



Experience supporting Belle II CDB server Infrastructure for Phase 3

Carlos Fernando Gamboa, Ruslan Mashinistov, Maxim Potekhin
Brookhaven National Laboratory

Benedikt Hegner, Brookhaven National Laboratory
Current affiliation: CERN

Martin Ritter
Ludwig-Maximilians-Universität München

Marko Bracko
University of Maribor and Institut Jozef Stefan

24th International Conference on Computing in High Energy and Nuclear Physics – CHEP 2019

BROOKHAVEN
NATIONAL LABORATORY



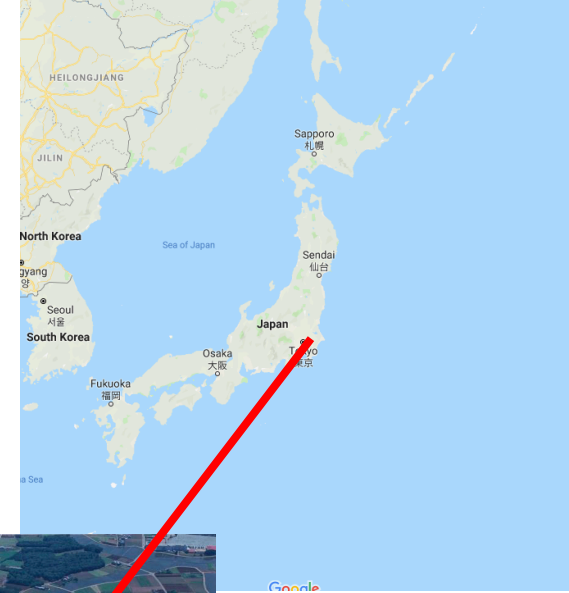
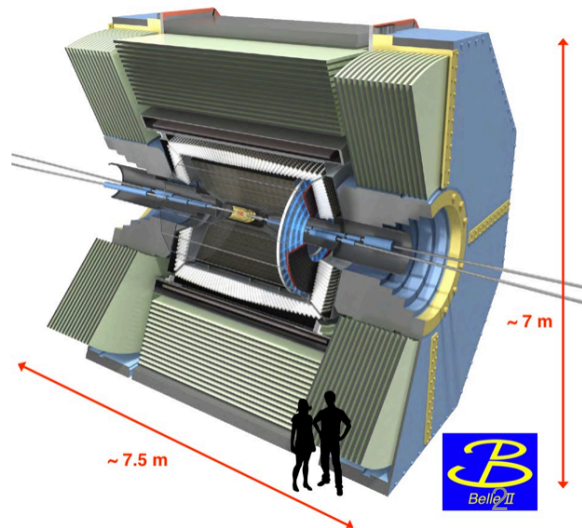
BROOKHAVEN SCIENCE ASSOCIATES

Belle II experiment hosted at KEK Tsukuba, Japan

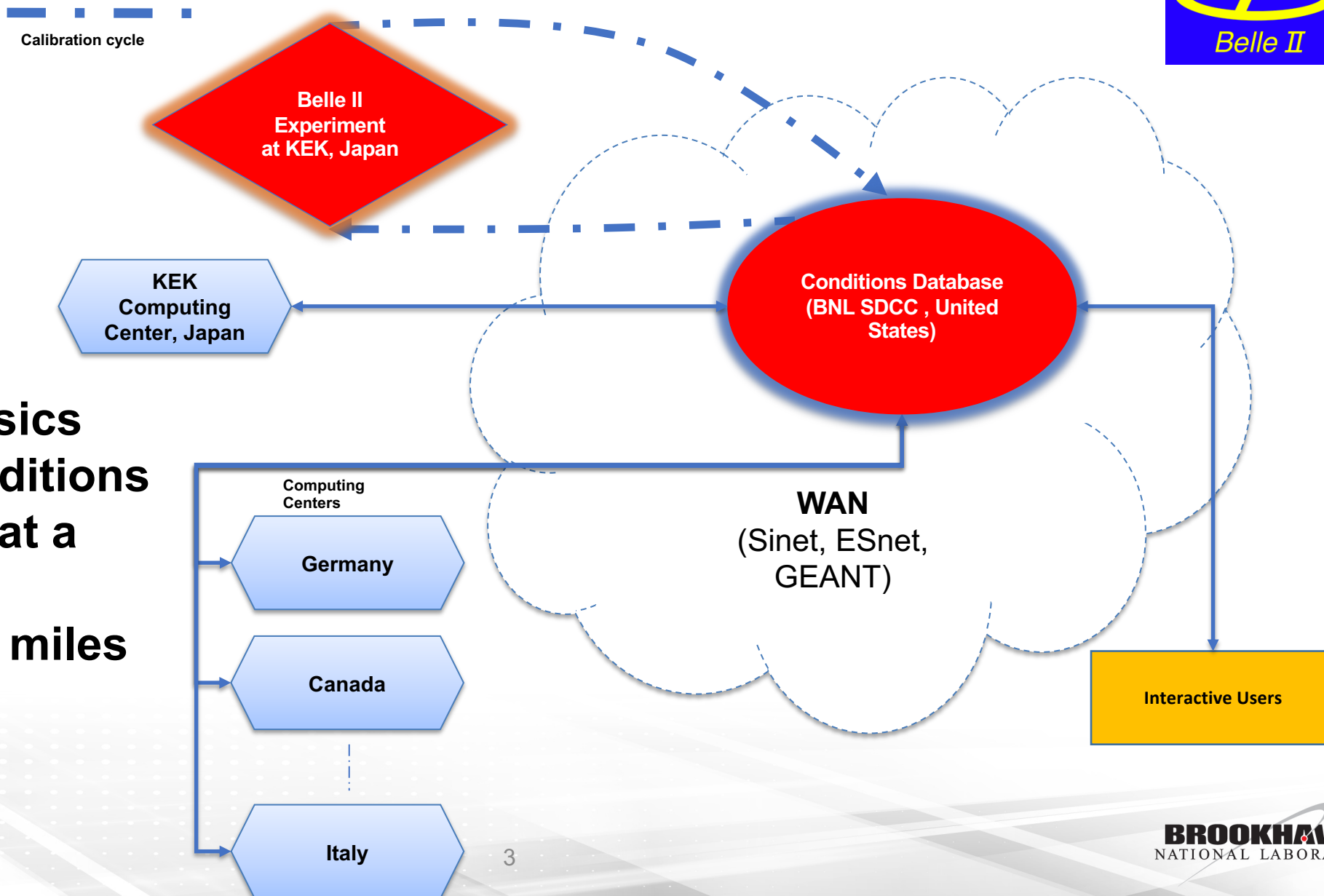
Belle II designed to find physics beyond the Standard Model of particle physics

