

First Production with the Belle II Distributed Computing System

CHEP

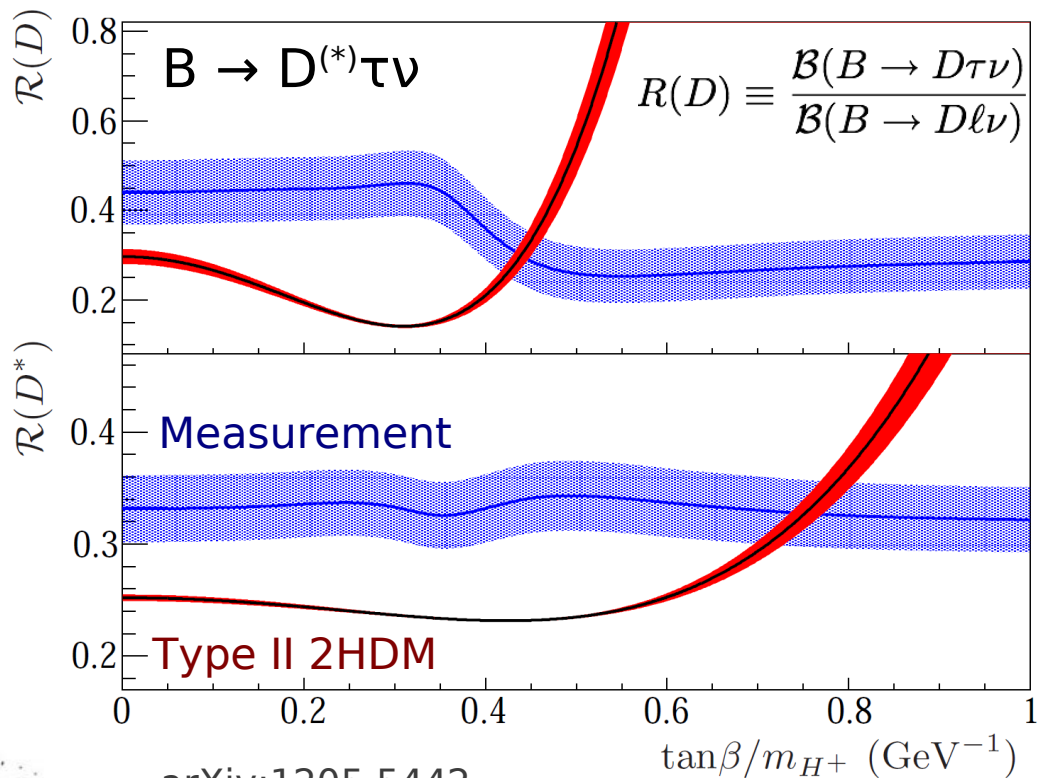
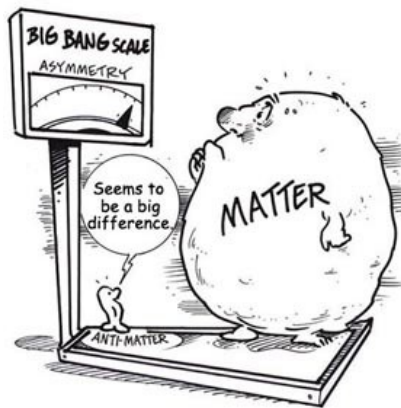
14.10.2013

Takanori Hara, Thomas Kuhr,
Hideki Miyake, Martin Sevier
for the Belle II Computing Group

Physics Objective of Belle and Belle II



- ✓ Confirmation of KM mechanism of \mathcal{CP} in the Standard Model
- ✗ \mathcal{CP} in the SM too small (by many orders of magnitude) to generate observed baryon asymmetry in the universe



