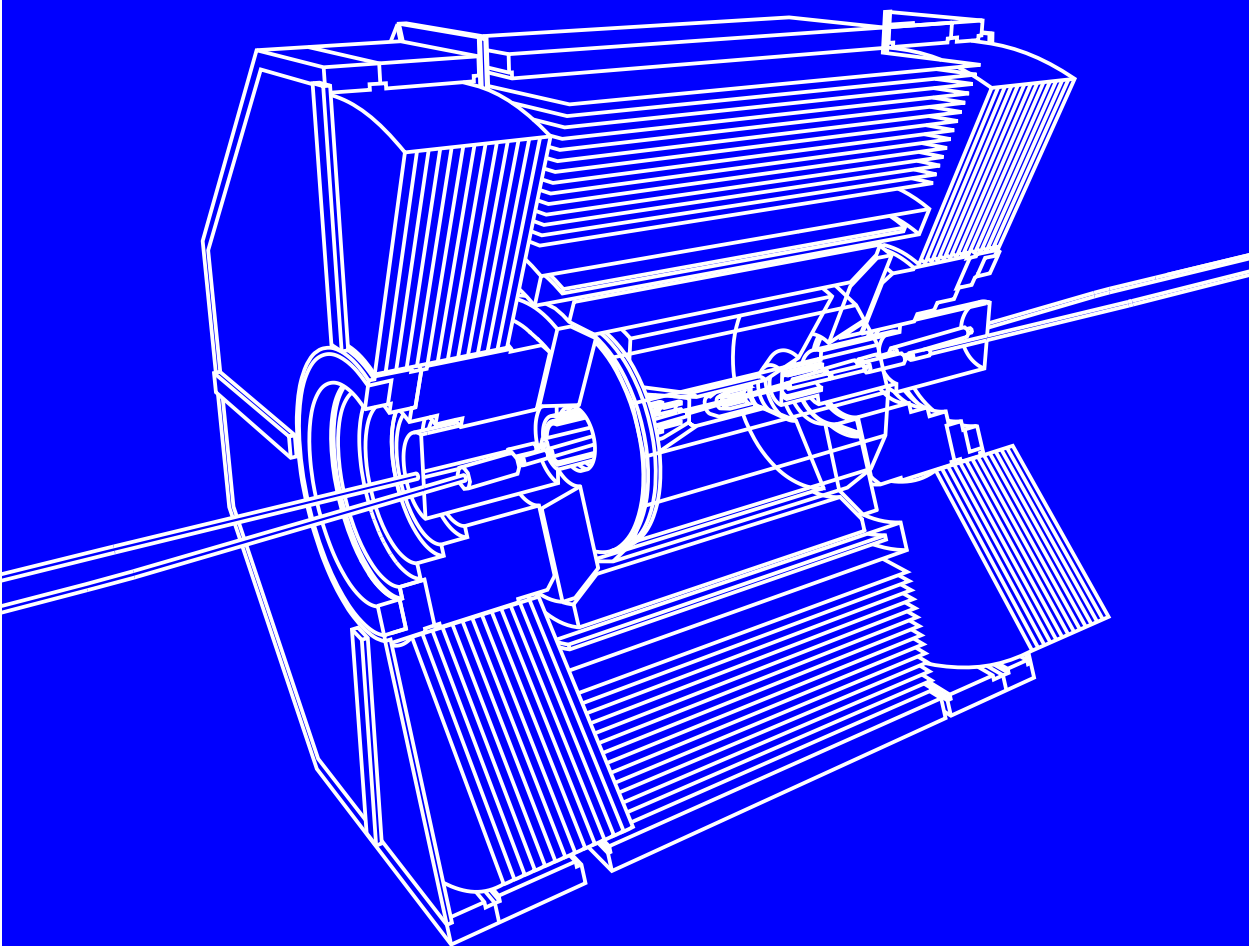


Computing at the Belle-II experiment



Takanori Hara (KEK)

takanori.hara@kek.jp

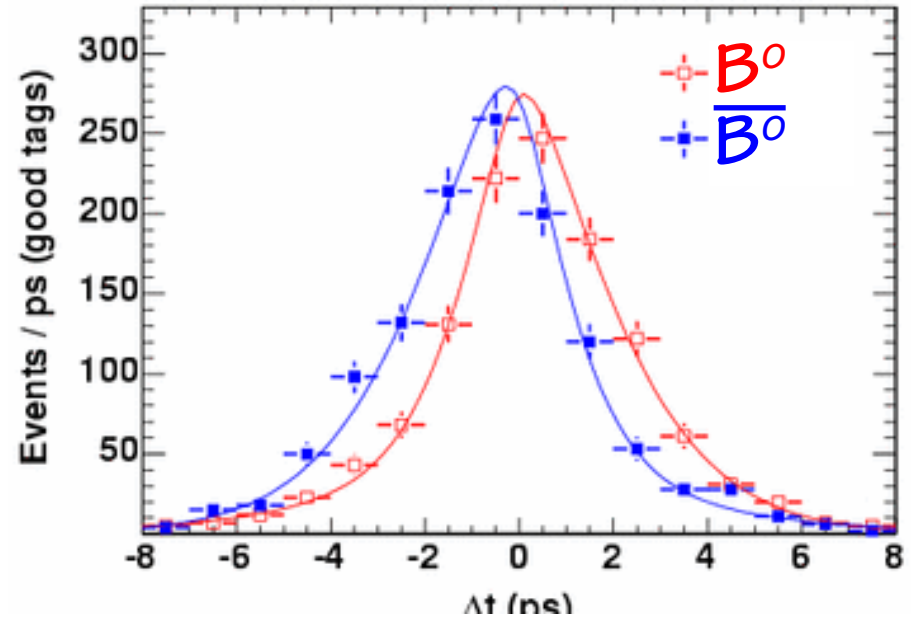
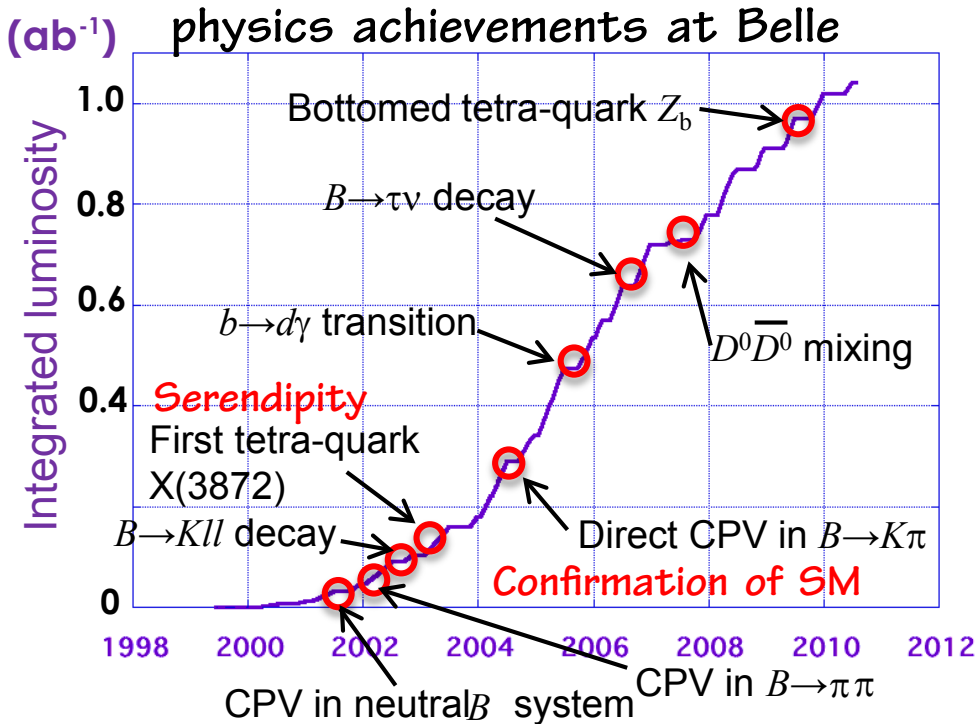
For the Belle II distributed computing group

13 April, 2015 @ CHEP2015 in Okinawa

From Belle to Belle II

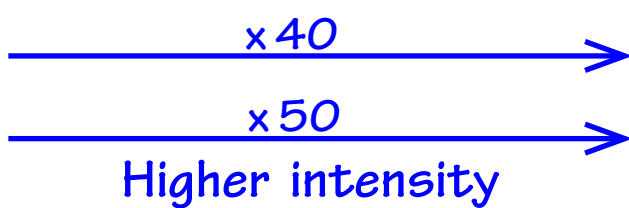
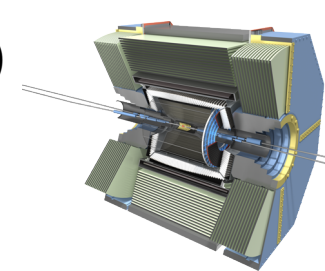
Accelerator	KEKB
Beam Energy (GeV)	3.5×8 ($\gamma = 0.425$)
CM energy, $Y(4S)$,
Luminosity ($\text{cm}^2 \text{s}^{-1}$)	2.1×10^{34}
Total data (ab^{-1})	1
	raw data: $\sim 1\text{PB}$
	mDST data/MC : 0.14/0.6 PB (for one version)

Computing one big center @ KEK (non-grid)

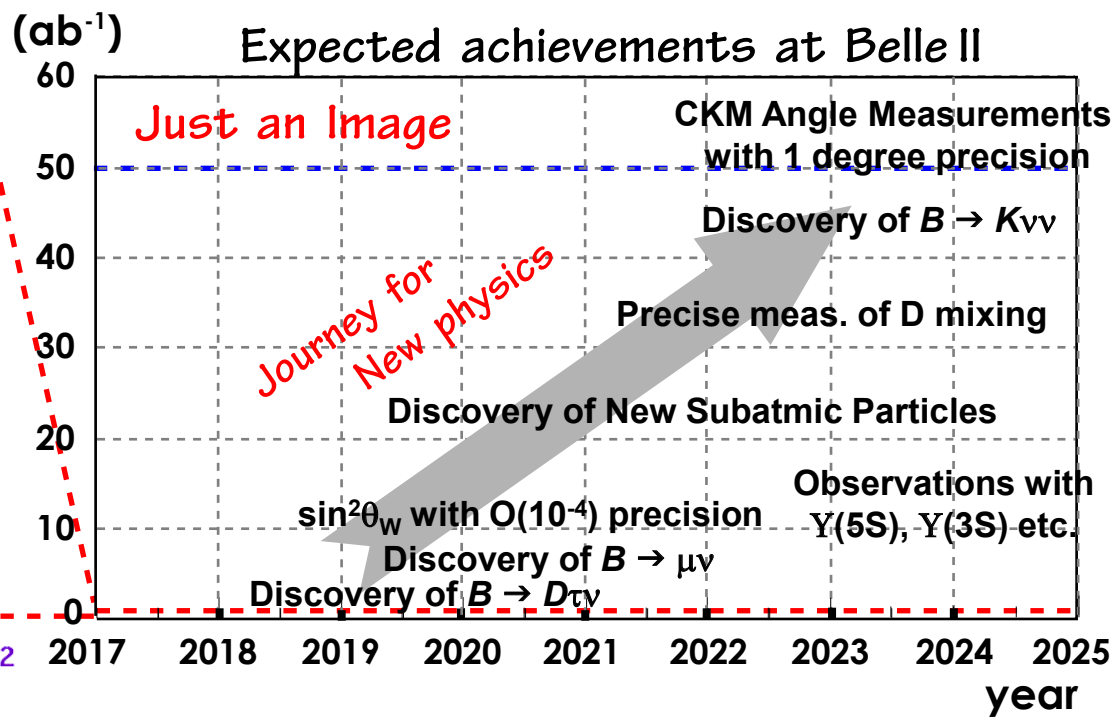
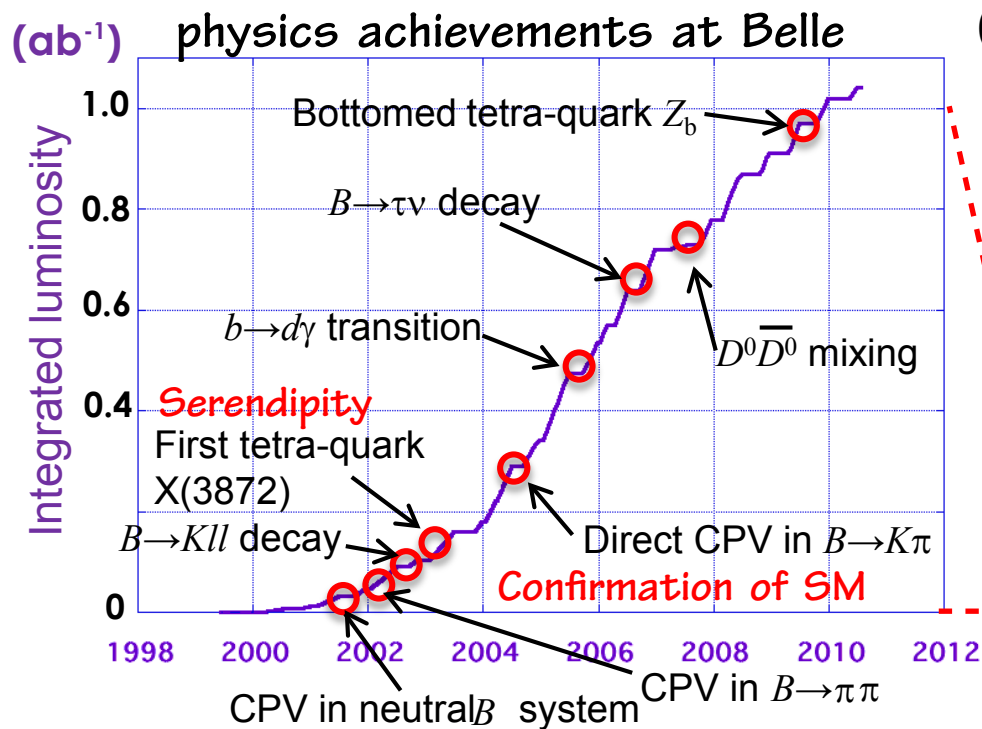


