



PNNL-SA-127779

Belle II Distributed Data Management and Networking

VIKAS BANSAL, MALACHI SCHRAM, ANTONIO LEDESMA

Pacific Northwest National Laboratory, Richland, WA

DPF '17, Fermilab, August 3, 2017

Belle II Collaboration



c.f. CERN Greybook July 2017
ATLAS: 39 countries, 217 inst., 7783 members
CMS: 49 countries, 208 inst., 6217 members
ALICE: 41 countries, 167 inst., 2799 members
LHCb: 17 countries, 74 inst., 1494 members

Computing requirement on par with LHC Run I

KEK Laboratory, Tsukuba, Japan



- ▶ The Belle II experiment is for Super B factory at KEK in Japan
- ▶ Complementary physics to the LHC based on precision measurements from high-intensity beams
- ▶ Total integrated luminosity : 50 ab^{-1}
- ▶ Collisions start in Early 2018
- ▶ Similar data rate as from LHC Run I

Raw Data : 100 KB / event

Detector mDST : 5 KB / event

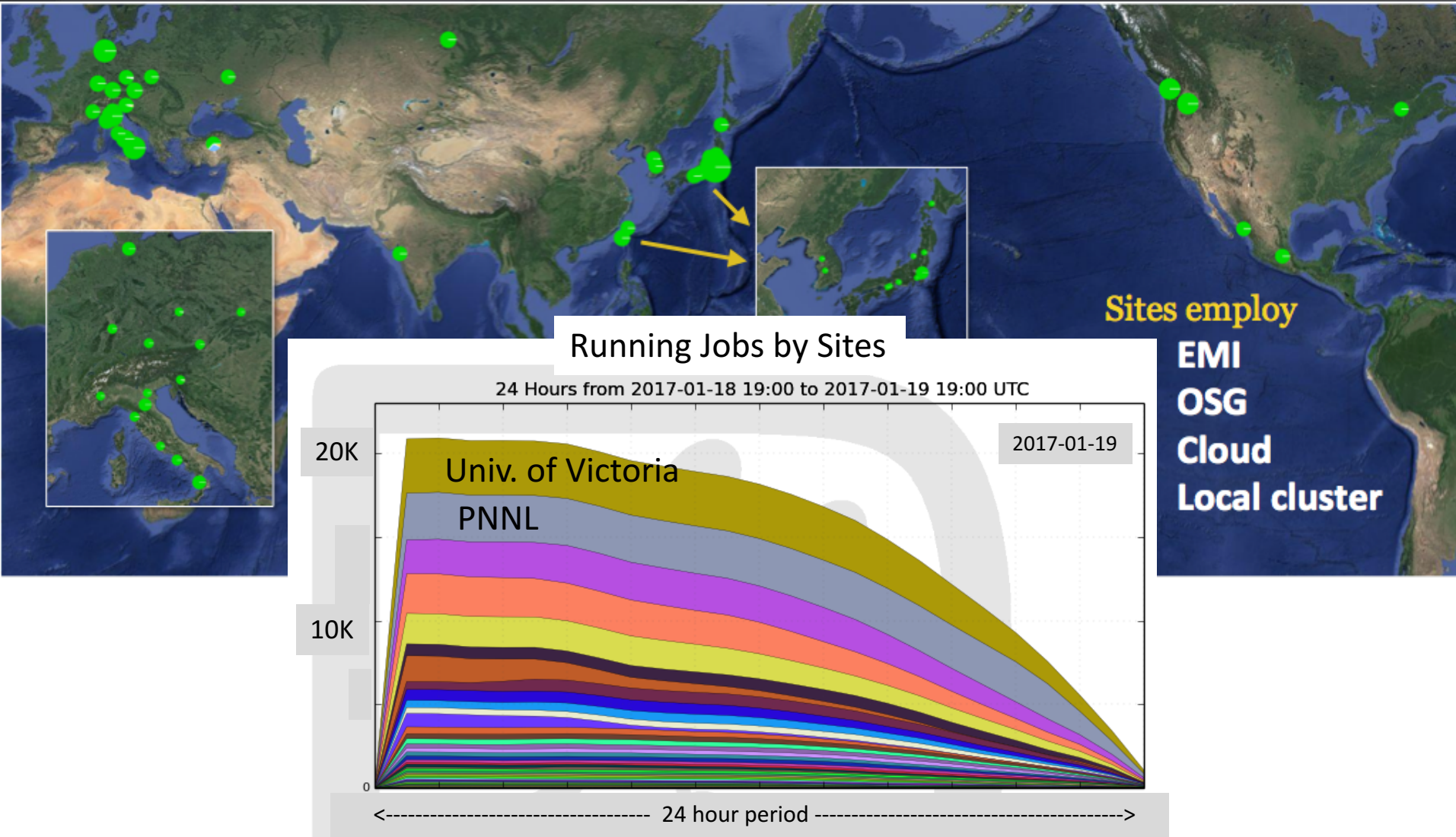
MC mDST : 6 KB / event

Reconstruction : 20 HEPSPC *s / event

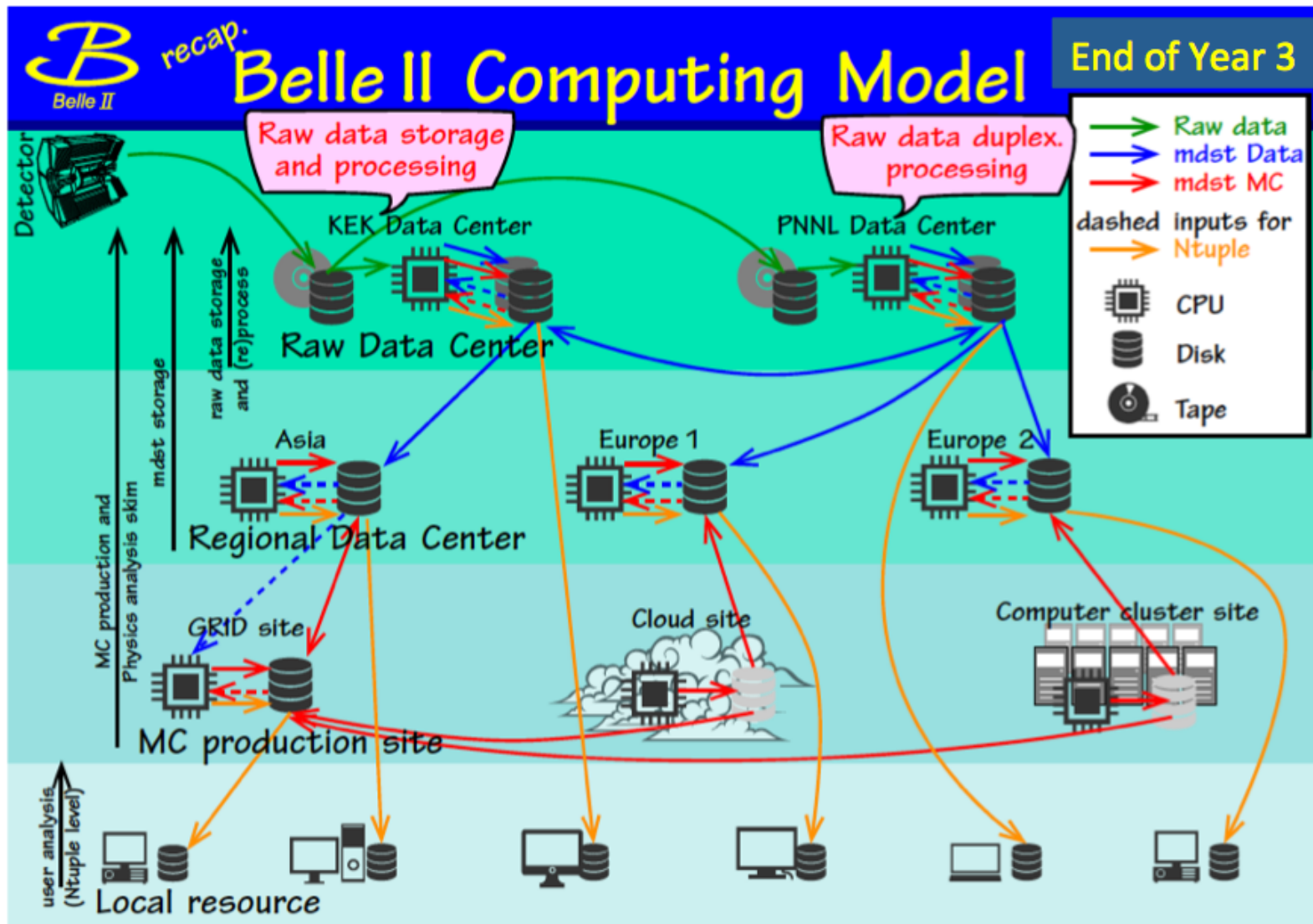
MC : 60 HEPSPC *s / event

	2017	2018	2019
Disk [TB]	3000	4500	11000
Tape [TB]	2000	2200	6500
CPU [KHEPSPC06]	210	400	480