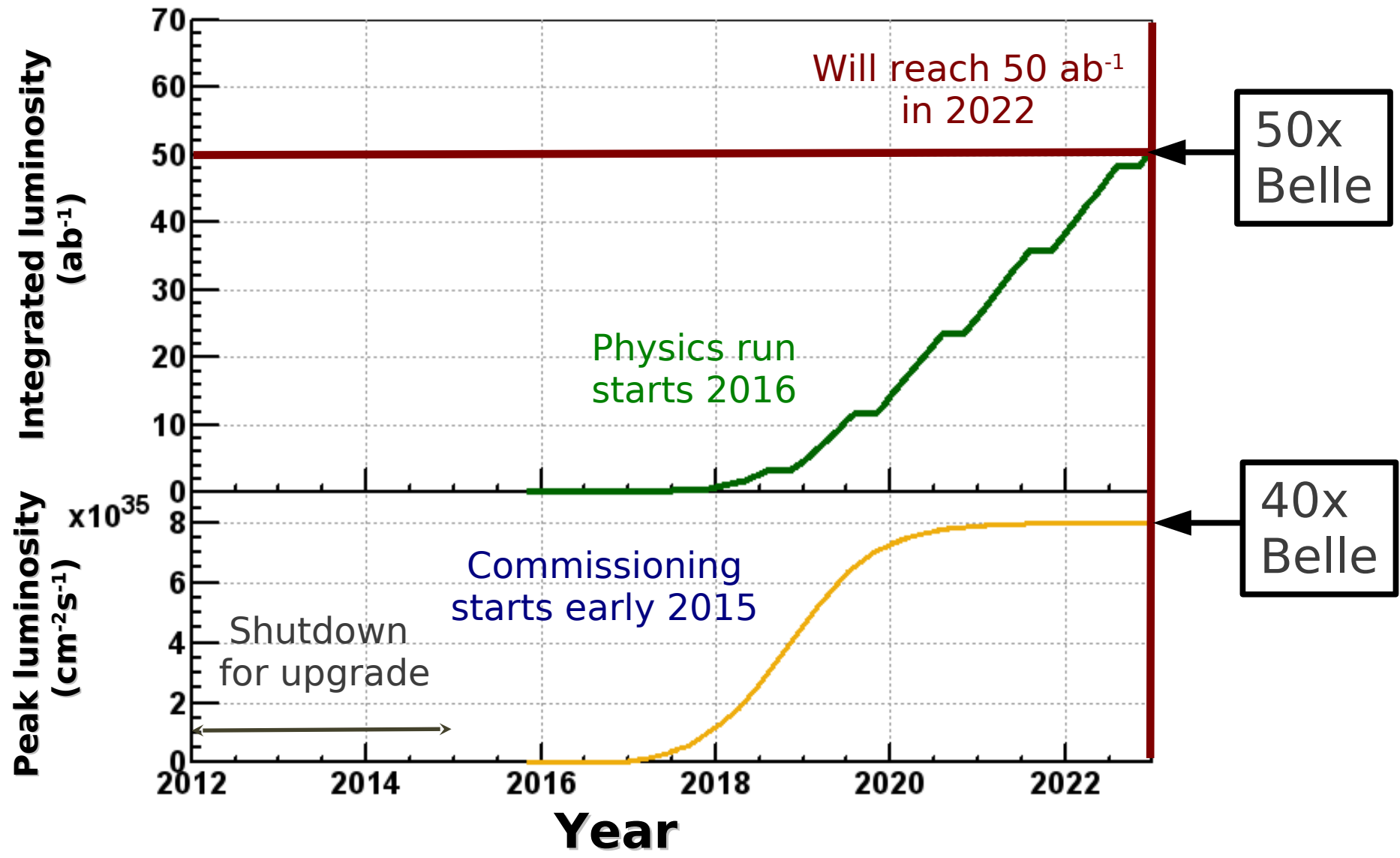
arXiv:1205.5442

- Need sources of \mathcal{CP} beyond the SM

→ Super B factory

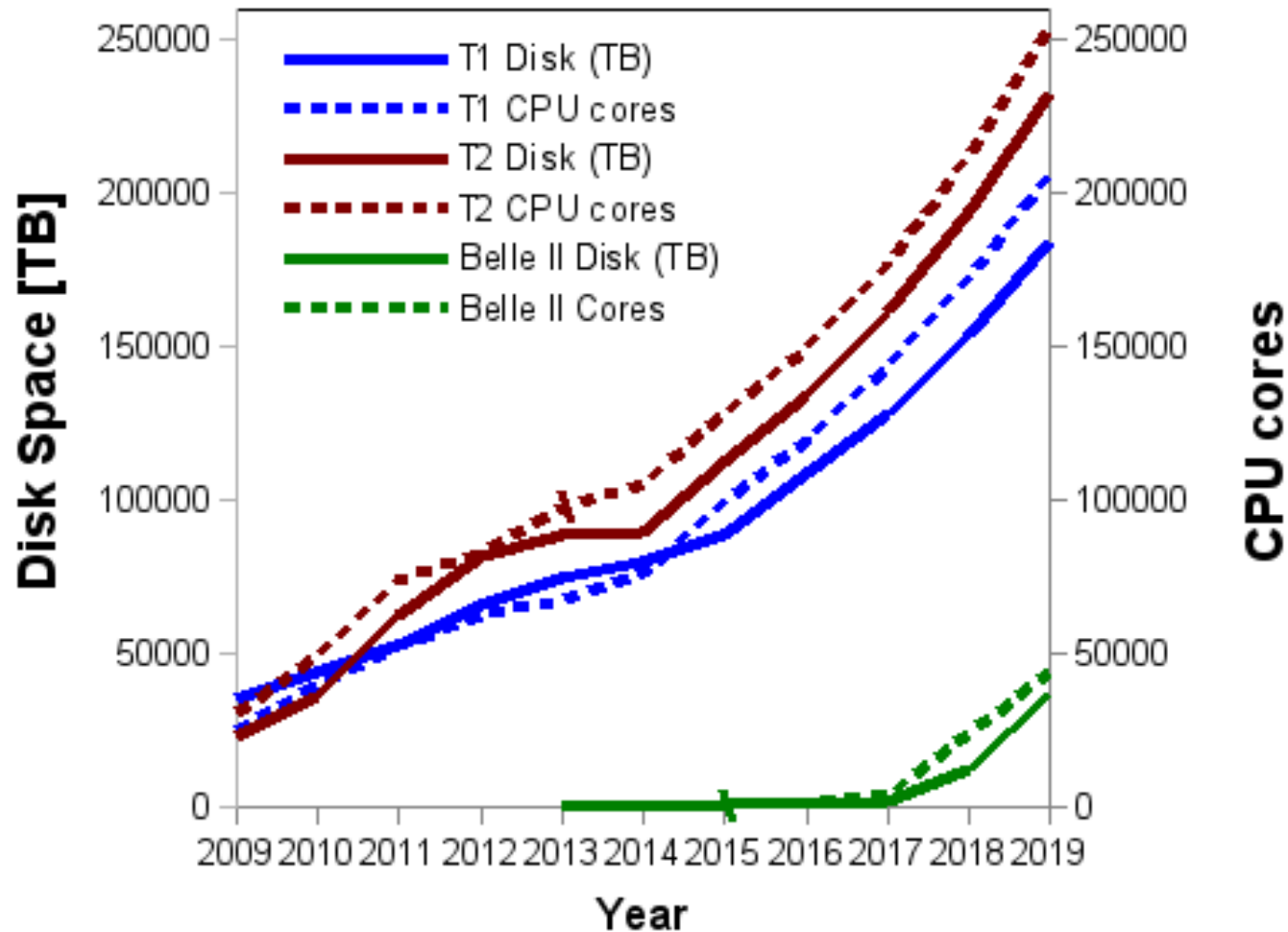
Complementary to LHCb

Projection of Luminosity at SuperKEKB



Resource Estimates

WLCG & Belle II Resources



→ Similar data rate as LHC experiments!