SuperKEKB accelerator designed to increase the instantaneous luminosity in 40x with respect to KEKB accelerator, and expected to deliver 50x more data to Belle II compared to its predecessor Belle

**“Phase 3” - operation with complete Belle II detector:
First run from March 2019 to June 2019, operation will continue until 2029**



Belle II Conditions Database accessibility overview



First particle physics experiment's conditions database hosted at a remote location, ~10,815 km/6,720 miles apart.

Belle II Conditions Database (CDB)



Events recorded by the Belle II detector are grouped in runs

- A run sets a data taking period with stable operating parameters
- The database manages conditions data on run granularity

Conditions data is composed by metadata and files, stored for persistency in a:

- Shared filesystem → **files or payloads**

Payload binary objects containing information like (calibration, Beam Parameters, etc.)

- Relational Database Management System (RDBMS) → **Metadata**

Allows identify payloads which are accessible in an external service

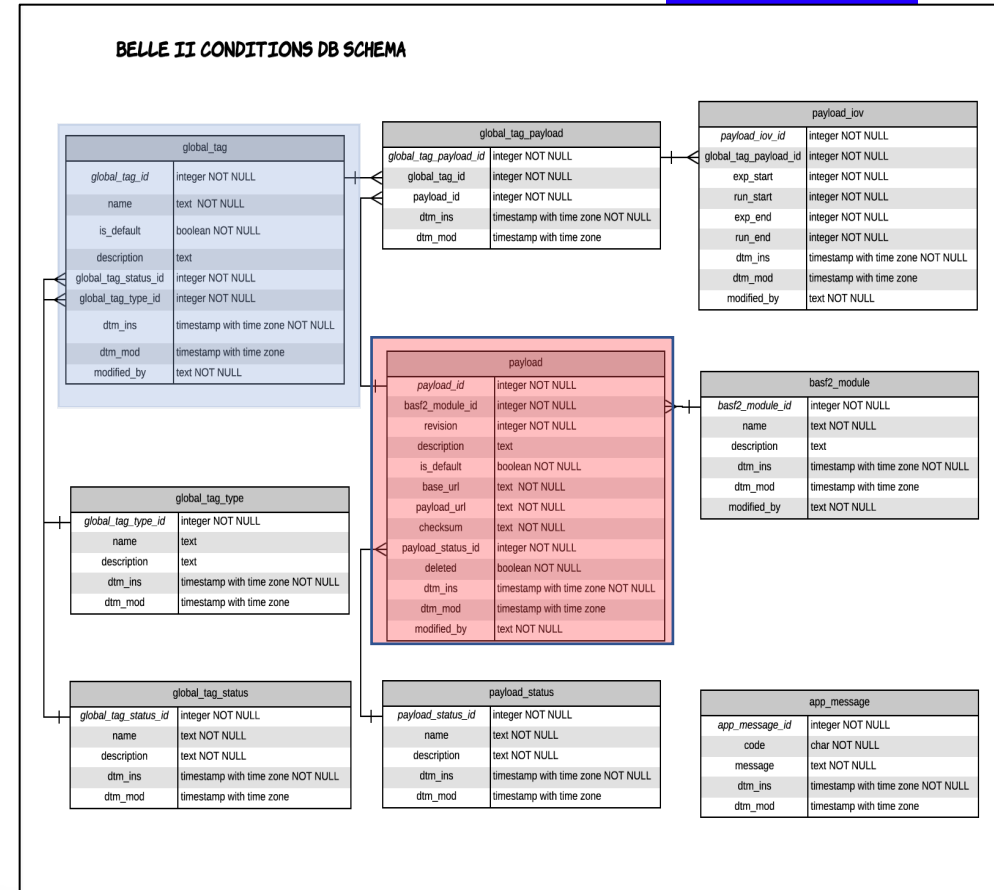
Belle II CDB data model relational database



Payloads contain a sample of conditions data (i.e. BeamParameters) stored in a file:
 File type is agnostic for the CDB server
 ROOT file format when restricted at client side
 CDB server metadata keeps track of the checksum of the file

Intervals of Validity (IOV) specify starting and ending experiments and runs for a given payload for a specific global tag. Can be a fixed run range (closed) or starting at a given run (open)

Global tag (GT) contain list of IOV-payload relationships and are used to select a complete set of conditions for a given reprocessing effort



Belle II CDB general service architecture



1. HAproxy

- Encryption/Decryption, Authorization and Authentication
- Load balancing and resiliency