From Belle to Belle II

Accelerator	KEKB	SuperKEKB
Beam Energy (GeV)	3.5 x 8 ($\gamma = 0.425$)	4 x 7 ($\gamma = 0.28$)
CM energy, Y(4S),, Y(4S),
Luminosity ($\text{cm}^2 \text{s}^{-1}$)	2.1×10^{34}	8×10^{35}
Total data (ab^{-1})	1	50
raw data	$\sim 1 \text{ PB}$	$\sim 100 \text{ PB}$
mDST data/MC	0.14/0.6 PB (for one version)	(another raw data copy outside KEK)



Computing one big center @ KEK (non-grid) \longrightarrow world-wide distributed computing



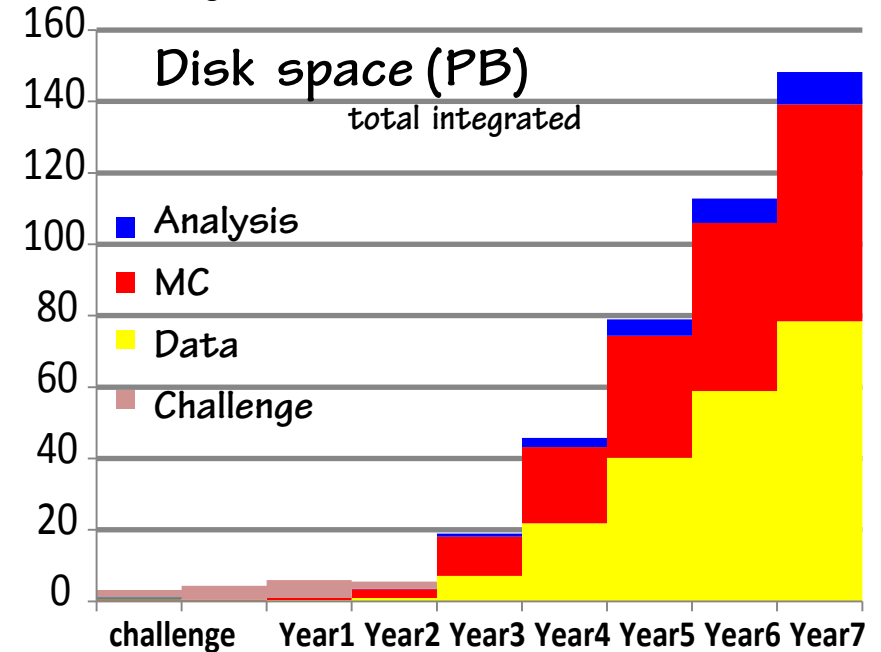
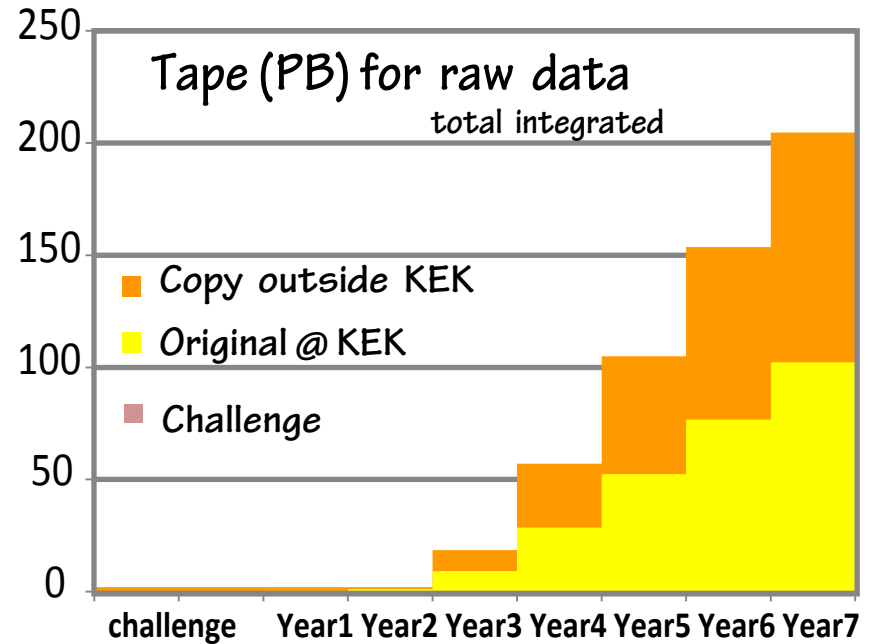
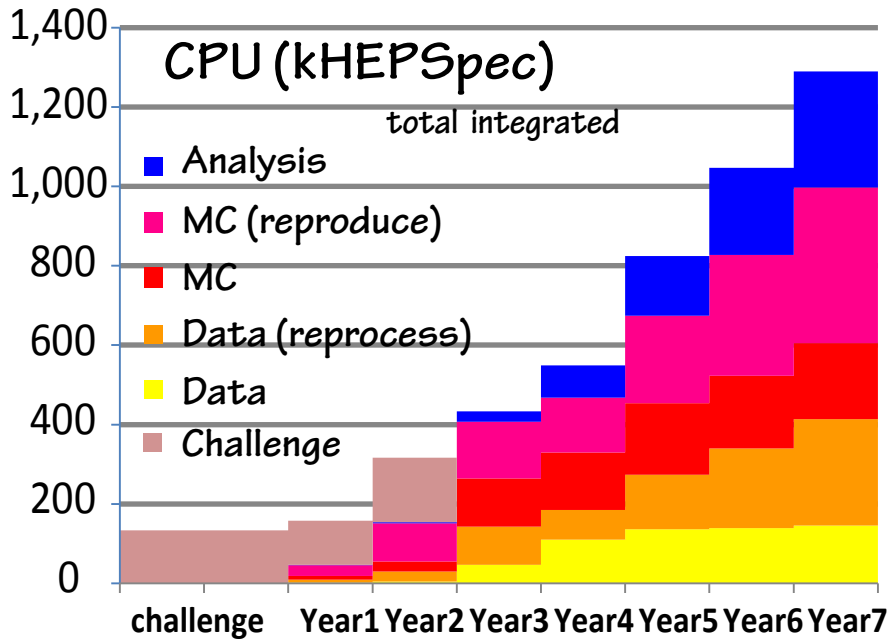
Hardware Resources for Belle II

version estimated in early 2014

- uncertainties
- Performance of accelerator
- beam background condition
- improvement of software

The yearly profile may change

The total at the last year should stay the same level

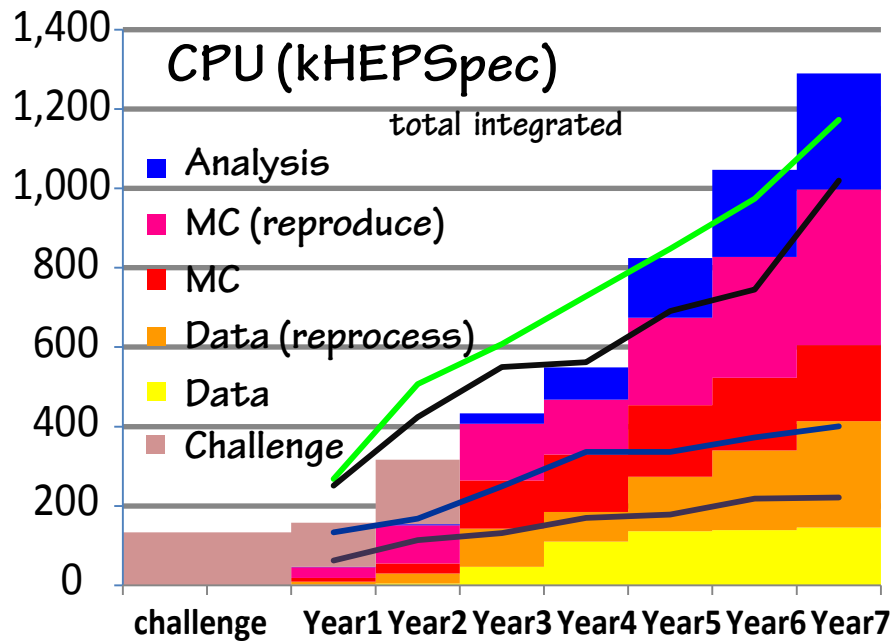


Hardware Resources for Belle II

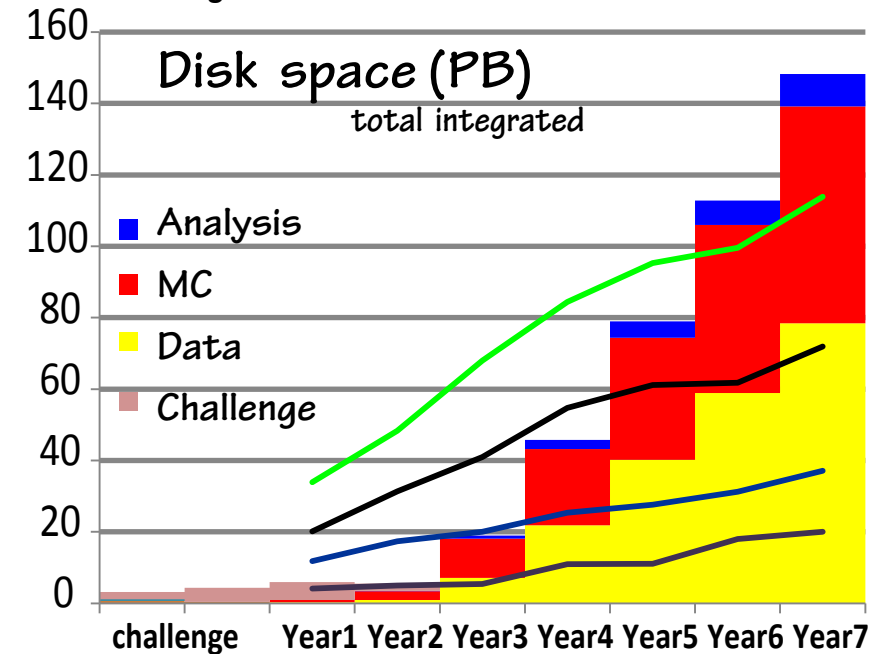
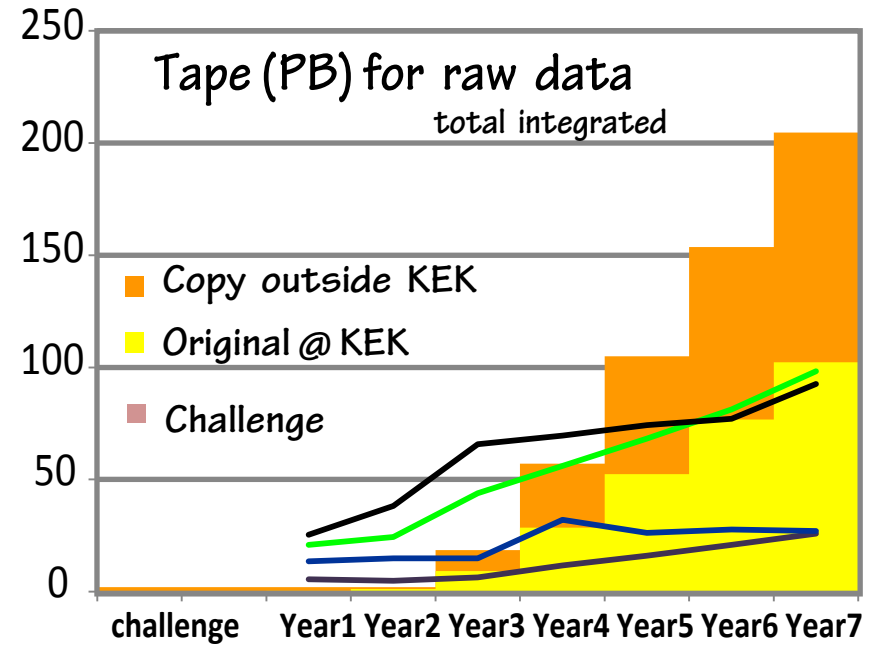
LHC resources based on the published pledges