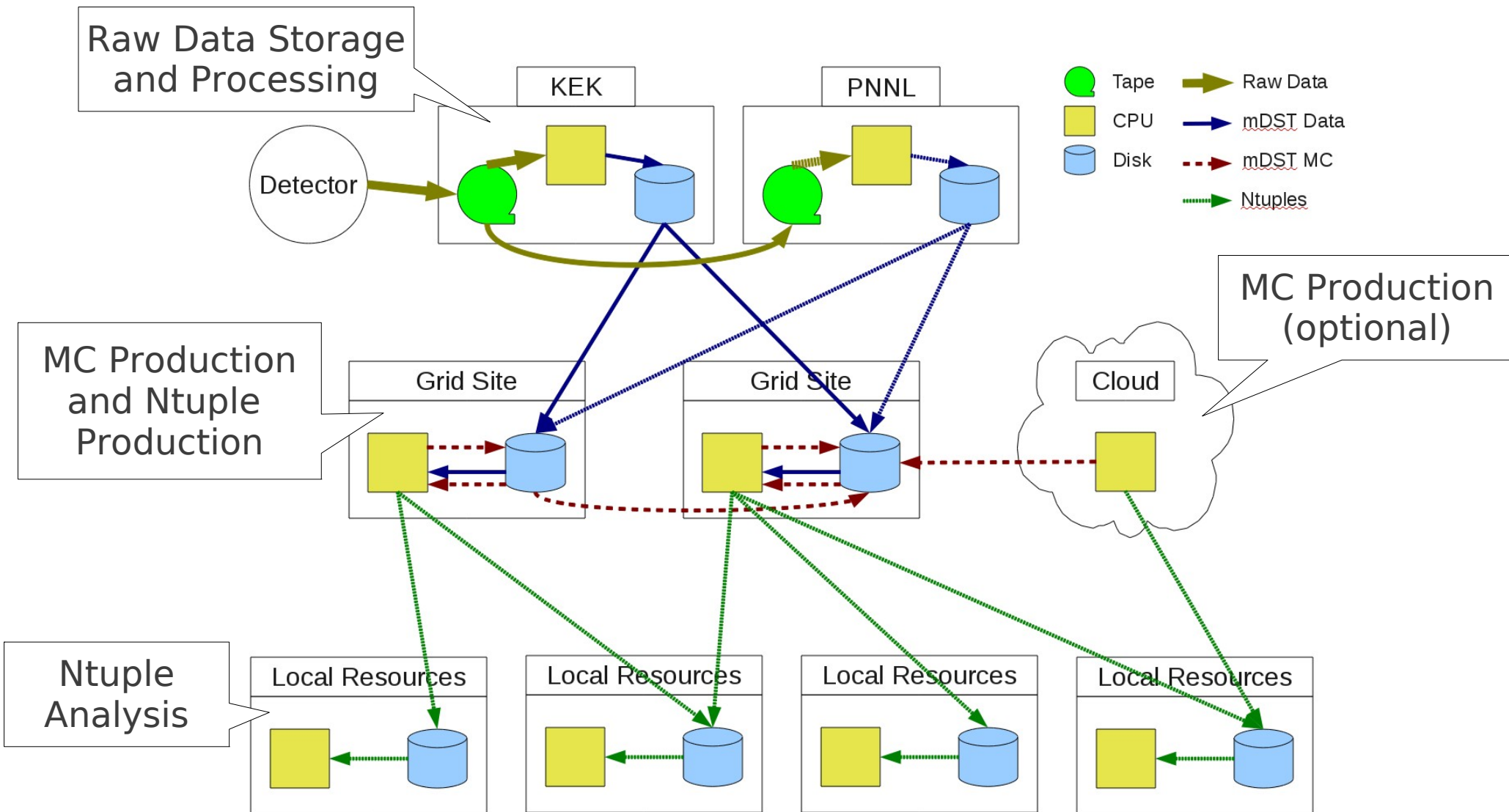
Belle II Collaboration



~530 members
94 institutions
from 23 countries



Computing Model



Distributed Computing System

- Based on existing, well-proven solutions plus extensions for Belle II
- DIRAC for job management
- AMGA for metadata

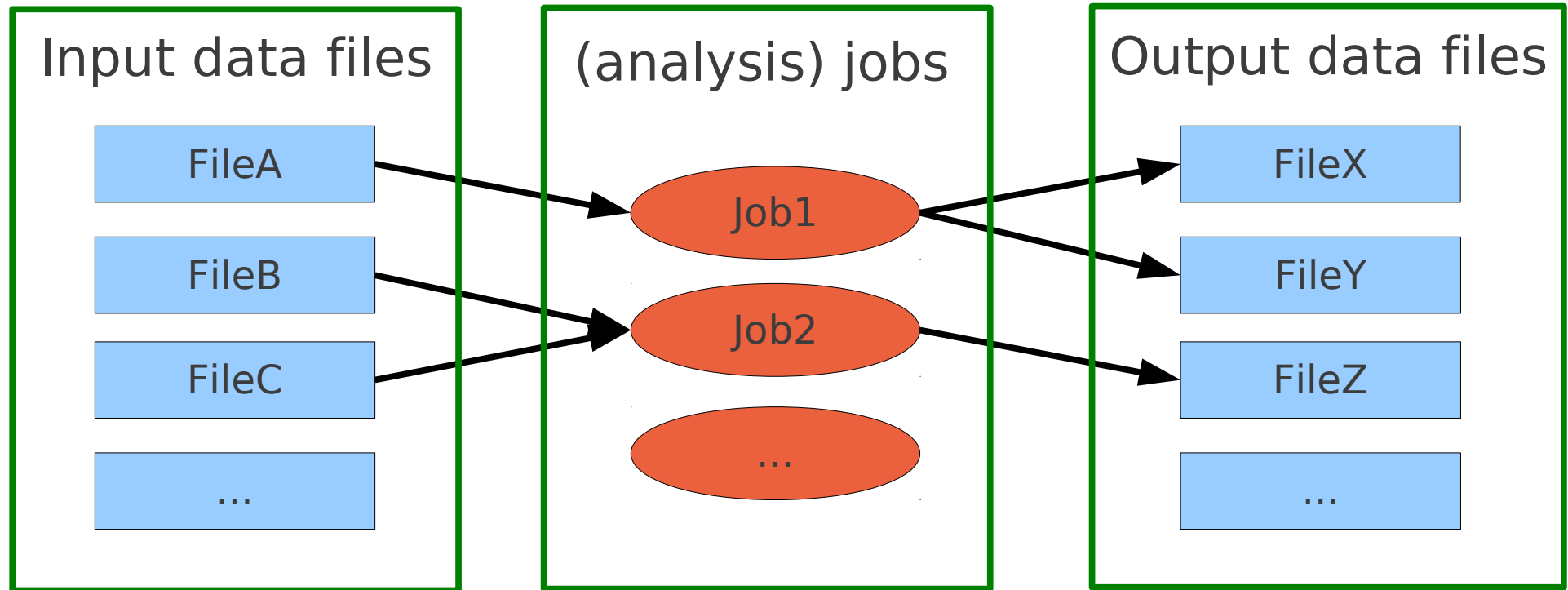


- CVMFS for software distribution
(thanks to CERN and Steve Traylen for providing the Stratum-0 server, and to GridKa for the stratum-1 server)

root/svn/trunk/grid/BelleDIRAC

FrameworkSystem/	4326 (9 months ago) by myco: basic sites management service for BelleDIRAC
Web/	5519 (4 months ago) by hideki: remove unused AMGA API
WorkloadManageme...	6098 (2 months ago) by hideki: fix a bug
gbasf2/	6647 (2 weeks ago) by hideki: fix unnecessary AMGA initialization
README	4325 (9 months ago) by myco: init files for BelleDIRAC distribution
__init__.py	6348 (6 weeks ago) by hideki: release for 2nd MC campaign

Workflow Abstraction



Input dataset

Project

Output dataset

- Don't deal with single files and jobs, but with datasets and projects

Projects

```
[ccx13] ~ $ gb2_project_summary -g belle_mcprod
Project      Owner      Status  Done/Fail/  Run/Wait  Submission  Time(UTC)  Duration
B2Kstargamma_BGx1_s1 tkuhr      Good    1000/  0/  0/  0  2013-08-14 14:41:57  06:47:32
B2Kstargamma_BGx0_s1 tkuhr      Good    1000/  0/  0/  0  2013-08-14 14:45:15  05:18:30
```

- Job submission

- `gbasf2 -r 1000 -s B2Kpi.py -p B2Kpi_s01`

- Job monitoring

- `gb2_project_summary`
- `gb2_project_analysis --Project B2Kpi_s01`
- `gb2_job_status --Project B2Kpi_s01 --Status=failed`

```
[ccx13] ~ $ gb2_project_analysis --Project testneb_b1
100 jobs are selected.
```

```
Project testneb_b1 summary:
Done (60)
  Execution Complete (60)
    Done (60)
      OSG.Nebraska.us: 60
Completed (25)
  Pending Requests (25)
    Done (25)
      OSG.Nebraska.us: 25
Failed (15)
  Application Finished With Errors (6)
    Exit Status 1 (6)
      OSG.Nebraska.us: 6
  Job stalled: pilot not running (9)
    Preparing to upload (1)
      OSG.Nebraska.us: 1
    Registering (1)
      OSG.Nebraska.us: 1
    Running (3)
      OSG.Nebraska.us: 3
    Selecting SE (1)
      OSG.Nebraska.us: 1
    Uploading (3)
      OSG.Nebraska.us: 3
```

