



# Running Rucio: A Facilities Perspective

John S. De Stefano Jr, of behalf of  
The BNL Belle II DDM Team, and  
Rucio Service Maintainers Worldwide

01 Mar 2022  
Software & Computing Round Table



# Contributors

# BNL Belle II DDM Team

- John S. De Stefano Jr (SDCC)
- Hironori Ito (SDCC)
- Paul Laycock (NPPS)
- Ruslan Mashinistov (NPPS)
- Cédric Serfon (NPPS)

# Rucio Instance Information Contributors

- ATLAS: CERN – Martin Barisits
- Belle II: BNL (\* *also see: Cédric's earlier talk*)
- CMS: CERN – Eric Vaandering
- DUNE, ICARUS: FNAL – Brandon White
  - Also: Rubin: SLAC
- EGI, GRIDPP, SWIFT-HEP: RAL – Tim Noble, Wenlong Yuan
- SKA: Rosie Bolton (\* *also see: Rosie's earlier talk*)
- sPHENIX: BNL – Matt Snyder

# Rucio Instances

# Rucio for ATLAS: CERN

(Almost) everything running in Kubernetes

- Exceptions (4 CPU, 8GB memory VMs):
  - Load balancers: 3 HAProxy VMs
  - Authentication: 3 VMs
- Server & daemon Kubernetes clusters:
  - 3 production clusters: x16 VMs each (4-Core, 8Gb)
  - 1 integration cluster: 6 VMs
- Database:
  - Oracle 19c cluster
  - 8TB in table & index space

# Rucio for CMS: CERN

## Kubernetes via OpenStack

- 2 k8s clusters:
  - 4 development machines
  - 8 production machines
    - Includes all Rucio functionality, except monitoring
- All CC7-based
- All accessible via Ixplus8
- Database: Oracle CMSR (offline)

# Rucio for DUNE, ICARUS: FNAL

## Kubernetes via OKD

- Cluster provided by FNAL Computing
- Containers:
  - Rucio 1.26.9 (LTS)
  - FNAL, OKD, policy customizations
- 2 production instances for replica management with different, custom FNAL metadata solutions:
  - DUNE: Metacat
  - ICARUS: Sequential Access via Metadata
- Database: FNAL Computing cluster
- Deployment framework
  - <https://github.com/bjwhite-fnal/rucio-fnal>

