# Belle II distributed computing

## Pavel Krokovny
Novosibirsk State University and
Budker Institute of Nuclear Physics
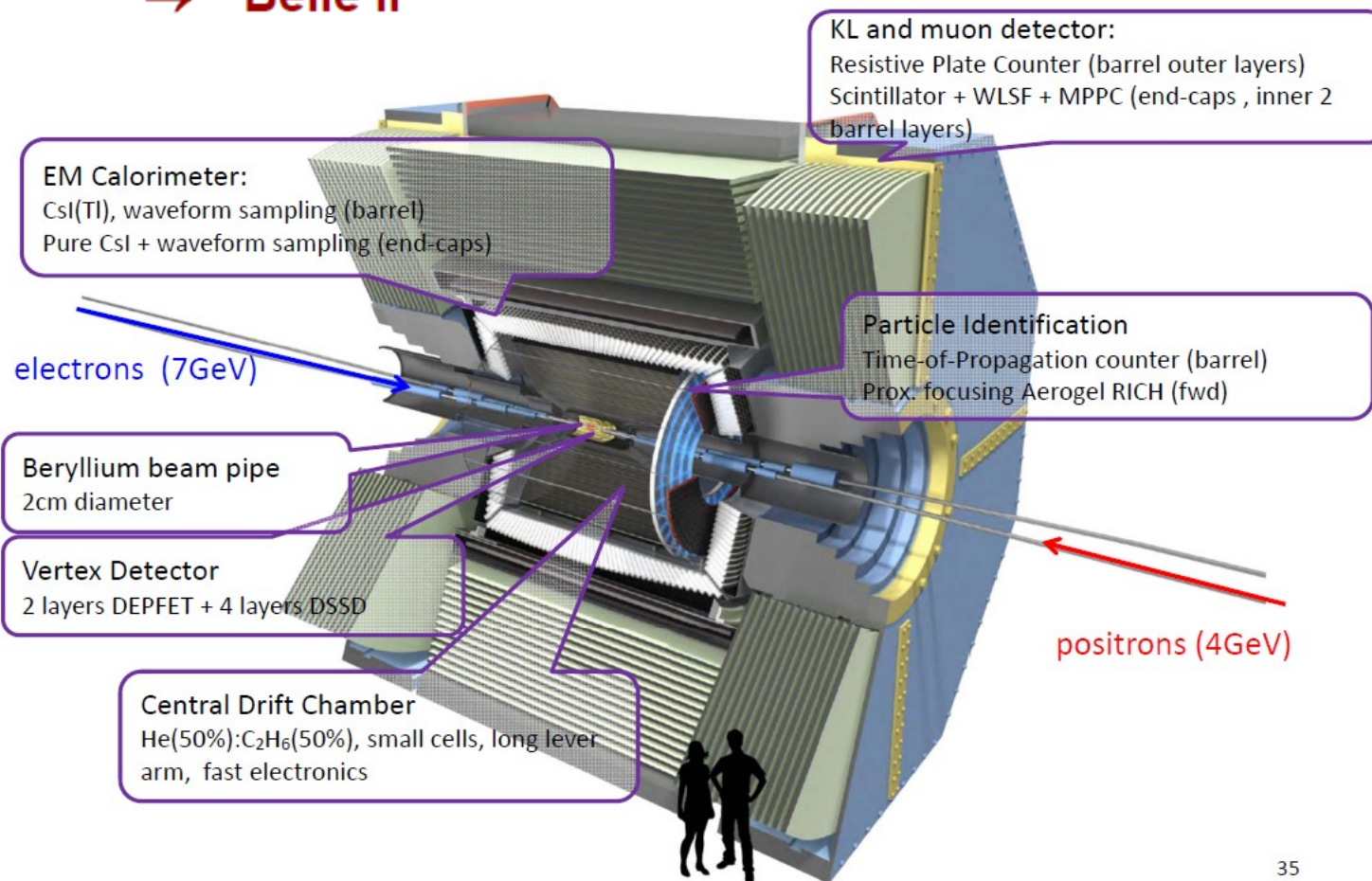
Outline:
- Belle II experiment
- Computing model
- Resource estimation
- Production system
- Computing in BINP / NSU

ACAT 2014
– bridging disciplines
1-5 September 2014
Prague, Czech Republic
Europe/Prague timezone

# The Super B factory project

|              | KEKB                  |     | SuperKEKB          |         |
|--------------|-----------------------|-----|--------------------|---------|
| Luminosity:  | $2.1 \times 10^{34}$  | →   | $8 \times 10^{35}$ | (x 40)  |
| Total Data:  | $1\ ab^{-1}$          | →   | $50\ ab^{-1}$      | (x 50)  |
| Detector:    | Belle                 | →   | Belle II           |         |

KL and muon detector:
Resistive Plate Counter (barrel outer layers)
Scintillator + WLSF + MPPC (end-caps , inner 2 barrel layers)

EM Calorimeter:
CsI(Tl), waveform sampling (barrel)
Pure CsI + waveform sampling (end-caps)

Particle Identification
Time-of-Propagation counter (barrel)
Prox. focusing Aerogel RICH (fwd)

electrons (7GeV)

Beryllium beam pipe
2cm diameter

Vertex Detector
2 layers DEPFET + 4 layers DSSD

positrons (4GeV)

Central Drift Chamber
He(50%):$C_2H_6$(50%), small cells, long lever arm, fast electronics

Belle II Collaboration:
· 26 countries,
· 97 institutes,
· ~600 scientists

# Luminosity prospect



SuperKEKB Commissioning starts in 2015

9 months/year
20 days/month

Belle
~1 ab$^{-1}$
$2.1 \times 10^{34}$/cm$^2$/s

target integrated luminosity
50 ab$^{-1}$ in 2022

target instantaneous luminosity
$8 \times 10^{35}$/cm$^2$/s

Physics run starts in 2017

| Experiment | Event size | Rate @ Storage | Rate @ Storage | |
|---|---|---|---|---|
| | [kB] | [event/sec] | [MB/sec] | |
| Belle II | 300 | 6,000 | 1,800 | (@ max. luminosity) |
| | | | | |
| ALICE (Pb-Pb) | 50,000 | 100 | 4,000 | |
| ALICE (p-p) | 2,000 | 100 | 200 | |
| ATLAS | 1,500 | 600 | 700 | |
| CMS | 1,500 | 150 | 225 (<~1000) | |
| LHCb | 55 | 4,500 | 250 | |

