirement for the complex cases**

- X We would need to see a benefit
- X We would need to agree on
  - A namespace synchronization protocol
    - ··· → On client request
    - ••• On client request plus subsequent lazy background fill (like Dynafed)
    - Message queues for selected namespace events (producer consumer)
  - ACL synchronization
  - Identity Mapping

#### **XDC** minimum viable product":

- read-only operation
- fetch data on miss with service credentials
  - -- data can be chunked or full-file
- manage cache residency, evicting data when necessary
- HTTP frontend with group-level authorization



# **Caching infrastructure with XROOTD**



### A Distributed XRootd Cache



Production Goal: Distributed cache that sustains 10k clients reading simultaneously from cache at up to 1MB/s/client without loss of ops robustness.

#### https://indico.cern.ch/event/729930/contributions/3021531/attachments/1661988/2663057/WLCG\_DOMA.pdf

2018-06-14

# **Caching infrastructure with XROOTD**



INFN

### A Distributed XRootd Cache



#### Distributed scenario with XROOTD CACHE redirector

Credits: D. Ciangottini, D.Spiga, T.Boccali – CMS and XDC

Geographically distributed cache

- The very same technology used on local scenario can be geodistributed
- Use ephemeral storages to enhance jobs efficiency
- Leverage high speed links to reduce the total amount of allocated space



Operations

D. Geshii / F. Fuhimanii - The DalaLake mouel - Infin GGK WS 2010 - Kimmi

# **Discussion on caching @DOMA**

X Caching will cost hardware

-----> caching will have to be used wisely

me making sure that one caches most of the reuse data.

# X We do not have today clear numbers that caching is really effective

facilities w/o storage (HPCs, CPU only sites)

facilities with storage, but dealing with non local data

We even do not have a forum to discuss this

X Caches need to offer the capability to be pre-filled (preplacement) or be passive caches (cache based on access)



# Where is the data lake in those scenarios ? 🕉

- X The complex scenarios above could either
  - → be inside a data lake
  - → connecting data lakes.
- X Optionally data lakes are connected by DM solutions (Rucio). See Vincent's presentations for the NorduGrid Meeting.



# Summary

🗙 Data Lake

Expose unique storage entry point – looks like one but it's (or can be) made of many

# Support different QoS for the internal components Hierarchical storage

- Caching at (every?) storage and CPU location
  Pre-placement vs access-based
- Intelligent management of data movement RUCIO?
- Unique namespace? (maybe)
- To be deployed in parallel to the current infrastructure
- ---- Cost improvement to be understood



