

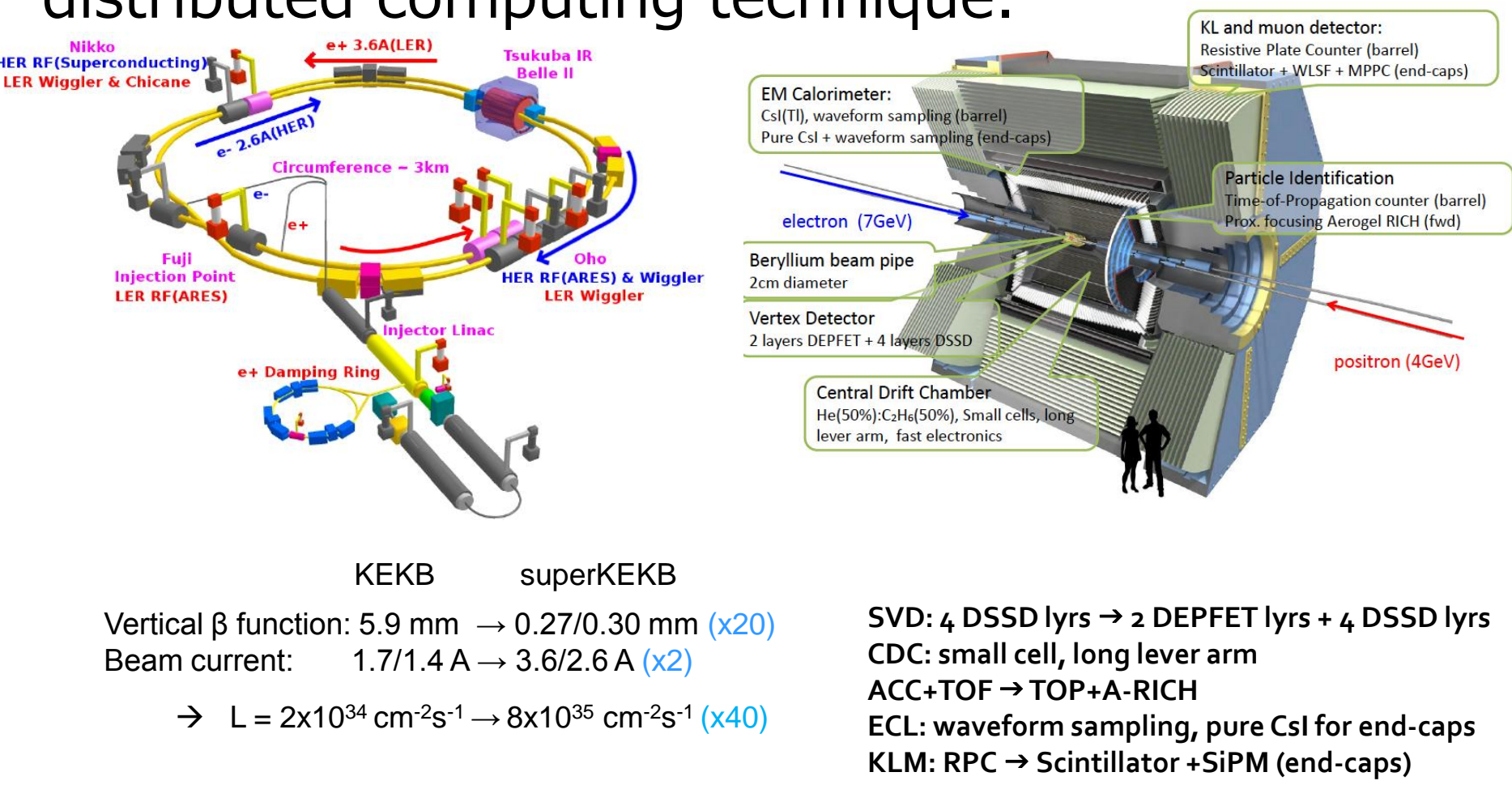
Monitoring system for the Belle II distributed computing (Poster ID: 314)