Belle II computing sites across the globe

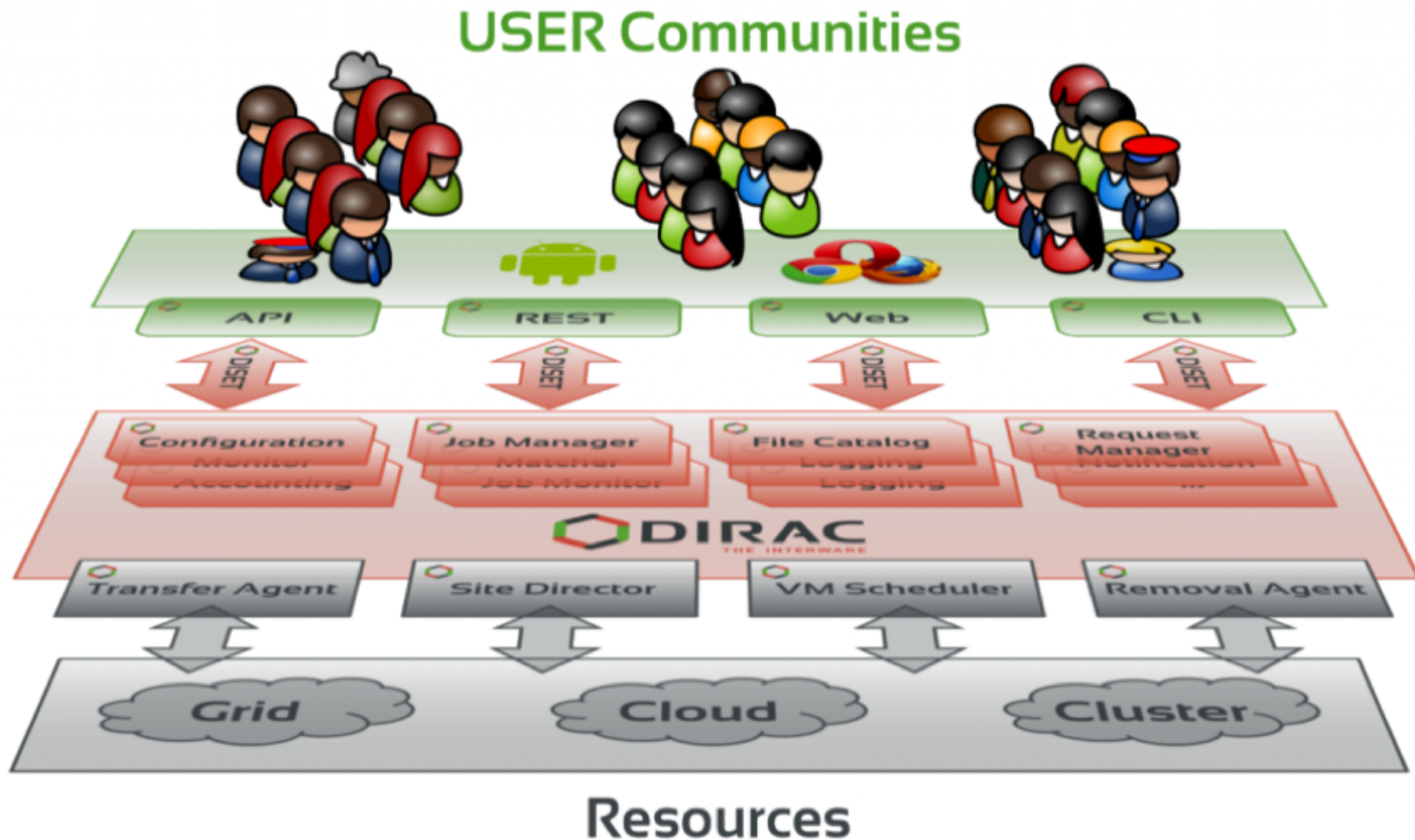


Belle II computing model until Year 3



Belle II distributed computing software

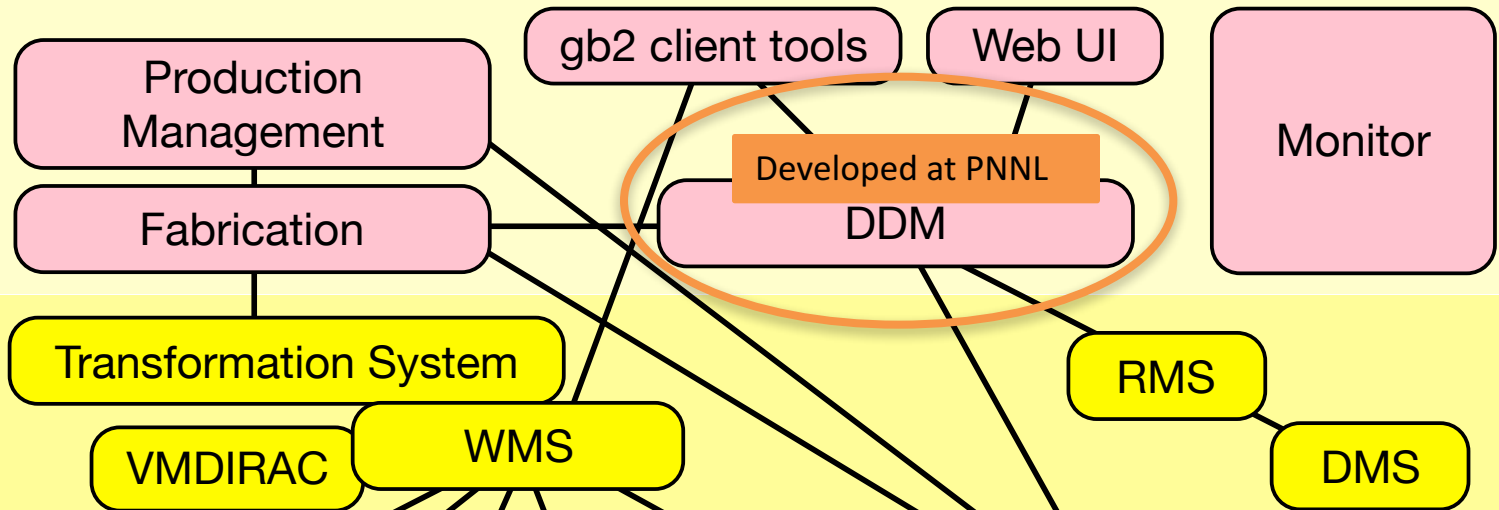
- ▶ DIRAC (Distributed Infrastructure with Remote Agent Control) as the solution of choice
- ▶ Successfully used by LHCb. Extended for Belle II.



