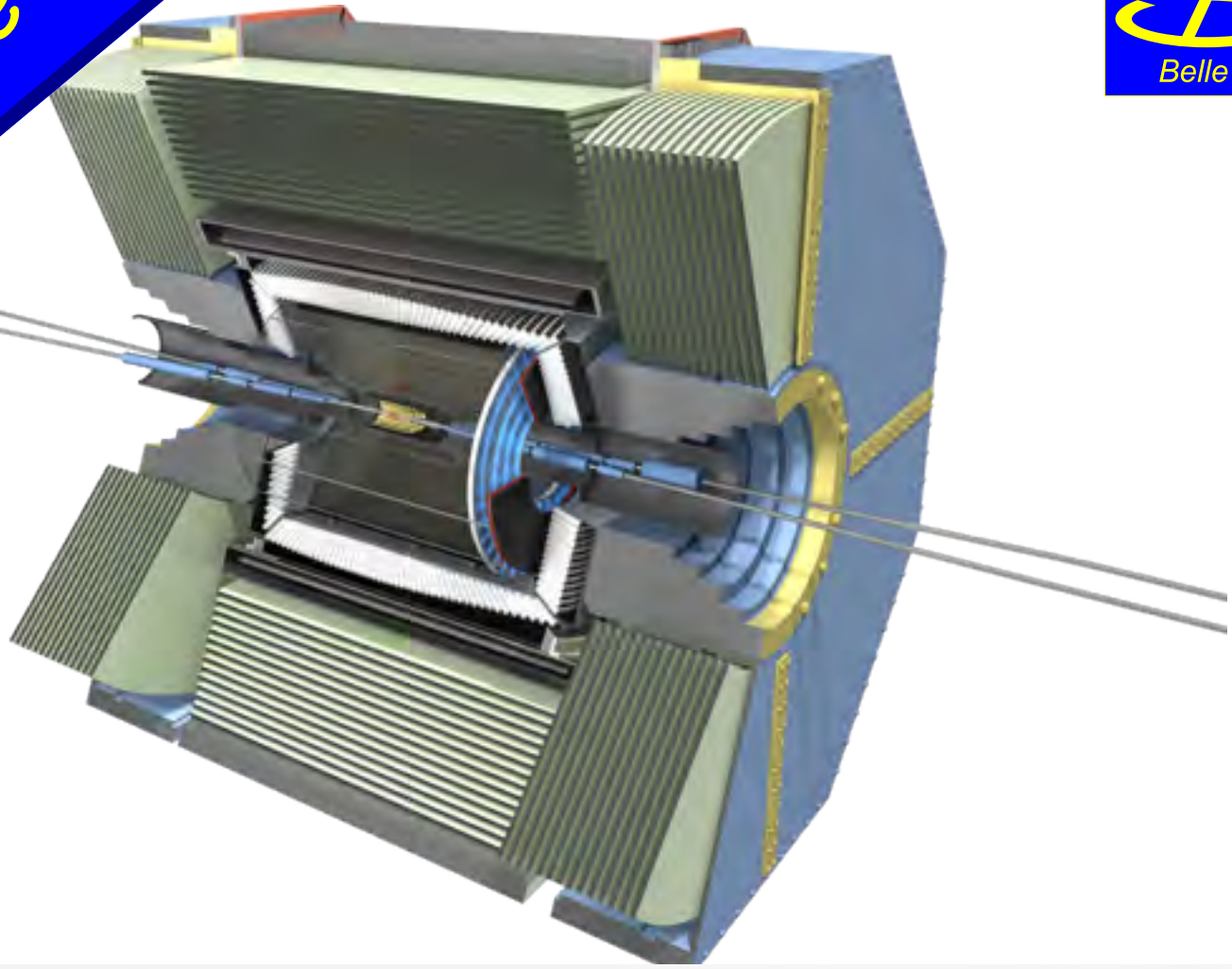


# Computing at Belle II



Takanori Hara (KEK)  
Mar. 2nd, 2012

ISGC @ Academia Sinica



## International Symposium on Grids & Clouds 2012

*Convergence, Collaboration, Innovation*

26 February - 2 March 2012, Academia Sinica, Taipei, Taiwan

# the Universe we live in now

To explain the dominance of matter in the universe

the Sakharov conditions

- . Departure from thermal equilibrium.
- . Baryon number violation.
- . CP-symmetry violation.