- Rescheduling of failed jobs

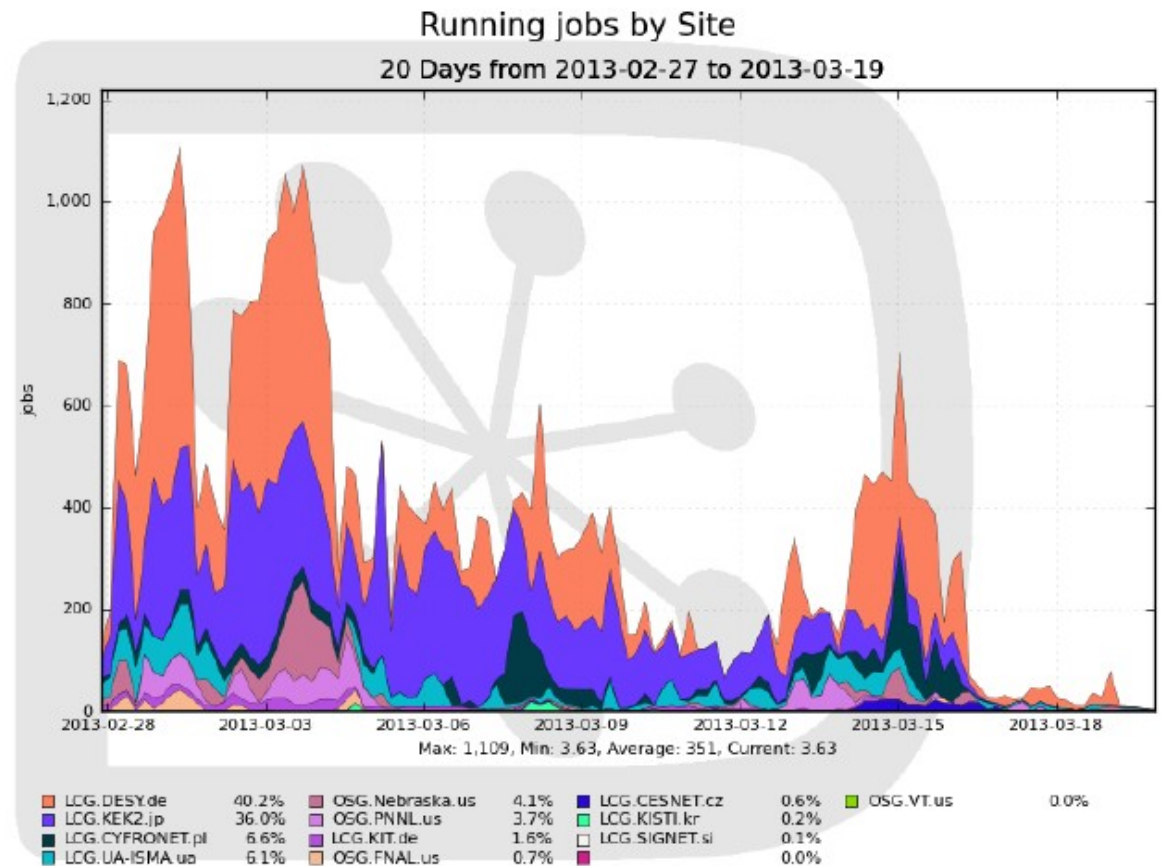
- `gb2_job_reschedule --Project B2Kpi_s01`

- Job output

- `gb2_job_output --Project B2Kpi_s01 --Status=failed`

First MC Production Campaign

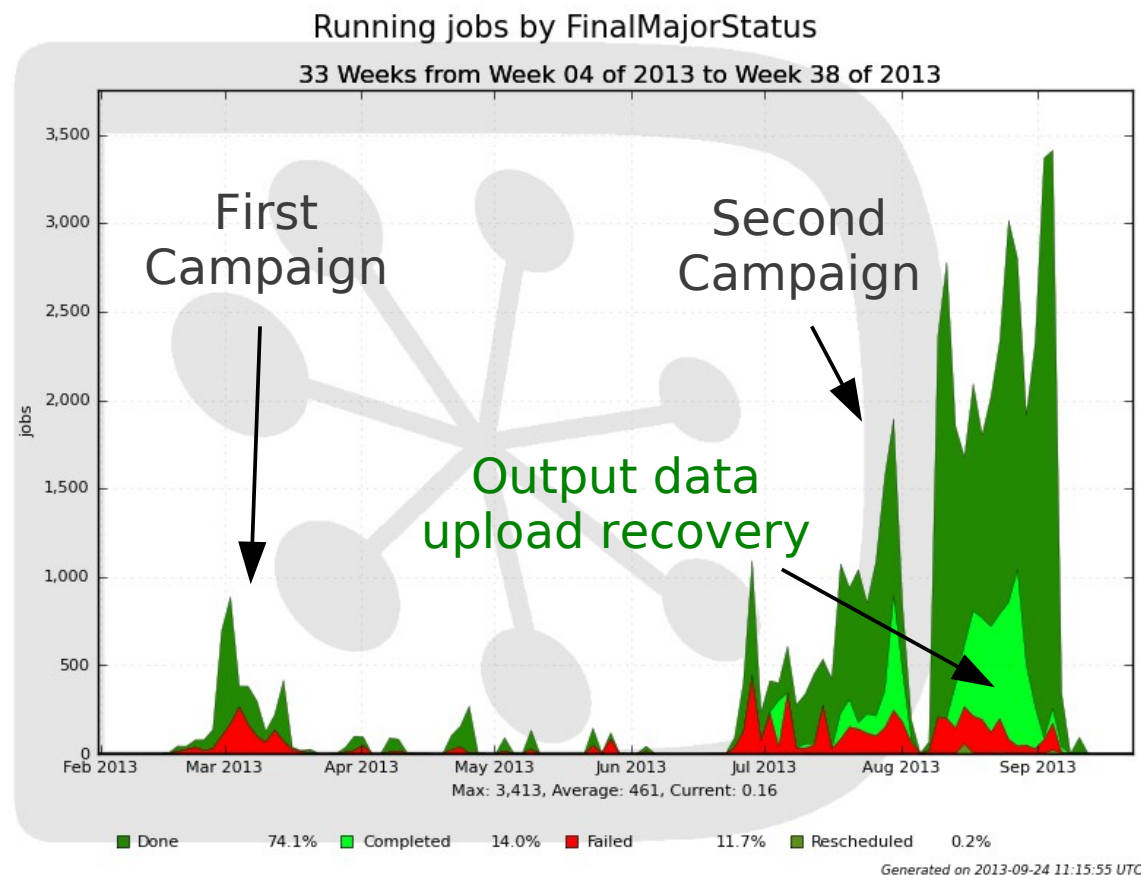
- February 28 – March 19, 2013
- 1st stage: event generation and detector simulation
→ raw data
- 2nd stage: reconstruction
 - 240k jobs,
40 kHS*days
 - 60M events,
190 TB of output data
- ~20% failure rate: metadata registration, input data download, application errors