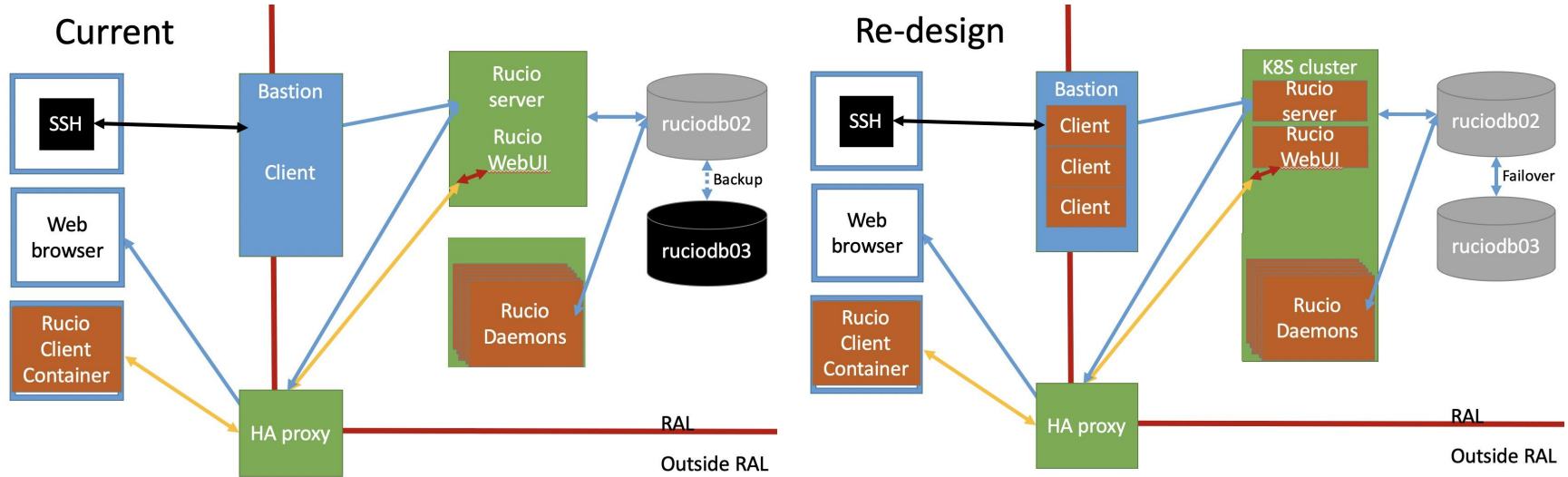


# Rucio for EGI, GRIDPP, SWIFT-HEP: RAL

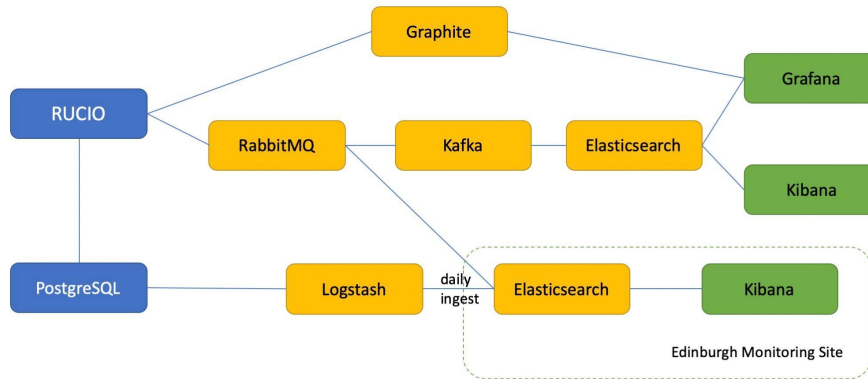
## Current deployment

- Amazon Cloud VMs
  - Servers:
    - Rucio 1.23.17
    - c1.xlarge VM (4 CPU, 16GB RAM)
  - Daemons:
    - Rucio 1.27.3
    - m2.medium VM (4 CPU, 4GB RAM)
- Future: Kubernetes deployment

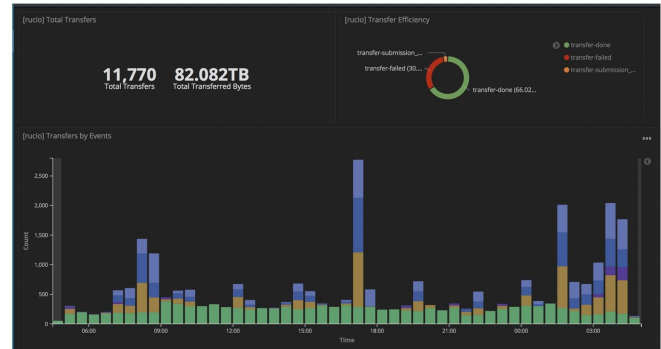
# Rucio for EGI, GRIDPP, SWIFT-HEP: RAL (2)



# Rucio for EGI, GRIDPP, SWIFT-HEP: RAL (3)



Shows **queued**, **failed**, **submitted**, **done** in real time



# Rucio for SKA: RAL

## Kubernetes deployment

- Maintained by SKA
- Deployed in UK STFC cloud
- CERN pilot FTS
- Production auth via x509
  - In development: tokens via [ESCAPE's IAM](#) at INFN