<http://wlcg-rebus.cern.ch/apps/pledges/summary/>

The real capacities and the usages can be different

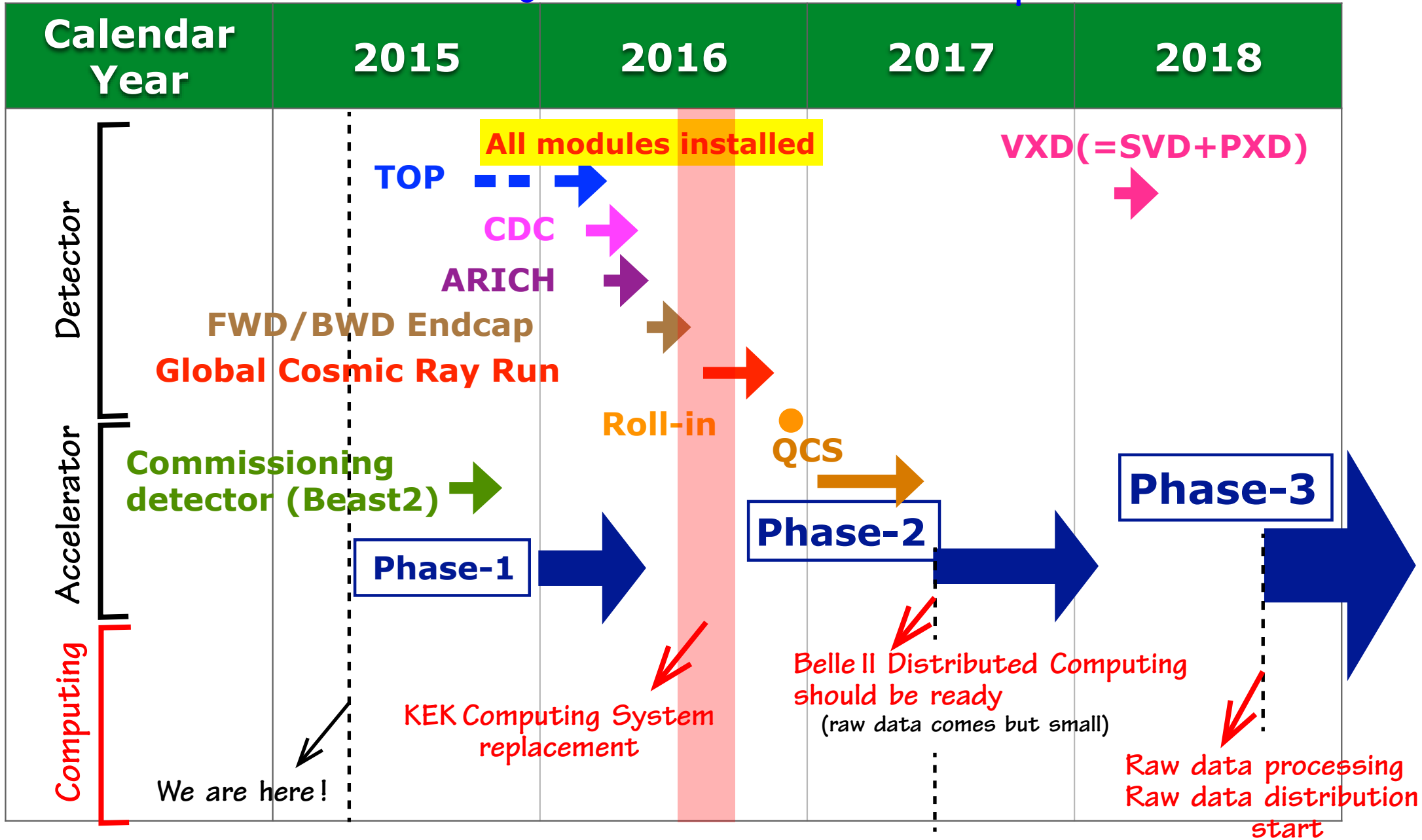


— ATLAS
— CMS
— ALICE
— LHCb



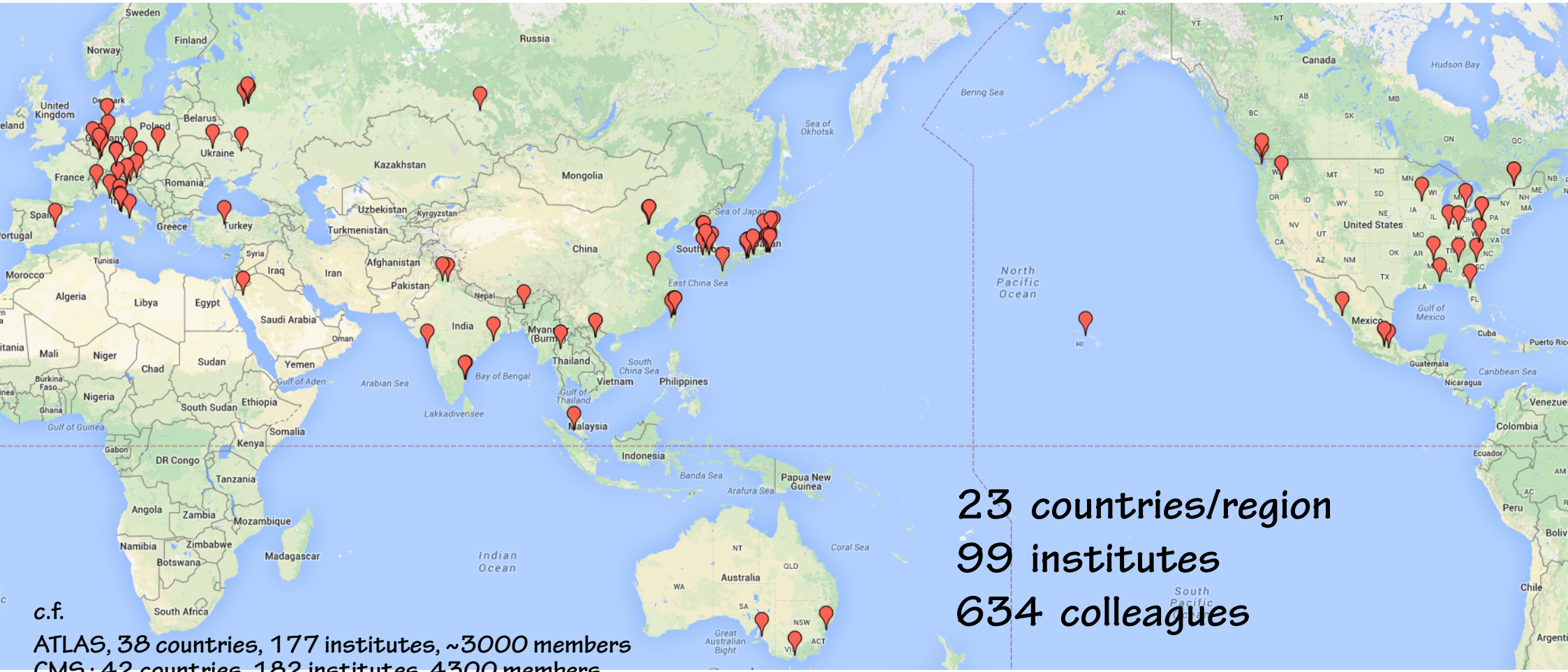
SuperKEKB/Belle II Time line

KEK is the hosting institute of the Belle II experiment





Belle II Collaboration



23 countries/region
 99 institutes
 634 colleagues

c.f.
 ATLAS, 38 countries, 177 institutes, ~3000 members
 CMS : 42 countries, 182 institutes, 4300 members
 ALICE : 36 countries, 131 institutes, 1200 members
 LHCb : 16 countries, 67 institutes, 1060 members

as of April 4, 2015

Asia : ~43%	N. America : ~17%	Europe : ~40%
Japan : 139	US : 78	Germany : 89
Korea : 37	Canada : 20	Italy : 62
Taiwan : 25	Mexico : 8	Russia : 40
India : 25		Slovenia : 17
China : 18		Austria : 14
Australia : 22		Poland : 11
		Czech rep. : 8

others : < 8 colleagues / country



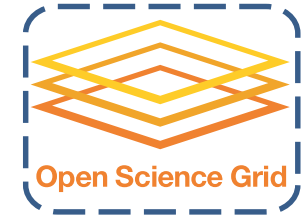


Distributed Computing Resources

KIT, CNAF, CESNET, SiGNET, HEPHY, UA-ISMA, ULAKBIM, CYFRONET,

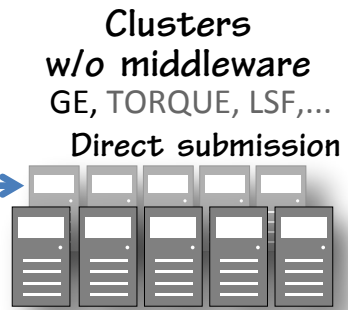


GRID Middlewares

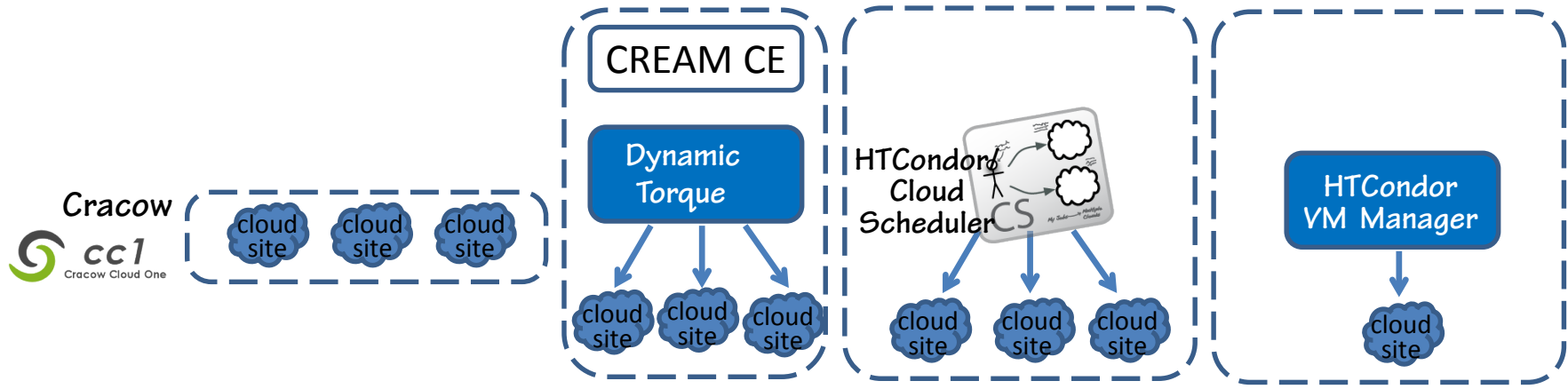


Belle II resources / infrastructure overlap with WLCG

batch system



BINP, NSU, universities in Japan, Korea



Melbourne
• Seen as a traditional CREAM CE site
• Installed in each cloud site

UVic
Academic clouds
Commercial clouds, Amazon EC2, etc

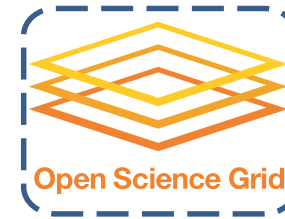
PNNL
HPC

Interoperability with DIRAC

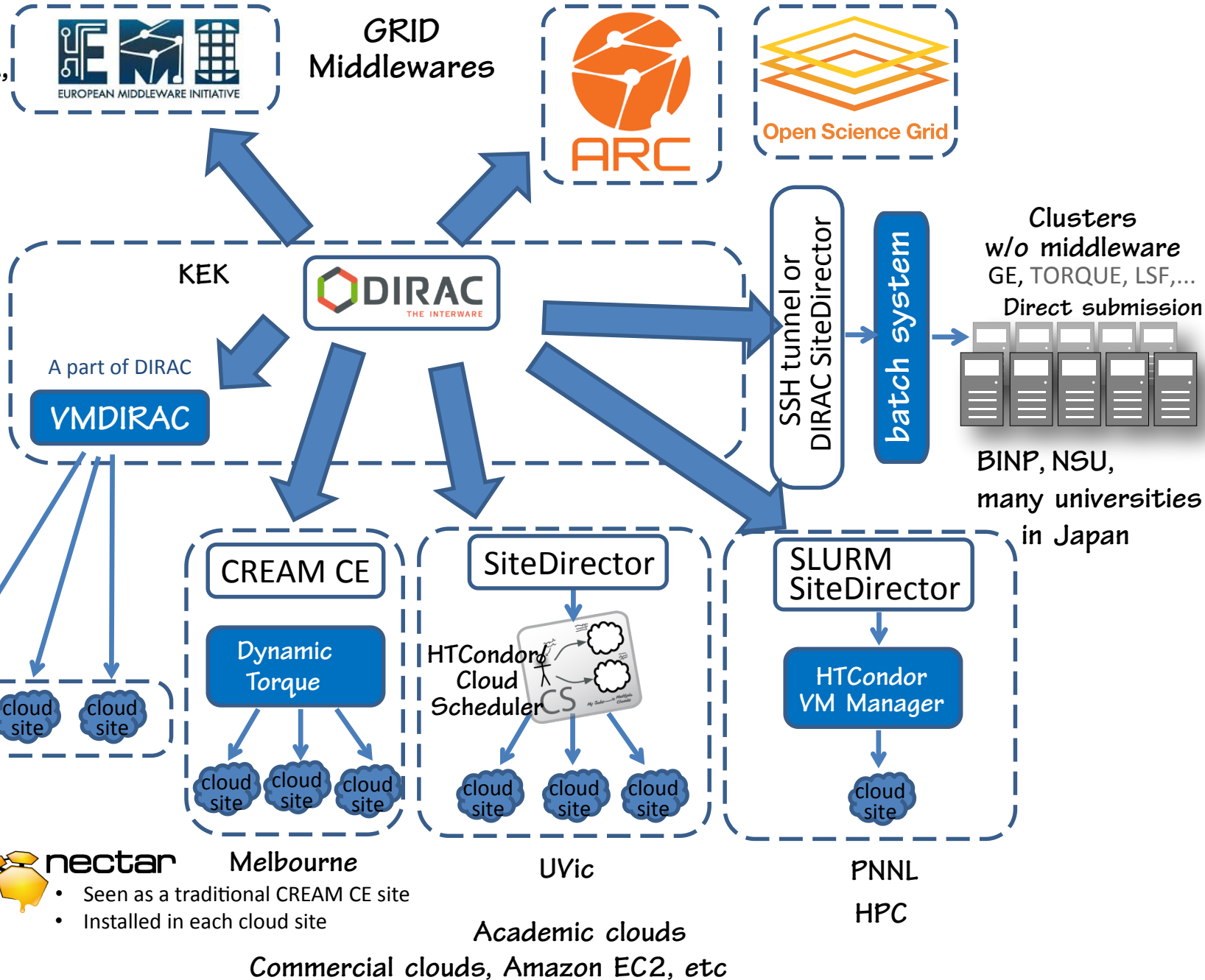
KIT, CNAF, CESNET, SiGNET, HEPHY, UA-ISMA, ULAKBIM, CYFRONET,



GRID Middlewares



Distributed
Infrasturcture with
Remote
Agent
Control
(originally developed for LHCb)