(LHC experiments : as seen in 2011/2012 runs)

# Software system

- A "framework" system with dynamic module loading, parallel processing, Python steering, and ROOT I/O
- Full detector simulation with Geant4
- C++ 11 and gcc 4.7
- Supporting major Linux distributions: SL, Fedora, Ubuntu, etc

- Code management: Subversion
- Formatting tool: astyle
- Building: scons and buildbot system
- Documentation: Doxygen, Twiki
- Issue tracking: Redmine

# Computing model

The BELLE II Computing model has to accomplish, in a distributed environment, the following main tasks:

- RAW data processing
- Monte Carlo Production
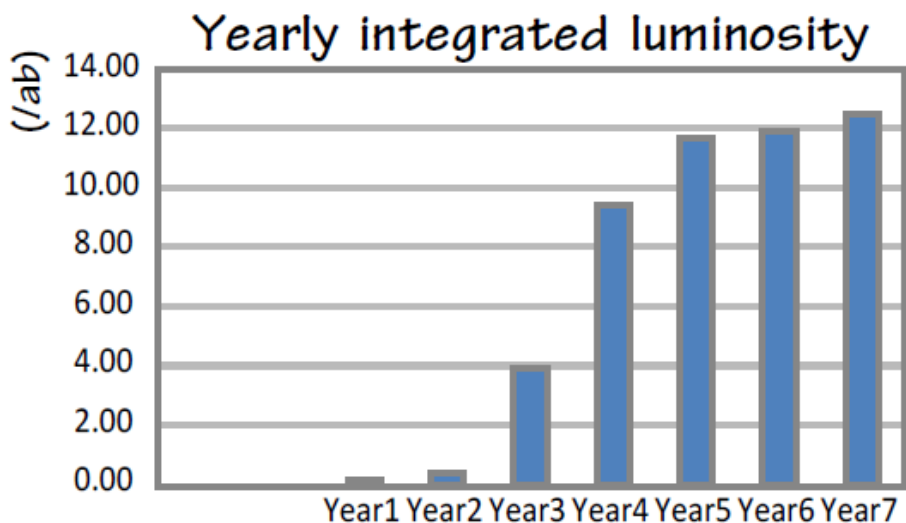- Physics analysis
- Data Storage and Data Archiving

Resource Estimation:

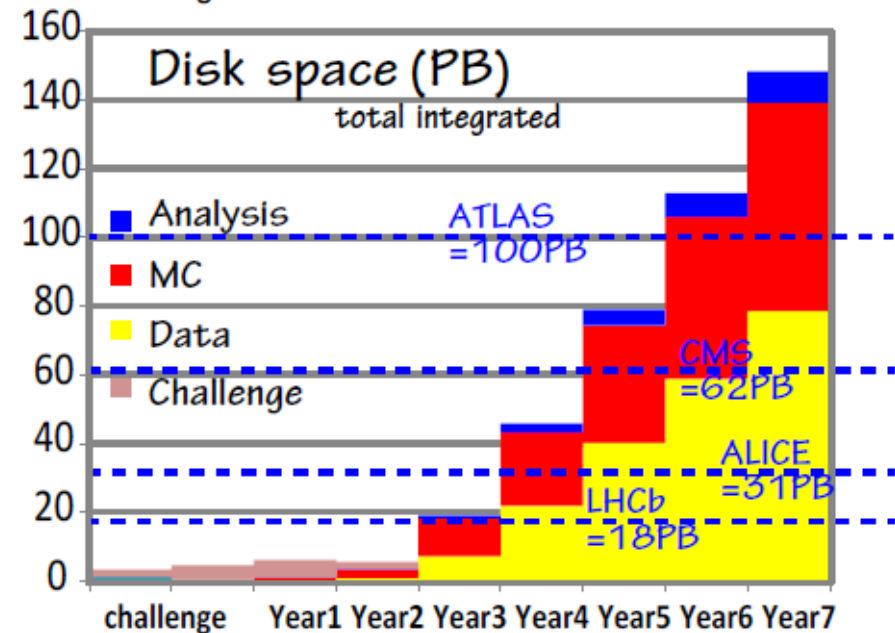- Event size for RAW (mDST) data is 300 (40) Kb
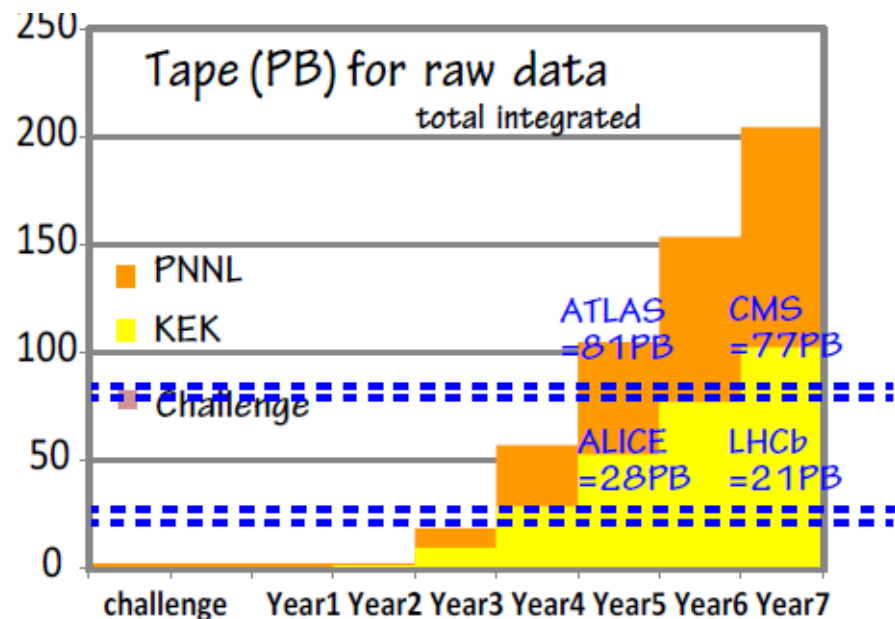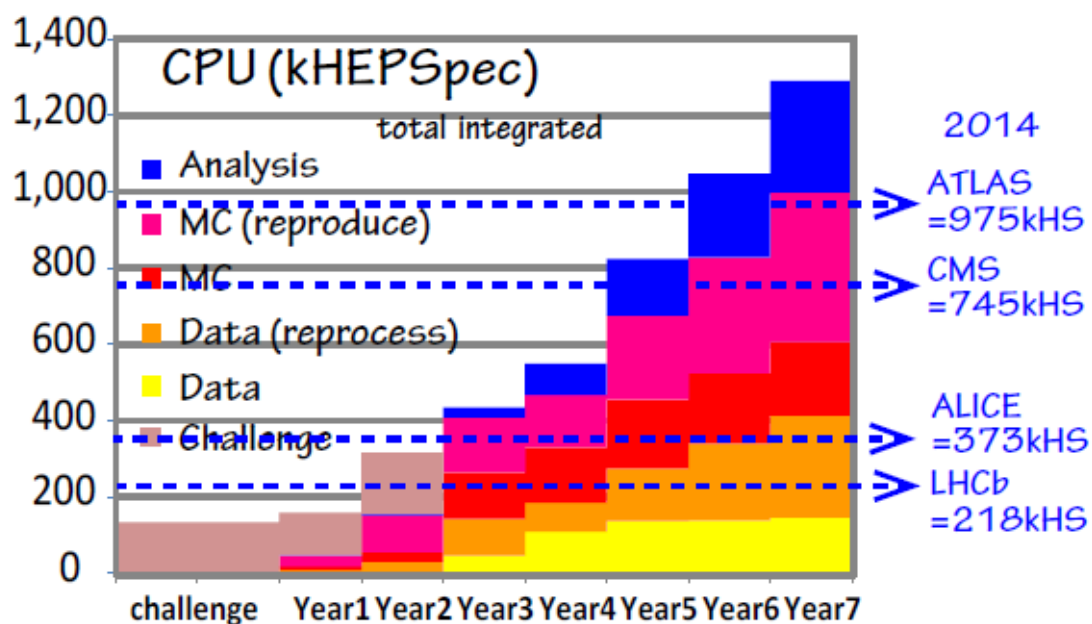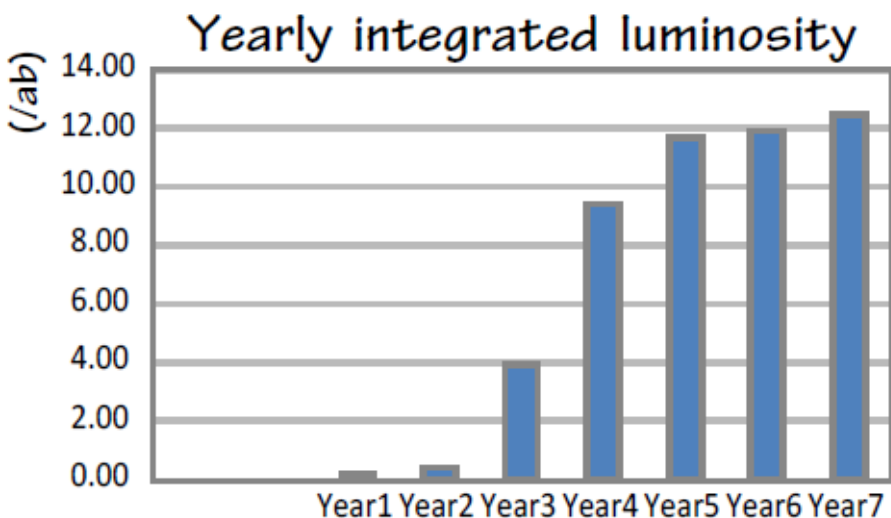- RAW data (MC) processing: 45 (90) HepSPEC*s / event

