# Data management
## Long term planning

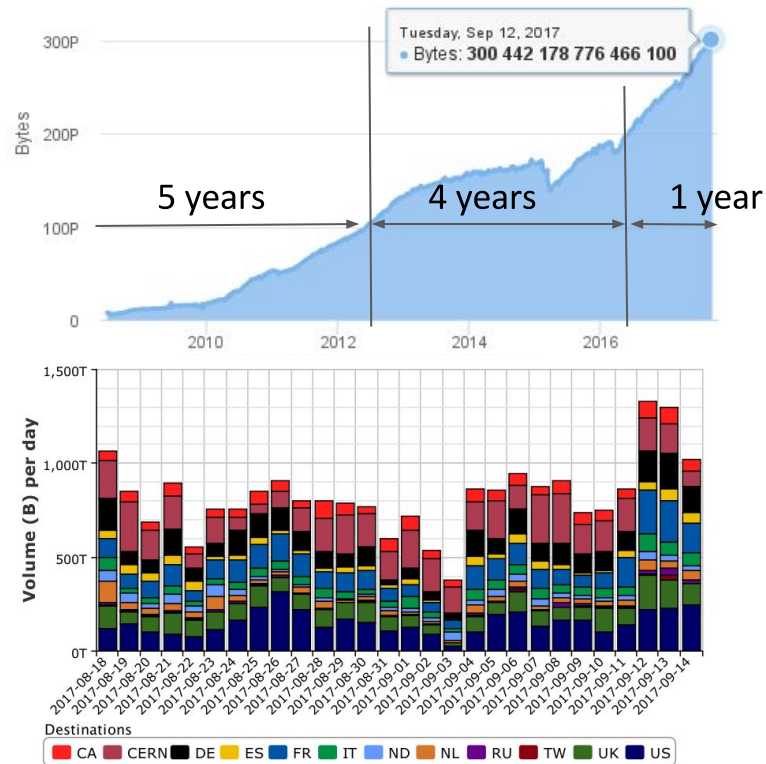Mario.Lassnig@cern.ch
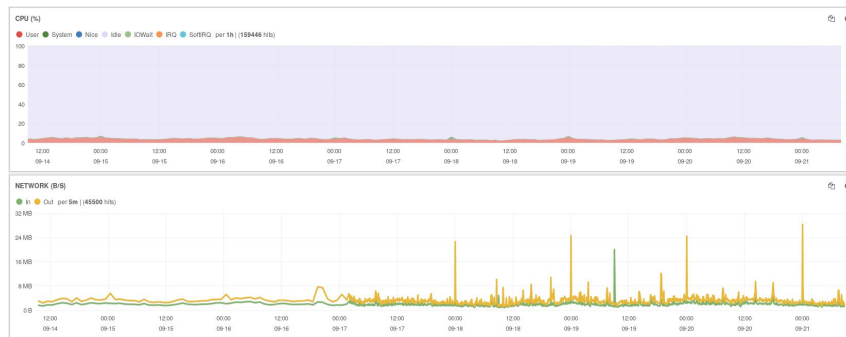Cedric.Serfon@cern.ch
and the data management crew

# Current situation

↳ We recently breached 300 Petabytes!

   ↳ 850M files, 200M containers, 130M datasets

↳ Many Rucio developments

   ↳ Release 1.13.0 ("Donkerine") is out

   ↳ 2 epics, 19 improvements, 27 features, 23 bug fixes

↳ Smooth operations

   ↳ Consolidation of resources / getting rid of small RSEs

   ↳ Validation of new monitoring

↳ Many new projects, some already in production

   ↳ XCache

   ↳ ECHO

   ↳ Cloud storage federations

   ↳ Network tuning and alarming

# Run-3 readiness and scalability 1/3

↳ Clean codebase (155'399 LOC)

  ↳ Patches/Features are tracked via JIRA & git

  ↳ PEP8 style-guide conformant, Flake code checker rates the code at ~8.0/10

  ↳ Thoroughly tested (>400 unit tests)

↳ Rucio is designed to be horizontally scalable

  ↳ Stateless and streaming requests handlers

  ↳ Elastic work sharding (can dynamically add/remove nodes or service instances)

↳ Frontend

  ↳ CPU utilisation ~5%

  ↳ Average network rate 2MB/sec

  ↳ Peaks to 25 MB/sec

  ↳ Stable memory behaviour
    (6GB/node, out of which 2-4GB are buffers)

