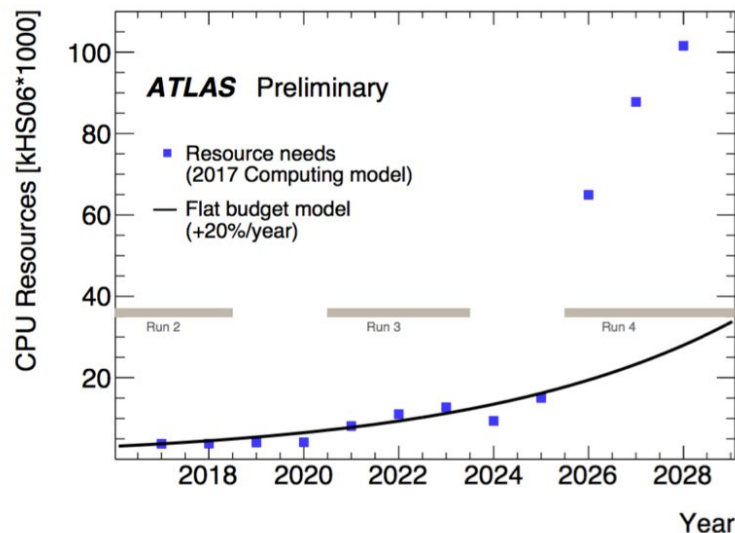
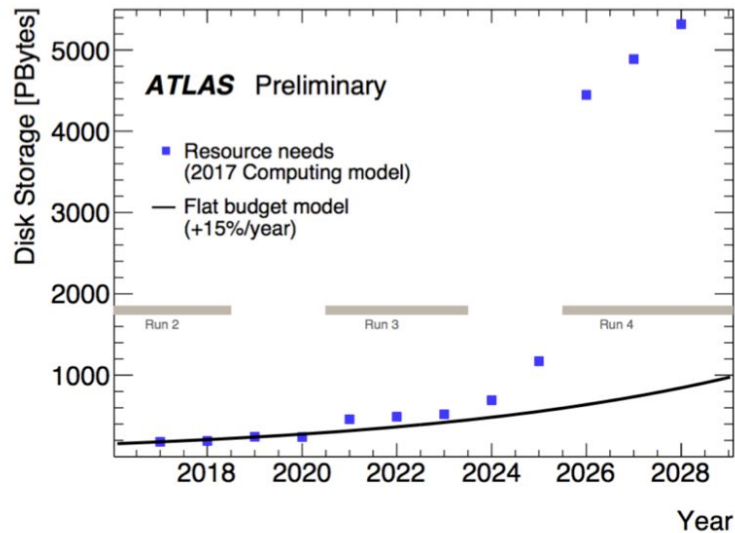


Understanding the performance of a prototype of a WLCG data lake for HL-LHC

Simone Campana, Xavier Espinal Currul, Maria Girone,
Ivan Kadochnikov, Gavin McCance, Jaroslava Schovancová
CERN IT