# Rucio for sPHENIX: BNL

Currently deployed on VMs in RHEV

- Testing on 1 front-end VM, 1 daemon VM
  - May scale to 2+ VMs each
- Rucio v1.27.4 via Python pip
- Separate PostgreSQL database VM on separate RHEV cluster
  - May move to HW hosting if necessary
- RHEL 8.4 VMs
- New deployment framework:
  - Foreman (v2.3.5)
  - Gitea (v1.14.2)
  - Puppet (v6.24)

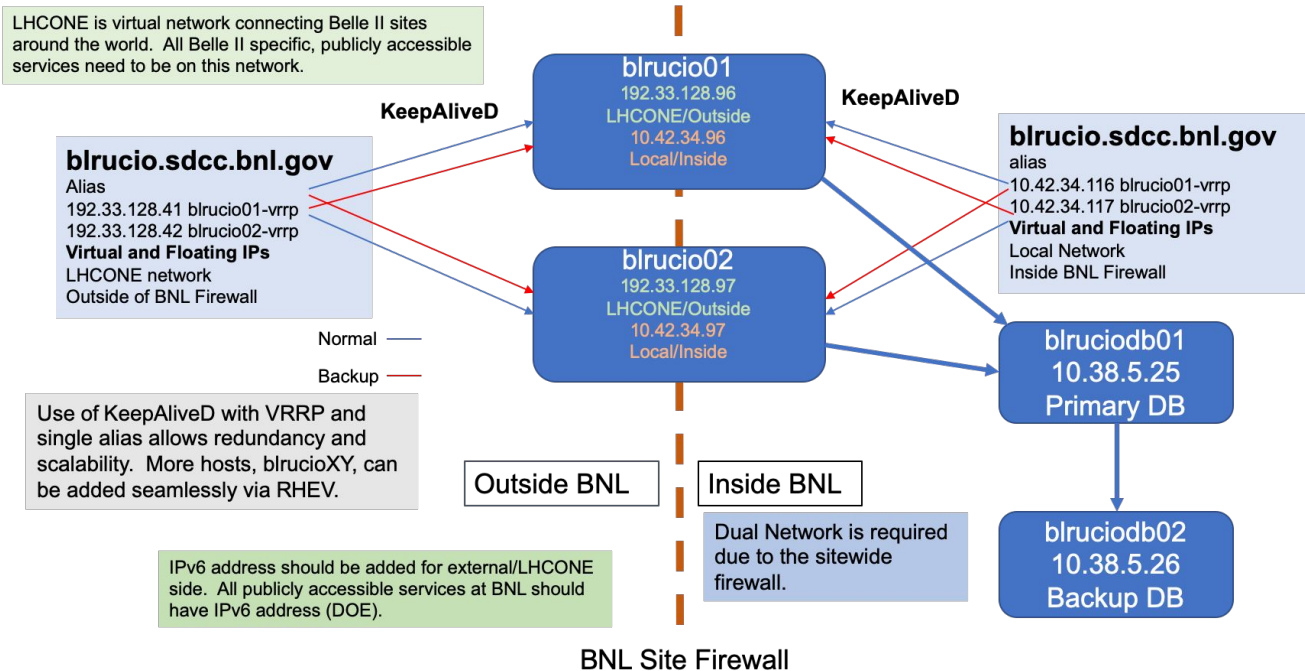