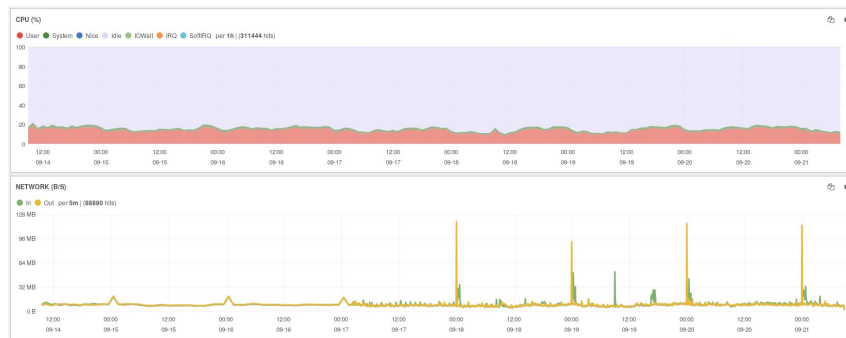# Run-3 readiness and scalability 2/3

↪ Clean codebase (155'399 LOC)

  ↪ Patches/Features are tracked via JIRA & git

  ↪ PEP8 style-guide conformant, Flake code checker rates the code at ~8.0/10

  ↪ Thoroughly tested (>400 unit tests)

↪ Rucio is designed to be horizontally scalable

  ↪ Stateless and streaming requests handlers

  ↪ Elastic work sharding (can dynamically add/remove nodes or service instances)

↪ Services/Daemons

  ↪ CPU utilisation ~20%

  ↪ Average network rate 11MB/sec

  ↪ Peaks to 120 MB/sec

  ↪ Zigzag memory behaviour
     (6GB/node, out of which 3GB are buffers)
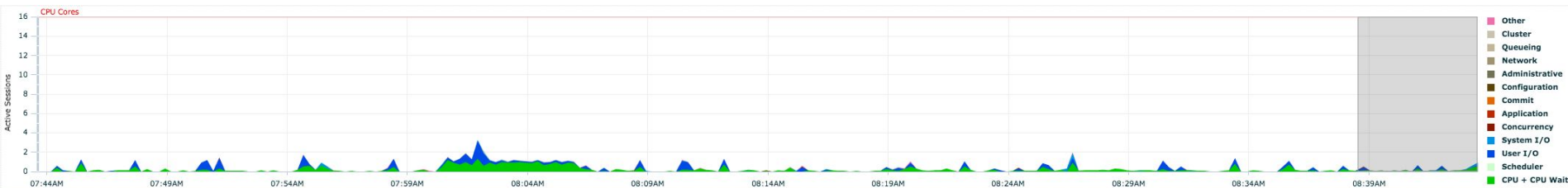
# Run-3 readiness and scalability 3/3

↪ Database

  ↪ Very low utilisation (compared to Oracle capabilities!), but very high usage

  ↪ Extremely optimised, in constant talks with DBAs

  ↪ No major problems for a very long time

  ↪ cf. Database Ops Slides

↪ However, lately we ran into some session problems (nr. of sessions exceeded)

  ↪ We artificially limit to 1000 reader and 1000 writer sessions to save Oracle resources

  ↪ Apache process handling seems to not terminate database sessions, even though process is gone

  ↪ Oracle-side process killer of sessions which are longer than Apache timeouts

# Of course, there are exceptions

↪ Rule evaluation service (a.k.a. judge) is limited by single node memory
  ↪ Rule evaluations are done atomically, thus one evaluation happens within one database transaction
  ↪ Daemon needs to load all files, replicas, etc. into memory
  ↪ **Largest single successful evaluation yet on 8GB node: 500'000 files**
  ↪ Instantaneous solution: high-memory nodes (but there are limits in what we can get)
  ↪ Long-term solution for infinite size: Partition the evaluation across multiple transactions

↪ Rucio is currently "manually" elastic
  ↪ Operations has to decide when to add new nodes or service instances
  ↪ Fully automate this process based on node health and service metrics

# Major developments in the next months

↪ Rucio development

   ↪ Transparent archive/ZIP support

   ↪ Conveyor scheme/protocol improvements     — protocols, clouds, and opportunistic resources

   ↪ ~~GlobusOnline transfertool~~     ~~— required for managed transfers to/from HPCs~~