Motivation

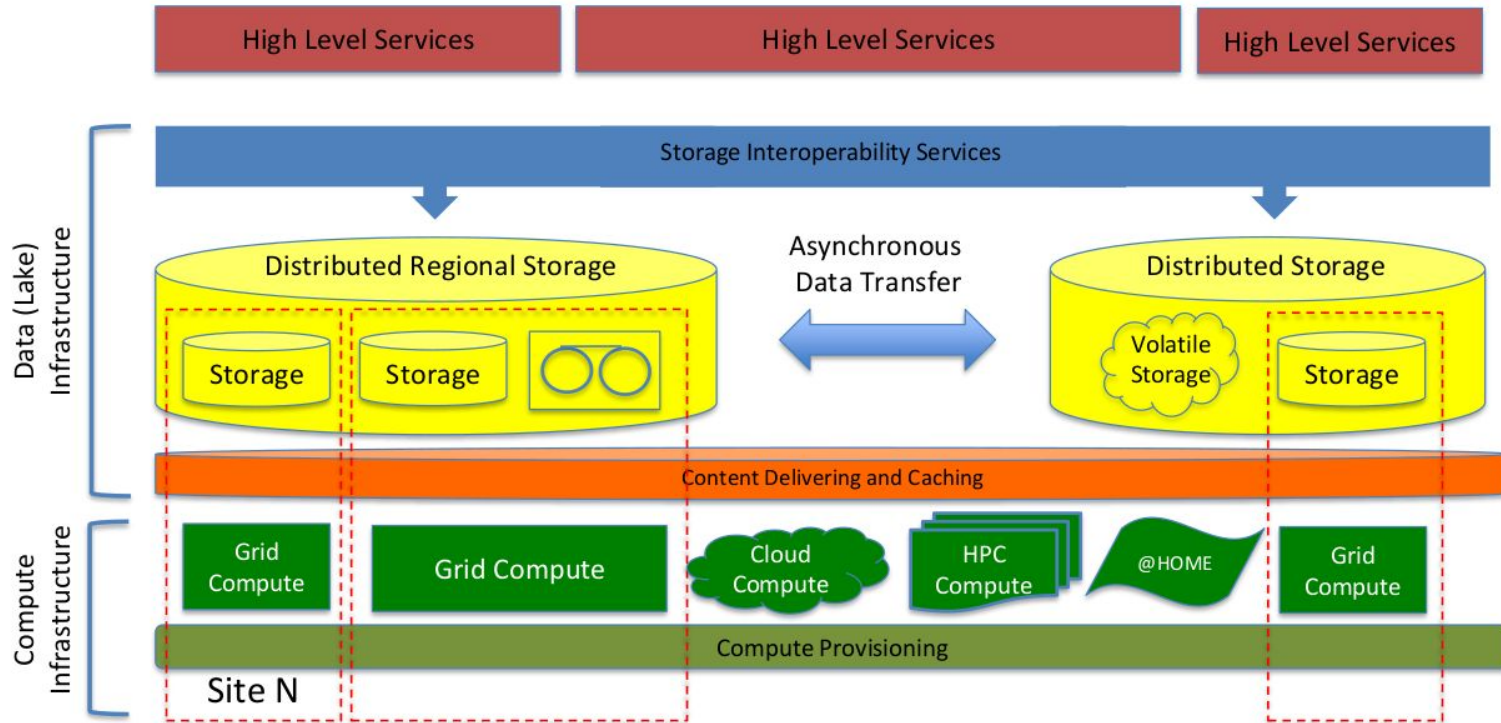
- HL-LHC storage needs are above the expected technology evolution (15%/yr) and funding (flat)
- We need to optimize HW usage and operational cost



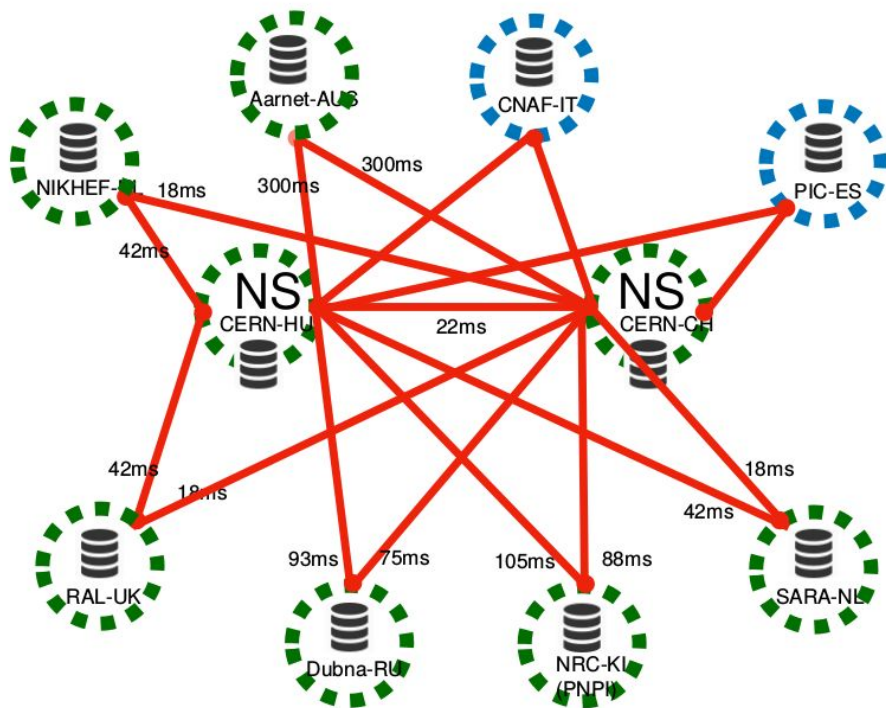
How to reduce cost???