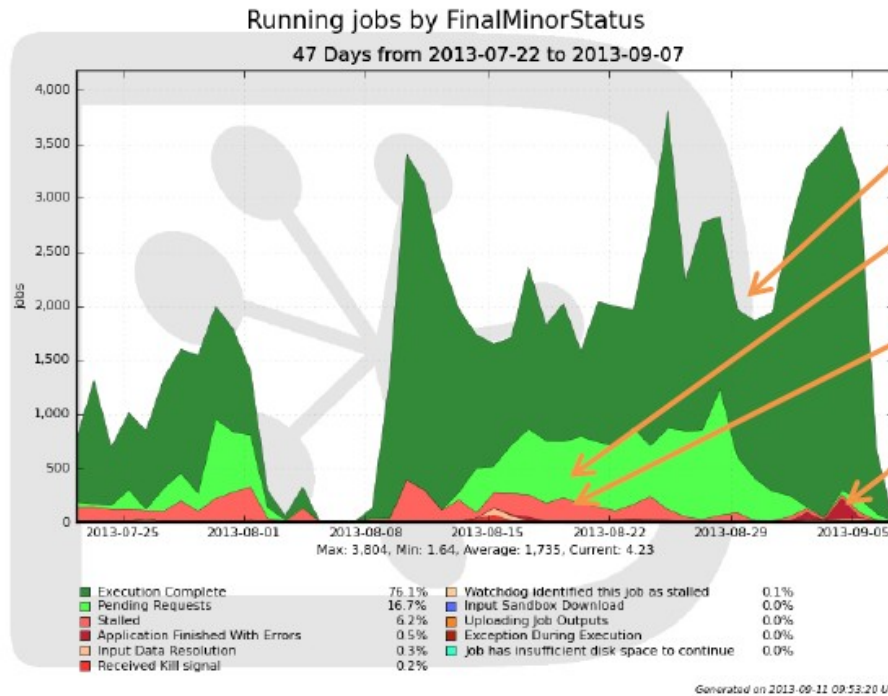
Generated on 2013-03-26 05:14:24 UTC

Second MC Production Campaign

- July 23 – September 8, 2013
- Simulation and reconstruction, with background mixing
→ mdst data
- 630k jobs,
700 kHS*days
- 560M events,
8.5 TB of output data
- ~10% → 1% failure rate:
site configuration/
downtime, proxy
expiration, server load,
human errors
- No crash of offline software



Issues and Solutions



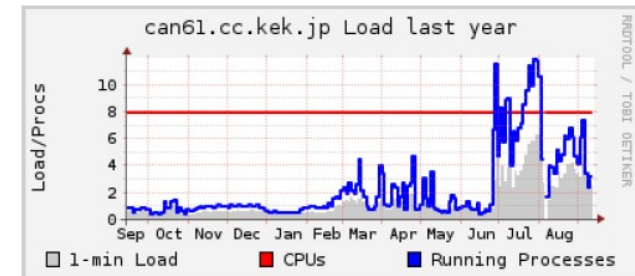
gLite bug

Failover (mainly SE problem)

Proxy expiration

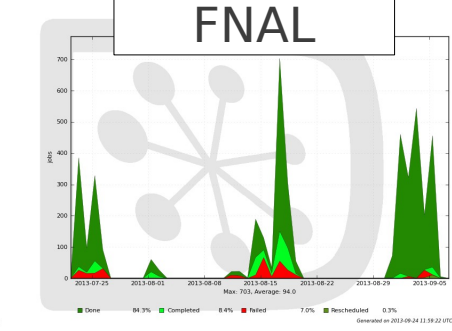
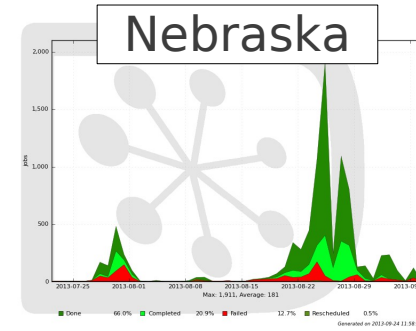
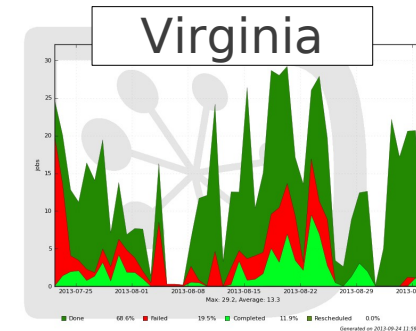
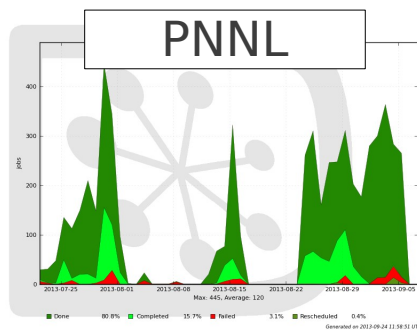
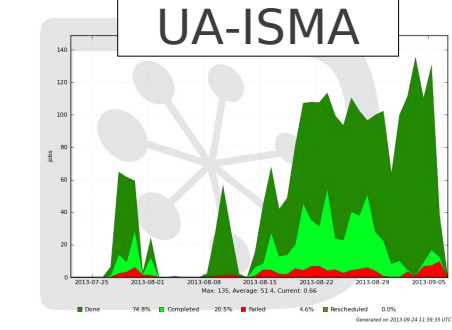
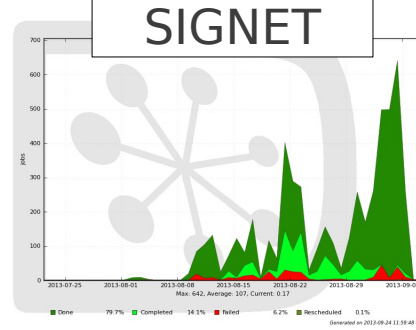
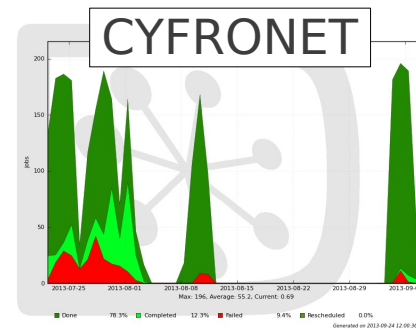
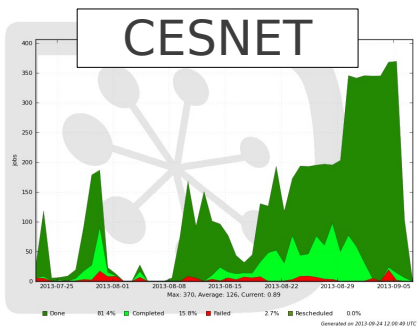
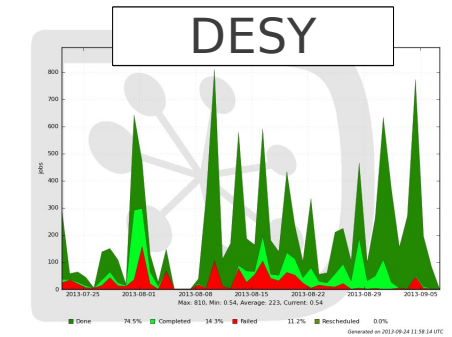
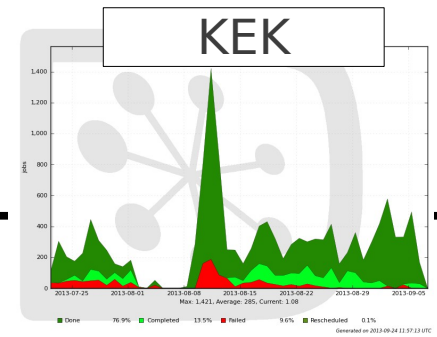
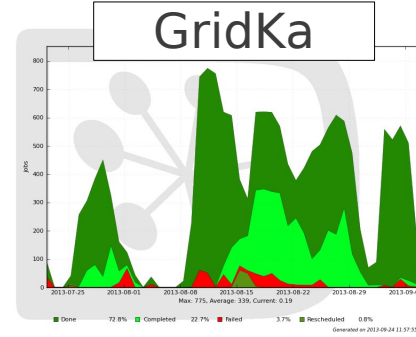
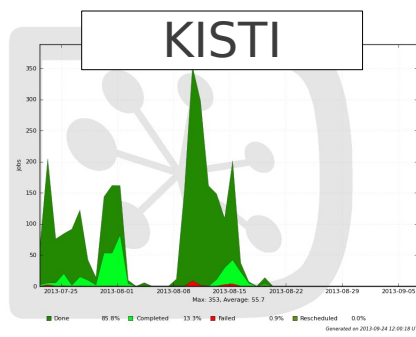
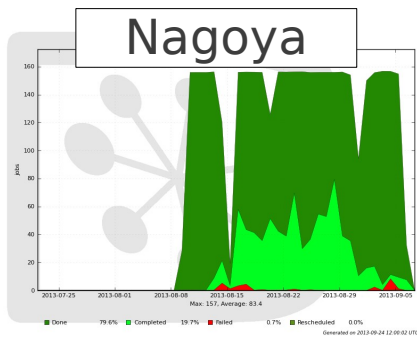
AMGA load