# Rucio for Belle II

# Belle II Rucio Infrastructure & Deployment

Mostly RHEV VMs

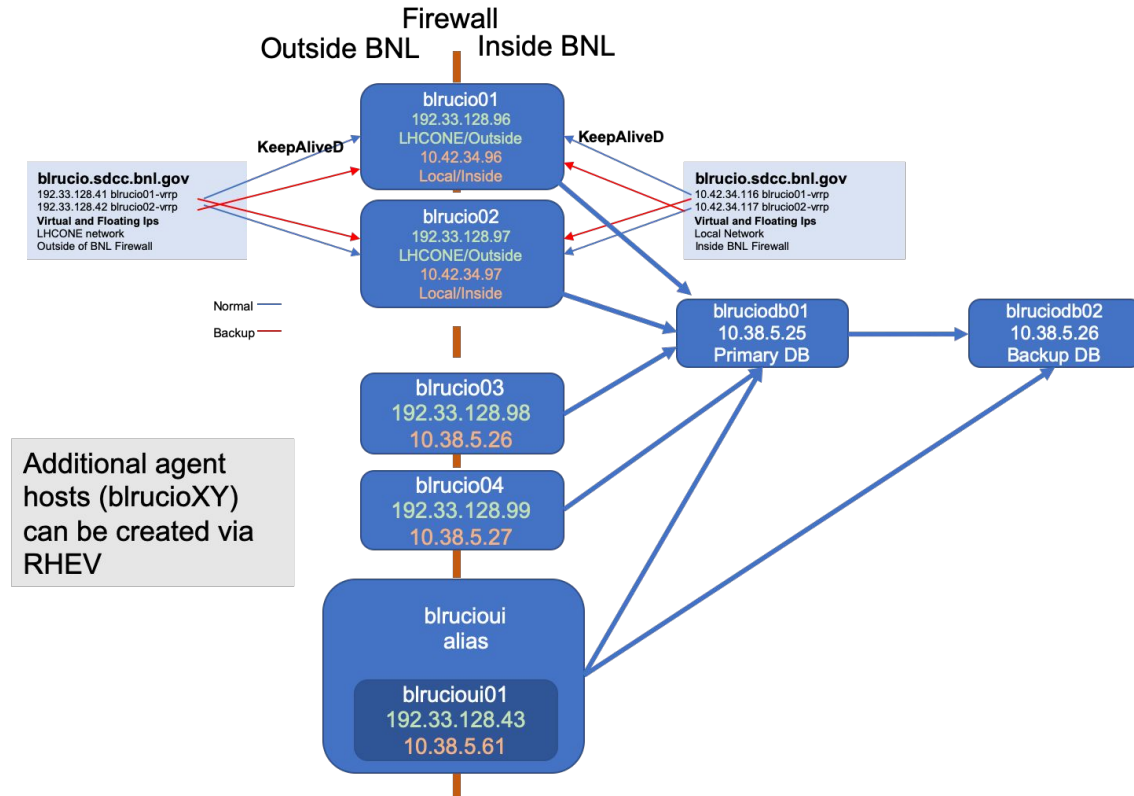
- 2 VM clusters:
  - 5 RHEL 7 production machines
    - 2 servers, 2 daemons, 1 UI
  - 4 integration machines
- Database: PostgreSQL (HW)
- Mostly Puppet via Git code
- Daemons fully integrated with systemd/journald
- Migrated from Python2 to Python3 (pip)