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

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

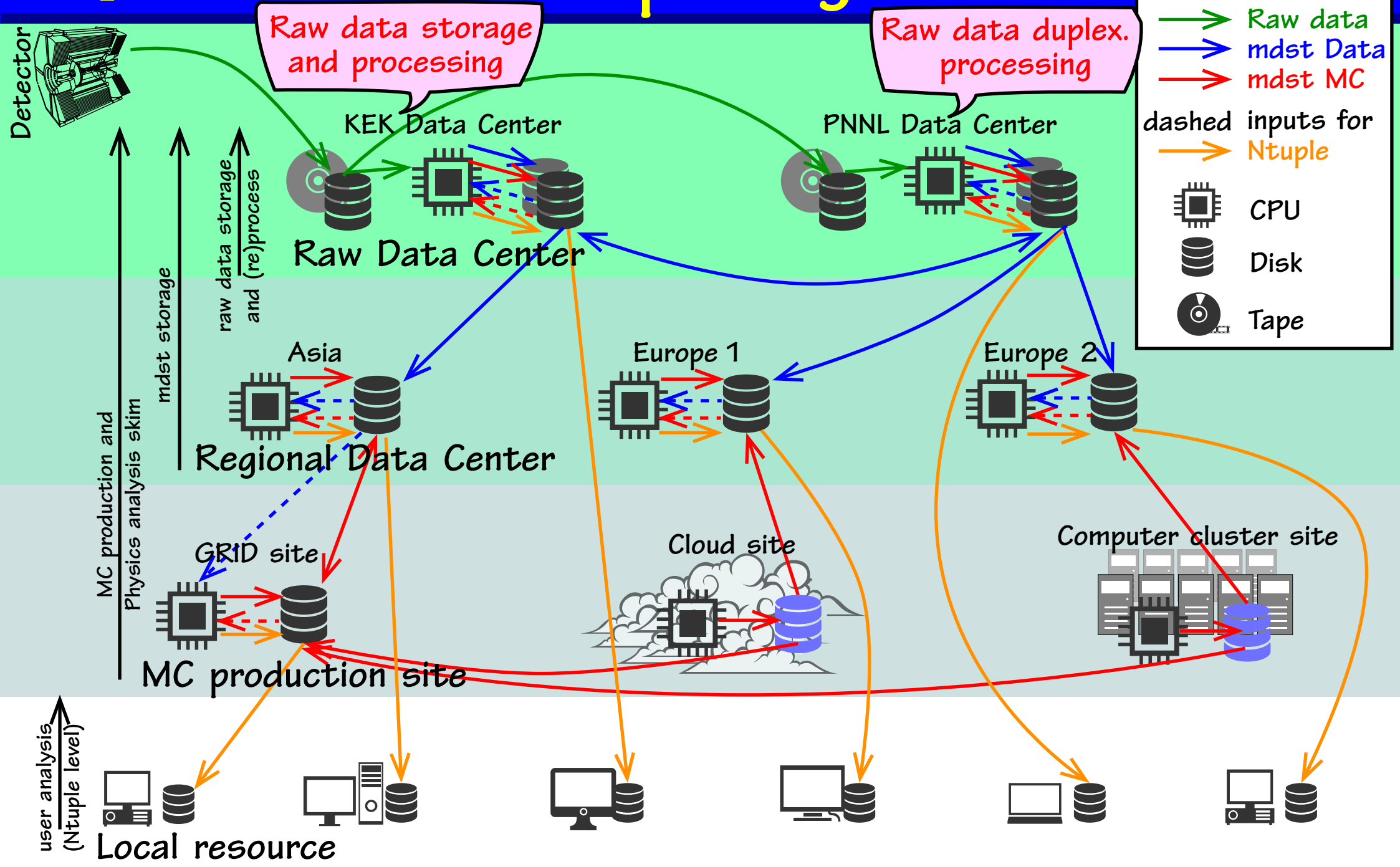
Academic clouds

Commercial clouds, Amazon EC2, etc



Belle II Computing Model

end of year 3





Belle II Distributed Computing system



DIRAC main servers @ KEK

DIRAC servers for test/development purpose at PNNL (USA), Cracow (Poland), etc.



VOMS @ KEK



AMGA
recent improvement

+



LFC : has been working well

Studies with DFC vs AMGA+LFC: not yet a stage to tell their scalabilities

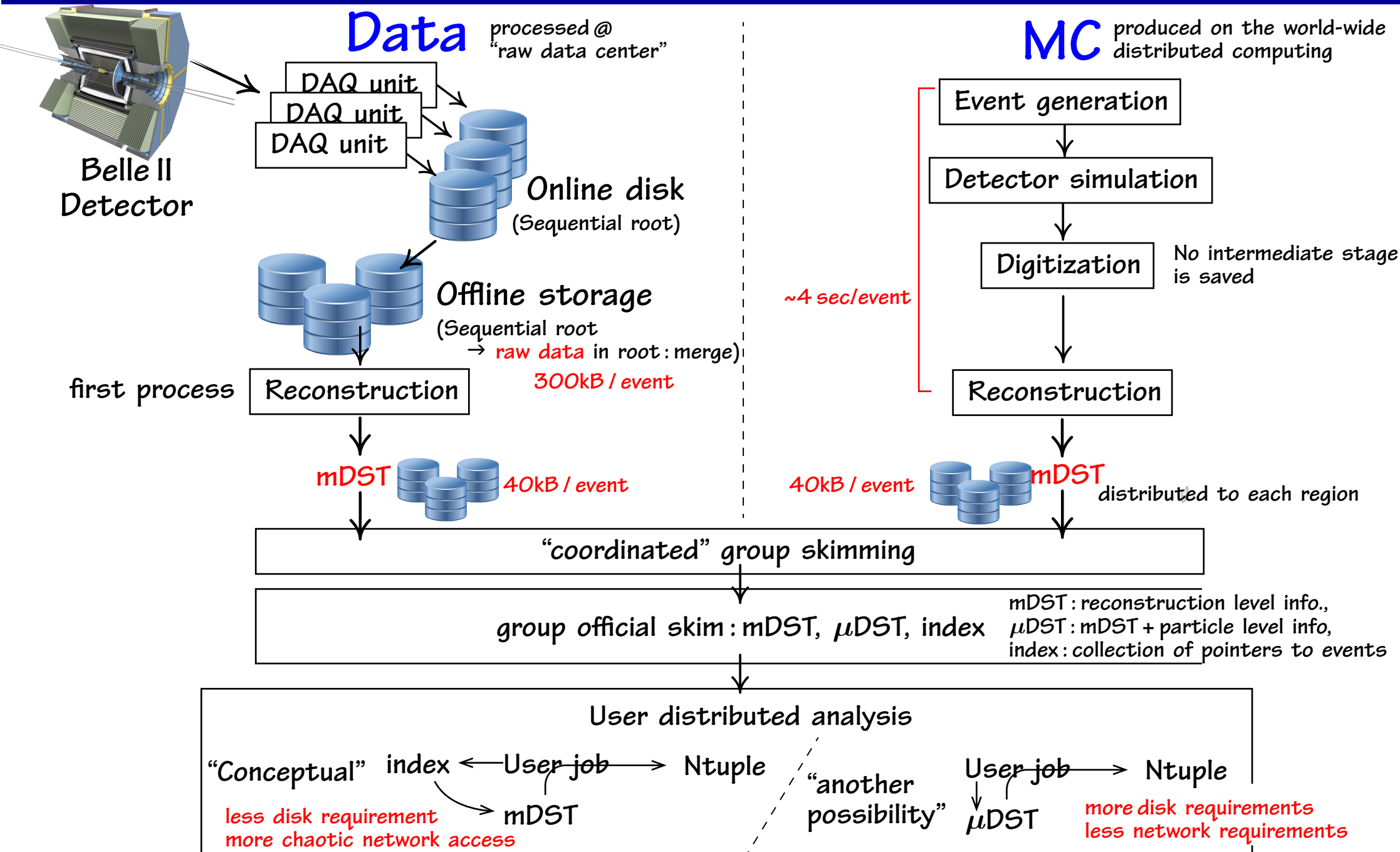


FTS3 : getting integrated



cvmfs is used for software distribution

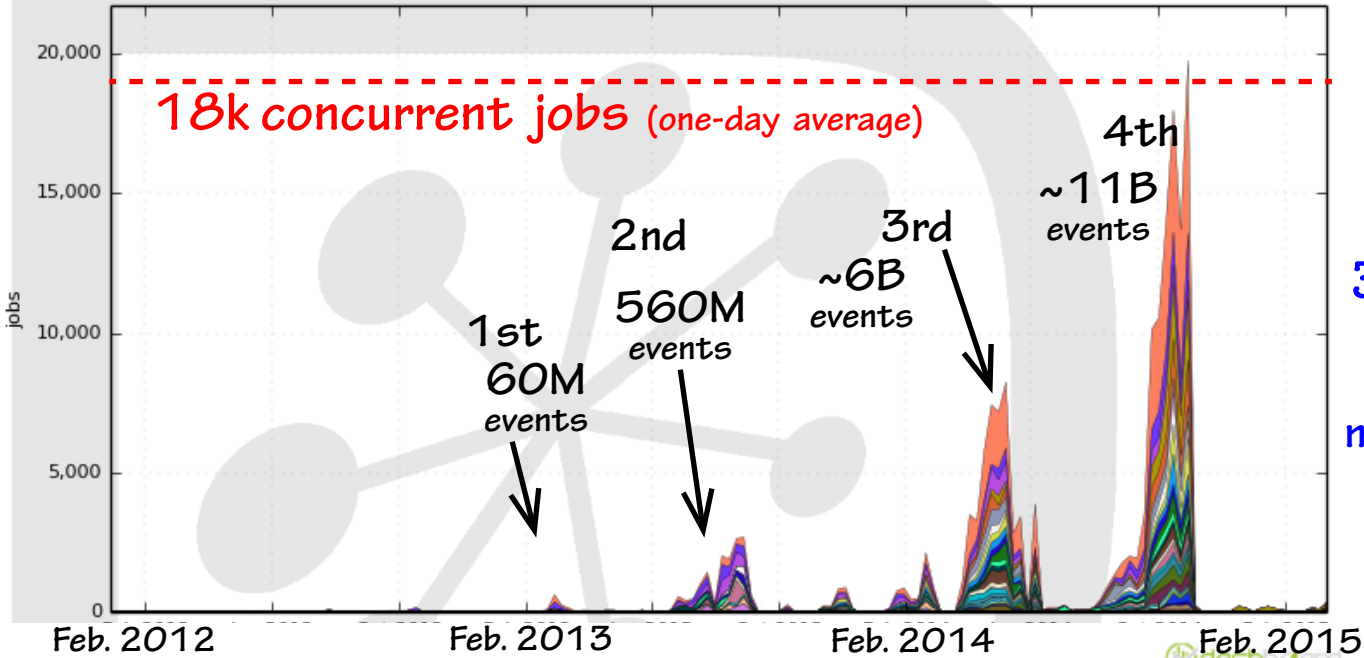
Data flow diagram



Operation at larger scale

Running jobs by Site

167 Weeks from Week 52 of 2011 to Week 10 of 2015



15 countries/regions

Australia, Austria, Canada, Czeck R., Germany, Italy, Japan, Korea, Poland, Russia, Slovenia, Taiwan, Turkey, Ukraine, USA

31 sites

GRID, Cloud, local cluster is available

more than $3ab^{-1}$ data produced in 2014

However, still a factor of x10 below requirements for full Belle II luminosity

LHCb (~120 sites)

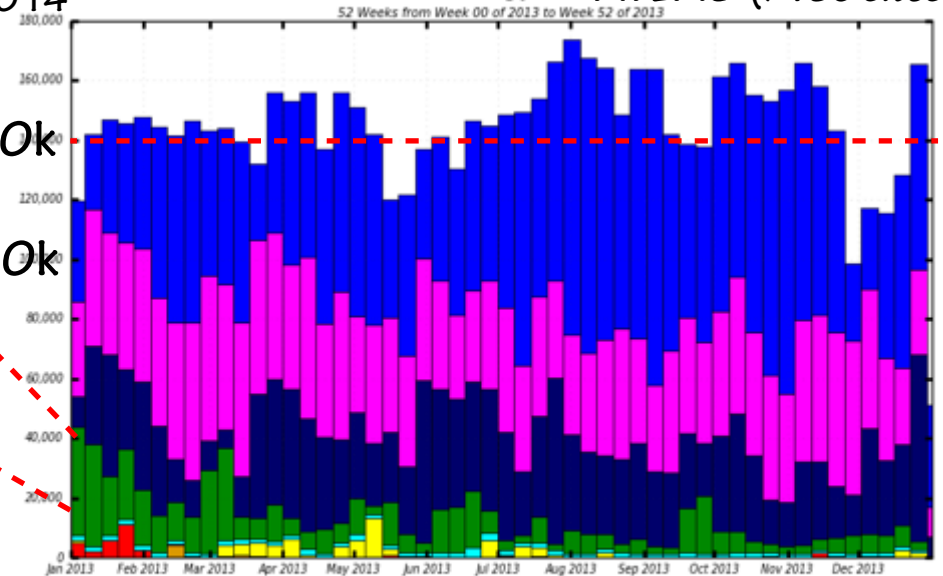
Running jobs by Site

252 Weeks from Week 52 of 2009 to Week 43 of 2014



Oct. 2014

Running jobs ATLAS (>130 sites)



Mar. 2010

Sep. 2012

Sep. 2014

Monitoring system

- DIRAC web interface can provide some monitors for job status etc.

*It is not sufficient for our operation.
e.g. pilot status*

→ We are developing our own monitoring system with web interface.

2-way monitoring

- Active way

Submit test job, test SE access, check port access

- Passive way

Keep statistics on the result of pilot jobs, Analyze log of pilot jobs



*Visualize monitoring result on the web page.
(enable to check it easily, anywhere)*



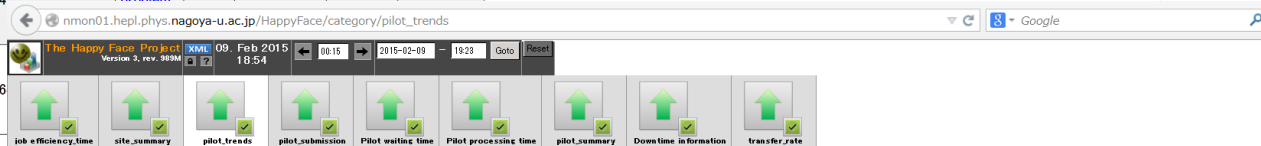
Monitoring system

Site status summary

site	worker node	CPU	#core	memory	OS	Kernel	rpm	cvmfs	releases	CPU Norm.	last updated
DIRAC.BINP.ru	gcf-110-1.belle2	QEMU Virtual CPU version 0.12.1	x2	1443MB/cores	Scientific Linux release 6.6 (Carbon)	2.6.32-504.1.3.el6.x86_64	OK	N.I.	OK (build-2014-10-18)	10.6 HS06	2015/01/25 09:34:33
DIRAC.Niigata.jp	ngtbel5.sc.niigata-u.ac.jp	Intel(R) Xeon(R) CPU E5-2660 0 @ 2.20GHz	x32	2013MB/cores	Scientific Linux SL release 5.5 (Boron)	2.6.18-371.6.1.el5	one problem found	Rev. 50	OK (build-2014-10-18)	11.2 HS06	2015/02/09 16:36:37
DIRAC.Osaka-CU.jp	cpu4	Intel(R) Celeron(R) CPU 550 @ 2.00GHz	x1	1868MB	Scientific Linux CERN SLC release 6.6 (Carbon)	2.6.32-504.3.3.el6.x86_64	OK	Rev. 50	OK (build-2014-10-18)	7.5 HS06	2015/02/09 16:38:18
DIRAC.PNNL-CASCADE.us	g1107	Intel(R) Xeon(R) CPU E5-2670 0 @ 2.60GHz	x16	8064MB/cores	Scientific Linux release 6.6 (Carbon)	2.6.32-504.3.3.el6.x86_64	one problem	Rev.	OK (build-	12.9	2015/01/21
DIRAC.PNNL.us	cwn1	AMD Opteron 62xx class CPU	x32	1950MB/cores	Scientific Linux release 6.5 (Carbon)	2.6.32-431.17.1.el6.x86_64					