K.Hayasaka (KMI, Nagoya U.), Y.Kato (KMI, Nagoya U.), T.Hara (KEK), H.Miyake (KEK), I.Ueda (U. Tokyo/KEK) for Belle II computing group

Introduction

- Belle II experiment is a next-generation B-factory at KEK in Japan, which will start for physics run w/o (with) vertex detector in 2017 (2018), where 50ab⁻¹ data sample will be collected for 10 years, which corresponds to about 5x10¹⁰ B \bar{B} -pair events.
- We roughly need to handle 1MHS06 CPU resources, 100PB storage for one set of raw data and 100PB one for MC/analysis data, finally.
- In order to utilize these huge resources, we adopt distributed computing technique.

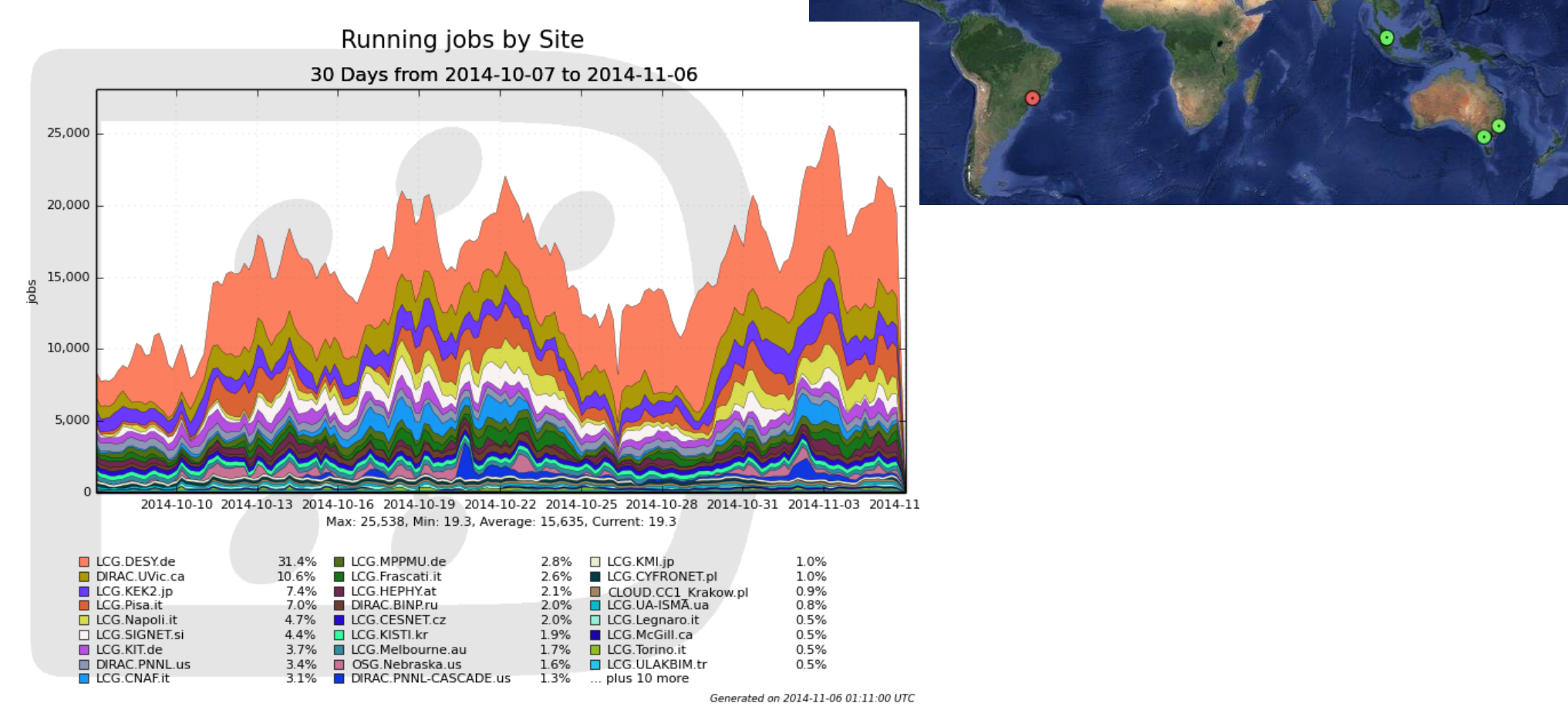


KEKB superKEKB
Vertical β function: 5.9 mm \rightarrow 0.27/0.30 mm (x20)
Beam current: 1.7/1.4 A \rightarrow 3.6/2.6 A (x2)
 $\rightarrow L = 2 \times 10^{34} \text{ cm}^{-2} \text{ s}^{-1} \rightarrow 8 \times 10^{35} \text{ cm}^{-2} \text{ s}^{-1}$ (x40)

SVD: 4 DSSD lyrs \rightarrow 2 DEPFET lyrs + 4 DSSD lyrs
CDC: small cell, long lever arm
ACC+TOP \rightarrow TOP+A-RICH
ECL: waveform sampling, pure CsI for end-caps
KLM: RPC \rightarrow Scintillator+SiPM (end-caps)

Belle II computing

- Belle II has adopted DIRAC as the distributing computing software framework, which can handle grid, cloud and local cluster resources. (<http://diracgrid.org/>)
- CVMFS is used to provide Belle II software and libraries.