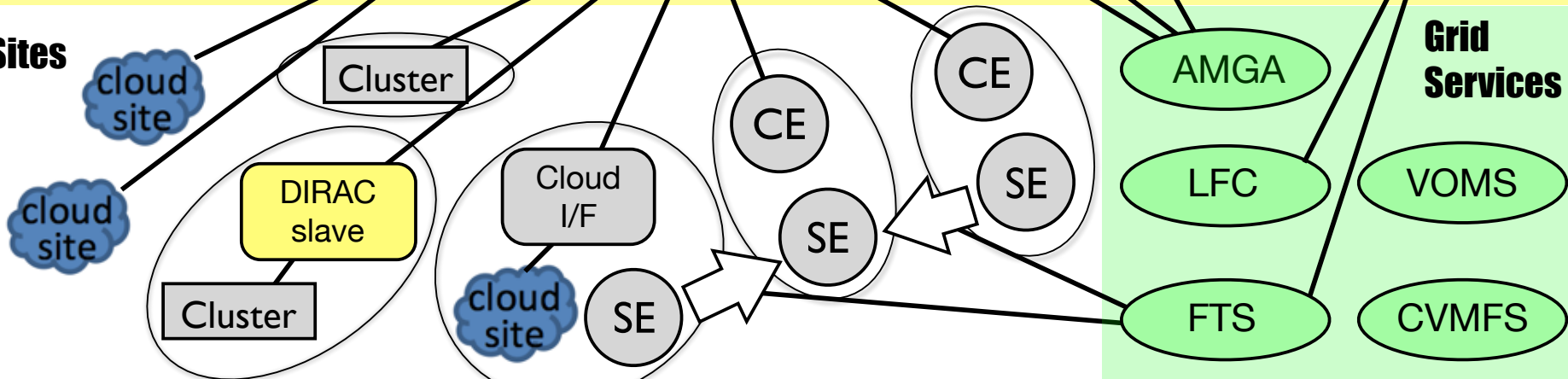
Belle II distributed computing layers

Production Manager Data Manager End Users Operations

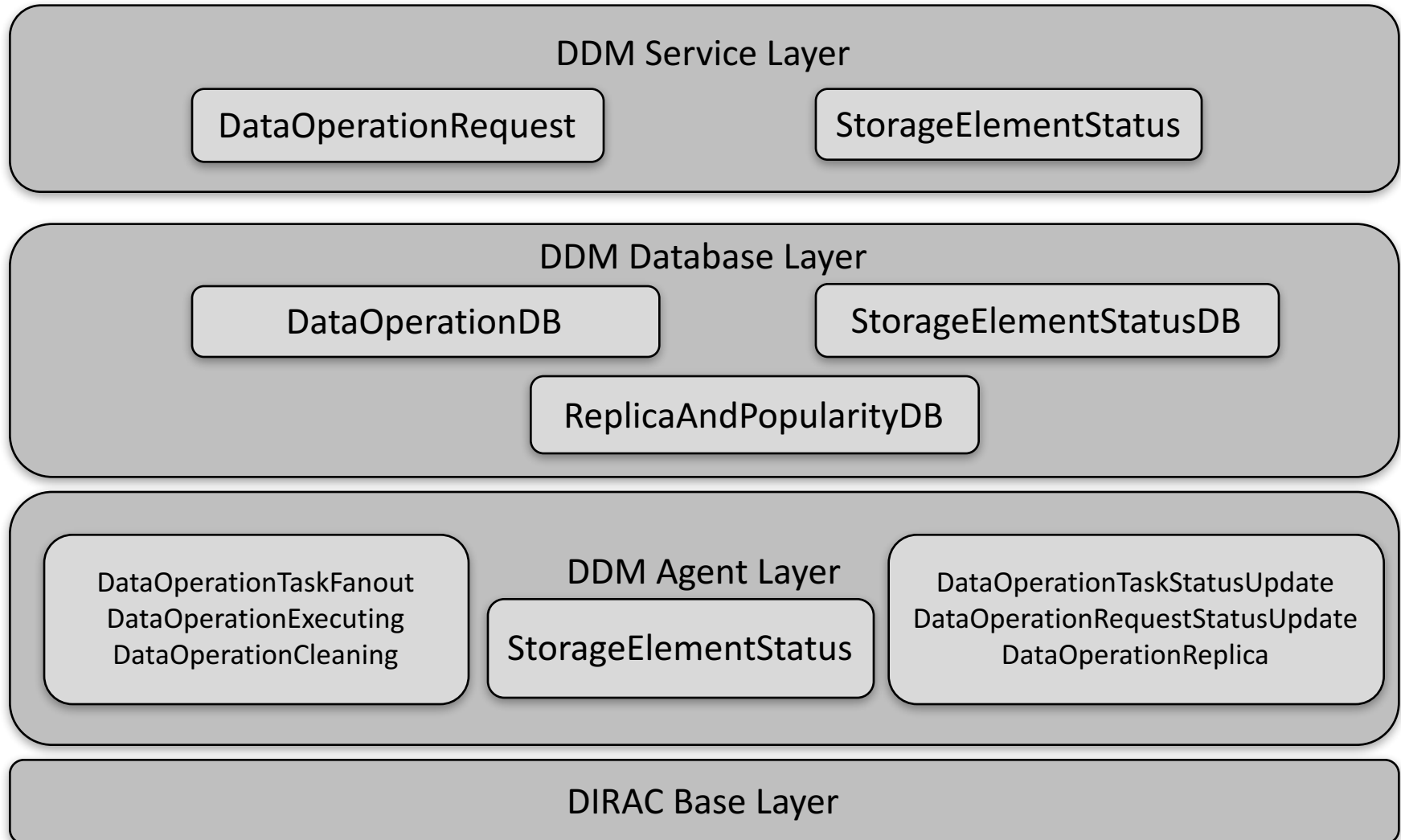
BelleDIRAC



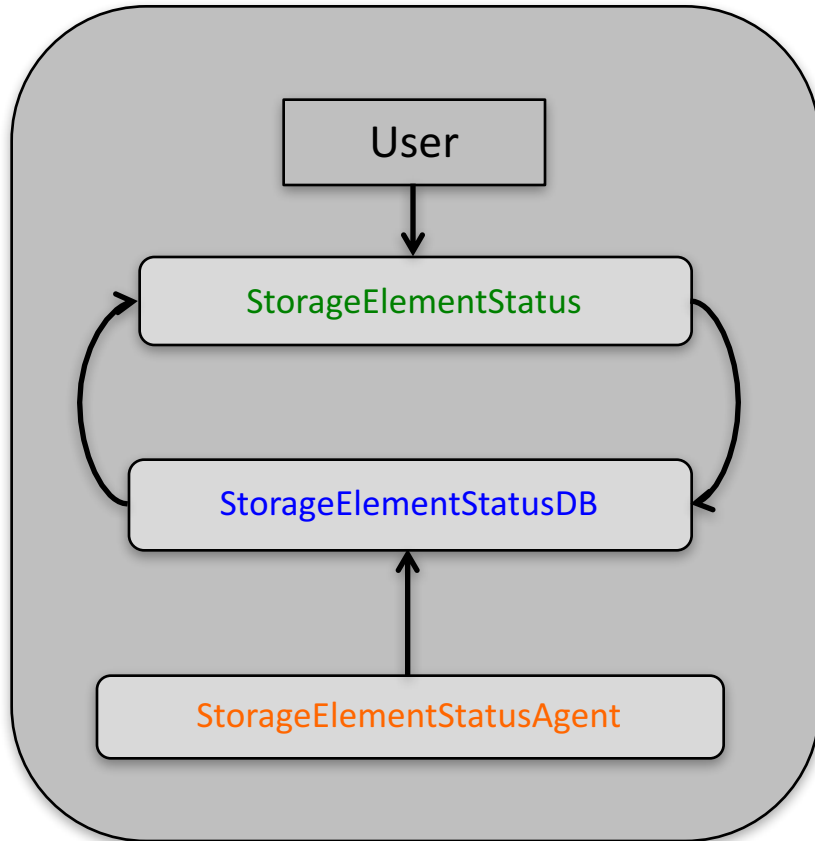
Sites



Distributed data management system overview



StorageElement Status Workflow



Purpose:

- Provide near real time storage elements information

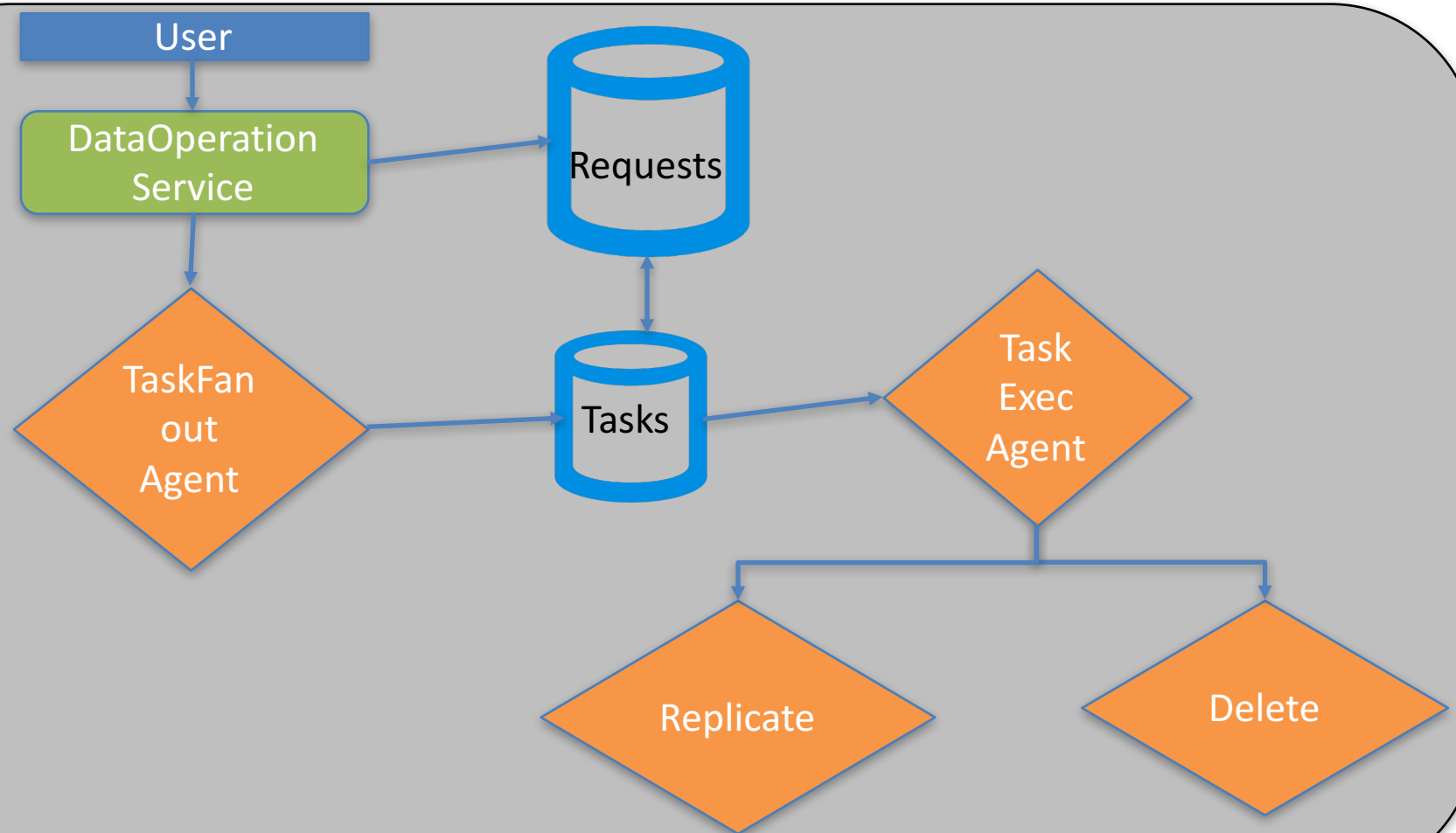
Current implementations:

- ① **StorageElementStatus** provide API layer for internal and external users.
- ② **StorageElementStatusAgentDB** provide persistified information
- ③ **StorageElementStatusAgent** :
 - Available space
 - Access rights (w/r) at file and director level
 - ADLER32 Checksum