↪ Data management operations

   ↪ Enable even more sites to use rucio-mover and different upload/download protocols

   ↪ Evaluate ROOT and WebDAV als third-party copy protocols

   ↪ Ongoing monitoring validation

   ↪ Small sites decommissioning / consolidation

   ↪ New deployment model based on containers/dockers under evaluation

↪ Storage & network technologies

   ↪ Deployment of XrootD 4.7     — most importantly VOMS authz and client-side caching

   ↪ Xcache phase-1 deployment     — SLAC, MWT2

   ↪ Network alarming & alerting     — low throughput, high packet loss, …

# What about long term? Our internal document

↪ Trace & document the progress of the different data management projects

↪ Long term planning tool (current milestones: end of Run-2, end of LS-2)

↪ 3 work packages (matching with ATLAS OTP)

   ↪ Rucio development

   ↪ Data management operations

   ↪ Network & storage technologies

↪ Follows the discussions, workshops, documents from relevant coordinators

   ↪ ATLAS Software & Computing Management

   ↪ HEP Software Foundation

   ↪ WLCG Data Steering Group

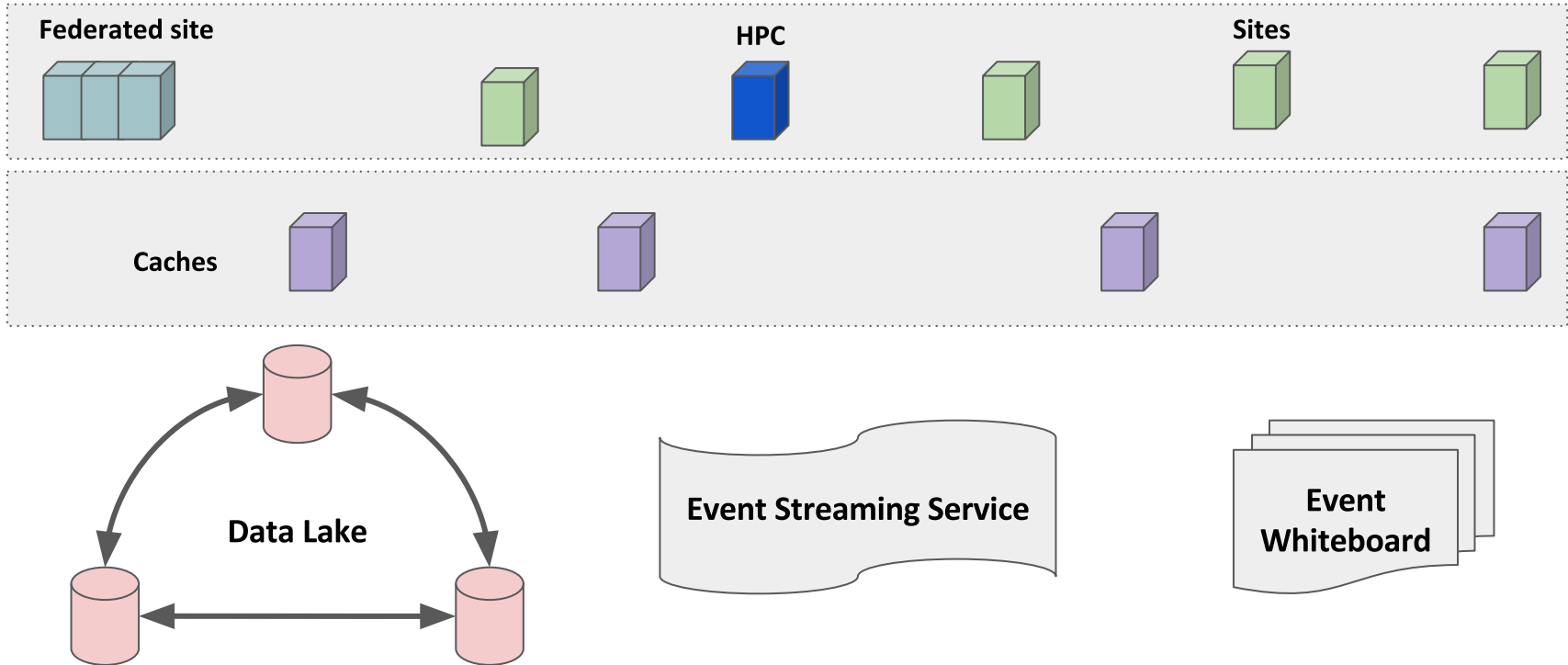   ↪ Non-ATLAS experiment feedback

# Collected recommendations

↪ WLCG Data Steering

  ↪ Working on straw-man proposals, nothing concrete yet

  ↪ Hot topics: authz without x509, increasing usage of TAPE, regional federations/caches