# Required hardware resources



Belle II

Yearly integrated luminosity

Tape (PB) for raw data — total integrated
- PNNL
- KEK
- Challenge

CPU (kHEPSpec) — total integrated
- Analysis
- MC (reproduce)
- MC
- Data (reprocess)
- Data
- Challenge

Disk space (PB) — total integrated
- Analysis
- MC
- Data
- Challenge

# Comparison of hardware resources

*Belle II*

## Normalized CPU usage by Site
### 76 Weeks from Week 52 of 2012 to Week 23 of 2014

70 kHS (100 kHS @ max)

1st
60M
events

2nd
560M
events

3rd
6200M
events

Max: 66,686, Average: 6,973, Current: 3.54

| | | | | | |
|---|---|---|---|---|---|
| LCG.DESY.de | 18.6% | LCG.CNAF.it | 3.6% | OSG.FNAL.us | 1.8% |
| LCG.KIT.de | 11.9% | DIRAC.PNNL.us | 2.9% | LCG.ULAKBIM.tr | 1.2% |
| LCG.KEK2.jp | 11.6% | LCG.SIGNET.si | 2.8% | OSG.PNNL.us | 1.0% |
| LCG.KMI.jp | 5.3% | LCG.Pisa.it | 2.8% | DIRAC.KrakowCloud.pl | 0.9% |
| LCG.UA-ISMA.ua | 4.8% | DIRAC.UVic.ca | 2.4% | LCG.McGill.ca | 0.6% |
| LCG.CYFRONET.pl | 4.4% | LCG.KISTI.kr | 2.3% | LCG.Torino.it | 0.6% |
| LCG.CESNET.cz | 4.3% | OSG.Nebraska.us | 2.3% | LCG.MPPMU.de | 0.5% |
| DIRAC.BINP.ru | 4.0% | LCG.Frascati.it | 2.2% | LCG.Legnaro.it | 0.5% |
| LCG.Napoli.it | 3.9% | LCG.Melbourne.au | 2.0% | plus 12 more | |

Generated on 2014-06-17 13:07:26 UTC

## 15 countries/regions
## 27 sites (+ 2 non-Belle II sites)

HEPHY (Vienna) and MPPMU (Munich)
joined recently
GRID, Cloud, local cluster
is available

## First official release of MC samples

BB generic decay/continuum
tau pair
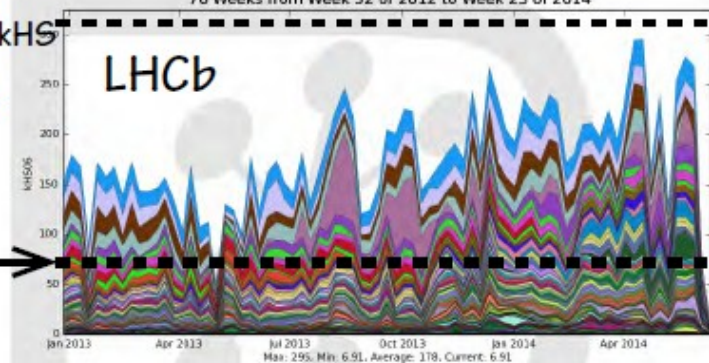(corresponding to $100fb^{-1}$ w/ and w/o BG)