# Belle II Rucio Front-End Service Network





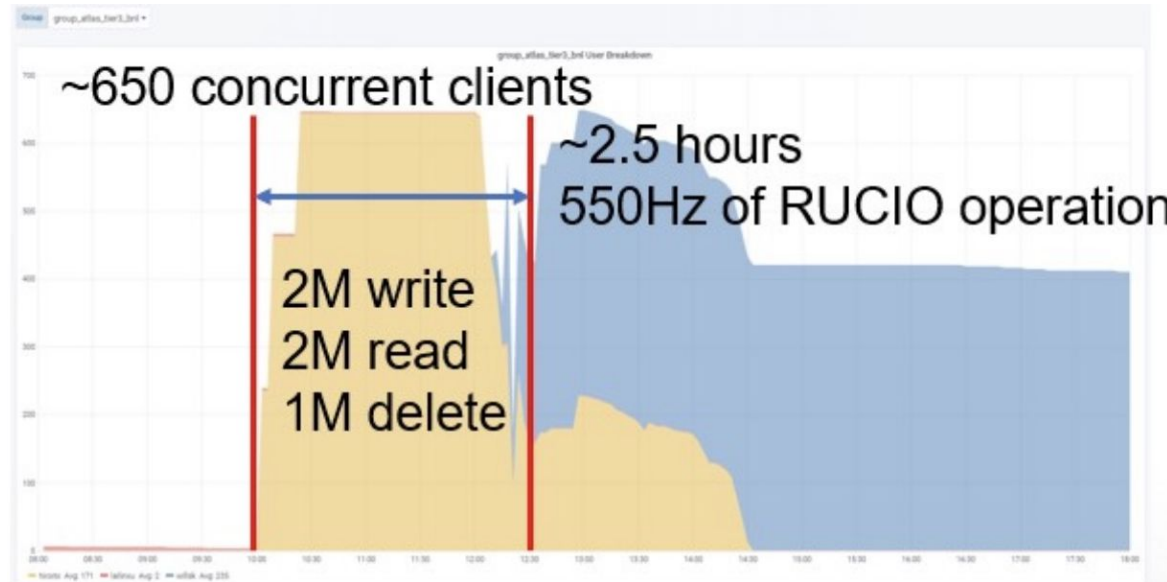
# Belle II Rucio Full Service Network



# Belle II Rucio Infrastructure Tests

Pre-production validation and scaling tests

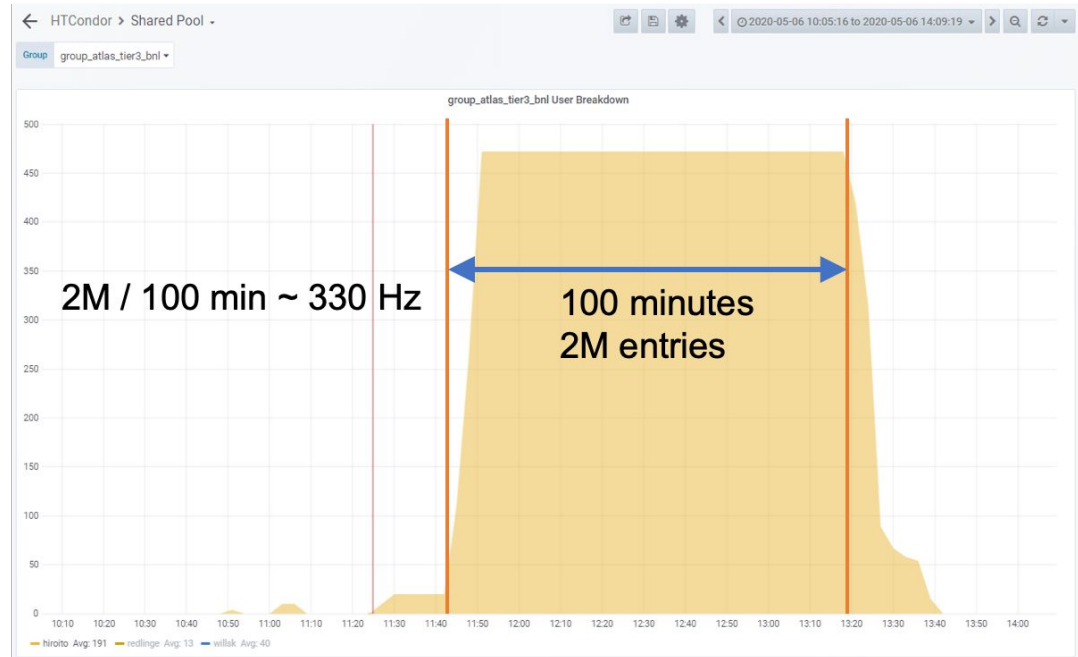
- 



# Belle II Rucio Infrastructure Tests (2)

- Submitting 1K jobs.
- ~ 470 clients total
- 2M entries in 100 min

470 job slots vs 700 slots earlier  
Via opportunistic queue

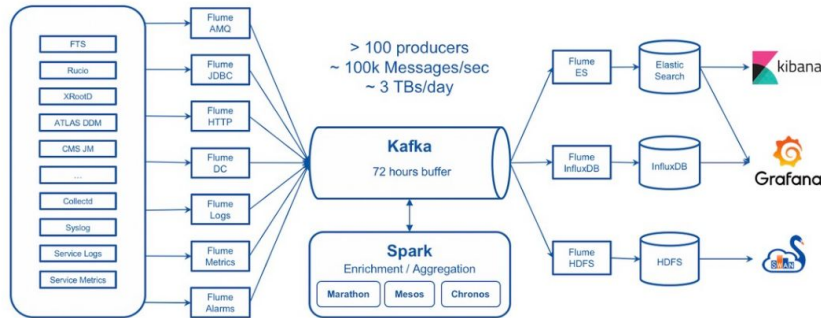


# Belle II Rucio Monitoring

## Consolidation and improvement

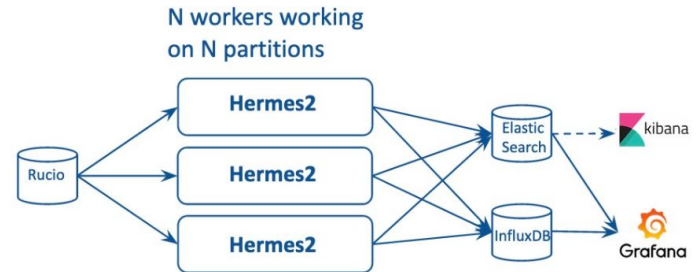
- Another development done was to simplify the monitoring stack :
  - The monitoring stack used by LHC experiments indeed relies on a complex machinery
  - Cost/benefit analysis was conducted to setup the same infrastructure for Belle II. Decision to simplify it
  - The simplification was done by introducing a new component within Rucio