- At the present, around 40 sites participate (LCG, OSG, HPC, cloud and traditional cluster) and 25K concurrent jobs are handled at peak.



Troubles during operation

- We observed different types of troubles so far:
 - CVMFS
 - CVMFS is not running.
 - CVMFS is running with some error.
 - When some file is accessed, I/O error appears.
 - Files on CVMFS are not up-to-date.
 - Computing Element (CE)
 - CE is down.
 - CE is alive but batch job system is down.
 - Server Certification on CE is expired.
 - Worker Node (WN)
 - WN has hardware trouble, mostly, HDD failure.
 - WN does not have enough free space for HDD.
 - Required package is not installed.
 - Clock is not adjusted.
 - Central DIRAC servers
 - Due to the heavy load, server can not reply.

Monitoring

- To maximize the operation efficiency of the resources = increase of the resources
- To do so, it is important to find any problem faster. \rightarrow Effective monitoring system is necessary!

Two kinds of monitoring:

Active way Where is a bug?

• Make some action and check the result. \rightarrow Submit test job to check WN condition, or to check CE reaction. Try to open the DIRAC ports.

Passive way Cloudy or sunny?

• Check the result of the action for the operation. \rightarrow Pilot/Job behavior, log of pilot/job or information inside DIRAC.

Of course, DIRAC equips some passive-way monitoring systems such as job accounting, etc...

We are developing the active-way monitor!

to check...