2. Caching technology

- Allow Global Tag caching and offload dependent services

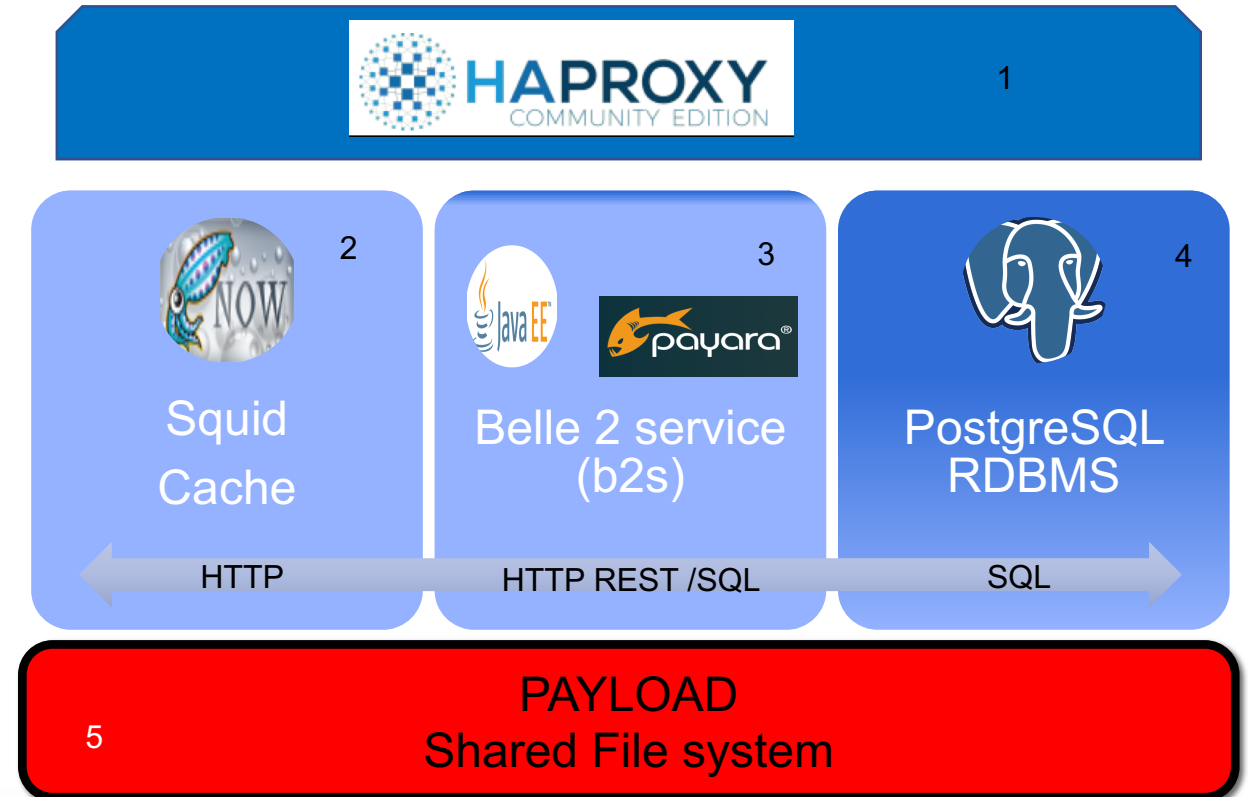
3. HTTP REST API presented by Belle 2 service (b2s) code

- HTTP REST API Swagger enabled

4. Metadata stored in a persistent database repository

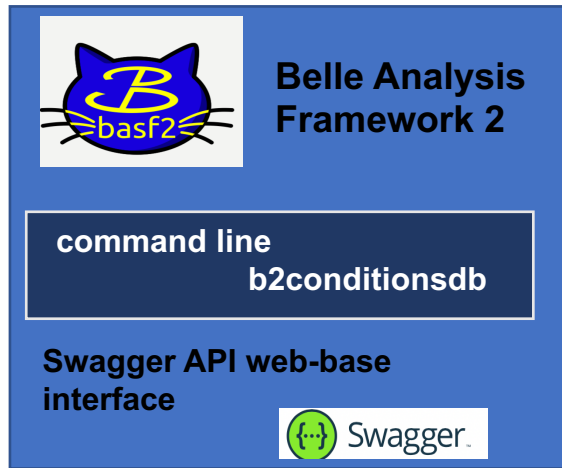
5. Payloads are stored in a shared file system

- General Parallel File System enabled for metadata and payload service (next slide)