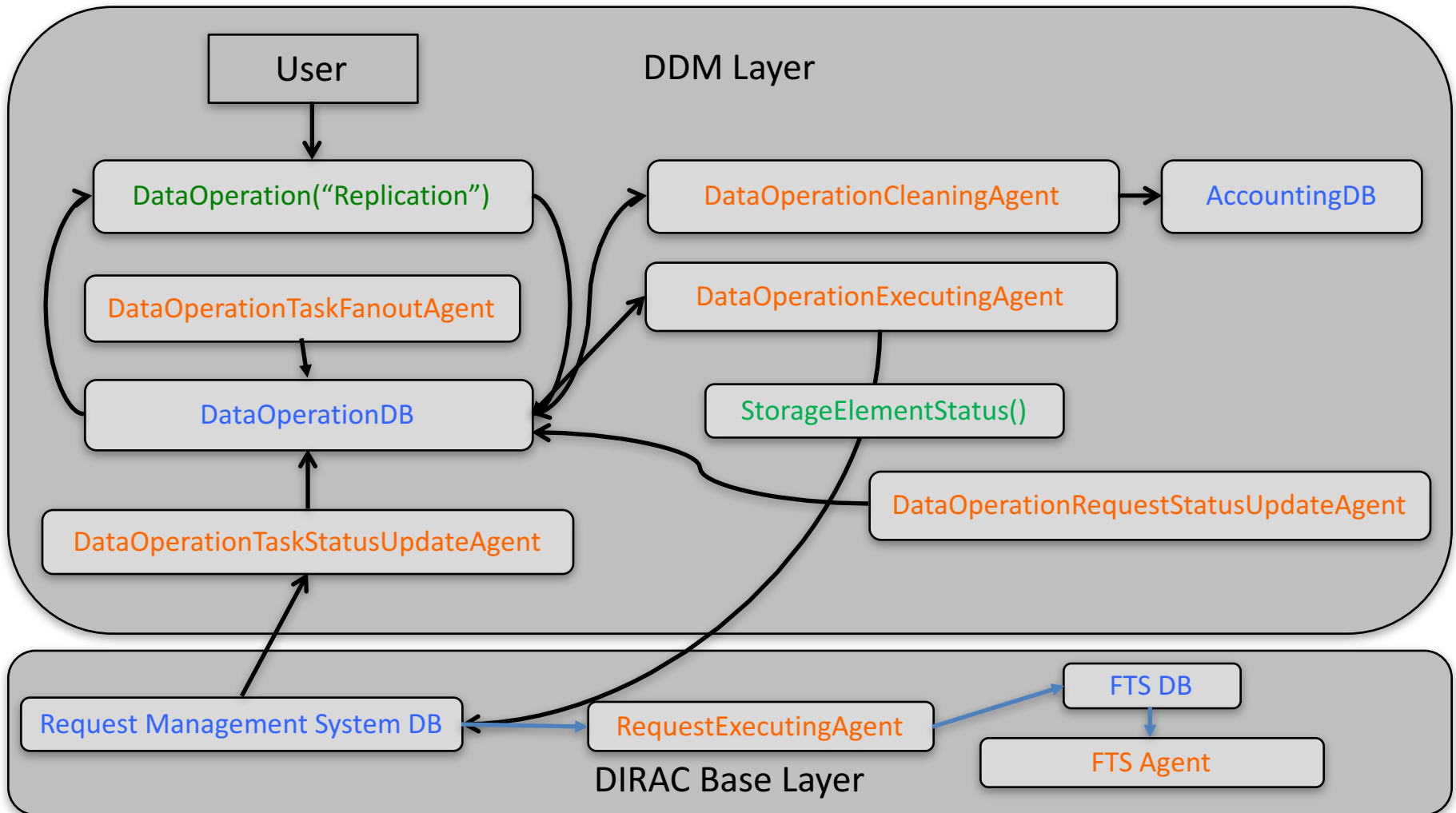
Extra features:

- Include full ACL

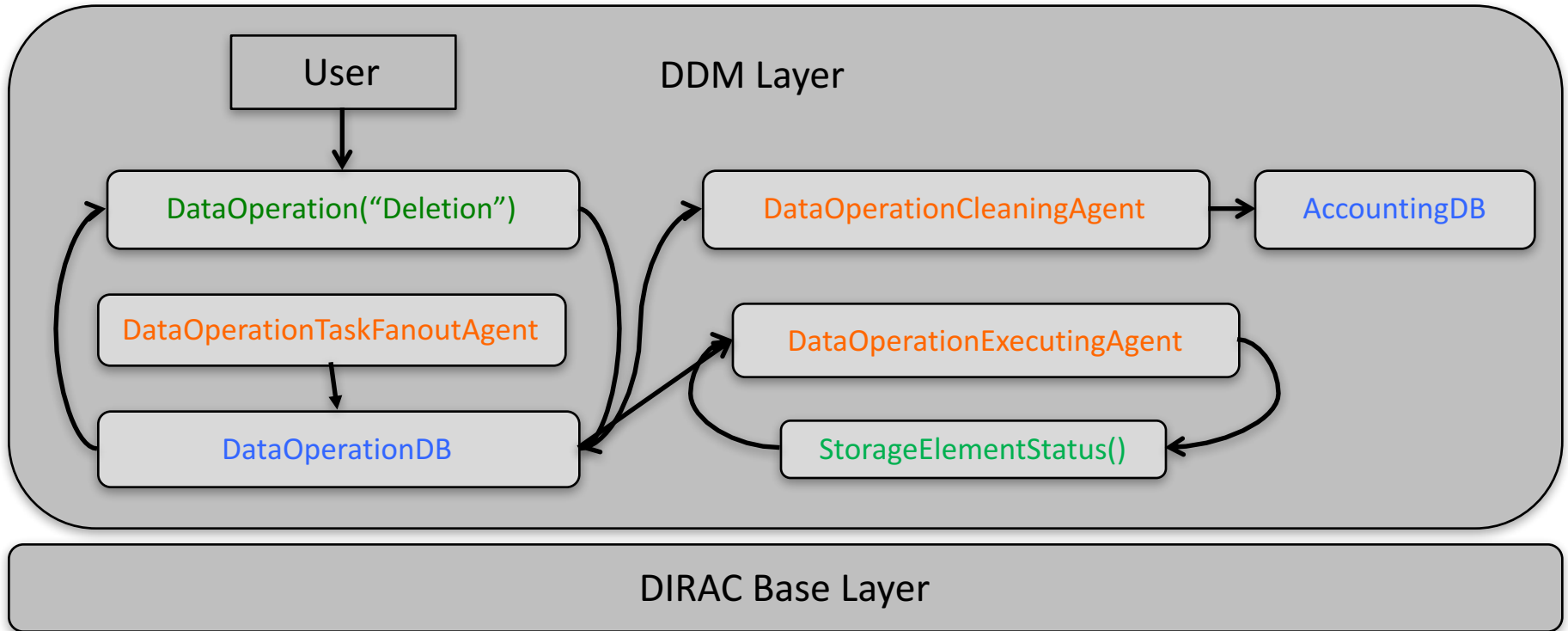
DataOperation : Workflow Overview



DataOperation : Replication Workflow



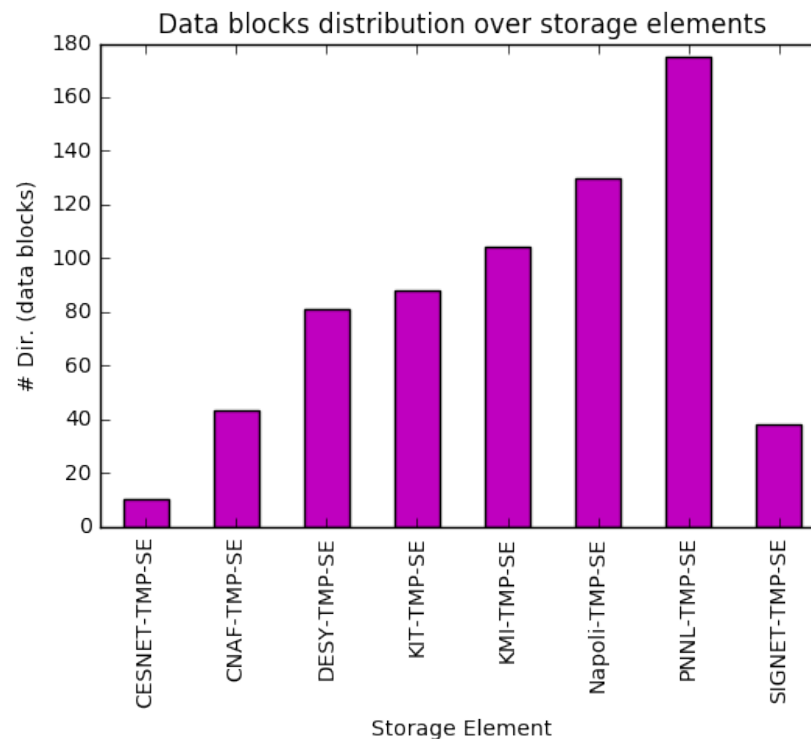
DataOperation : Deletion Workflow



- ▶ DDM by design does not delegate deletion of files to base DIRAC to avoid race conditions on data files
 - In addition Belle II utilizes AMGA (metadata service) that is not available in base DIRAC