- Many places where we can reduce cost.
Here we **focus on storage** which is one of the bigger contributors.
- **Reduce HW cost:** introduce the concept of Quality of Service (QoS)
 - we store more than we think today!
 - EOS: 2 copies
 - CEPH: 3 copies
 - dCache: Raid-N
- **Reduce Ops cost:** deploy fewer (larger) storage services
- **Co-location of data and compute** not guaranteed

Data and Compute Infrastructures



Data Lake Prototype



- Goal: testbed to test and demonstrate some of the ideas
- Deployed a Distributed Storage prototype, based on EOS
 - distributed storage
 - network links: latency, bandwidth
 - storage media: disk/cache/tape
 - evolving data access protocols: driven by the changes in networks
 - evolving inter-storage communication

The core metric: event throughput

- Compute side of things \Rightarrow boils down to the **event throughput at the same cost**
 - \Rightarrow **Are we able to support the same or even better event throughput at the same cost with the evolving storage configuration?**
- **Easier said than done!**
 - Which events? Which SW? How much I/O? How much memory? ...
 - How to measure job performance? Storage performance?
 - How to benchmark?
 - What to take into account for the storage configuration?
 - Topology of resources? its transparency?
 - (Co-)location of data vs. compute resources?
 - Types of storage media vs. access policies?
 - Direct vs. remote access to data?
 - How to evolve tools to support the core mission

Measurements

- Methodology, how to measure and benchmark
- What to measure: event throughput
 - I/O rate
 - Stage-in / Stage-out time
 - SW init time
 - Time spent in the event loop
- Production and Analysis workflows
- Core count preferences: MSCORE (production) vs. SCORE (analysis)
- Local vs. remote data access

Benchmark

- Resources: standard storage vs. distributed storage
 - can compare these flavors of resources
 - in different configurations of the **distributed storage**
 - hot/warm/cold storage
 - caching
 - local vs. remote access
 - data replication policies/stripping
 - downtime/recovery of subset of storage resources
 - benchmarking per resources, VM

⇒ study and benchmark both

- **job performance**, and
- **distributed storage performance**, at once

Workflows types - ATLAS

- G4 simulation
 - CPU intensive, not so much RAM demanding, not much I/O intensive
 - ttbar full simul, reference workflow to compare HS06
- Digi+reco
 - some I/O (not that much IOwaits for jobs), RAM-demanding, sensitive to latency
 - Event mixing, digitization, trigger, trigger reconstruction
 - 50 GB in
- Production derivation
 - More I/O intensive
 - Skim, slim, ...
 - 5 GB in
- Analysis - focusing on analysis derivation

Workflows types - CMS

- Understanding the equivalents
 - G4 simulation: quick
 - Reco takes more time
 - Premixed pile-up
 - CMS pre-mixes min bias \Rightarrow huge files, less copies. Perhaps lower I/O?
 - ATLAS does not pre-mix min bias \Rightarrow smaller files, more copies
 - No derivations
 - Analysis
- Production workflows in CMS: leverage the “1-chain” job <https://doi.org/10.1007/s41781-017-0001-9>
 - Generation - Simulation - Digitization - Reconstruction steps in 1 job, to save data stage-out and stage-in among jobs
 - \Rightarrow very small input and 1 output of the full chain

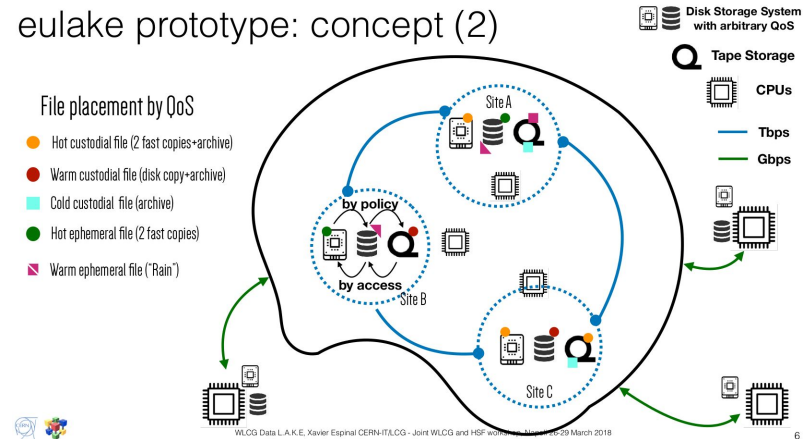
Data access modes

- ATLAS: copy to scratch vs. directIO from co-located storage vs. read over WAN
- CMS: remote read

ATLAS

storage vs. compute	Data access mode	Standard storage	eulake
co-located	copy to scratch	✓	✓
	directIO	✓	✓
not co-located	copy to scratch	?	✓
	directIO	?	✓

eulake prototype: concept (2)



by Xavier Espinal

CMS: investigation of data access modes ongoing

Data Lake Prototype in use...