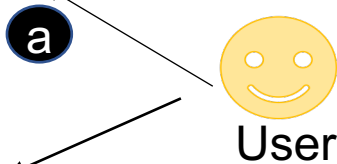
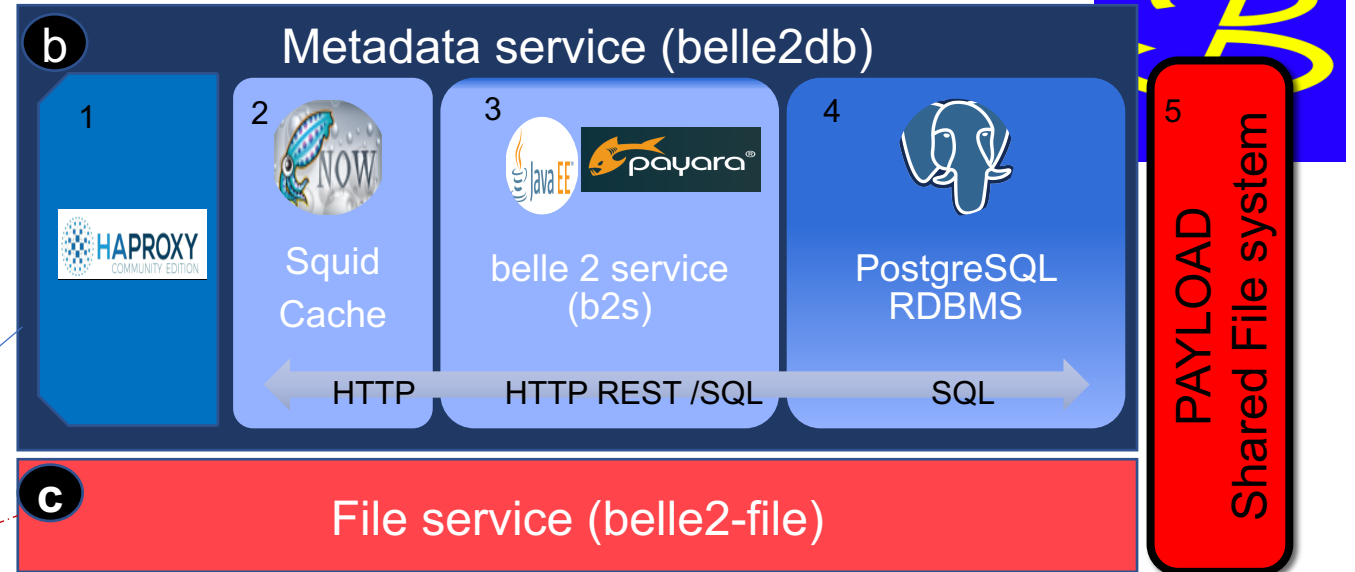


Belle II CDB general service data flow

B2 CLIENTS

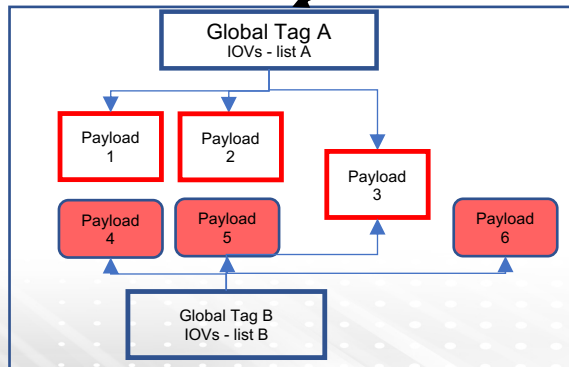


User access
published
Global tags



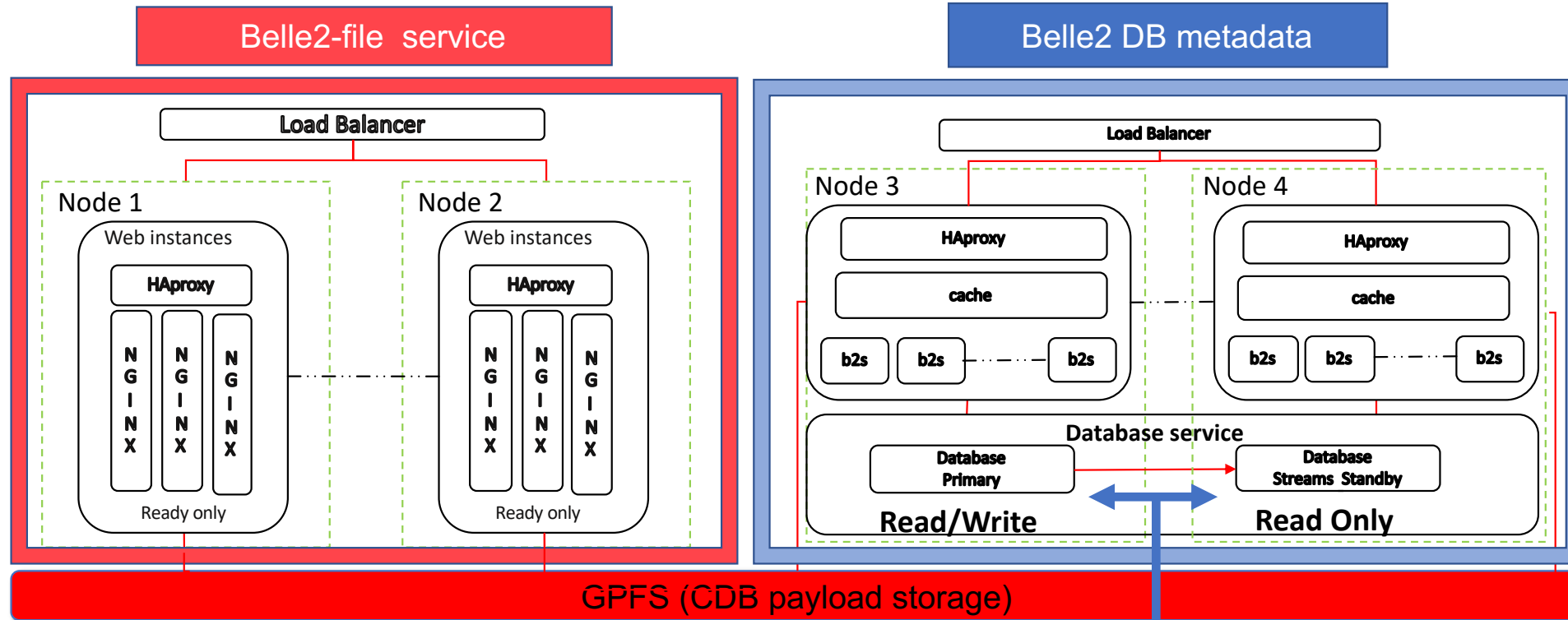
User needs to retrieve Global Tags:

- Uses a client to interact with the CDB metadata and file services
- Metadata service provides information about the payloads and related IOVs of Global Tags:
 - Request is Authenticated/Authorized and/or encrypted/decrypted
 - The cache provide the request if possible
 - The request is translated from HTTP to SQL (or vice versa) by the **b2s** service
 - Database provide the metadata associated to the request
- User retrieve payloads identified by the metadata on **b** via Belle2-files service



Belle II CDB server architecture (under the hood)

Orchestration of individual containerized components



Streams standby replication technology enabled

- Database service open in READ only mode, mainly serving the **CDB web monitor** application
- Technology improves lag synchronization compared to the initial deployment
 - Initially hot standby was deployed for database disaster and recovery

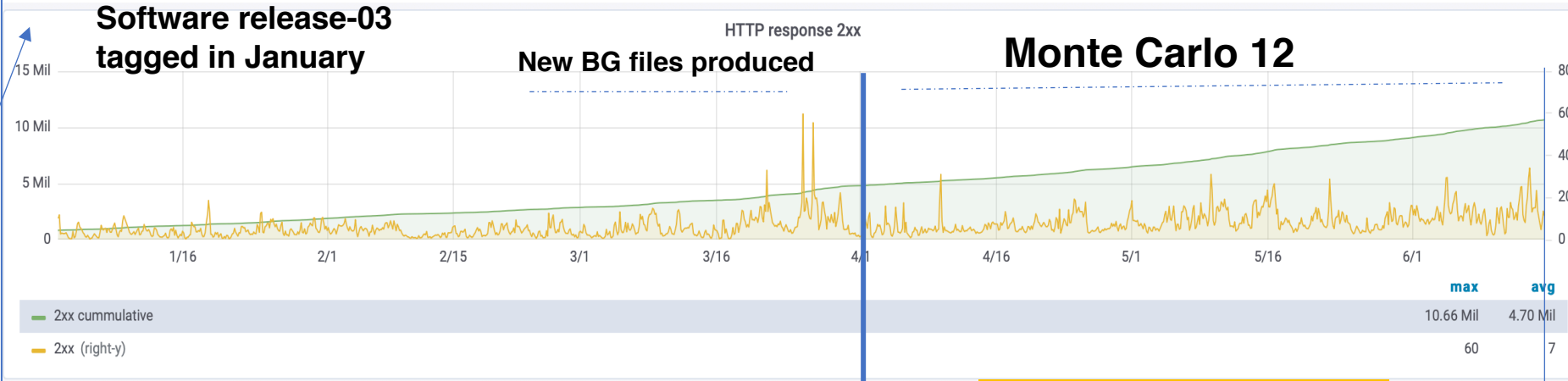
Performance Belle II Conditions service CDB January-June 2019

all service response in terms of HTTP response code



4.7M of HTTP 2XX succeeded

Belle II Conditions metadata service

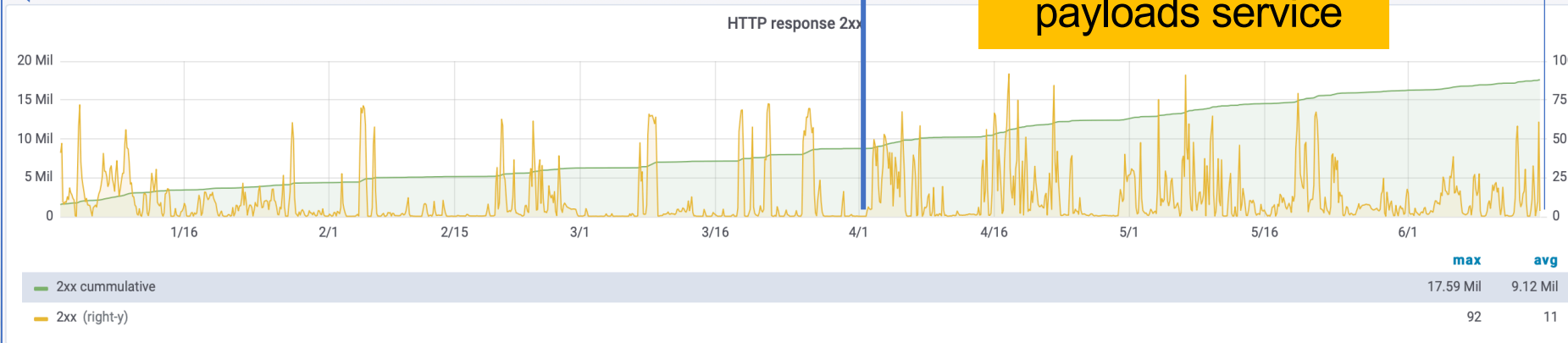


CDB service resource usability within the capacity of resources provisioned. See extra slides for capacity details.