↪ SW&C Management

  ↪ Monday plenary

  ↪ ESS (event level processing), CDN (efficient caching and placement), DKB/Whiteboard (metadata)

  ↪ Storage Performance Improvement Team (SPIT?) — also build on existing work of our CompSci PhDs

  ↪ Reassess the use of object stores — more closely involve objectstore experts

↪ HSF CWPs which touched data management

  ↪ Computing Models        — equivalent to SW&C stance

  ↪ Data Access             — slightly more detailed on cataloguing, CDNs, caching, and IO patterns

  ↪ Machine Learning        — concerned about I/O throughput, file-types (serial/random IO), filesystems
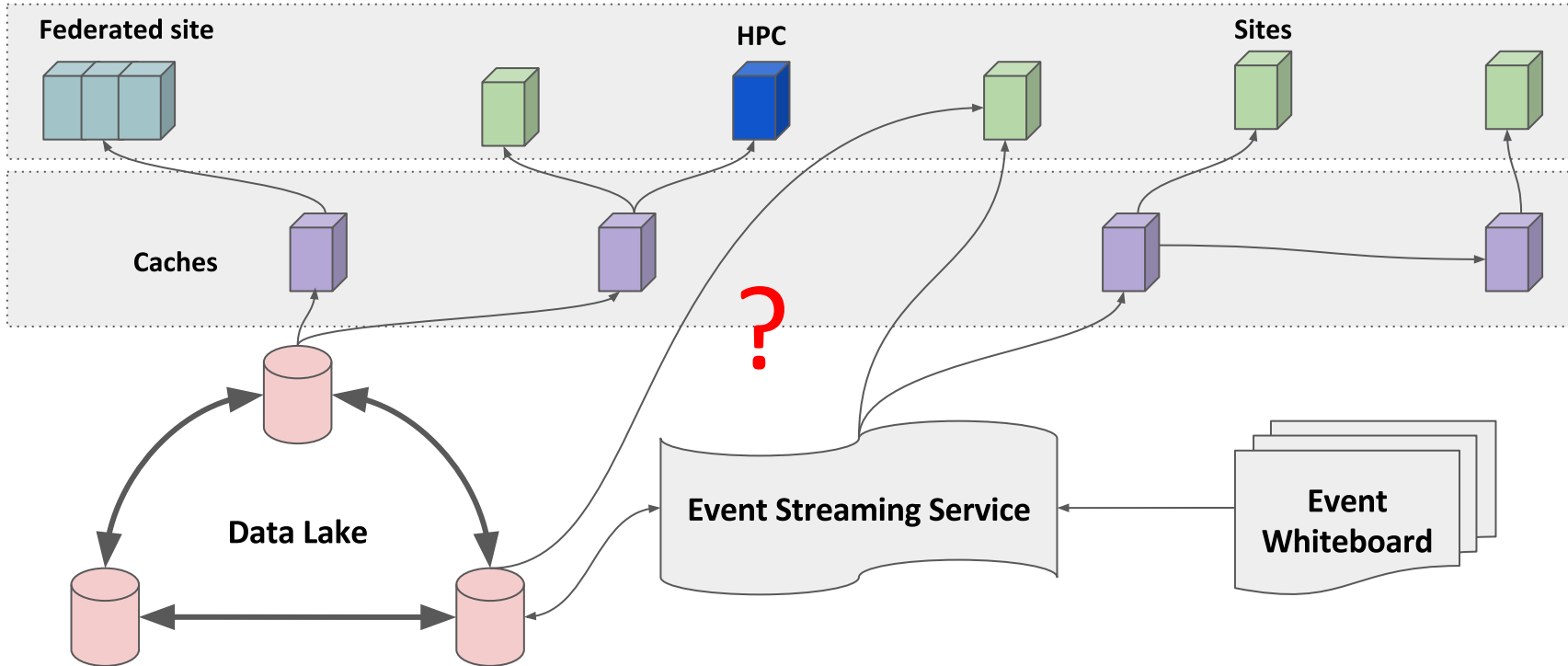
# Some projects under discussion