← Check basic requirements for each site by test jobs. (libraries, free disk space, etc)

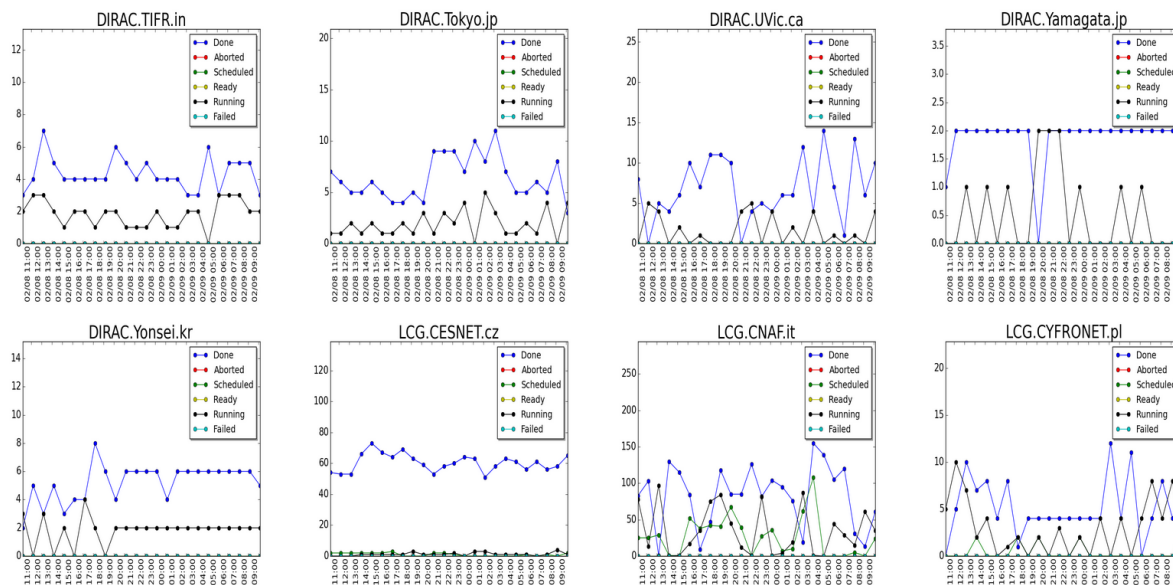
↓ Check pilot status as well as processing time, heartbeat time and so on.



We can quickly find many kinds of problems by our monitoring system now.

As a next step, more sophisticated system is being developed:

more precise diagnosis, automated notification, etc.



→ made the MC production shift easier

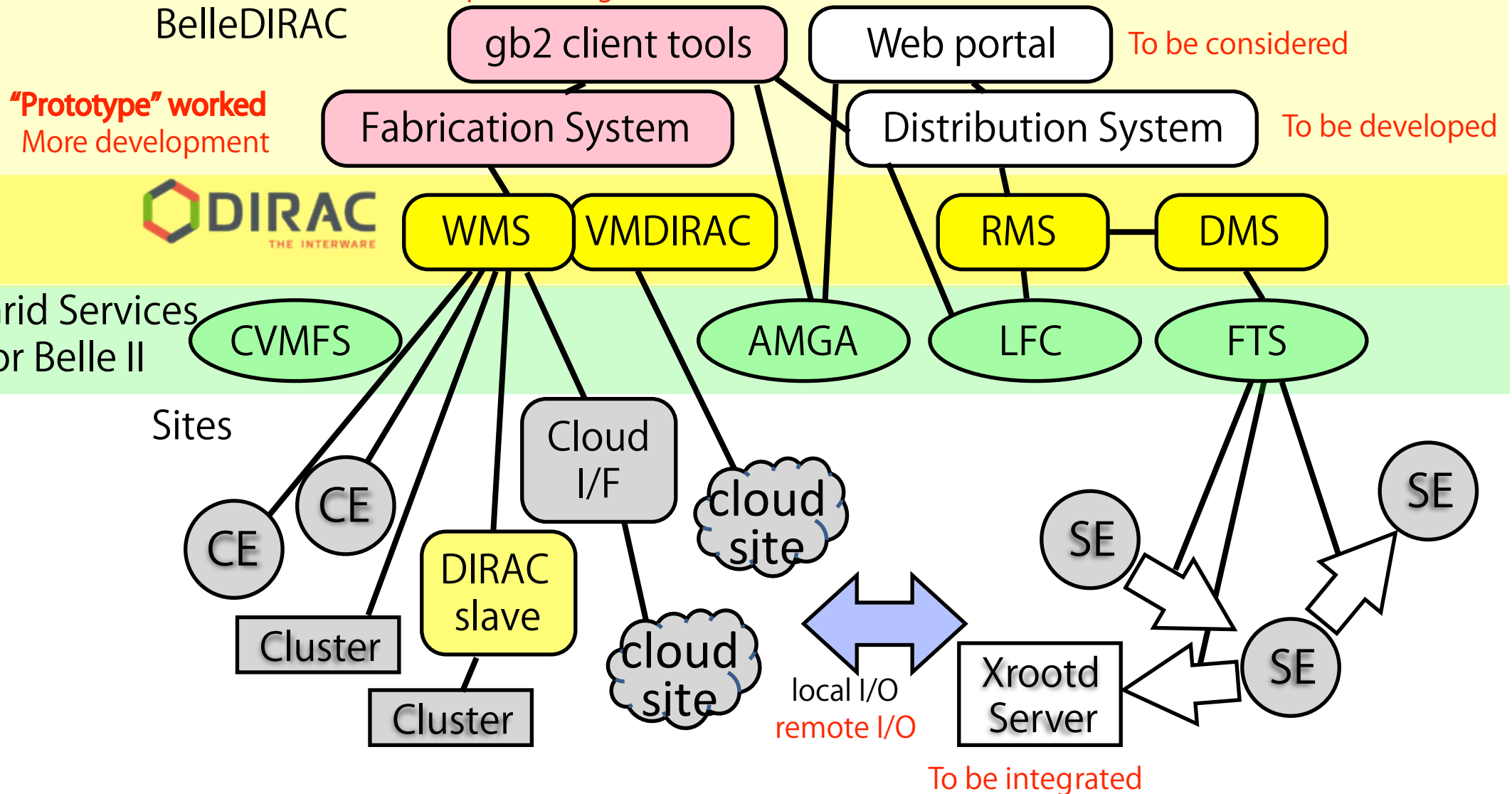
Distributed Computing in future

Production Manager

Data Manager

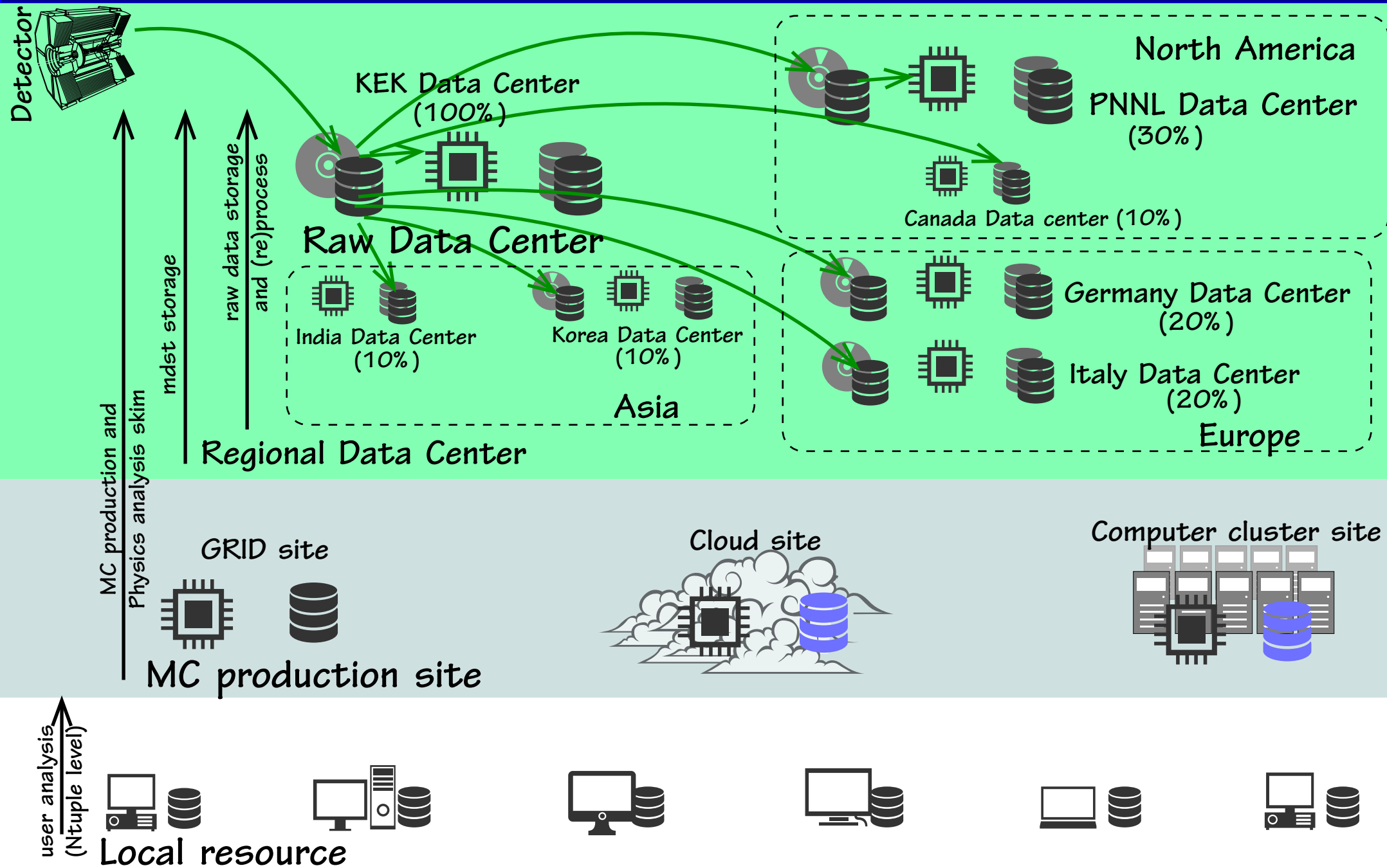
End Users

Keeps evolving



Belle II Computing Model

start year 4
(raw data part)



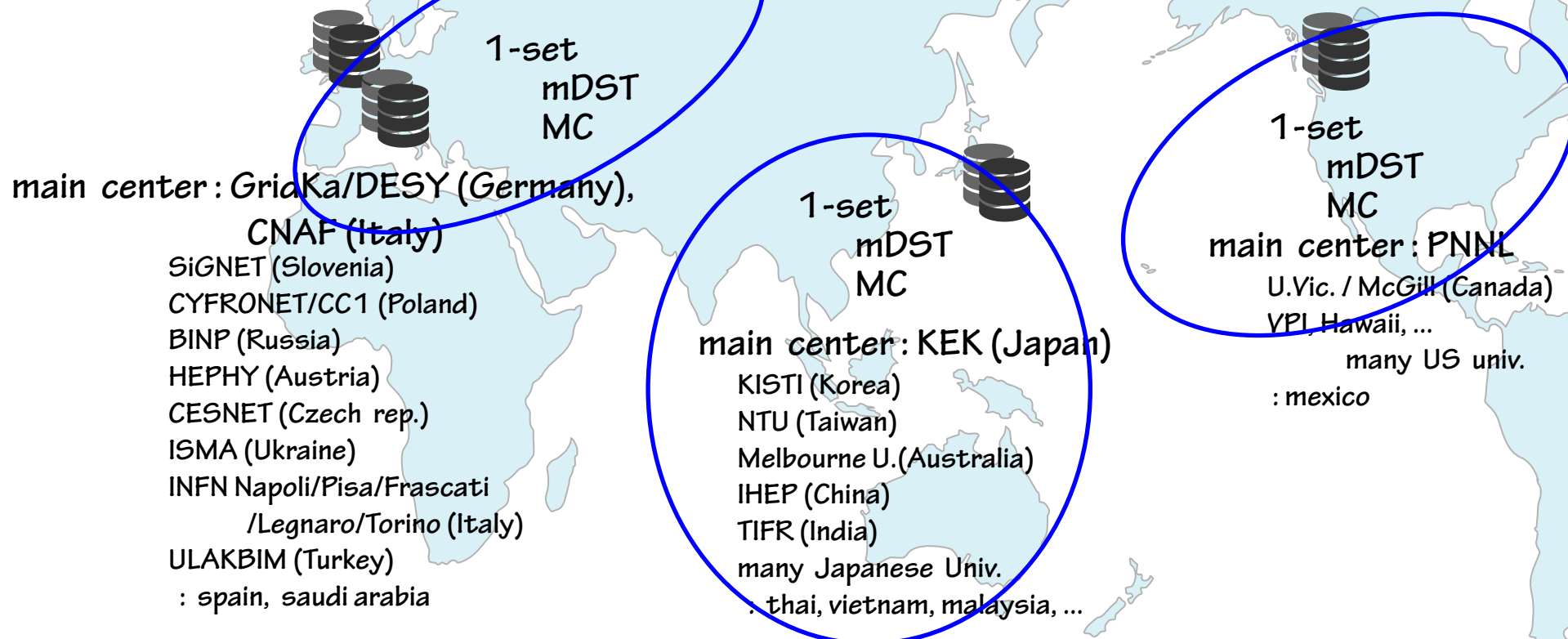
Processed Data Distribution

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

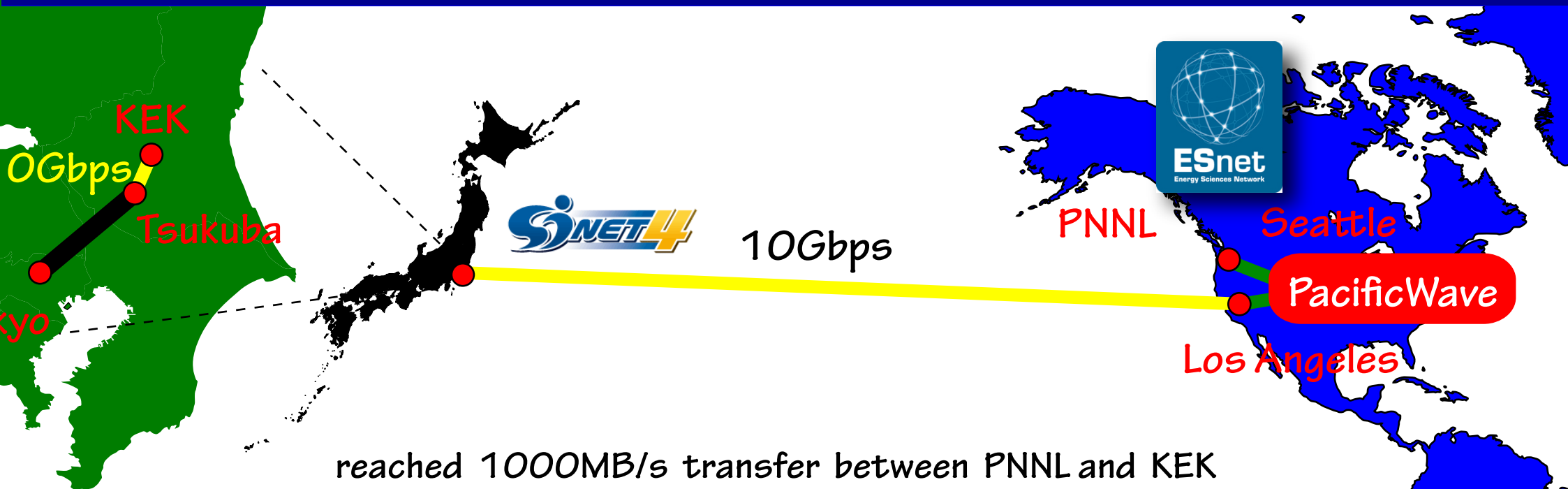
For the MC data seems to be natural to be the similar structure

better network? in each region
 completeness of the dataset in each region
 easier maintenance?