- First, **integrate** it with the Experiment's Distributed Data Management and Workload Management Systems
 - ATLAS
 - ✓ DLP exposed as a storage endpoint to ATLAS DDM (Rucio)
 - ✓ Data can be transferred from any ATLAS site into the DLP end.
 - ✓ Integrated with ATLAS WMS (PanDA)
 - CMS
 - ✓ DPL exposed as a storage endpoint to CMS DDM
 - ✓ Integrated with CMS CRAB3

Data Lake and HammerCloud

- ✓ We integrated the Data Lake Prototype with HammerCloud
 - We can test real workflows and data access patterns of ATLAS and CMS

Initial focus on ATLAS

(Data is copied from storage to WN)

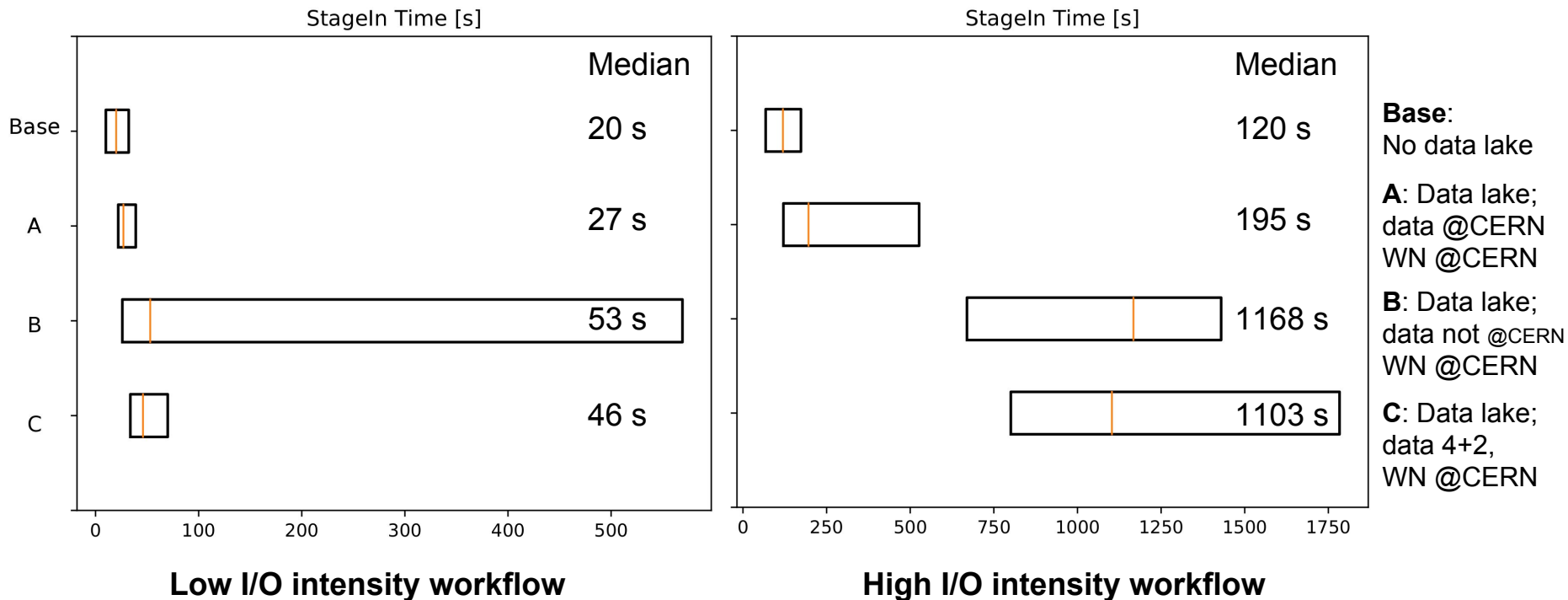
4 test scenarios, stage-in from

1. Base: Local access (no data lake)
2. A: DLP, data @CERN, WN @CERN
3. B: DLP, data NOT @CERN, WN @CERN
4. C: DLP, 4+2 stripes, WN @CERN