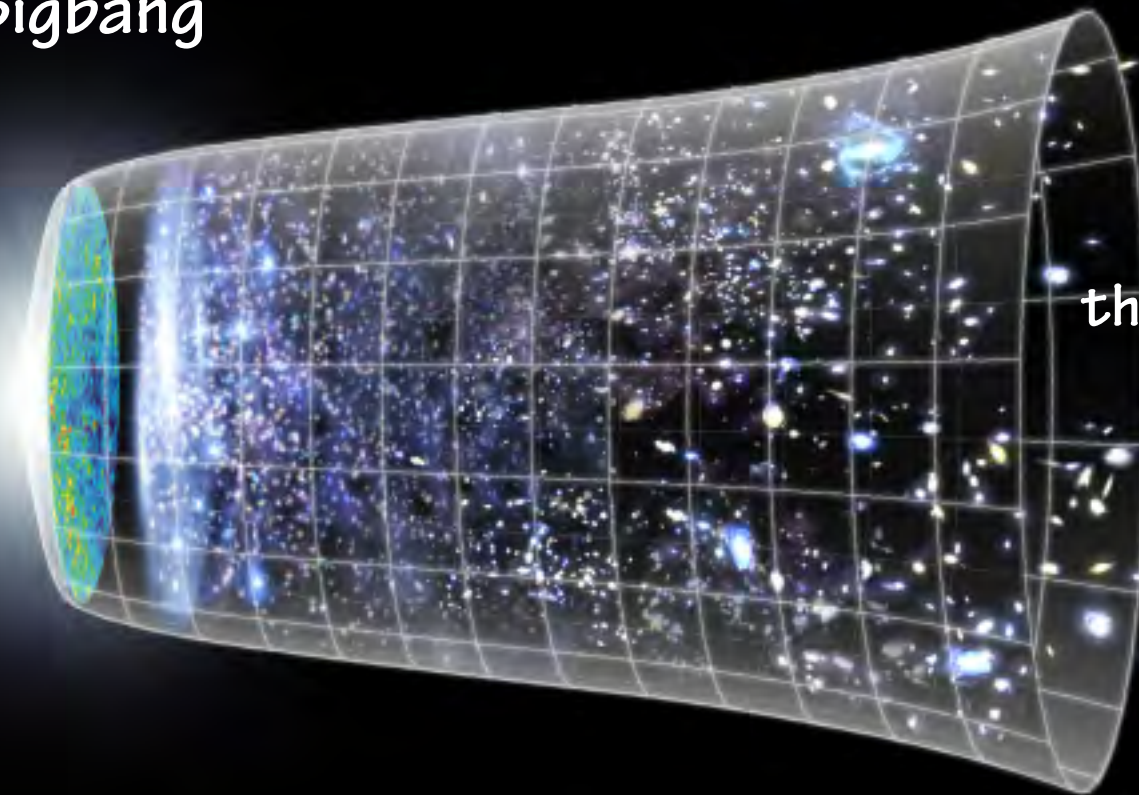
= an asymmetry between  
the behavior of matter and antimatter