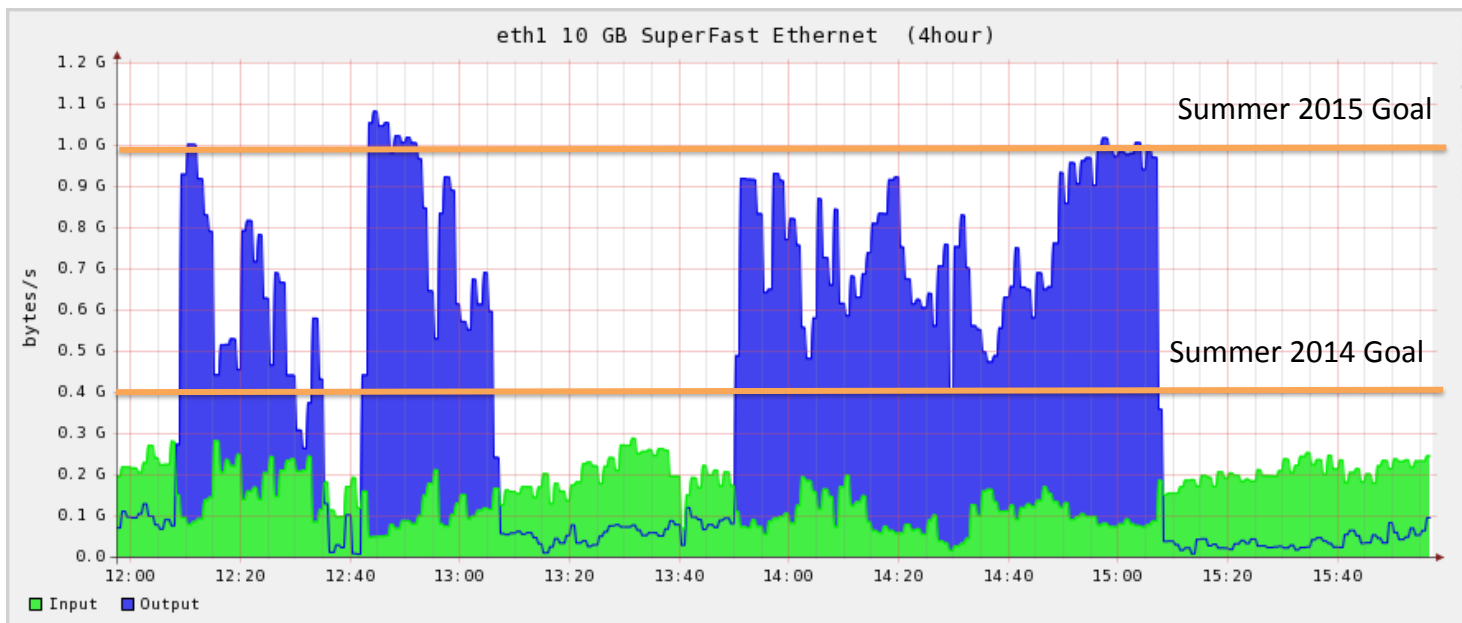
unbalance of resources
 data copy between three regions



Trans-Pacific data challenge

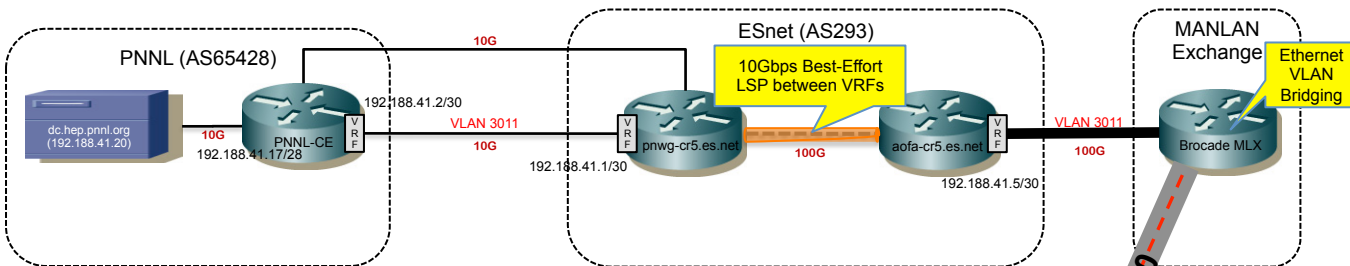


reached 1000MB/s transfer between PNNL and KEK



Trans-Atlantic data challenge

US side

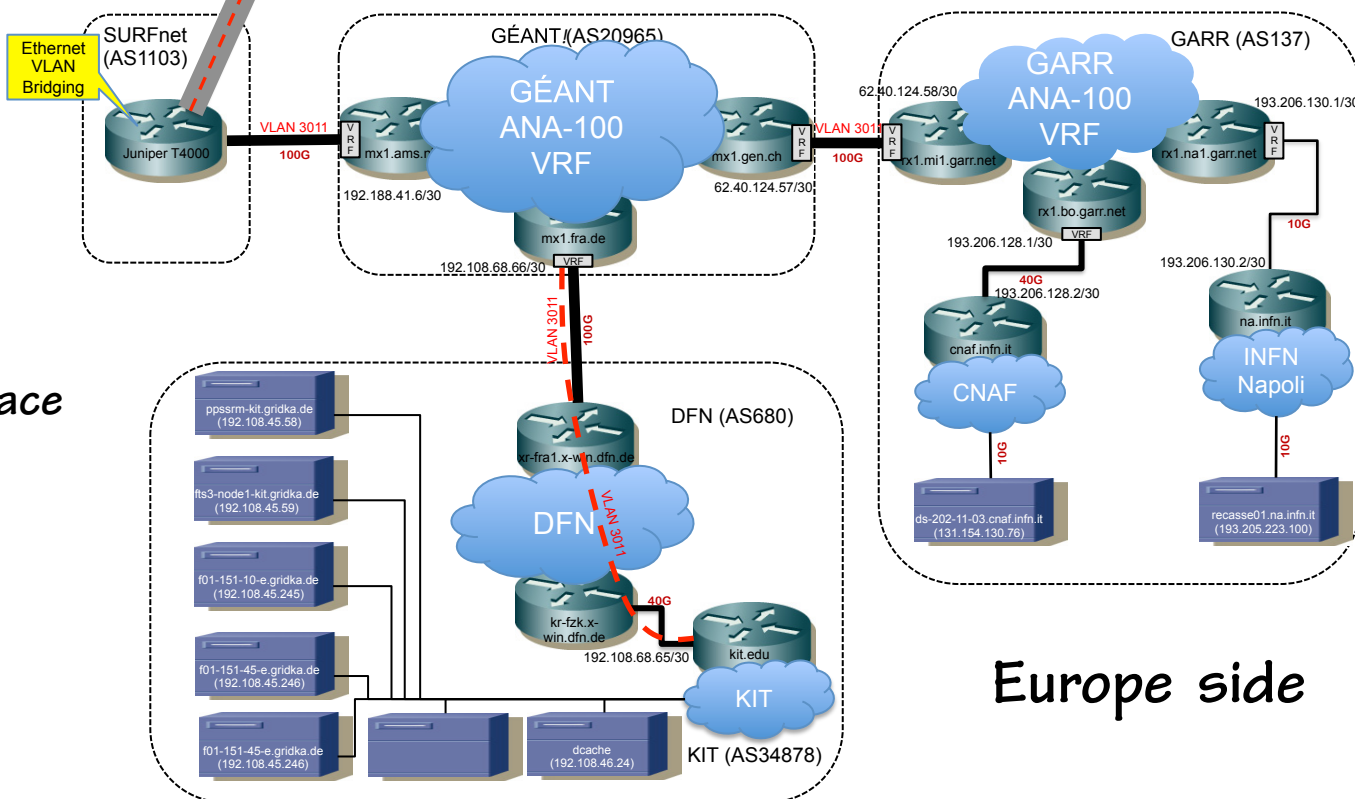


- “traceroute” was used to confirm the routing to each SE
- “iperf” was used to do initial network transfer rate test
- FTS3 server at GridKa was used to schedule data transfers

Dedicated 10G link between PNNL SE and ESNet
10G best-effort Label Switched Path in ESNet backbone

Test was done in May/June 2014

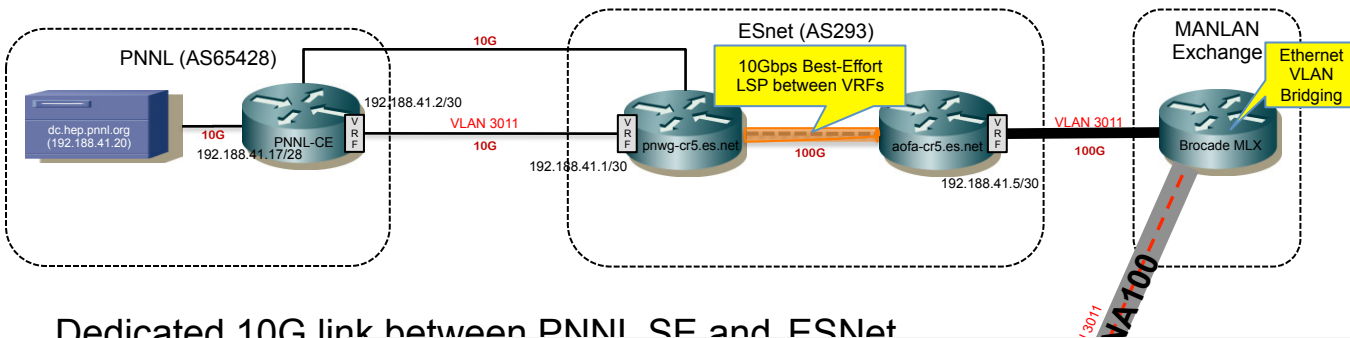
- Network providers setup the VLAN
- Local network providers and sites coordinated final configurations
- Sites must configure hardware interface to match destinations



Europe side

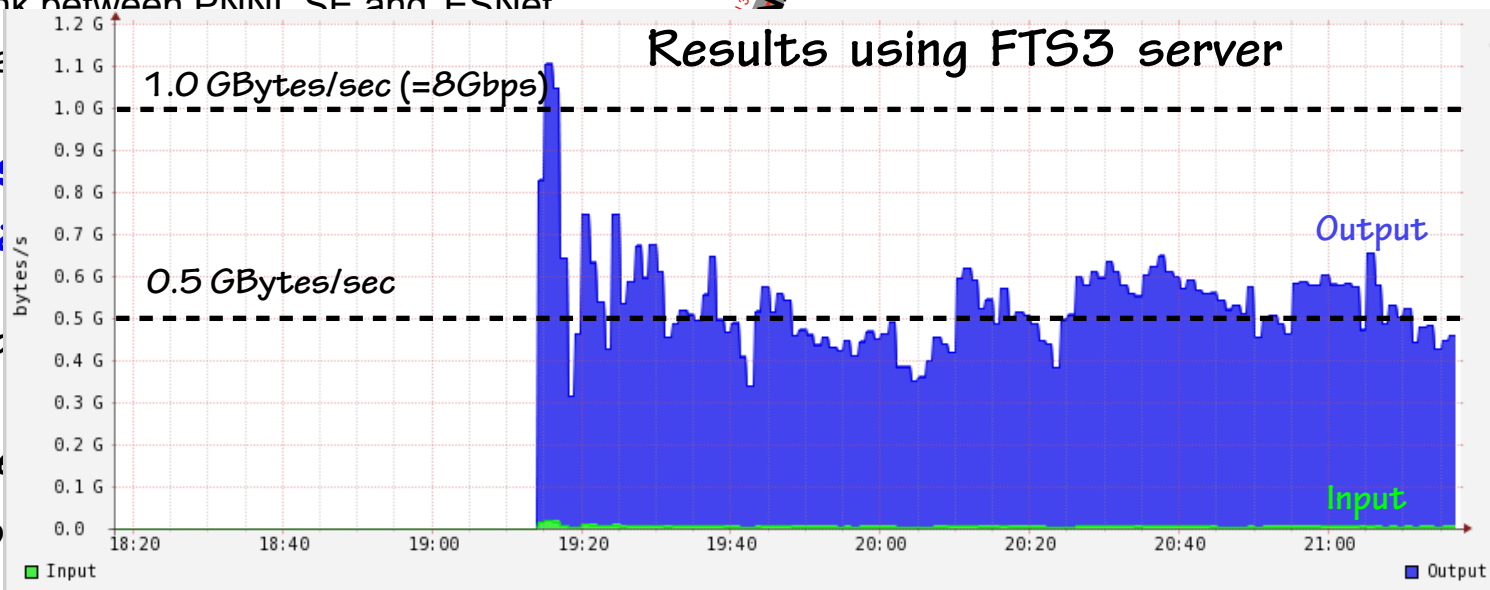
Trans-Atlantic data challenge

US side



- “traceroute” was used to confirm the routing to each SE
- “iperf” was used to do initial network transfer rate test
- FTS3 server at GridKa was used for transfers