## Trans-pacific / trans-atlantic network data tranfer challenge
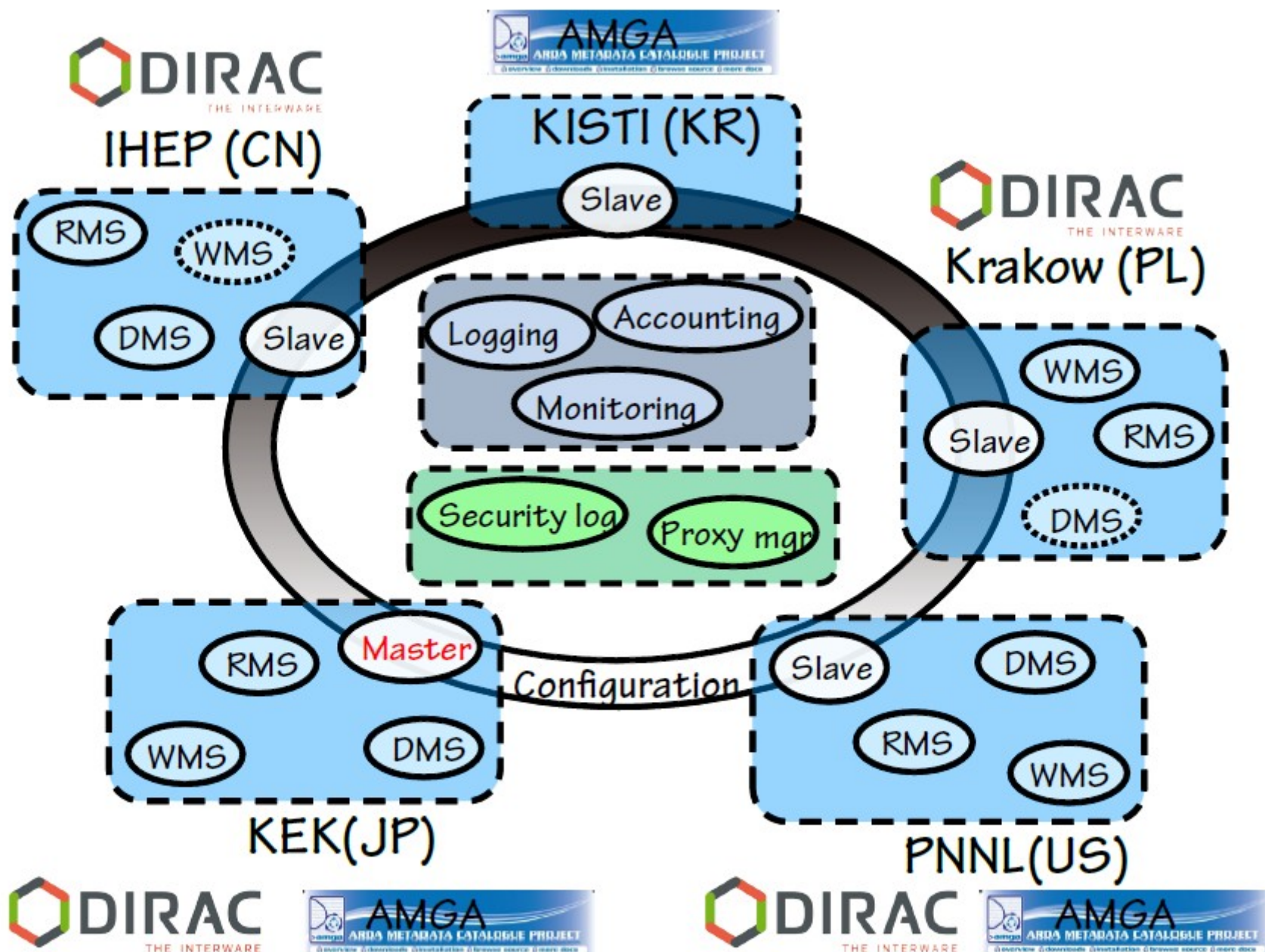
300 kHS

120 sites

Belle II now
70 kHS

### Normalized CPU usage by Site
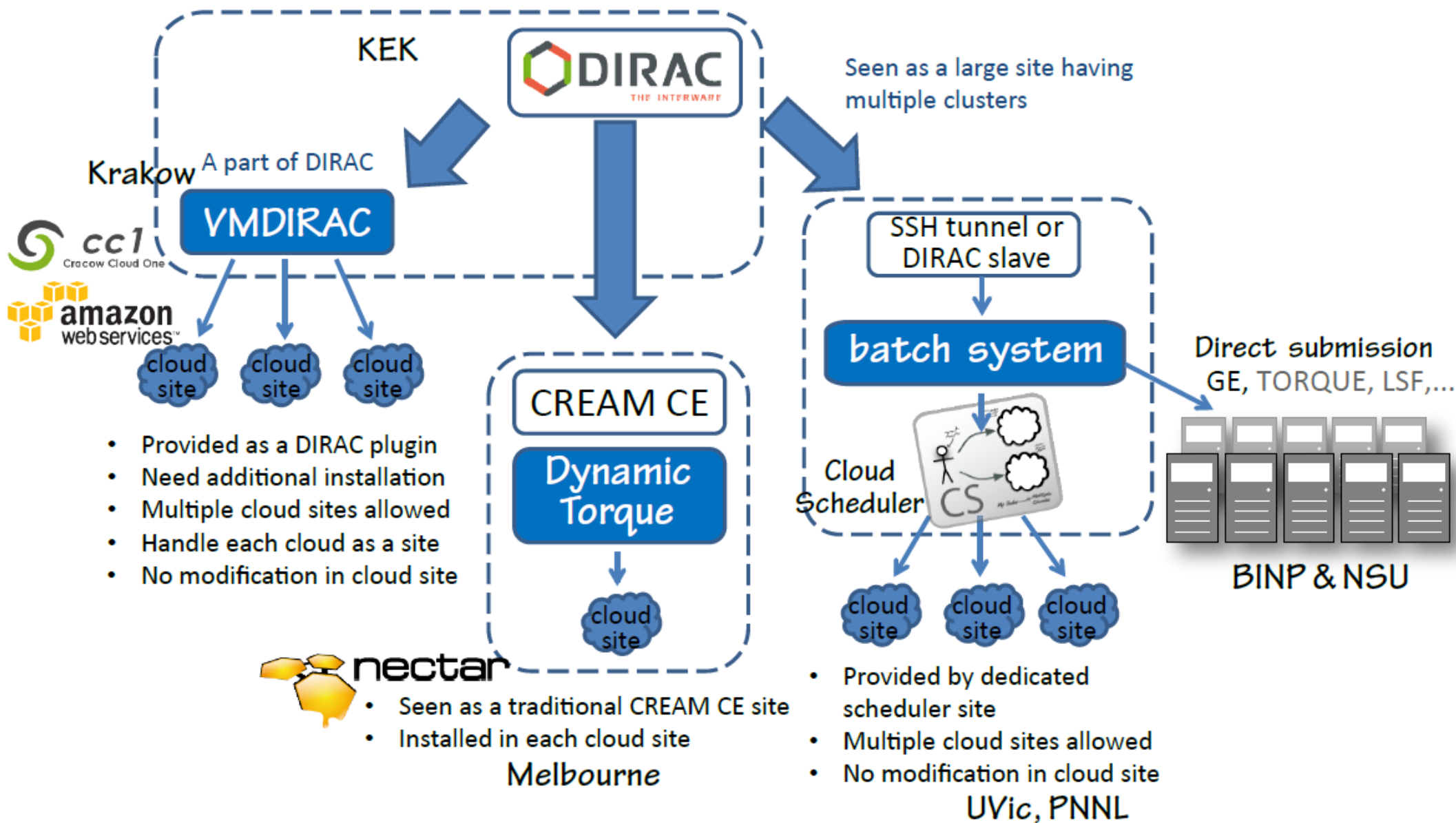76 Weeks from Week 52 of 2012 to Week 23 of 2014