Kobayashi-Maskawa theory

# How to observe

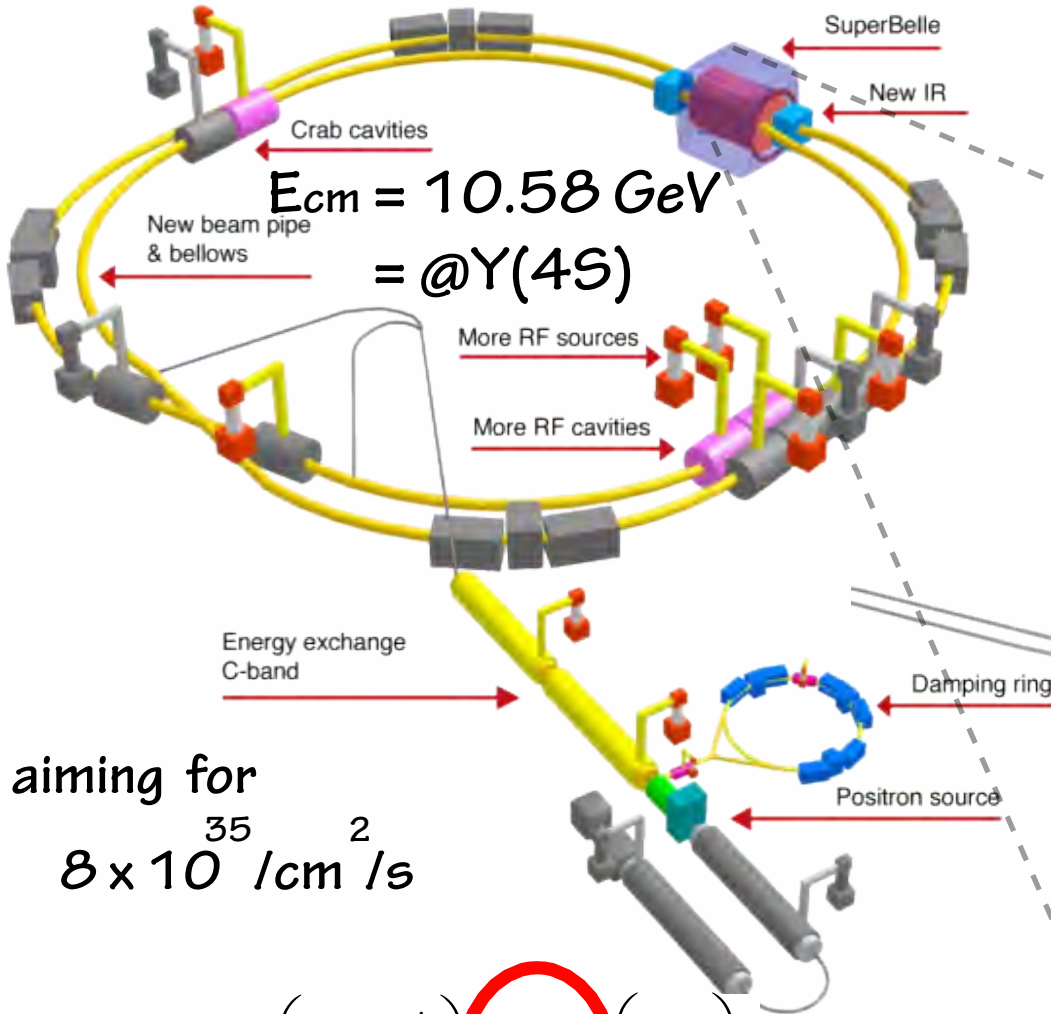
the Bigbang

a 13.7 billion years later  
(=now)



Now  
the Matter-dominated  
Universe

- . using telescopes (Hubble telescope) → after 0.5 billion years
- . checking the Cosmic Microwave Background → after 0.4 million years