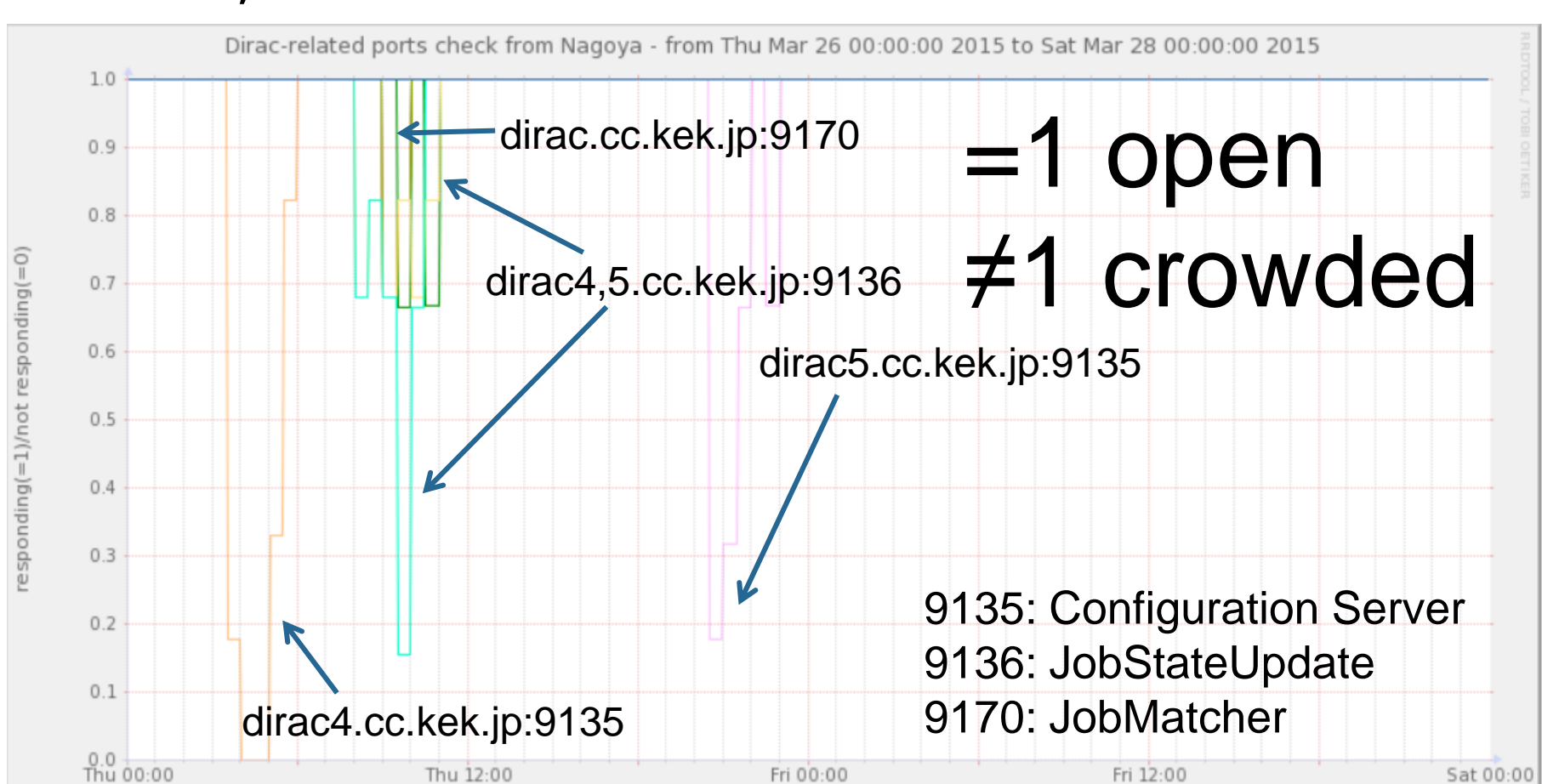
- Worker node status by submitting diagnosing job.
- CE health by test submission.
- DIRAC load by port check.

These results are stored in DB and summarized/visualized in the web pages.

We are also developing the passive-way monitor. \rightarrow See Poster#337.

DIRAC port checker

- When DIRAC server experiences a heavy load, it sometimes gets down. But, before it, we can observe foretastes.
- One is port accessibility. DIRAC assigns ports from 9130 to 9200 function by function. When some port is crowded, it becomes difficult to open it.
- We periodically open the ports and check if the port is accessible or not. This can be job failure reason, too.



Even if the server gets down, we can reboot quickly since mail notification system to the operator is equipped.