↪   Grab bag of ideas to automate our data management

↪   Service to automate archiving of small files to TAPE

↪   Self-healing rules / automatic handling of suspicious files

↪   Xcache Phase-2: Add more sites and test client-level caching

↪   Smarter lifetime model with finer granularity

↪   Testing of SDNs with network function virtualisation

  ↪   The R&E networks will become heavily shared once more WLCG-scale experiments come online

  ↪   Ensure that we have proper network shaping and control in place

↪   Usage of TAPE without SRM

↪   Downloads via Rucio WebUI

↪   … and many more

# Streaming Content Delivery? Is it the new model?

Federated site

HPC

Sites

Caches

Data Lake

Event Streaming Service

Event Whiteboard

# Streaming Content Delivery? Is it the new model?

# Streaming Content Delivery? Is it the new model?

↪ At first glance it looks like the tiered MONARC model all over again.
    What are the practical implications?

↪ Objective seems to be to trade CPU for storage

   ↪ Secondarise data on all sites for quick eviction, keep primaries only on data-lake
   ↪ Recreate data products at fine granularity when necessary
   ↪ Exploit small but fast caches everywhere for these fine granularity data products

↪ For many features mentioned in the CWP we have solutions either
    already in place or in some stage of development

   ↪ e.g., rule-based data distribution, Xcache, network monitoring, asynchronous prefetch, …

↪ Will have to come up with a very detailed plan for each step

   ↪ We risk low utilisation/throughput at the beginning, is it acceptable?

# Data management beyond ATLAS

↪ AMS and Xenon1t are using Rucio in production

↪ COMPASS and LSST are evaluating it

↪ More experiments expressed interest, also from non-HEP domains


↪ Engage interested experiments —> **Rucio Community Workshop** early 2018!

↪ Ensure community-friendly Rucio development and deployment

   ↪ move to Github, out-of-the-box deployment, out-of-the-box development, revise documentation

↪ Establish Rucio as complete, though modular, data management solution

   ↪ e.g., experiments who would like to use Xcache, Object Stores, Event Streaming Service, ~~GlobusOnline~~

   ↪ and ensure smooth cooperation with our WMFS system PanDA

# Licence and copyright

↪ Rucio redistribution licence is Apache Licence 2.0

  ↪ Open source and free software (OSI approved)

  ↪ Compatible with GNU General Public Licence 3.0 (but neither 1.0 nor 2.0)

↪ CERN holds Rucio copyright 'for the benefit' of ATLAS to permit the widest possible adoption and reuse (http://legal.web.cern.ch/licensing/software)

↪ Will seek licence clarification where we know/assume tight cooperation  to ensure smooth further development and deployment

  ↪ E.g.,  with external/non-CERN/non-ATLAS institutes

  ↪ HSF Note: http://hepsoftwarefoundation.org/notes/HSF-TN-2016-01.pdf

  ↪ We do not want to end up in a situation where a critical piece of software is in one way or the other licence-incompatible for a complete data management system

    ↪ E.g., FTS, DynaFed, GFAL, XrootD/Xcache, ~~GlobusOnline~~, …

# Current effort allocation

↳ Rucio development

  ↳ 3.2 FTE across 6 persons

  ↳ Only 3 full-time (>66% time) developers though, reduction by 1 FTE in 2018 foreseen

  ↳ Partial help (<= 20% time) currently zero-sum due to overhead for existing developers

    → But very useful to attract new people who then want to stay for longer!

  ↳ No outside-ATLAS Rucio contributions yet, but actively pursuing

↳ Data management operations

  ↳ 1.8 FTE across 2 persons

  ↳ DDM/WFMS review recommended at least 2 FTE, actively trying to involve new people

↳ Storage & network technologies

  ↳ 2.2 FTE across 7 persons

  ↳ + unaccounted behind-the-scenes efforts from non-ATLAS persons

# Summary

↪ Data management is in a very good shape

  ↪ For now, no major redesigns or rewrites required for Run-3

  ↪ Lots of ideas and potential improvements to make our lives easier

↪ Recommendations from WLCG, SW&C, HSF CWP are the foundation for long-term

  ↪ Event-level processing

  ↪ Distributed caching

  ↪ Metadata

↪ Paving the way for non-ATLAS deployments and contributions

↪ Effort allocation needs to be improved & dedicate larger fractions of time

↪ (and of course finish writing that journal article)