- . using accelerator (KEKB, LHC, ...) → at around  $10^{-12}$  sec



aiming for  
 $8 \times 10^{35} / \text{cm}^2 / \text{s}$

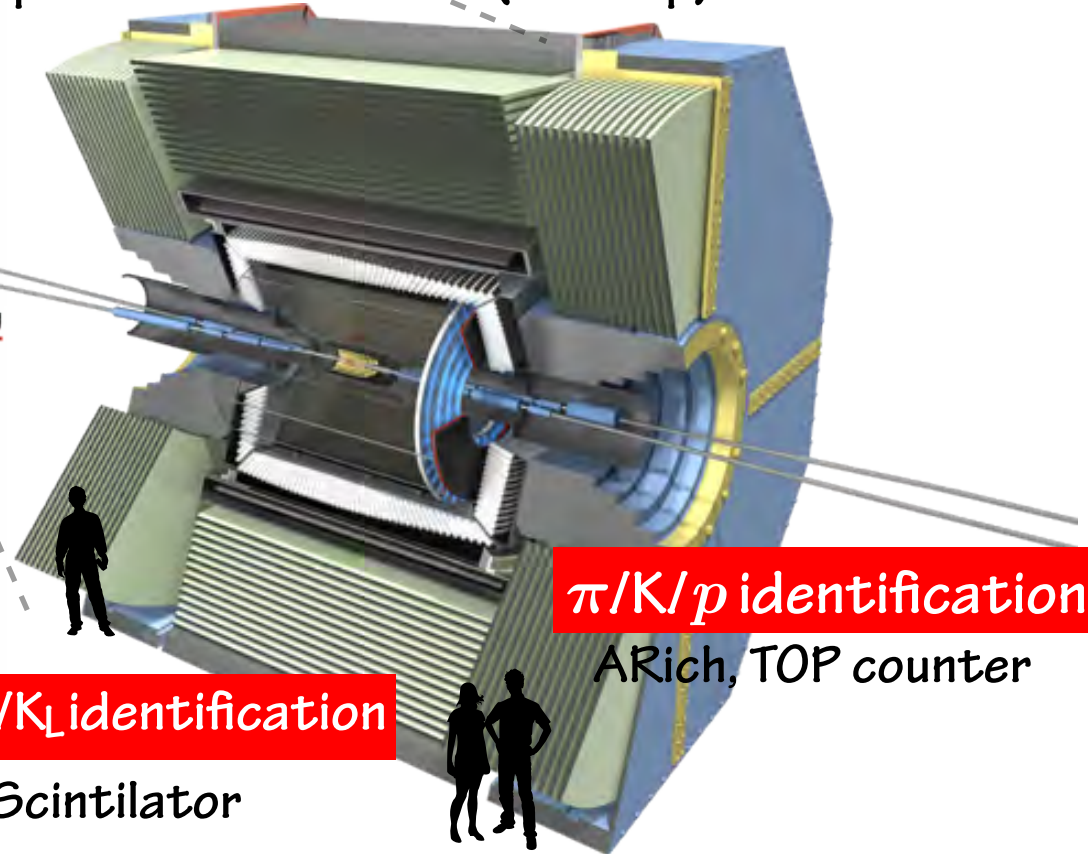
$$L = \frac{\gamma_{\pm}}{2e r_e} \left( 1 + \frac{\sigma_y^*}{\sigma_x^*} \frac{I_{\pm} \xi_{\pm y}}{\beta_y^*} \left( \frac{R_L}{R_y} \right) \right)$$