Dedicated 10G link between PNNL SE and ESNet
 10G best-effort LSP



Test was done in March

Network provided by GARR
 Local network coordinators
 Sites must coordinate to match

transfers



Vincenzo Capone,
 Aleksandr Kurbatov, Mian Usman



Chin Guok

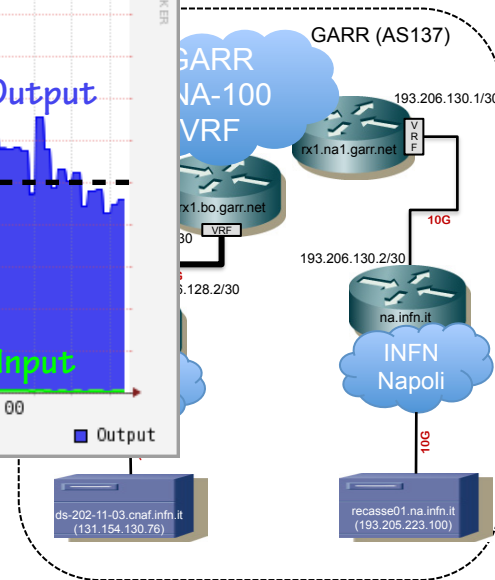
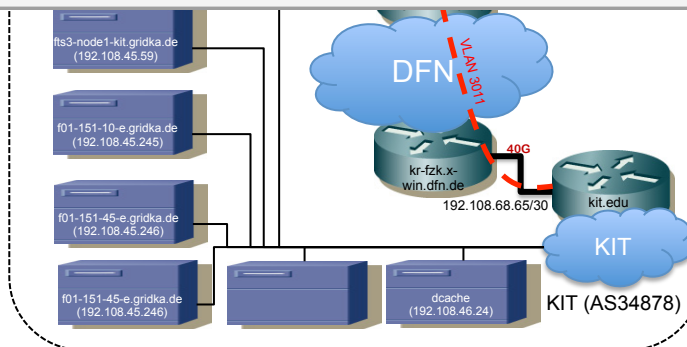


Thomas Schmid, Hubert Weibel

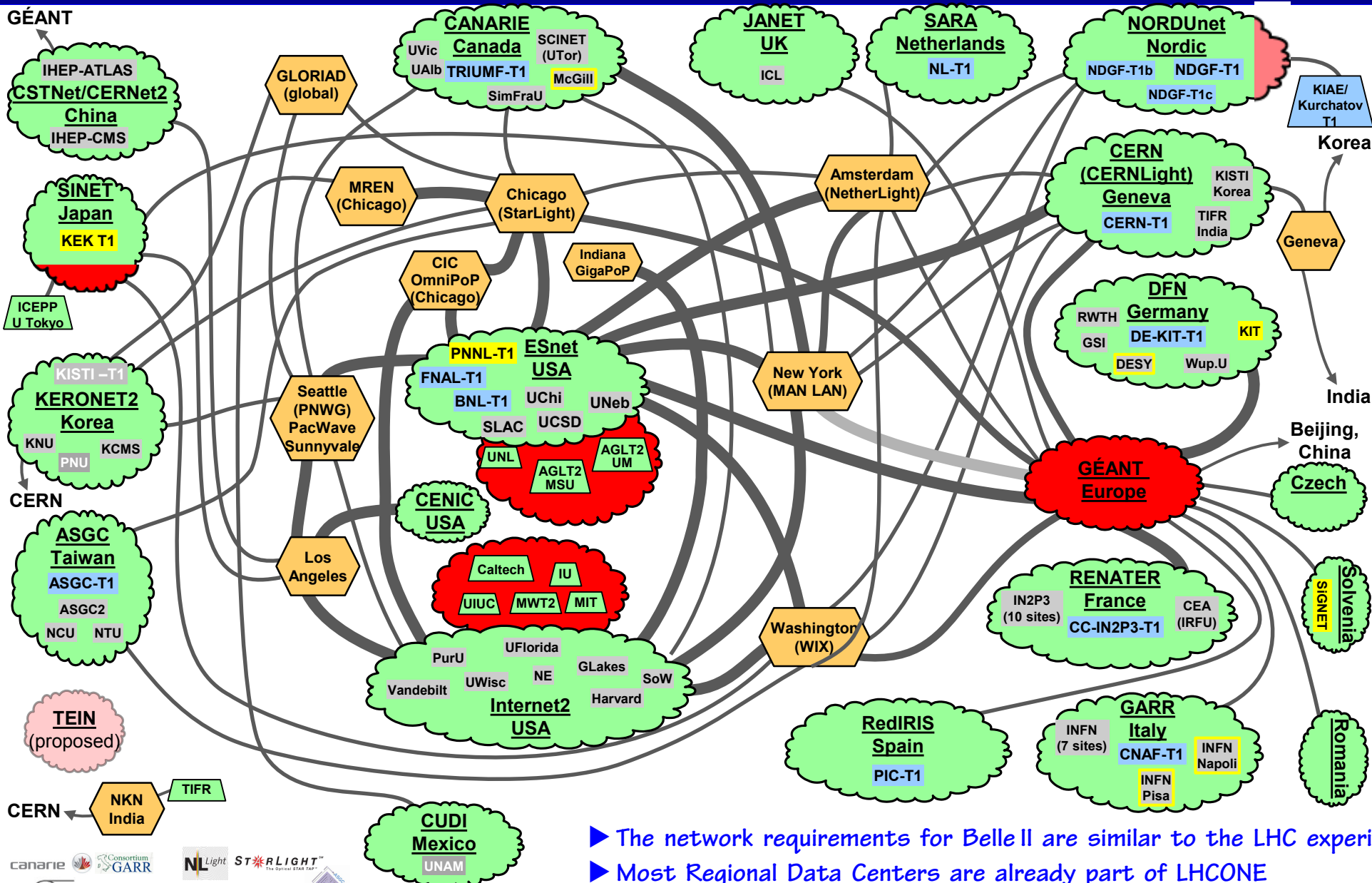
M. Schram (PNNL)



Marco Marletta



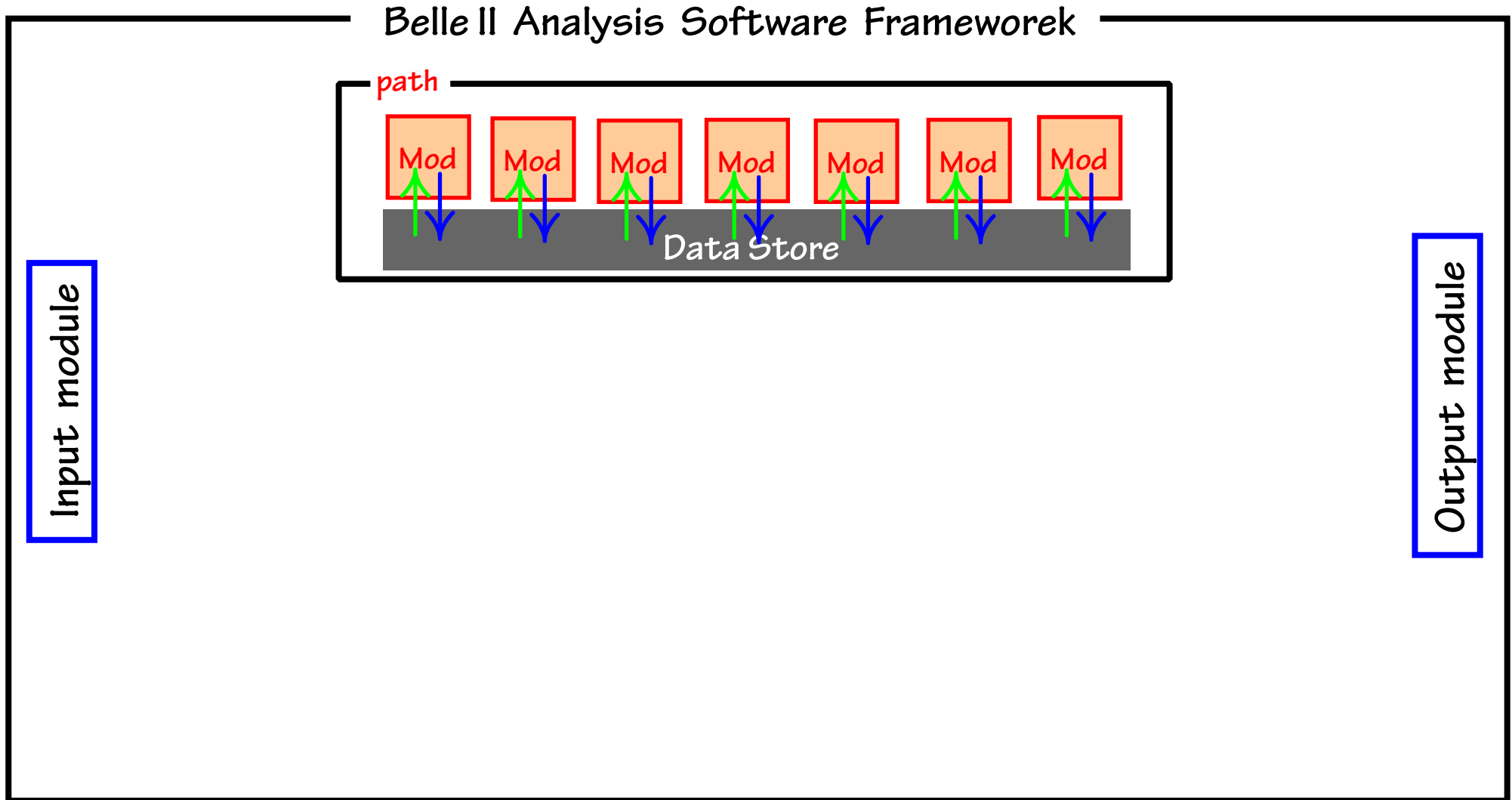
Europe side



- ▶ The network requirements for Belle II are similar to the LHC experiments.
- ▶ Most Regional Data Centers are already part of LHCONE
→ LHCONE has been extended to include Belle II