Average of HTTP requests

9.5M of HTTP 2XX succeeded

Belle II Conditions payloads service



Operational Experience

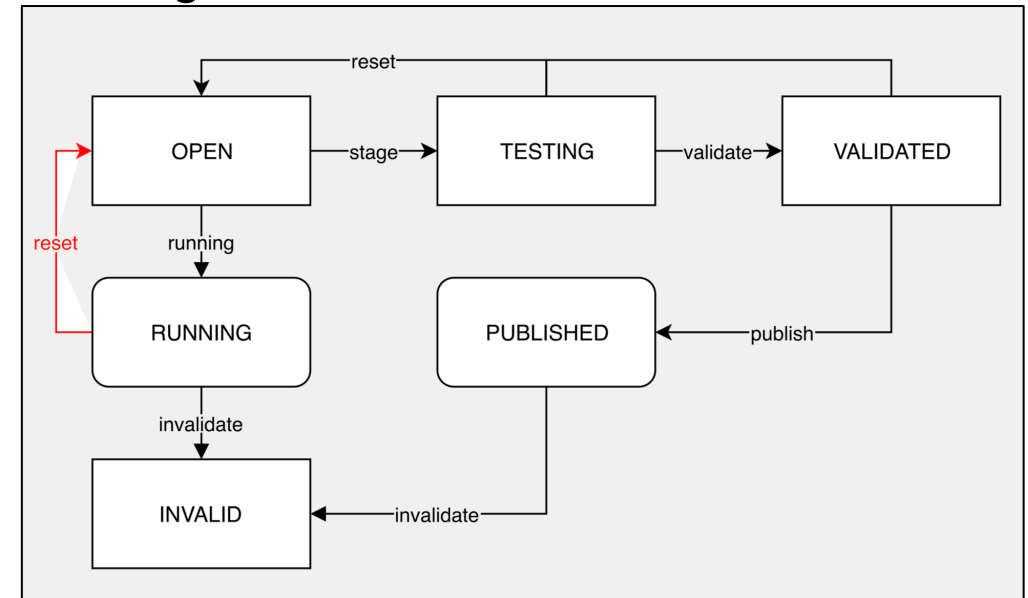


Global tag state machine implementation into the CDB server

The initial design of the state machine had to be adjusted to changes in Belle II workflows

Changes included:

- Database schema modification
- New API HTTP REST endpoint
- Transition restrictions are implemented
- New functional tests server side
- Enhancement of Authorization schemes
 - Only privilege user can execute a change in the Global Tag state i.e. Open to Running



This required effort from different experts from software and database group.

Deployed in July 1st 2019
All of this implementation deployed while minimize service disruption

Previously Global Tag states only supported three states, NEW, PUBLISHED, INVALID

New developments: CDBweb browser



- Allows user to visualize database
- Conditions database content accessed via web browser.

The Global Tag page

Questions? Write to gotekhin@bnl.gov

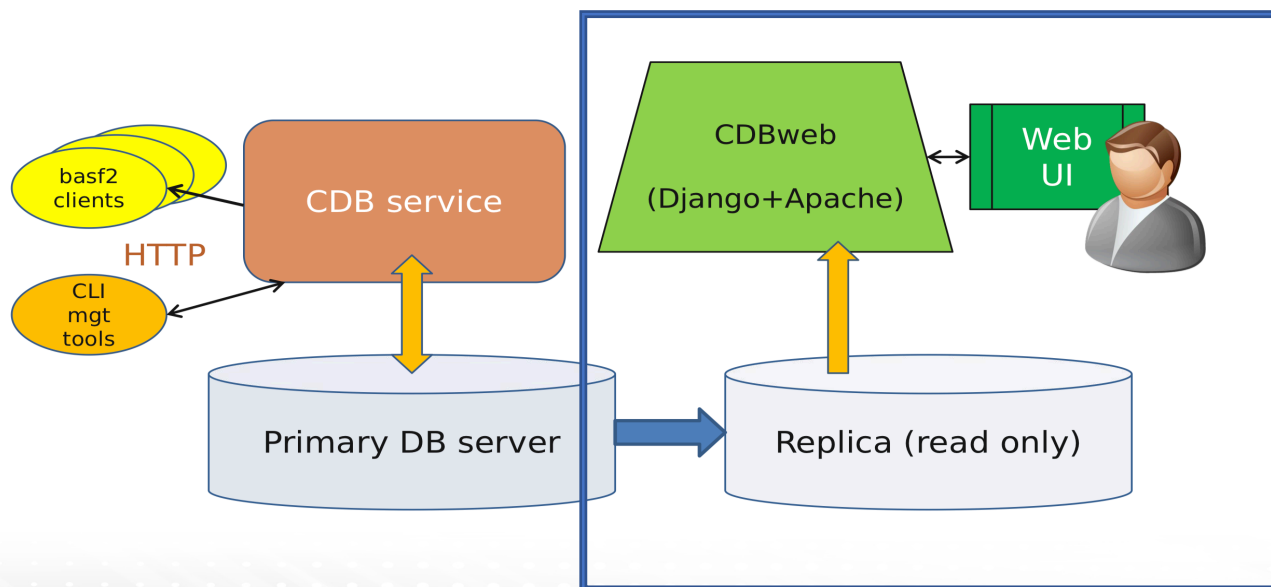
Home GlobalTag Payload Types of Payload Global Tag Comparison

GlobalTag: 534 items found [Click on items for more details](#)