**tracking/vertexing**

- small-cell Drift Chamber
- Silicon Strip det.
- longer lever arm
- Pixel detector

**$e/\gamma$  detection**

pure-CsI calorimeter (end-cap)



**$\pi/K/p$  identification**

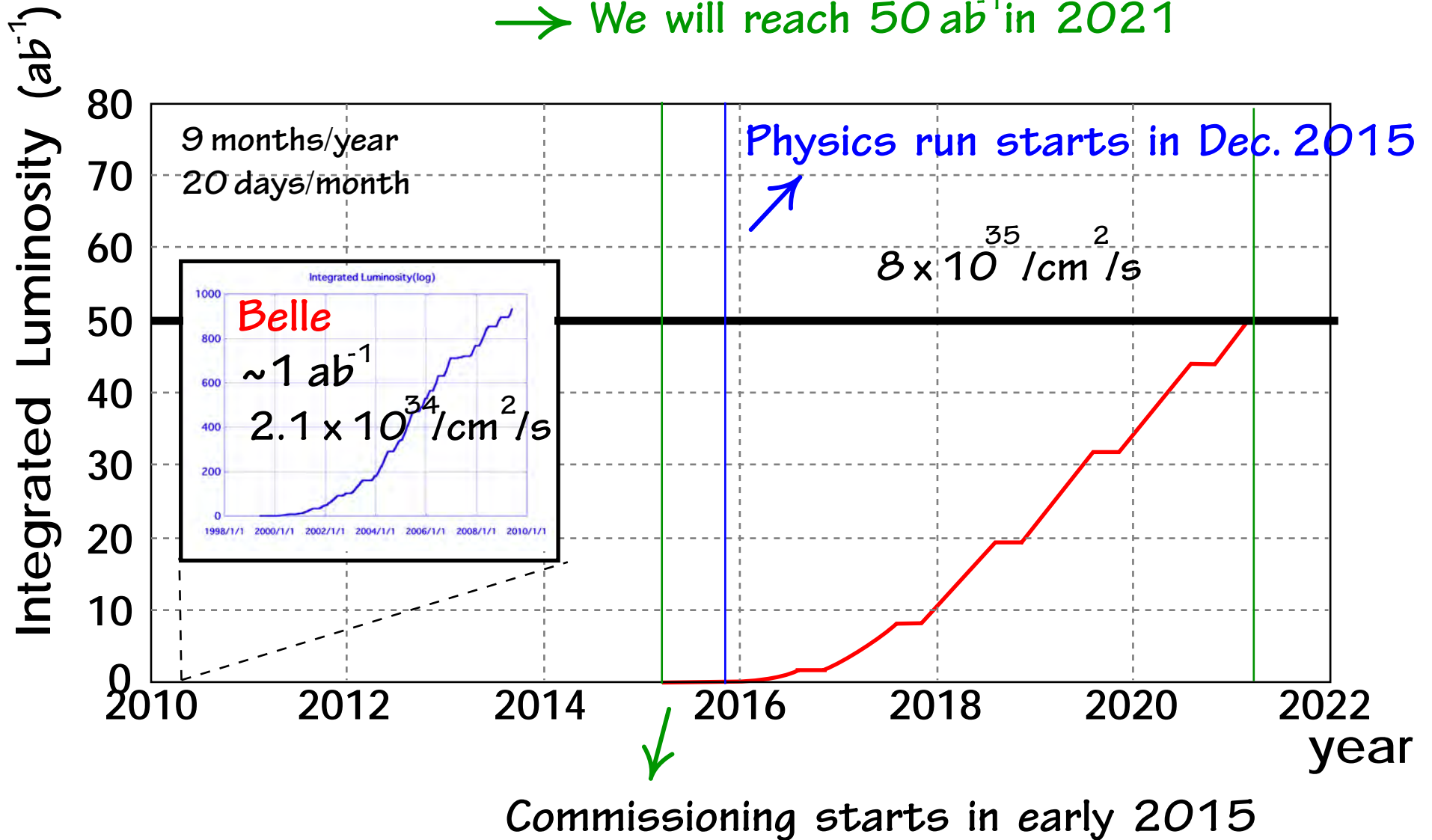
ARich, TOP counter

**$\mu/K_L$  identification**

Scintillator

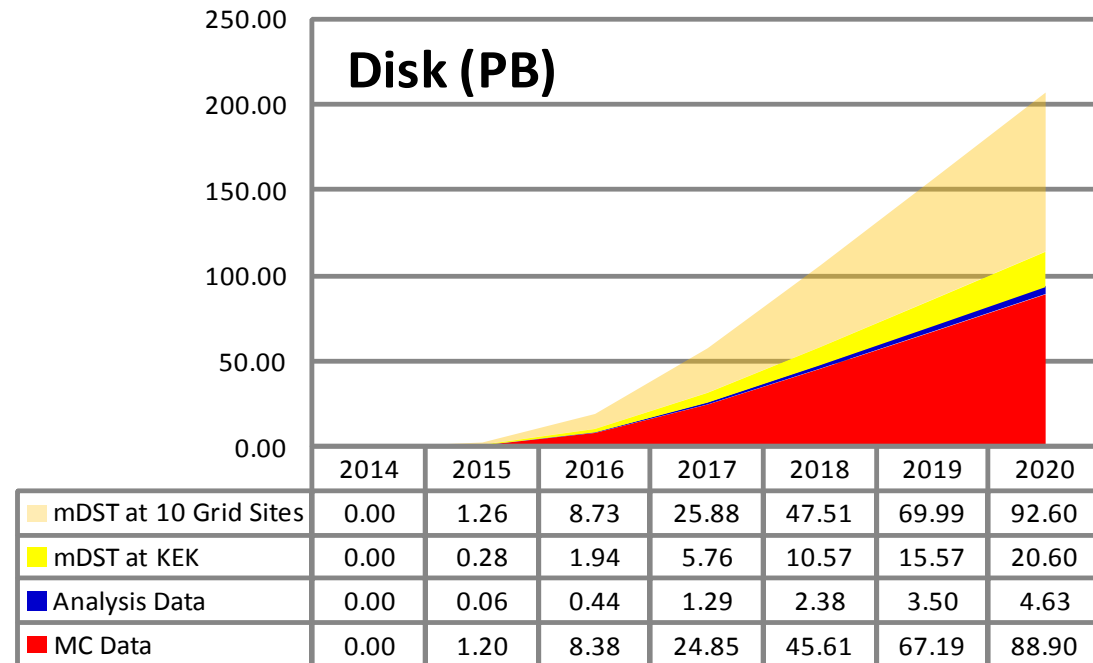
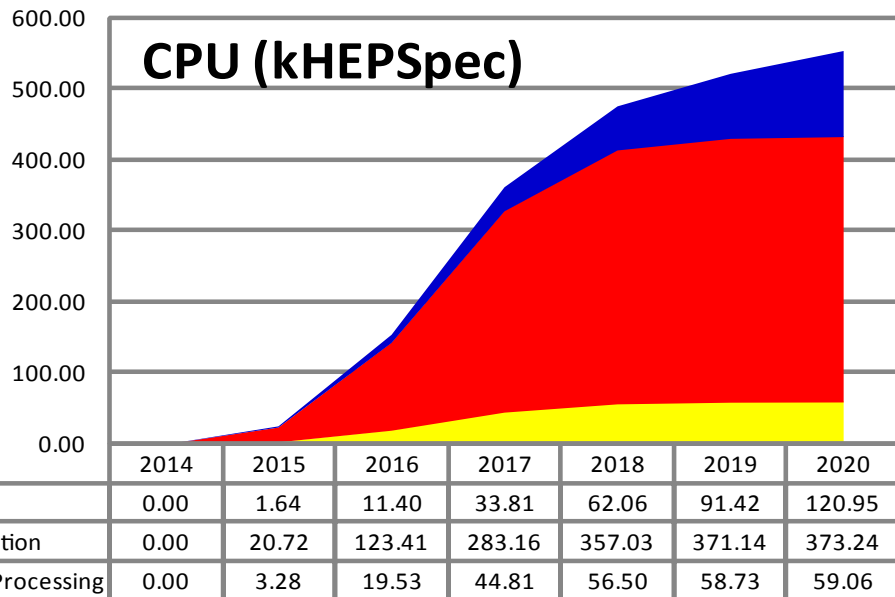
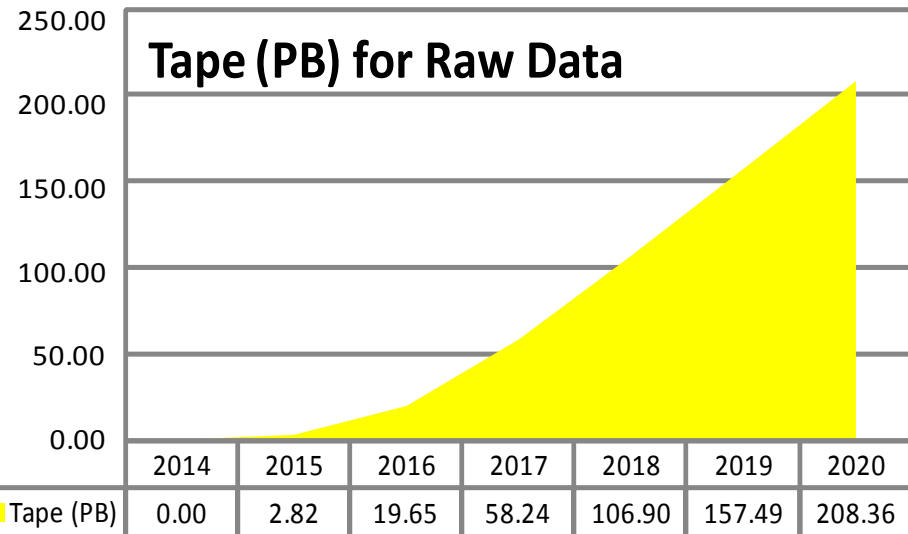
50  $ab^{-1}$  by the end of 2020JFY =  $\times 50$  present

→ We will reach 50  $ab^{-1}$  in 2021



Preliminary estimates depend on many unknown parameters

- . accelerator performance
- . data reduction
- . performance of simulation/reconstruction
- . analysis requirements, ...



Experiment	Event Size [kB]	Rate [Hz]	Rate [MB/s]
<i>High rate scenario for Belle II DAQ:</i>			
Belle II	300	6,000	1,800
<i>LCG TDR (2005):</i>			
ALICE (HI)	12,500	100	1,250
ALICE (pp)	1,000	100	100
ATLAS	1,600	200	320
CMS	1,500	150	225
LHCb	25	2,000	50

# Belle II Collaboration



Japan : 129  
Asia : 214 (=55%)