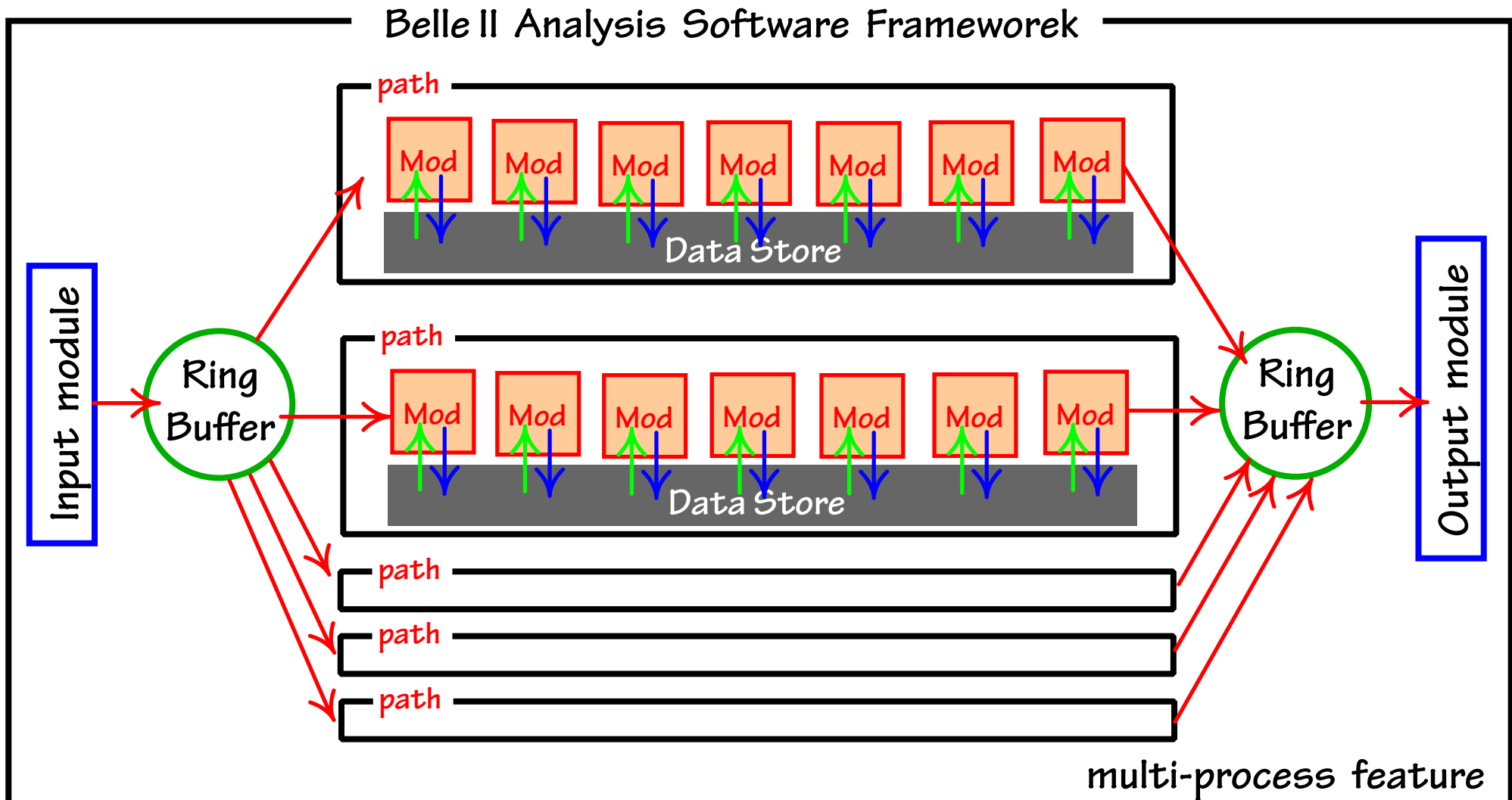


Common Software Framework



steering file is written in python

Common Software Framework



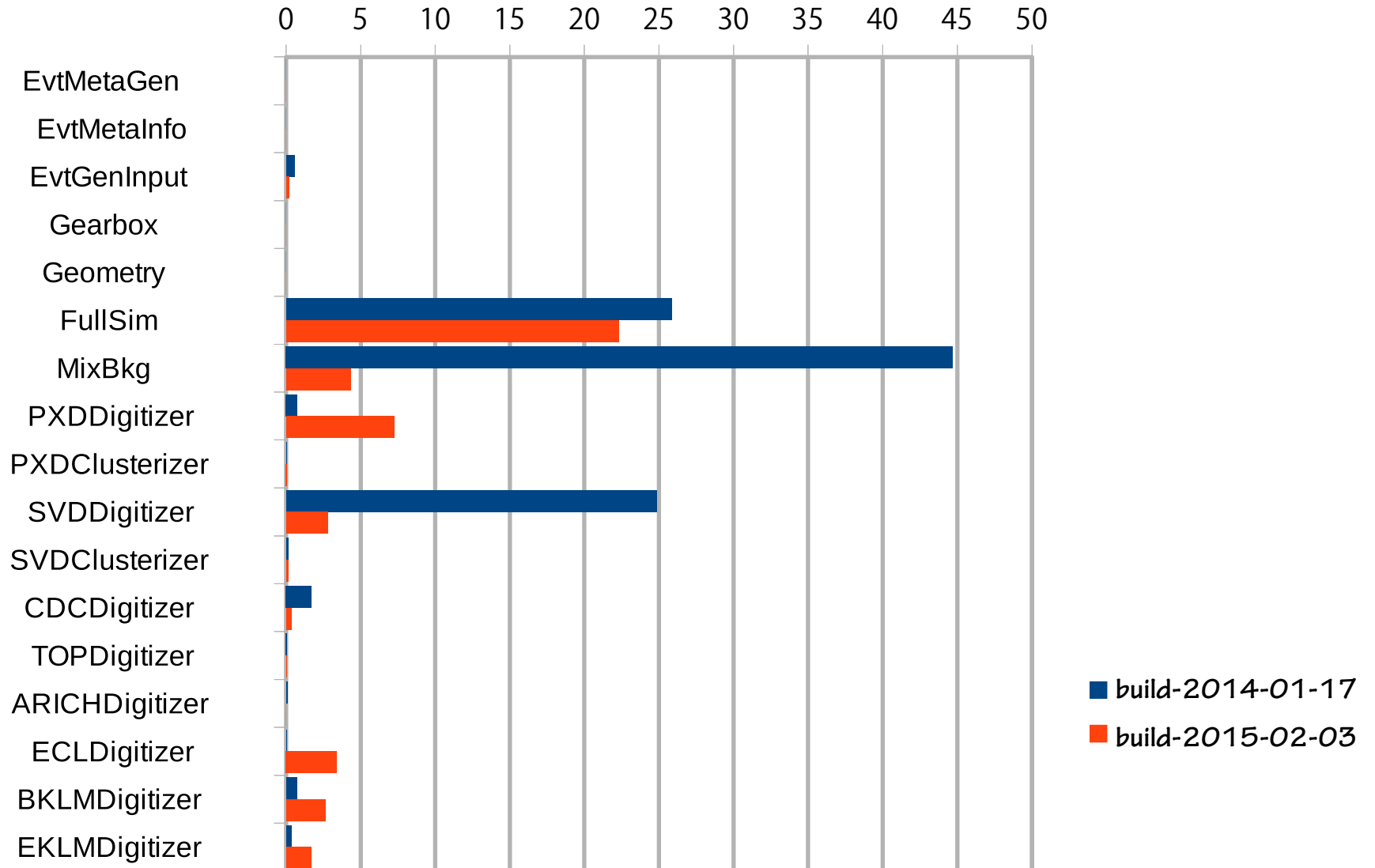
Geometry initialization
Event Generator

Detector simulation, Digitization,
Reconstruction (track finding/fitting, π^0/γ clustering)
Particle identification

steering file is written in python

Software update/improvement

CPU time for simulation [HEPSpec06*s]



100 HS06 * s/event → ~40 HS06 * s/event

Coordinated group skimming

From experience of Belle

- ▷ Many users cannot use resources effectively in skimming process
 - iterate the skimming with different selection
- ▷ Not so many users pay attention
 - whether or not the submitted jobs successfully finished
 - consideration of exp/run-dependences, log file check, ...
- ▷ Users' interests can be easily concentrated on certain datasets
 - in particular, the early stage of the experiment
 - data taken under good accelerator/detector condition
 - new datasets (new exp#/run#)
 - same input but many different outputs
- ▷ Usually, it takes too much time to carry out the skimming process and to finish by users

Job failure rate ↓

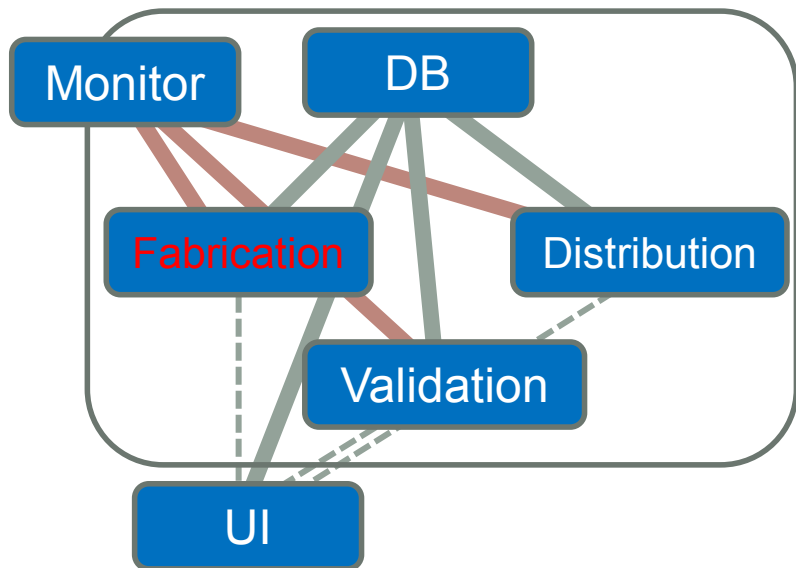
“coordinated” group skimming must make the physics analysis faster.

Production system

MC in early 2014 : no “Production System”

2.5 M jobs : manually submitted by the shifters

→ some mistakes e.g. same job submitted twice
wrong destination SE ...

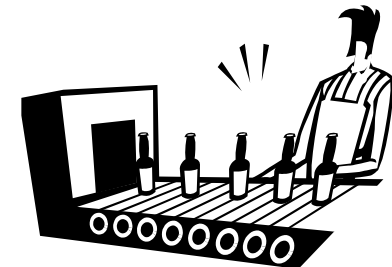


Production system, what we need...

Once a production manager defines a production parameters
software release, simulation channel, background, number of events, ...
↓
A system takes care of the rest and produces the output

MC in late 2014 : with “Prototype Production System”

4.7 M jobs : controlled by a single person



“Full Production System” yet to be developed

Summary

▶ Belle II time line

SuperKEKB accelerator commissioning (phase1) starts in early 2016

Phase2 (w/o VXD) starts in 2017

→ Belle II Distributed Computing should be ready

Phase3 (w/ Full detector) starts in 2018

→ raw data distribution starts

▶ MC mass production & Data transfer challenge

the basic concept of the Belle II computing model was proven

but still we have many things to do (e.g. Data distribution, user distributed analysis, etc.)

▶ Efforts to utilize the limited hardware resources AMAP are on-going toward the physics run.

multicore jobs (memory size ↓)

Tuning/Optimization of software (memory size ↓, shorter CPU time)

Coordinated group skimming (less human-error)

Production system

Belle II Oral/Poster Presentations

Software

Oliver FROST @ Track2, April 13
Cellular Automaton based Track Finding
for the Central Drift Chamber of Belle II

Marko BRACKO @ Track3, April 13
The Belle II Conditions Database

Marko STARIC @ Track2, April 14
Physics Analysis Software Framework for Belle II

Thomas HAUTH @ Poster A, 233
Software Development at Belle II

Tobias SCHLÜTER @ Poster A, 469
The GENFIT Library for Track Fitting and its Performance in Belle II

Leo PIILONEN @ Poster B, 320
The Simulation Library of the Belle II Software System

Tadeas BILKA @ Poster B, 469
Alignment and calibration of Belle II tracking detectors

Computing

Hideki MIYAKE @ Track4, April 13
Belle II production system

Randy SOBIE @ Track7, April 14
Utilizing cloud computing resources for Belle II

Yuji KATO @ Poster A, 337
Job monitoring on DIRAC for Belle II distributed computing

Chia-Ling HSU @ Poster A, 468
The Belle II analysis on Grid

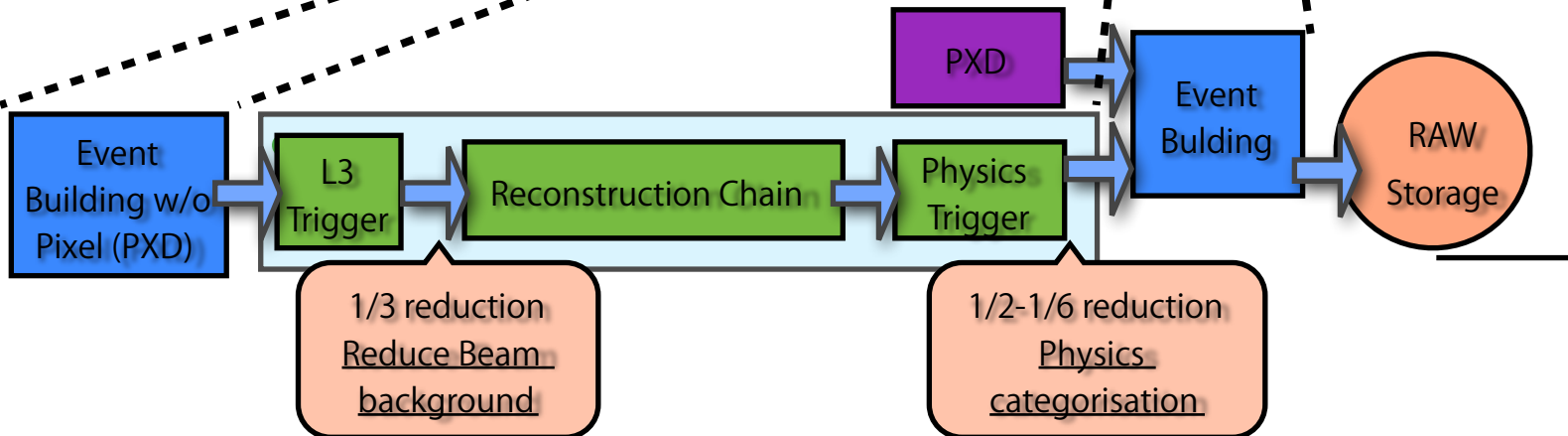
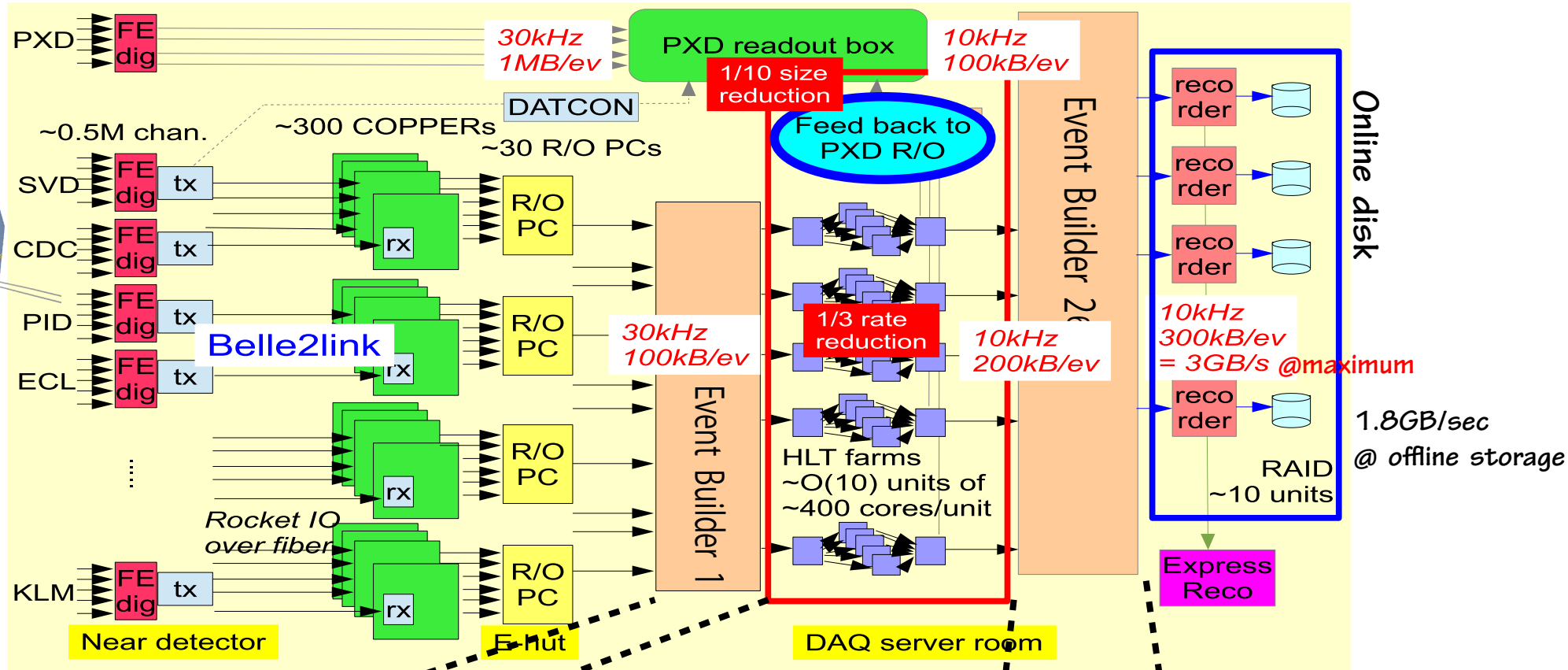
Rafal Zbigniew GRZYMKOWSKI @ Poster B, 342
Belle II public and private clouds management in VMDIRAC system.

Kiyoshi HAYASAKA @ Poster B, 314
Monitoring system for the Belle II distributed computing

Jae-Hyuck KWAK @ Poster B, 466
Improvement of AMGA Python Client Library for the Belle II Experiment

Geun Chul PARK @ Poster B, 313
Directory Search Performance Optimization of AMGA
for the Belle II Experiment

DAQ + physics trigger



AMGA metadata catalogue

ARDA Metadata Grid Application
 – Metadata server for GRID environment
 (EMI product)



Metadata : data of data
 LFN, run range, software version...