ID: Name (can be partial): Status: Type: Modified by: items per page: 25

| Name | ID | Default? | Description | Status | Type | Modified | Mod. by | Total Payloads | Distinct Payloads |
|--|-----|----------|---|-----------|---------|-----------------------|-----------|----------------|-------------------|
| Temp_GT_testing | 433 | • | GT only for tool testing | NEW | DEV | 06/16/2019 9:35 a.m. | pkumar | 0 | 0 |
| val_GTooltesting | 632 | • | for software validation on data | NEW | DEV | 06/15/2019 9:26 a.m. | pkumar | 257 | 192 |
| data_reprocessing_procd_default_payloads | 631 | • | This GT contains the default payloads used in release-03-02, but with IPv6 variables for data reconstruction. It has to be used only for TESTING or as a fallback option | NEW | DEV | 06/15/2019 12:24 a.m. | temponi | 309 | 202 |
| staging_master_nbraun_20190614-003411 | 630 | • | Removed prescaling from shabtra | NEW | DEV | 06/14/2019 12:36 a.m. | nbraun | 1 | 1 |
| procd_phase2 | 629 | • | GT containing final set of Phase2 calibration | NEW | DEV | 06/12/2019 9:01 a.m. | sakaliam | 2337 | 3 |
| ARCH_pre_procd | 628 | • | GT to gather all the ARCH calibrations that are required for procd. | NEW | DEV | 06/11/2019 10:07 a.m. | lika | 1322 | 10 |
| test_overlap_fix | 627 | • | A test GT to try to automatically solve the old overlaps | NEW | DEV | 06/10/2019 3:38 p.m. | temponi | 761 | 1 |
| staging_master_nbraun_20190610-142259 | 626 | • | Added skim BelleHECL and raster to skim menu | NEW | DEV | 06/10/2019 2:23 p.m. | nbraun | 3 | 3 |
| mc_production_MC12c_rev1 | 625 | • | GT for payloads necessary for simulation as part of the MC12c campaign | PUBLISHED | DEV | 06/09/2019 9:53 p.m. | ghenett | 20 | 17 |
| procd_geometry | 624 | • | GT to store PYD geometry payloads | NEW | DEV | 06/08/2019 12:29 p.m. | sakaliam | 5 | 1 |
| procd_calibration_scratch | 623 | • | GT to prepare procd calibration payloads for official upload | NEW | DEV | 06/08/2019 12:26 p.m. | sakaliam | 4433 | 4 |
| ECL_pre_procd | 622 | • | GT to gather all the ECL calibrations that are required for procd. Procd will start from an empty global tag. | NEW | DEV | 06/07/2019 10:38 a.m. | ehill | 718 | 49 |
| release-03-02-00_rev2 | 621 | • | Software Development and Testing Tag This global tag contains necessary payloads for testing of the software and experiment independent MC. It cannot be used to analyze data and only contains payloads for the following intervals of validity - exp 0 nominal Belle I Configuration with full PYD (full phase 3) - exp 1002 Phase 2 Configuration with minimal vertexing detectors and additional background detectors - exp 1003 Belle II Configuration with partial PYD (early phase 3) Previous global tag: release-03-02-00_rev1 Changes are collected in Full request #4122, see https://clouds.desy.de/projects/B2/traces/softwarebuild-requests/4122 | PUBLISHED | RELEASE | 06/07/2019 7:15 a.m. | hupp_ware | 313 | 202 |
| master_2019-06-07 | 620 | • | Software Development and Testing Tag This global tag contains necessary payloads for testing of the software and experiment independent MC. It cannot be used to analyze data and only contains payloads for the following intervals of validity - exp 0 nominal Belle I Configuration with full PYD (full phase 3) - exp 1002 Phase 2 Configuration with minimal vertexing detectors and additional background detectors - exp 1003 Belle II Configuration with partial PYD (early phase 3) Previous global tag: | PUBLISHED | RELEASE | 06/11/2019 12:59 p.m. | hupp_ware | 328 | 307 |

Read Only Mode



Jason Web Token (JWT) authorization



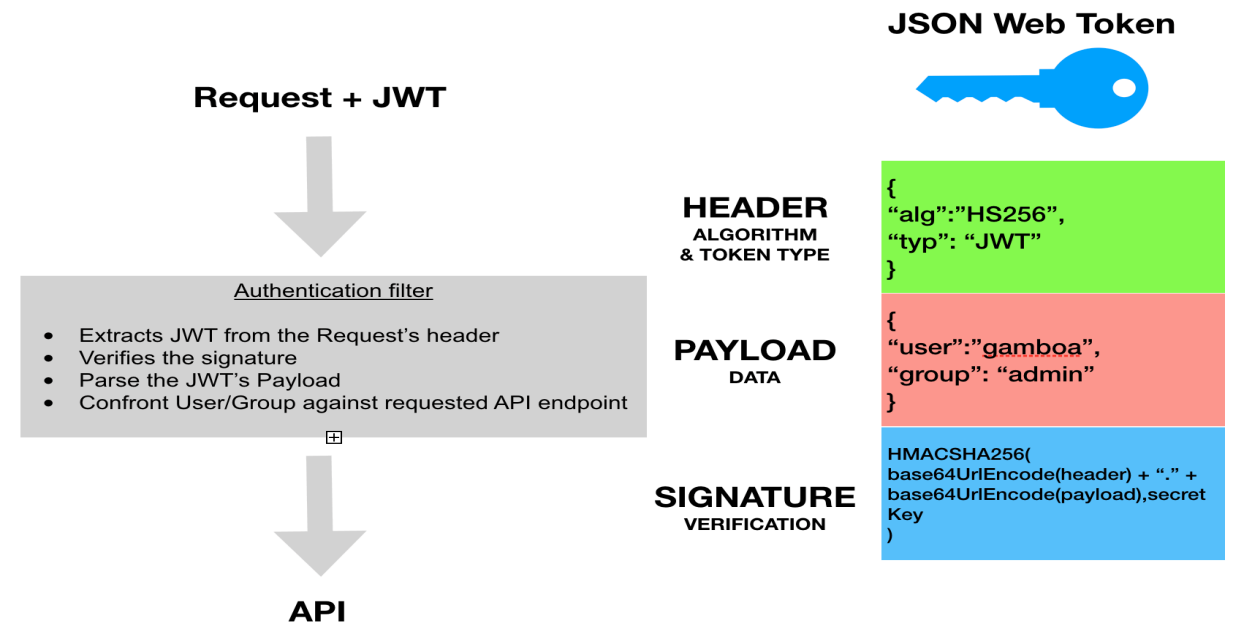
Belle II CDB is enhancing the authorization layer to enable fine-grained accessibility to serve conditions data.