Sources > Transport > (Processing) > Storage > Access



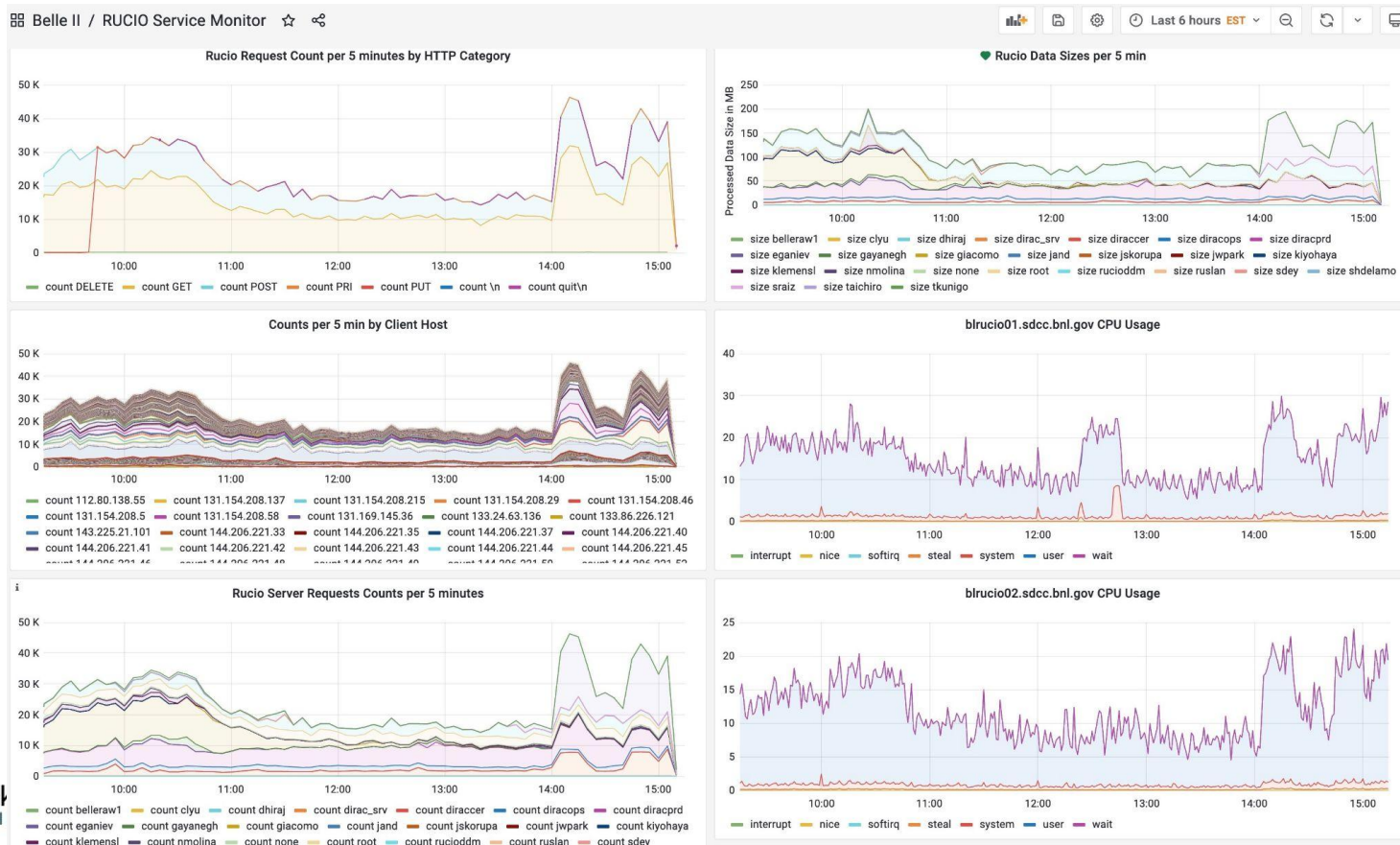
ATLAS monitoring stack

Sources > Transport > (Processing) > Storage > Access



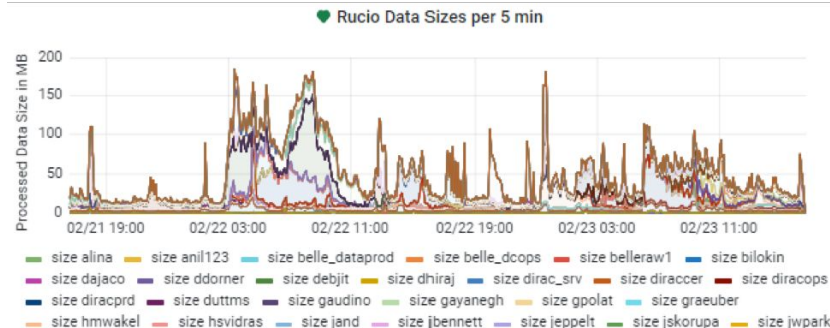
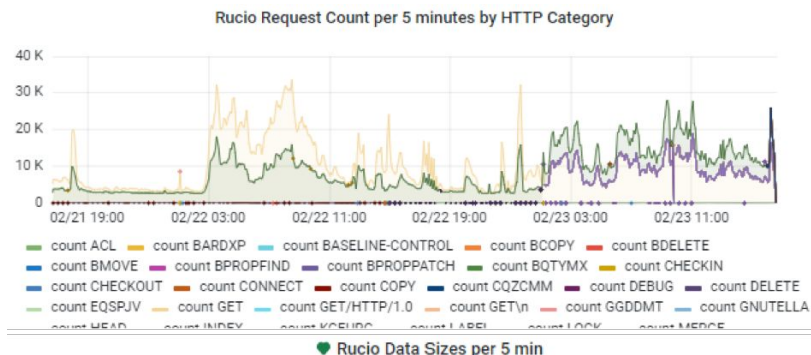
Belle II monitoring stack