DIRAC load

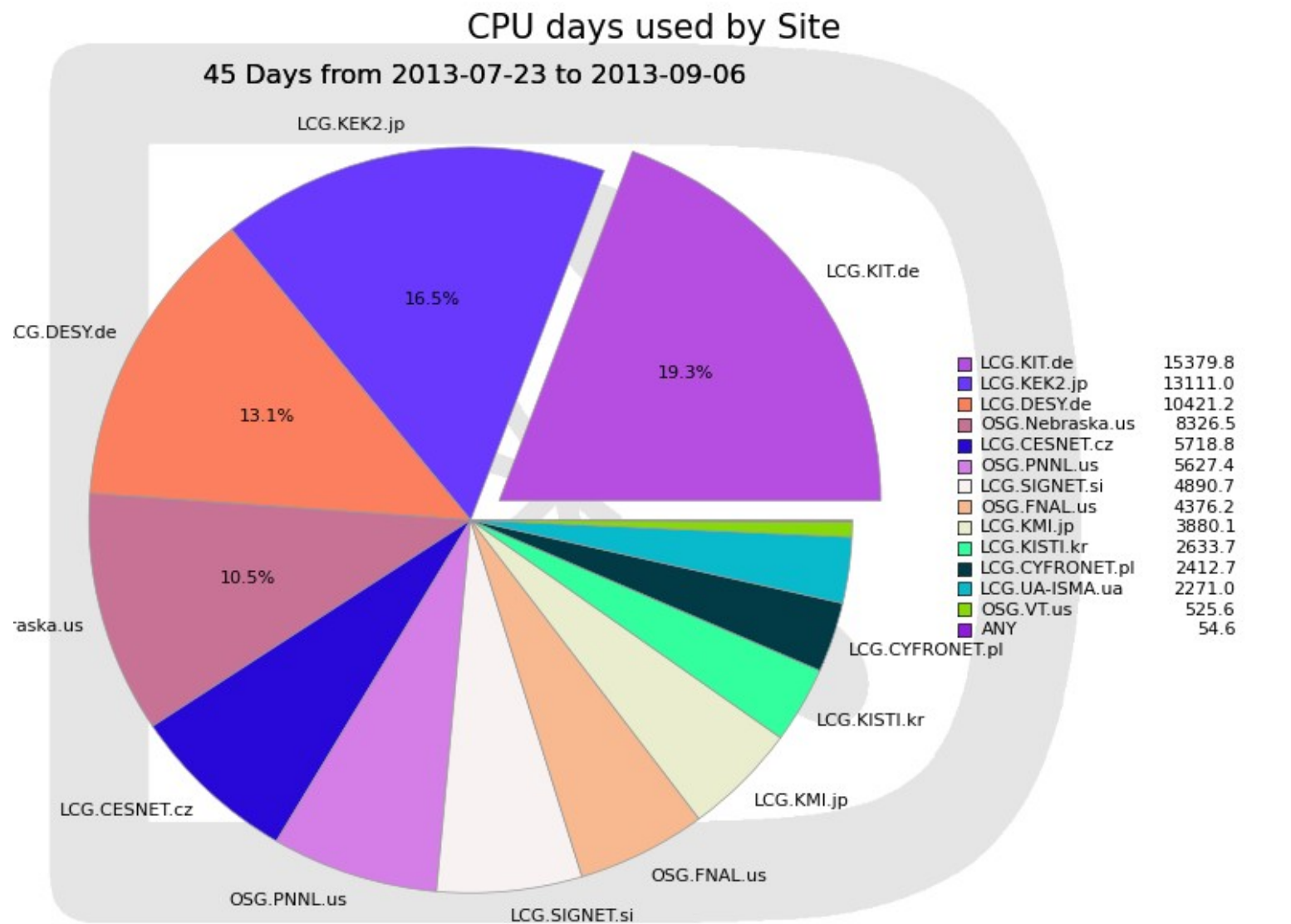


- Failover mechanism for output storage (Increased number of pool accounts on KEK SE)
- Proxy lifetime extended to 168 hours
- Communication frequency with DIRAC decreased and services distributed over more nodes