Two possible options:

- Public Key Infrastructure (X.509 certificates)
- Token based infrastructure, JWT
 - Industry and HEP community are leaning towards token based technology as an alternative to provide Authorization as opposed to Public Key Infrastructure

Proof of concept of JWT for authorization implemented

Basic functionality tested via RESTful CDB API



Not a trivial task as the implementation involved troubleshooting different system used in the CDB framework

Future developments



HTTP REST API metadata

Augment current HTTP API to enhance support for Global Tag state machine

Authentication and Authorization

Extend JWT proof of concept by integrating an external token granting system and Belle 2 database clients

Summary



- Overview of the Belle II CDB database server presented
- **CDB service successfully delivered conditions data for the first run of the experiment in phase 3**
- Global tag state granularity extended
- PostgreSQL Streams standby replication of CDB metadata successfully enabled
- Alternative mechanism to visualize CDB content developed

Extra slides



Hardware capacity installed for Belle II CDB services



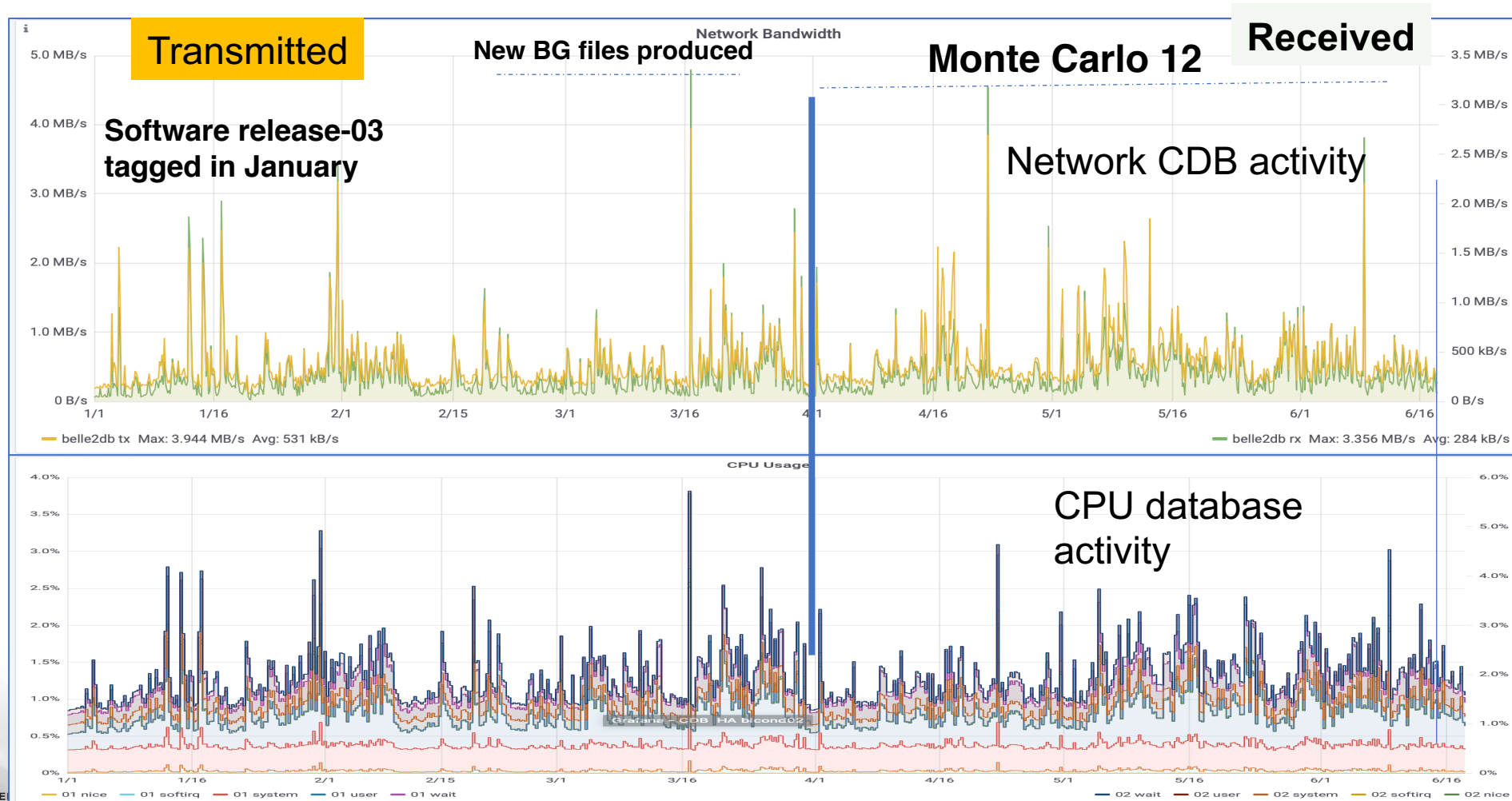
Belle2db service (metadata)

2 Nodes Dell R730xd
Two Intel(R) Xeon(R) CPU E5-2667 v4 @
3.20GHz -> total CPU thread 32
Memory 256 GB
Disk for Database Intel DC P3700
NVMe SSD on PCI Express 800 GB
20Gb/s Channel Bounded connectivity

Belle2db-files service (payloads)

2 Nodes Dell R430
Two Intel(R) Xeon(R) CPU E5-2650 v4
@ 2.20GHz total CPU threads 48
Memory 256GB
20Gb/s Channel Bounded connectivity

Performance CDB service (belle2db metadata service) CDB January-June 2019



Performance CDB service (belle2db-file services) from January-June 2019



CDB service resource usability within the capacity of resources provisioned despite of the increase in the usage of the belle2db-file service

Software release-03 tagged in January

New BG files produced

Monte Carlo 12



Belle2db-file services files performance



Despite of the change in the access to the Belle2 DB file service, system remained stable and not major issues observed.