- ▶ Data operation executing agent only acts on one operation on physical file at any given time

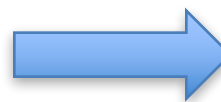
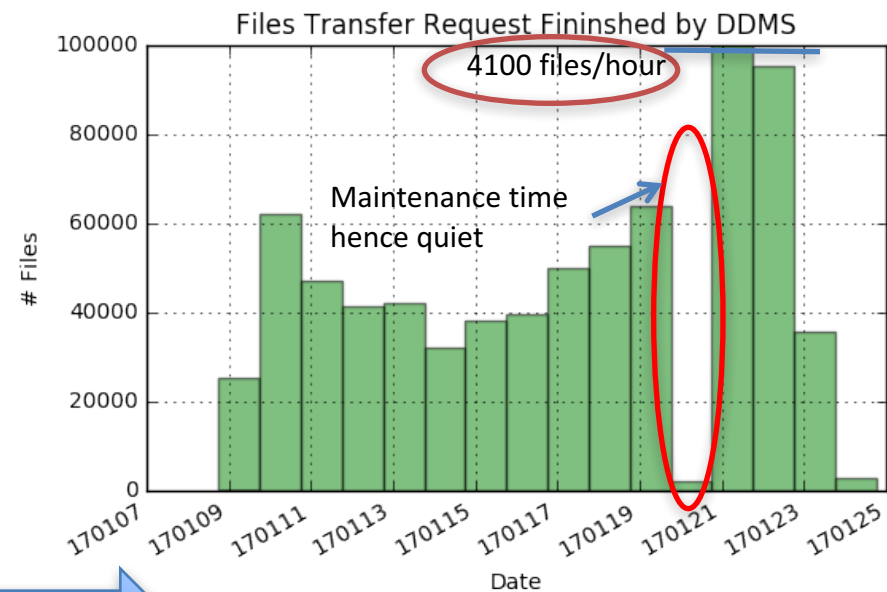
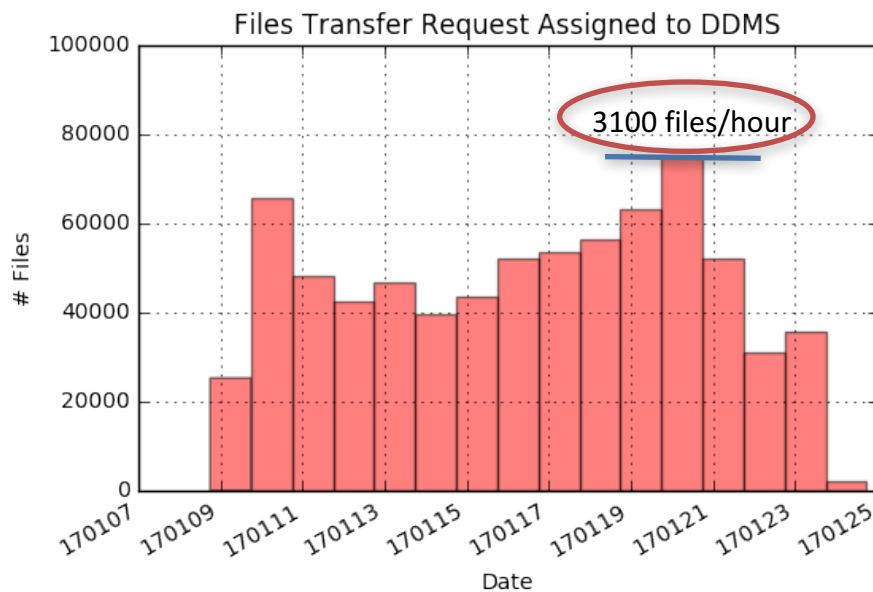
Distributed data management system in production

- ▶ Distributed Data Management System (DDMS) is successfully used in Monte Carlo (MC) samples production
- ▶ Production system hand overs data files to DDMS
 - DDMS decides to transfer the files/data to specific sites as per policy.