# Belle II Rucio Monitoring (2)



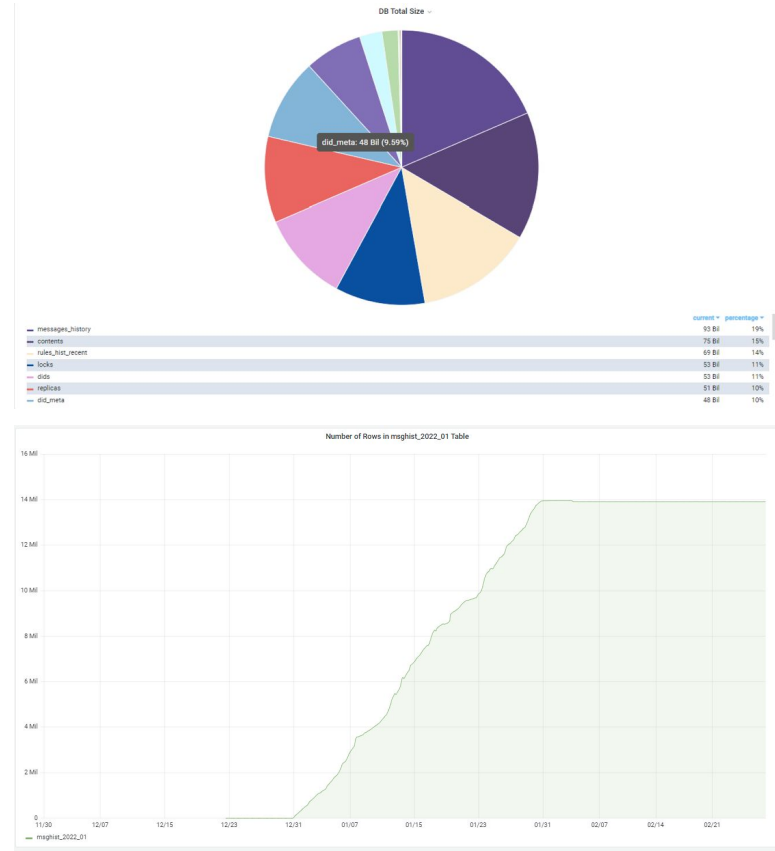
# Belle II Rucio Service Request Monitor

- Log parsing, input to PostgreSQL TimescaleDB for analysis
- Logs all HTTPs requests
- Sophisticated alarm to identify problematic requests



# Belle II Rucio Database

- PostgreSQL
- History table grows rapidly
- Requires **monthly** partitioning



# Belle II Rucio Issues & Challenges

Rucio was originally built for one experiment (ATLAS) and one deployment (CERN)

- Server requires a core config file in a non-default path
  - Expects file in `/opt/rucio/etc/``
  - Installs to `/usr/rucio/etc/`` (and `/usr/local/etc/``, and also `/usr/local/rucio/etc/`` for UI)
- UI had “ATLAS” burned into header (despite config var)
  - Since changed to “Rucio UI”
    - `class="name"` in `web/ui/templates/base.html``
- UI requires direct DB access for read-only data
- Daemons don't all support “supported” options



# Belle II Rucio Issues & Challenges (2)

Things are improving – especially when we communicate and report them!

- “Server certificate expired” client error for valid server certificate ([#5021](#))
  - Related: Rucio client error on server CRL refresh delay/hang
- Server test tool dependency issue ([#5023](#))
- Daemons missing sleep option ([#3987](#))

# Next Steps

- Continue to improve monitoring
- Continue to improve Belle II DDM integration (DIRAC)
- Improve deployment automation
  - Puppet-ize more service components
- Test, implement token-based auth

# Thank you

*jd at bnl dot gov*