Contributing Sites



Contributing Sites



Generated on 2013-09-24 11:38:20 UTC

Shifters

Asia

JST 09:00-17:00

CEST 02:00-10:00

PDT 17:00-01:00

Europe

JST 17:00-01:00

CEST 10:00-18:00

PDT 01:00-09:00

USA

JST 01:00-09:00

CEST 18:00-02:00

PDT 09:00-17:00

24x7 production relied on shifters in three time slots

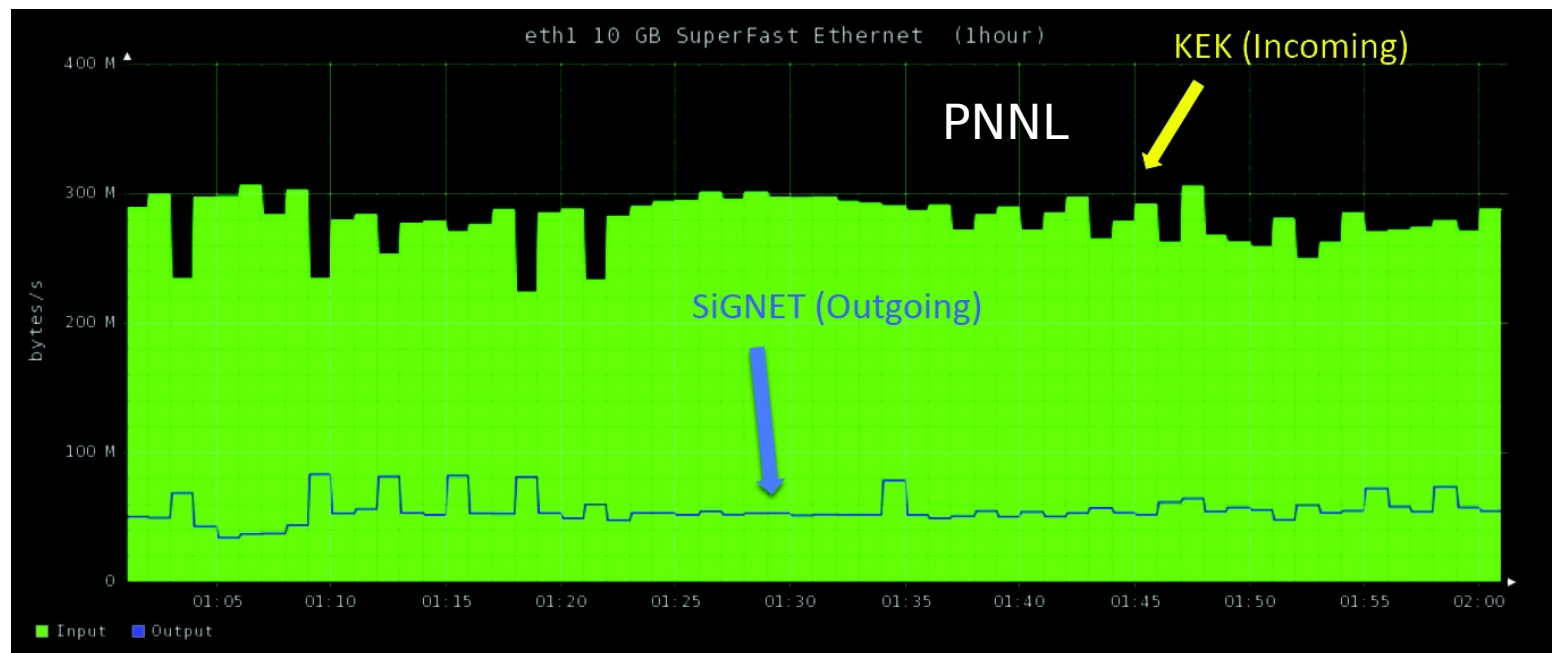
- Submit jobs, monitor progress, resubmit failed jobs
- Keep record in log book
- Hand over with SeeVogh/EVO

- Sergey Barannik, Matt Barrett, Craig Bookwalter, Marko Bracko, Kihyeon Cho, Kamal Dutta, Rafal Grzymkowski, Yanliang Han, Takanori Hara, Kiyoshi Hayasaka, Andreas Heller, Nam Gyu Kim, Peter Kodys, TK, Radek Ludacka, Hideki Miyake, Kamal Jyoti Nath, Geunchul Park, Malachi Schram, Martin Sevier, Oleksandr Sobolev, Michael Steder, Wenjing Wu



















Data Challenge

- Network connection between sites is essential
 - Raw data from KEK to PNNL
 - Mdst data and MC between sites world wide
- ➔ Transfer tests between different sites in May 2013 with FTS2 server at GridKa

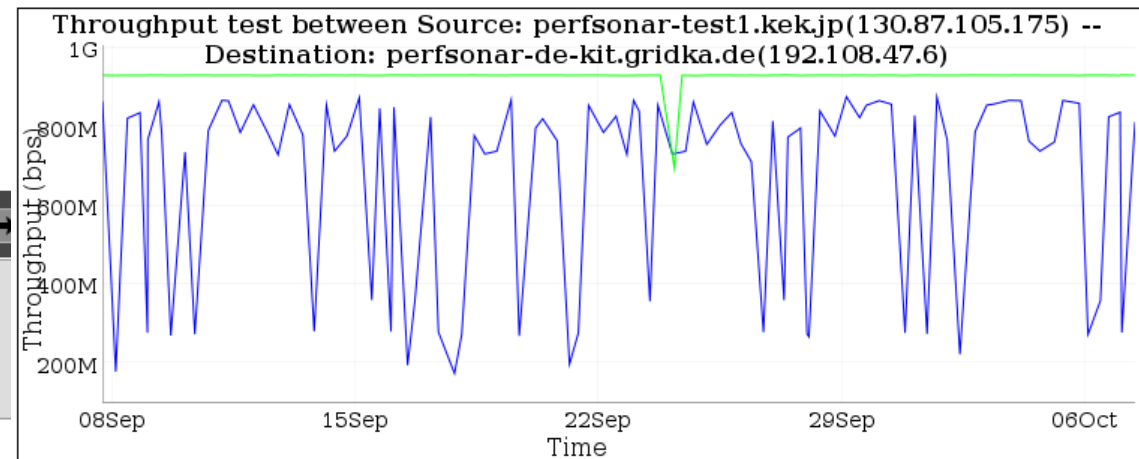
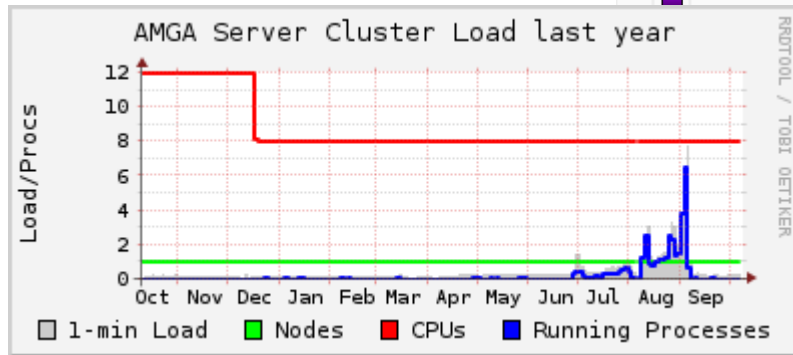


Monitoring

BelleII AMGA Server (belle2-amga)

Host	Status	Services	Actions
amga.pnl.gov	UP	1 OK	   
belledh.kisti.re.kr	UP	1 OK	   
can68.cc.kek.jp	UP	8 OK	   
can69.cc.kek.jp	UP	4 OK	   