LHCb

Max: 295, Min: 6.91, Average: 178, Current: 6.91

# Production system (core services)

# Production system (cloud usage)

Belle II

**KEK**

DIRAC THE INTERWARE

Seen as a large site having multiple clusters

**Krakow** — A part of DIRAC

cc1 Cracow Cloud One

**VMDIRAC**

amazon web services

cloud site  cloud site  cloud site

- Provided as a DIRAC plugin
- Need additional installation
- Multiple cloud sites allowed
- Handle each cloud as a site
- No modification in cloud site

**CREAM CE**

**Dynamic Torque**

cloud site

nectar

- Seen as a traditional CREAM CE site
- Installed in each cloud site

**Melbourne**

SSH tunnel or DIRAC slave

**batch system**

Cloud Scheduler  CS

Direct submission
GE, TORQUE, LSF,...

cloud site  cloud site  cloud site

BINP & NSU

- Provided by dedicated scheduler site
- Multiple cloud sites allowed
- No modification in cloud site

**UVic, PNNL**

# Computing in BINP & NSU

Experiments
- BINP: CMD-3, SND, KEDR
- International collaborations: Belle II, ATLAS, LHCb

Computing resources
- BINP: 512 CPU cores, 6 Tflops
- Novosibirsk State University: 2432 cores, 29 Tflops
- Supercomputer Center: 30 + 85 Tflops (CPU + GPU)

We want:
- keep the specific computing environment and user's experience
- be like a normal SC user

The solution is:
- run HEP tasks inside virtual machines,
- run VMs inside supercomputer's batch system jobs.

# Key features of the virtualized infrastructure

- Virtualization by KVM
  - included in modern Linux distributions,
  - quite stable,
  - does not require modified Linux kernel.

- VM disk images are located on SC's file system and accessible via InfiniBand.
- Local snapshots of VM images are used on physical nodes thus leaving master images unchanged.
- Input/output data are located at BINP and accessed by VMs via NFS.
- VMs are just regular batch tasks at a supercomputer.
- VMs are started automatically on user's demands.