Distributed data management system in production

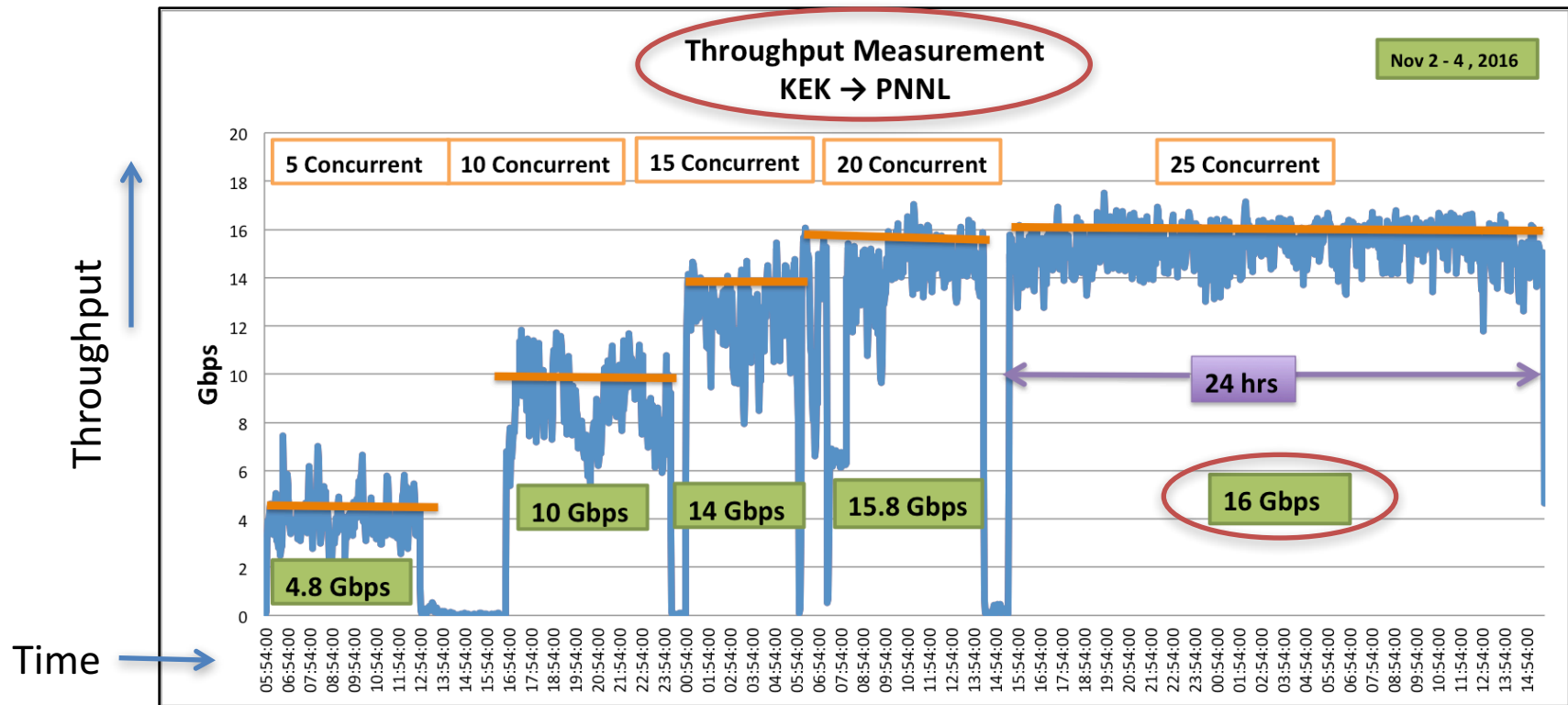
- ▶ Production system hand overs data files to DDMS
 - DDMS check for site health before starting any transfers. Replicates through DIRAC using File Transfer Server and reports back finished transfers



-
- The diagram illustrates the DDMS Network Component architecture. It features a central component box labeled "DDMS Network Component" containing three stacked modules: "Network Accounting Service", "Network Accounting DB", and "Network Accounting Agent". To the left is a "WebAppDIRAC" box, and to the right is a "MaDDash API" box. Arrows indicate the data flow: from WebAppDIRAC to the Network Accounting Service; from the Network Accounting Service to the Network Accounting DB; from the Network Accounting DB to the Network Accounting Agent; from the Network Accounting Agent to the MaDDash API; and from the MaDDash API back to the Network Accounting Service. A heatmap titled "Average Percentage of Dropped Packets Between Sites" is positioned next to the WebAppDIRAC box, showing network performance data between various sites.

WAN data challenge KEK ↔ PNNL

- ▶ Estimated network bandwidth for peak outgoing traffic from KEK : 9 Gbps
- ▶ KEK outgoing traffic measured at 16 Gbps
- ▶ DDM can be made to assess various network paths for optimal data transfers



Summary and Plan Moving Forward

- ▶ Belle II computing needs on par with LHC Run I
- ▶ Belle II distributed computing software is written in DIRAC framework
- ▶ Distributed data management system handles data transfers and deletion and is currently deployed in production mode
- ▶ We are actively working on next version of DDM that will directly schedule transfers to FTS to avoid DIRAC's RMS layer
- ▶ We foresee to add network health and monitoring in distributed data management system

Detailed Deletion Workflow