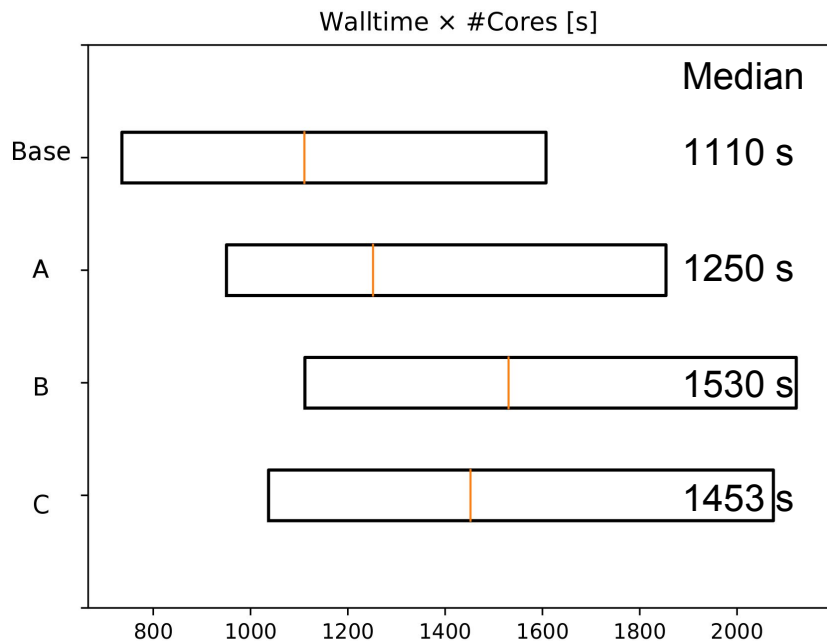
PoC with a CMS “1-chain job” running.

Running Tests backed by the WLCG Data Lake												
State	Id	Host	Template	Start (Europe/Zurich)	End (Europe/Zurich)	Sites	subm jobs	run jobs	comp jobs	fail jobs	tot jobs	
running	20126028	hammercloud-ai-12	1005: PF.T. mc16 Sim_tf 21.0.16 - WLCG Data Lakes - local data clone.989 EULAKE folder CERN	13/Sep, 11:42	14/Sep, 11:05	CERN-PROD_DATALAKES, CERN-PROD_DATALAKES_MCORE, CERN-PROD_DATALAKES_TESTA, 3 more...	2	3	84	16	15	107
running	20126030	hammercloud-ai-12	1006: benchmark derivation AthDerivation/21.2.8.0 1k events - WLCG Data Lakes - local data clone.977 EULAKE folder CERN	13/Sep, 12:08	14/Sep, 12:11	CERN-PROD_DATALAKES, CERN-PROD_DATALAKES_MCORE, CERN-PROD_DATALAKES_TESTA, 3 more...	1	4	43	6	11	55
running	20126032	hammercloud-ai-12	1012: A.F.T. AtlasDerivation 20.7.6.4 clone.808 clone.845 EULAKE folder CERN	13/Sep, 12:36	14/Sep, 13:51	ANALY_CERN-PROD_DATALAKES, ANALY_CERN-PROD_DATALAKES_TESTA, 2 more...	5	0	0	0	0	5
running	20126035	hammercloud-ai-12	1007: benchmark digi+reco derivation Athena/21.0.53 5 events - WLCG Data Lakes - local data clone.987 EULAKE folder CERN	13/Sep, 14:30	14/Sep, 13:11	CERN-PROD_DATALAKES, CERN-PROD_DATALAKES_MCORE, CERN-PROD_DATALAKES_TESTA, 3 more...	1	4	23	15	34	44
Running Tests backed by the standard storages, copy-to-scratch												
State	Id	Host	Template	Start (Europe/Zurich)	End (Europe/Zurich)	Sites	subm jobs	run jobs	comp jobs	fail jobs	tot jobs	
running	20126021	hammercloud-ai-73	845: AFT AtlasDerivation 20.7.6.4 clone.808	12/Sep, 20:42	13/Sep, 21:19	ANALY_ARNES, ANALY_ARNES_DIRECT, ANALY_AUSTRALIA, 142 more...	263	231	11967	1848	13	14338
running	20126036	hammercloud-ai-12	977: benchmark derivation AthDerivation/21.2.8.0 1k events - WLCG Data Lakes - local data	13/Sep, 14:46	14/Sep, 13:32	NIKHEF-ELPROD, SARA-MATRIX, BNL_PROD, 5 more...	2	7	36	0	0	45
running	20126040	hammercloud-ai-12	989: PF.T. mc16 Sim_tf 21.0.16 - WLCG Data Lakes - local data	13/Sep, 15:40	14/Sep, 14:57	NIKHEF-ELPROD, SARA-MATRIX, BNL_PROD, 5 more...	3	4	32	1	3	40
running	20126046	hammercloud-ai-12	987: benchmark digi+reco derivation Athena/21.0.53 5 events - WLCG Data Lakes - local data	13/Sep, 19:12	14/Sep, 18:10	NIKHEF-ELPROD, SARA-MATRIX, BNL_PROD, 5 more...	1	4	9	2	13	16
running	65532	hammercloud-ai-34	195: functional T3_CH_CERN_DOMA	12/Sep, 10:16	14/Sep, 8:15	CERN Tier-0 T3_CH_CERN_DOMA	24	3	415	0	0	442