# Batch System Integration Mechanisms

**NSU**

**NUSC**

PBS Pro Batch System

*STAGE 1*
*Job submission*
*and automated*
*VM group*
*deployment*
*sequence*

**NSC/SCN**

**BINP/GCF**

**BINP**

Orchestration Services

**SGE Batch System**

**Belle II Detector User Group**

Orchestration Services

**SGE Batch System**

SND Detector User Group

Orchestration Services

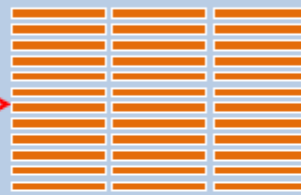**SGE Batch System**

ATLAS Data Analysis Group

# Batch System Integration Mechanisms

**NSU**

**NUSC**

PBS Pro Batch System

*STAGE 1
Job submission
and automated
VM group
deployment
sequence*

Group of VMs
of particular
type created on
demand

*Set of computing
nodes with KVM,
IPoIB & HT support
enabled on
demand*

**NSC/SCN**

**BINP/GCF**

**BINP**

Orchestration
Services

**SGE Batch
System**

**Belle II Detector
User Group**

Orchestration
Services

**SGE Batch
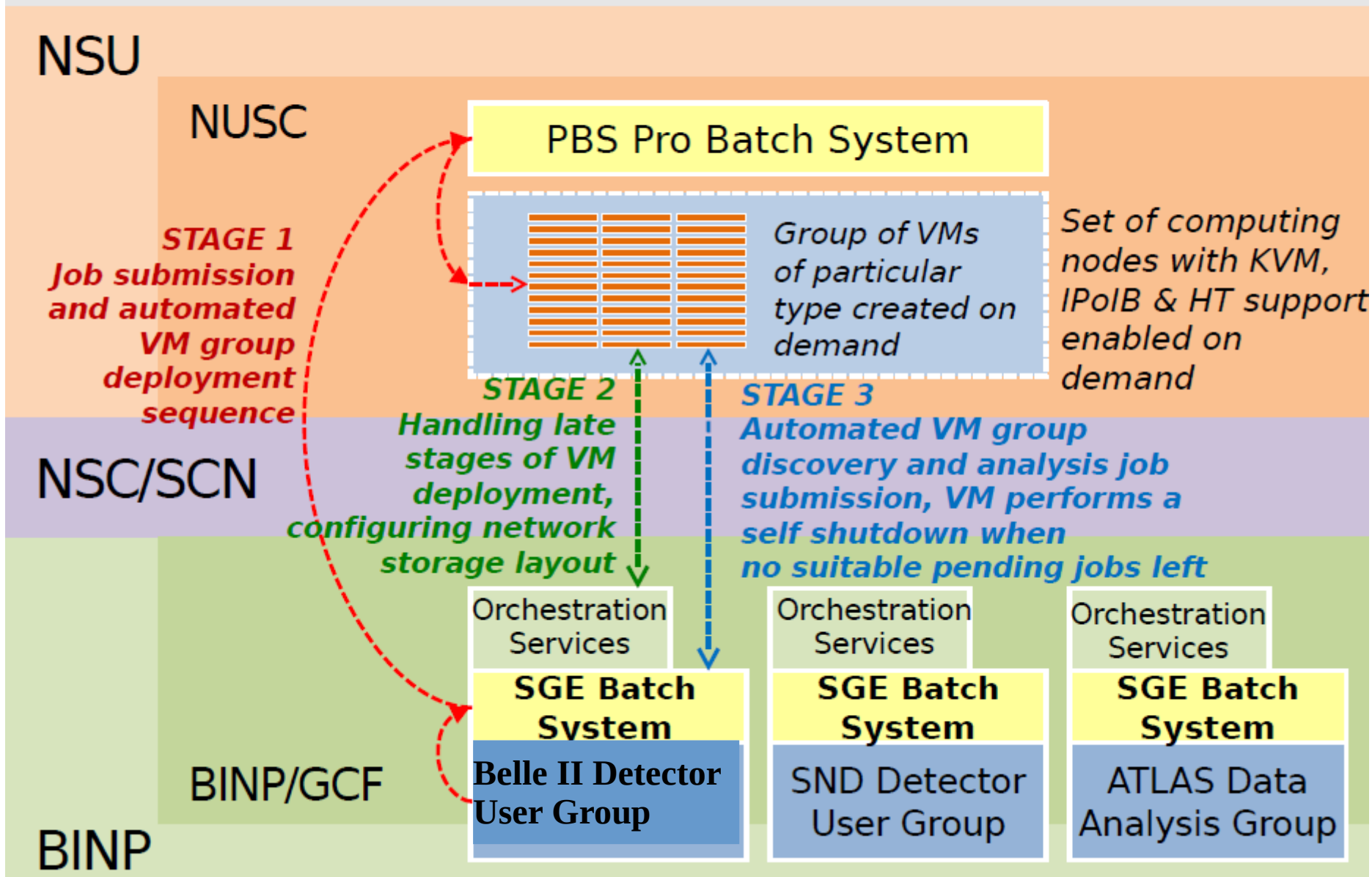System**

SND Detector
User Group

Orchestration
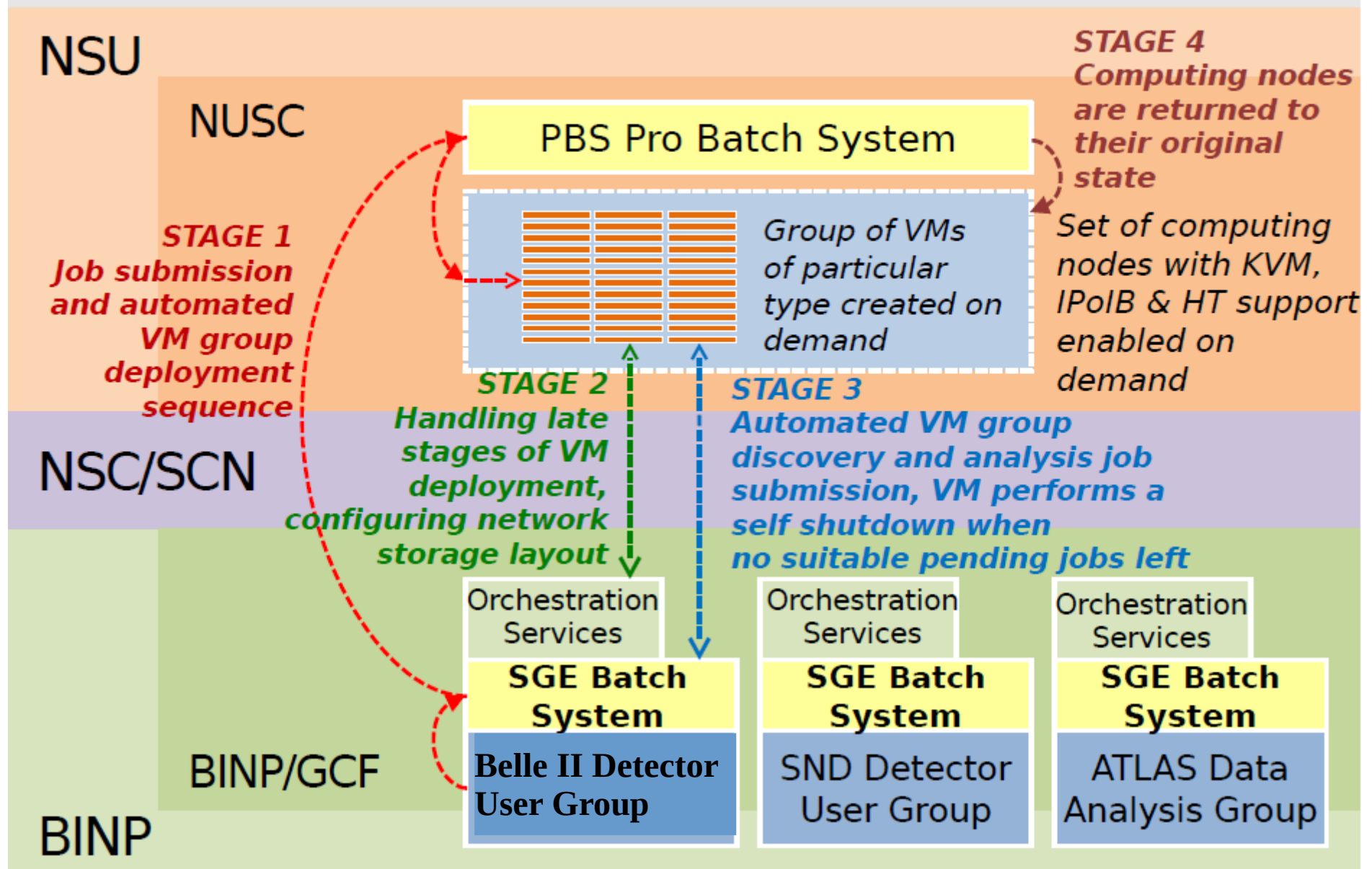Services

**SGE Batch
System**

ATLAS Data
Analysis Group

# Batch System Integration Mechanisms

**NSU**

**NUSC**

PBS Pro Batch System



Group of VMs of particular type created on demand

Set of computing nodes with KVM, IPoIB & HT support enabled on demand

*STAGE 1*
**Job submission and automated VM group deployment sequence**

*STAGE 2*
**Handling late stages of VM deployment, configuring network storage layout**