SiteCrawler

- We have developed diagnosing job submitting system to check CPU, size of memory, OS, kernel version, installed rpm packages, clock adjustment, free disk space, http proxy accessibility, cvmfs status, Belle II lib. files' md5sum and make a test execution for Belle II software on a WN.
- Every one hour, one job is submitted to each site.

click! You can see information for the checked WNs.

site	hostname	OS	kernel	mem	disk	release	CPU usage	last updated
LCG-INFN01	inf01-01-101-01	Scientific Linux release 6.3 (Carbon)	2.6.32-0504.3.4.el6.x86_64	4096	160	6.10-0.0	100%	2015/04/02 02:55:45 UTC
LCG-INFN02	inf02-01-101-01	Scientific Linux release 6.3 (Carbon)	2.6.32-0504.3.4.el6.x86_64	4096	160	6.10-0.0	100%	2015/04/02 02:55:45 UTC

WN status summary for LCG.UAKBIM.tr

worker node	CPU	mem	OS	kernel	disk	release	CPU usage	last updated
LCG-INFN01	AMD Opteron(TM) Processor 6174	4096	Scientific Linux release 6.3 (Carbon)	2.6.32-0504.3.4.el6.x86_64	160	6.10-0.0	100%	2015/04/02 02:55:45 UTC
LCG-INFN02	AMD Opteron(TM) Processor 6174	4096	Scientific Linux release 6.3 (Carbon)	2.6.32-0504.3.4.el6.x86_64	160	6.10-0.0	100%	2015/04/02 02:55:45 UTC

Before submitting cordial jobs, we can check if Belle II software can be actually executed or not, from various viewpoints. In particular, trouble relating cvmfs is often observed.

Summary

- Belle II experiment requires the huge computing resources to process its huge data sample.
 - Belle II operates the distributed computing system using DIRAC.
 - But, having DIRAC is not an end of the story.
 - Through the operation, many kinds of troubles have been experienced and they reduce the effective amount of the CPU resources.
 - To reduce the "dead time" of the CPU resources, it is important to find any problem faster and we need a good monitoring system.
 - Not only a standard monitoring system equipped in DIRAC but also our own system optimized to detect the trouble we faced are necessary.
 - We have been developing the monitoring system for Belle II computing.
 - Here, we introduce the "active-way" monitorings.
 - They help to find problems on WN, CE and central DIRAC servers.
- ("passive-way" monitoring is discussed at Poster#337)

CE test job submitter

- When DIRAC fails 70% of job submission to some site or SiteCrawler does not run for 6 hours, we submit a simple job to the site, outside DIRAC, i.e., using glite-ce-job-submit, by cron.
- The job executes just "cat /proc/cpuinfo".
- CE's reaction is recorded in DB.

CE Job Submission test result

site name	CE	queue	status	last updated time
LCG-INFN01	ce08-01-101-01	cream-1f-belle	IDLE	2015/04/02 03:40:14 UTC
LCG-INFN02	ce08-01-101-01	cream-1f-belle	IDLE	2015/04/02 03:40:14 UTC
LCG-INFN03	ce08-01-101-01	cream-1f-belle	IDLE	2015/04/02 03:40:14 UTC

CE Job Submission test result on LCG.Pisa.it

site name	CE	queue	status	job id	last updated time
LCG-INFN01	ce08-01-101-01	cream-1f-belle	submission failed	None	2015/04/02 02:55:45 UTC
LCG-INFN02	ce08-01-101-01	cream-1f-belle	submission failed	None	2015/04/02 02:55:45 UTC
LCG-INFN03	ce08-01-101-01	cream-1f-belle	submission failed	None	2015/04/02 02:55:45 UTC

You can get the failure reason by clicking "log". In this case: FATAL - Received NULL fault; the error is due to another cause: FaultString=[java.lang.NullPointerException] - FaultCode=[SOAP-ENV:Server] - FaultSubCode=[SOAP-ENV:Server]

When pilots/jobs do not run on some site, we can judge whether it is due to DIRAC or CE.

Future Plan

- More sophisticated/automated monitoring system
 - Combining various information collected by sub-monitoring systems and diagnose what happens automatically.
 - Provide "human-readable" operation summary:
 - Automated notification system:
 - To LCG site: submit GGUS tickets
 - To other site: similar system (redmine?)
 - Also monitoring system for SE
 - accessibility
 - read/write
 - checksum
- Belle II monitoring system pursues the most efficient operation for Belle II computing by quickly finding the problems.