Site	CE	Status	PilotJobEff (%)	PilotsPerJob	Waiting	Scheduled	Running	Done	Aborted
+		Multiple	Bad	0	0	0	0	0	20
+	LCG.CESNET.cz	Multiple	Fair	62.54	1	0	150	42	115
+	LCG.CYFRON...	Multiple	Good	88.31	1	0	0	1133	150
+	LCG.DESY.de	Multiple	Fair	76.73	1	0	1051	225	387
+	LCG.KEK2.jp	Multiple	Good	85.09	1	155	224	1128	264
+	LCG.KISTI.kr	Multiple	Poor	41.68	1	64	16	4005	5716
+	LCG.KIT.de	Multiple	Good	99.29	1.01	2042	558	4292	50
	LCG.KMI.jp	ncream01.hepl...	Good	100	1	80	159	1808	0
	LCG.SIGNET.si	creamce.ijs.si	Poor	26.55	1	100	53	1	426
	LCG.UA-ISMA.ua	gl-ce.isma.khar...	Bad	14.85	1.12	4	91	174	1543
		Multiple	Fair	74.63	1.01	133	864	612	1677
		Multiple	Poor	53.92	1	1484	491	916	13949
	osg.pnl.gov		Bad	21	1	0	0	30	15281
	valiant.phys.vt.e...		Poor	51.57	1	66	80	58	417







The Happy Face Project
Version 3, rev. 913M


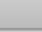
XML 07. Oct 2013 17:35


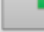
00:15

Banned sites  

Batch System  

Job CPU efficiency  

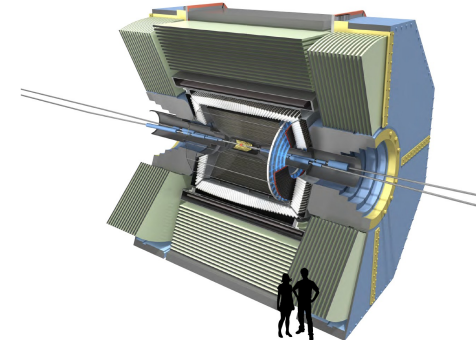
Ganglia  

Infrastructure  

Summary



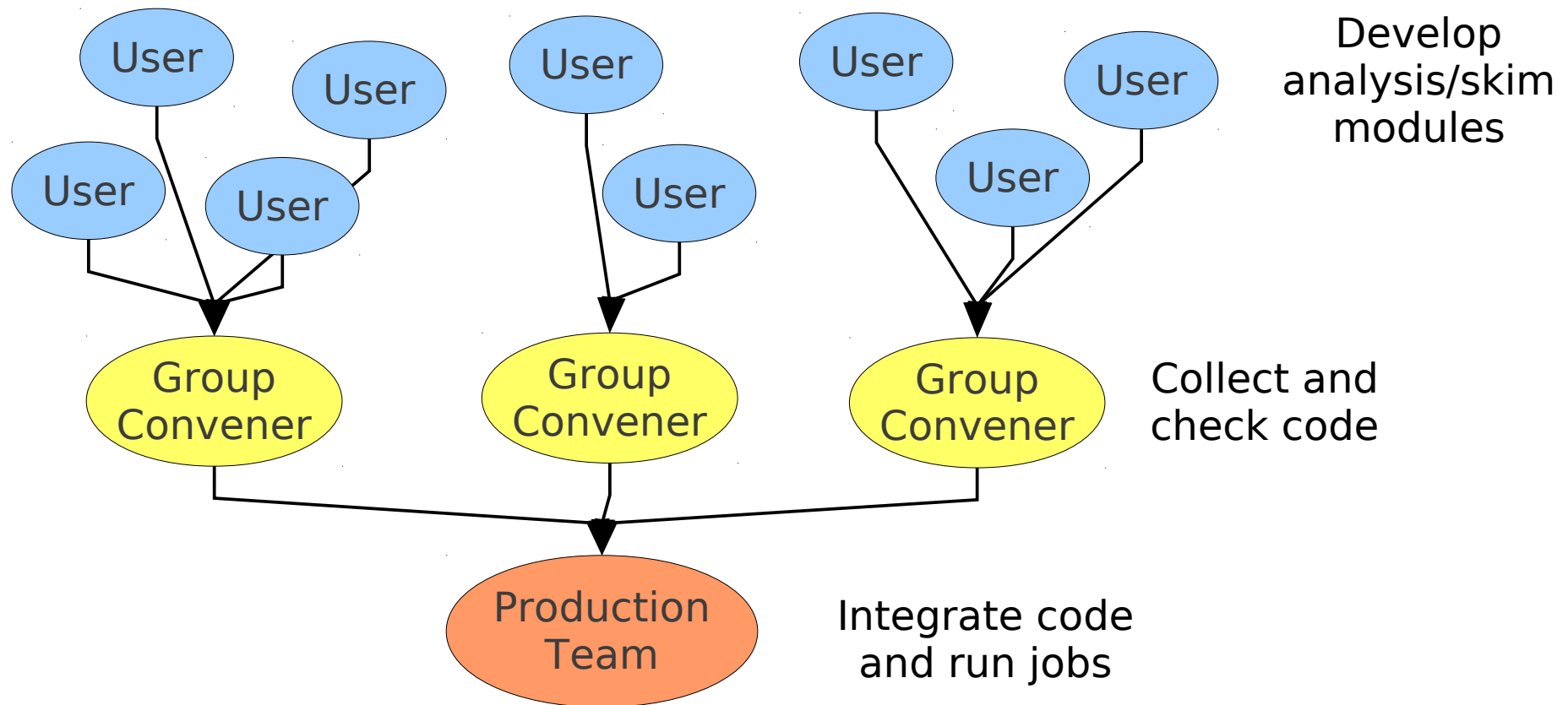
- Belle II will search for New Physics with $O(50)$ times more data than current B factories
- ➔ Huge data volume is a challenge for the computing
 - Distributed computing system based on existing technologies and infrastructures
 - Workflow abstraction with projects and datasets
- First two MC production campaigns this year
 - ✓ Belle II distributed computing system works!
 - ✓ Bottlenecks and issues identified
 - ➔ Many thanks to technology and resource providers!
- ➔ Next steps:
 - MC campaign with more (cloud) sites
 - Further automatize and harden the system
 - Exercise user analysis on the grid



Backup

Organized Analysis

- Problem: inefficient resource usage by many users
- ➔ Limit resources per user, but maintain free access to data
- Offer high-performance organized analysis as a service



Tasks of Computing Facilities

Non-grid Sites	Grid Sites	KEK	
		Storage and Processing of Raw Data	Main Center
	Experiment-specific Services	Experiment-specific Services	
	Monte-Carlo Production	Monte-Carlo Production	Grid
	Data Analysis	Data Analysis	
Ntuple-level Analysis	Ntuple-level Analysis	Ntuple-level Analysis	Local Resources
User Interface	User Interface	User Interface	

(Commercial) Cloud Computing

- Resource demands vary with time
- Fair-share can solve this issue only to some extent
- Cloud computing allows to buy resources on demand
 - Well suited to absorb peaks in varying resource demand