*STAGE 3*
**Automated VM group discovery and analysis job submission, VM performs a self shutdown when no suitable pending jobs left**

**NSC/SCN**

Orchestration Services

Orchestration Services

Orchestration Services

**SGE Batch System**

**SGE Batch System**

**SGE Batch System**

**Belle II Detector User Group**

SND Detector User Group

ATLAS Data Analysis Group

**BINP/GCF**

**BINP**

# Batch System Integration Mechanisms

**NSU**

**NUSC**

PBS Pro Batch System

*STAGE 1*
**Job submission and automated VM group deployment sequence**

Group of VMs of particular type created on demand

*STAGE 4*
**Computing nodes are returned to their original state**

Set of computing nodes with KVM, IPoIB & HT support enabled on demand

*STAGE 2*
**Handling late stages of VM deployment, configuring network storage layout**

*STAGE 3*
**Automated VM group discovery and analysis job submission, VM performs a self shutdown when no suitable pending jobs left**

**NSC/SCN**

Orchestration Services

Orchestration Services

Orchestration Services

**SGE Batch System**

**SGE Batch System**

**SGE Batch System**

**BINP/GCF**

**Belle II Detector User Group**

SND Detector User Group

ATLAS Data Analysis Group

**BINP**

# Summary

- Belle II experiment has successful experience in distributed computing using GRID, Cloud and local clusters.

- Current CPU power is about 70 kHS (25% of LHCb) and used for MC production.

- BINP successfully uses Virtual Machines technique to run up to 1K jobs simultaneously for 6 experiments on 3 supercomputer sites.
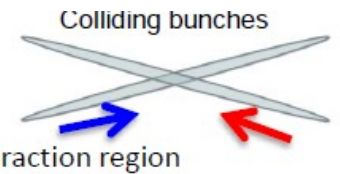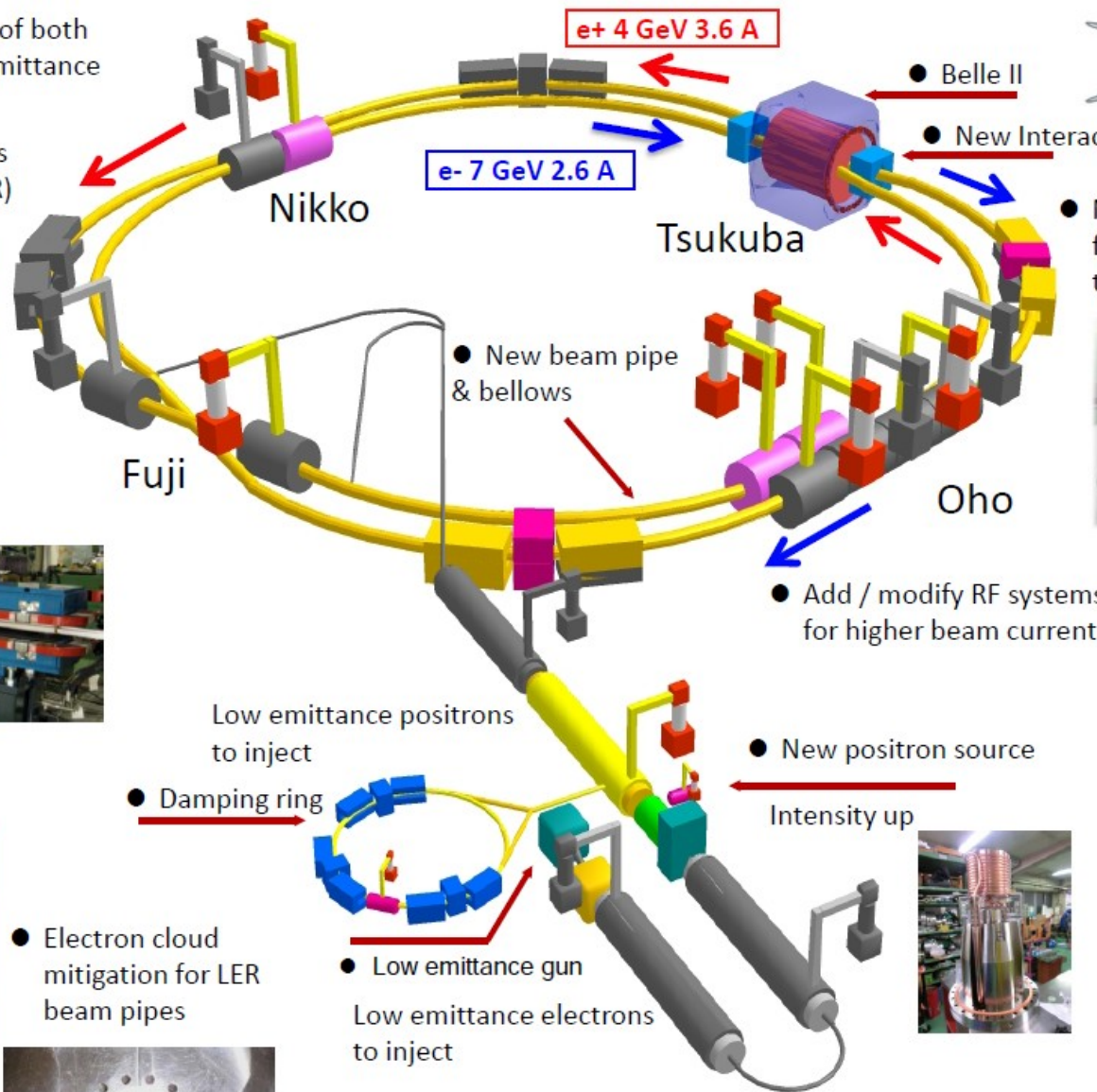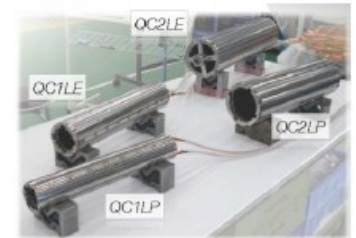
# Backup