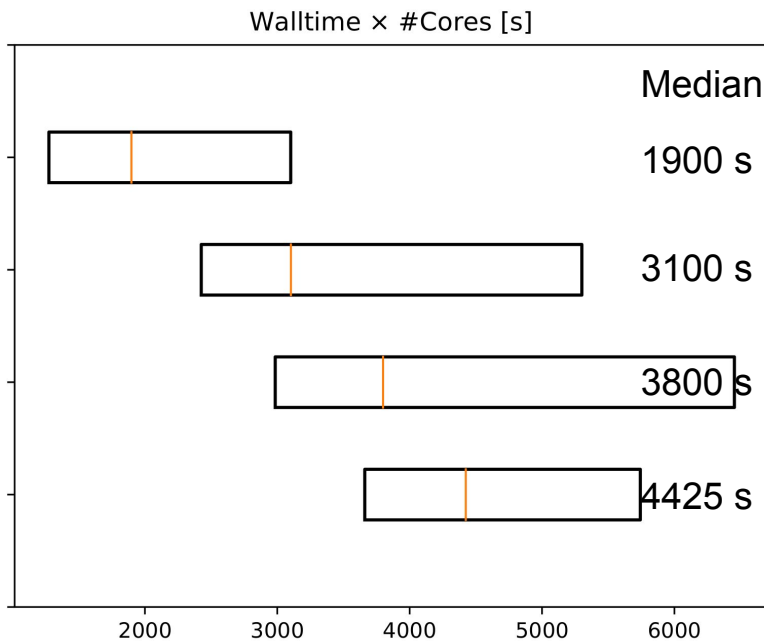
Data Lake, Stage-in Time



Data Lake, WallTime x cores



Low I/O intensity workflow



High I/O intensity workflow

Base:
No data lake

A: Data lake;
data @CERN
WN @CERN

B: Data lake;
data not @CERN
WN @CERN

C: Data lake;
data 4+2,
WN @CERN

WLCG DOMA Activities

- **Content Delivery and Caching**
 - data access performance, content delivery and caching
- **Third Party Copy**
 - investigate, commission & deploy alternative TPC protocols to gridFTP; prototype token-based auth in TPC
- **QoS**
 - at the storage level: define, implement & expose different classes based on performance/reliability need and affordability; integrate the notion of the storage classes up
- **DOMA and Related Network activities**
 - network R&Ds; focus on data transfer: DTNs, low level transfer protocols, bandwidth on demand, P2P channels, SDNs, ...
- **DOMA and AAI**
 - prototyping an architecture; x509 free, based on Jason Web Tokens
- N.B.: HEP Community White Paper Roadmap [arXiv:1712.06982](https://arxiv.org/abs/1712.06982)

Performance metrics and measurements in the Data Lake mode



- Trying to understand if distributed storage saves cost
 - With any distributed storage, we can study, measure, and benchmark
 - jobs and distributed storage performance
 - with different workflows
 - w.r.t. different data access modes
- ⇒ *Can we hide latency and average out bandwidth so that the data location becomes irrelevant?*

Simone Campana, Xavier Espinal Currul, Maria Girone,
Ivan Kadochnikov, Gavin McCance, Jaroslava Schovancová