# SuperKEKB accelerator

# Raw data distribution

# mDST & MC data distribution

mDST (data) is copied in Asia, Europe, and USA

For the MC data seems to be natural to be
  the similar structure
    better network? in each region
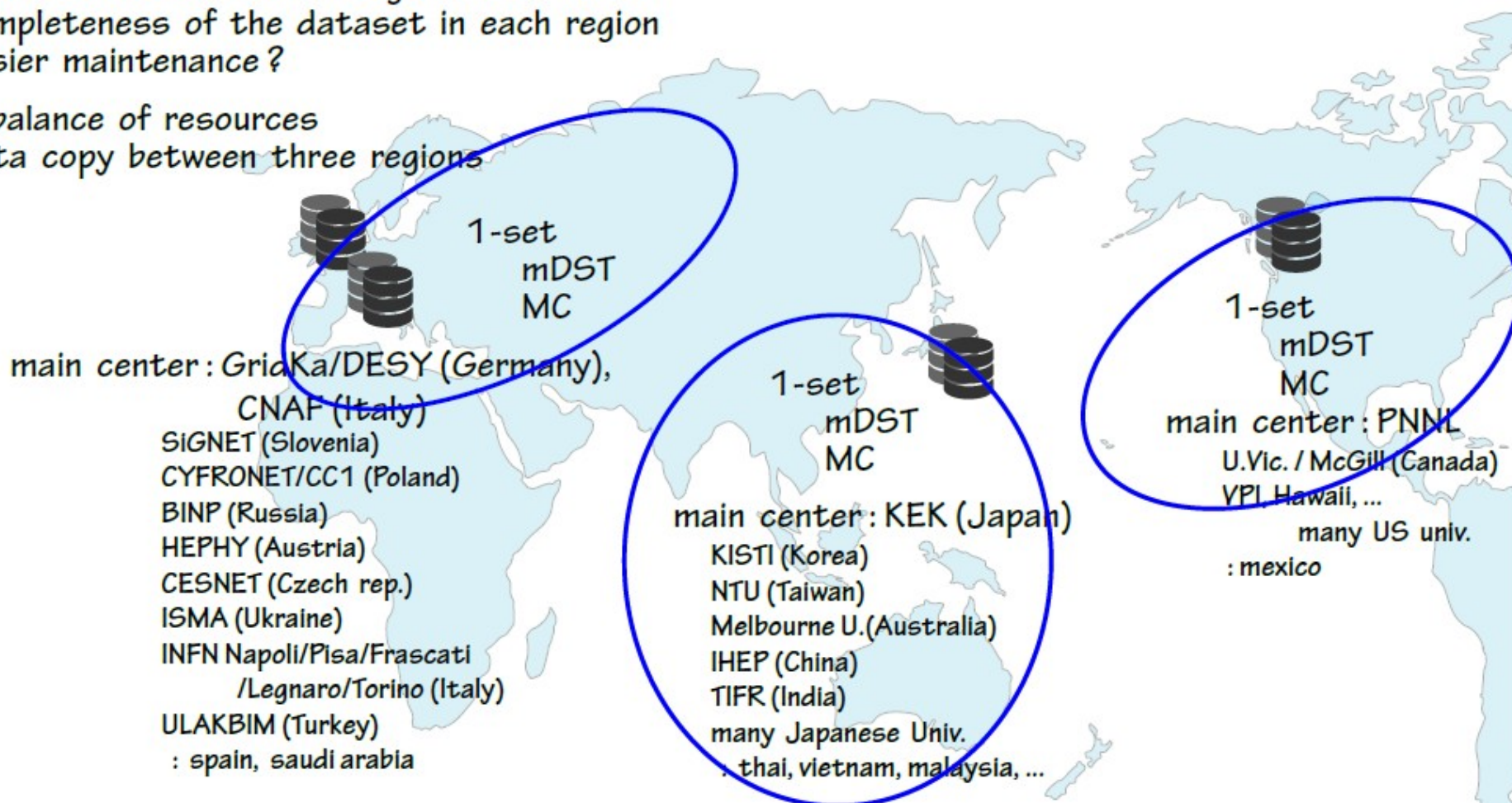    completeness of the dataset in each region
    easier maintenance?

    unbalance of resources
    data copy between three regions

1-set
mDST
MC

main center : GridKa/DESY (Germany),
         CNAF (Italy)
SiGNET (Slovenia)
CYFRONET/CC1 (Poland)
BINP (Russia)
HEPHY (Austria)
CESNET (Czech rep.)
ISMA (Ukraine)
INFN Napoli/Pisa/Frascati
        /Legnaro/Torino (Italy)
ULAKBIM (Turkey)
  : spain, saudi arabia

1-set
mDST
MC

main center : KEK (Japan)
KISTI (Korea)
NTU (Taiwan)
Melbourne U.(Australia)
IHEP (China)
TIFR (India)
many Japanese Univ.
  : thai, vietnam, malaysia, ...

1-set
mDST
MC
main center : PNNL
U.Vic. / McGill (Canada)
VPI, Hawaii, ...
        many US univ.
  : mexico

# Network

## Current Connectivity

### Trans-Pacific
10G : Tokyo - LA
10G : Tokyo - NY
10G : Osaka -Washington

### Trans-Atlantic
3 x 10G : NY - Amsterdam
3 x 10G : Washington - Frankfurt
ANA-100G NY - Amsterdam

### Trans-Asia
2.5G : Madrid-Mumbai
2.5G : Singapore-Mumbai
10G : Japan-Singapore

## "Planned" Connectivity

### Trans-Pacific
SINET5
100G link to US
  in 2016

### Trans-Atlantic
EEX (ESNet Extension to Europe)
2 x 100G : NY - London
100G : Washington - Geneva
40G : Boston - Amsterdam

### Trans-Asia
10G : Mumbai - GEANT
SINET ?

# Deployment of the virtualized computing infrastructure

Deployment stages:

- Virtualizing of an experimental group's computing environment.

- Tests of virtual machines locally on BINP resources.

- Transferring the VMs to a supercomputer and running them under remote batch system's control.

- Integration of local and remote batch systems.

Finally we have dynamical virtualized computing cluster. Physicists use computing resources in a conventional way.