China	: 15	Malaysia	: 1
India	: 14	Taiwan	: 20
Korea	: 32	Viet Nam	: 3





# Belle II GRID sites

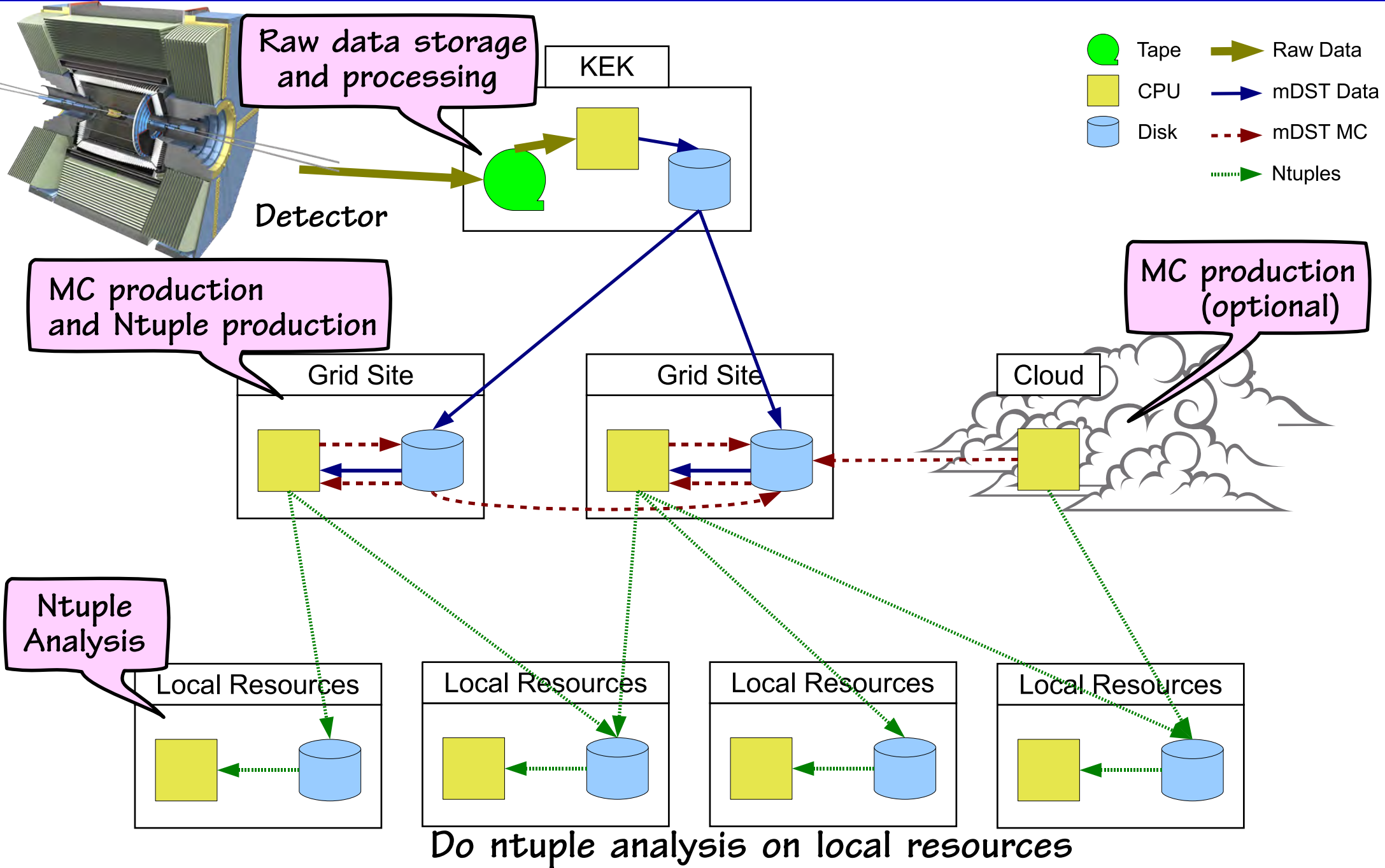
GRID  
middleware

gLite

OSG

- ◆ Australia : LHC Tier2/3, Belle VO, Cloud system
- ◆ Austria : LHC Tier2
- ◆ China (IHEP) : LHC Tier2, DIRAC server
- ◆ Czech Republic : LHC Tier2, Belle VO
- ◆ Germany : LHC Tier1/2, Belle VO
- ◆ India : LHC Tier2, Belle II data center planned
- ◆ Japan (KEK) : Belle VO
- ◆ Korea (KISTI) : LHC Tier2, Belle VO
- ◆ Poland : LHC Tier2/3, Belle VO, Cloud system
- ◆ Russia : LHC Tier2
- ◆ Slovenia : LHC Tier2, Belle VO
- ◆ Taiwan : LHC Tier1/2
- ◆ USA : OSG @ PNNL planned, Belle VO @ OSG exists

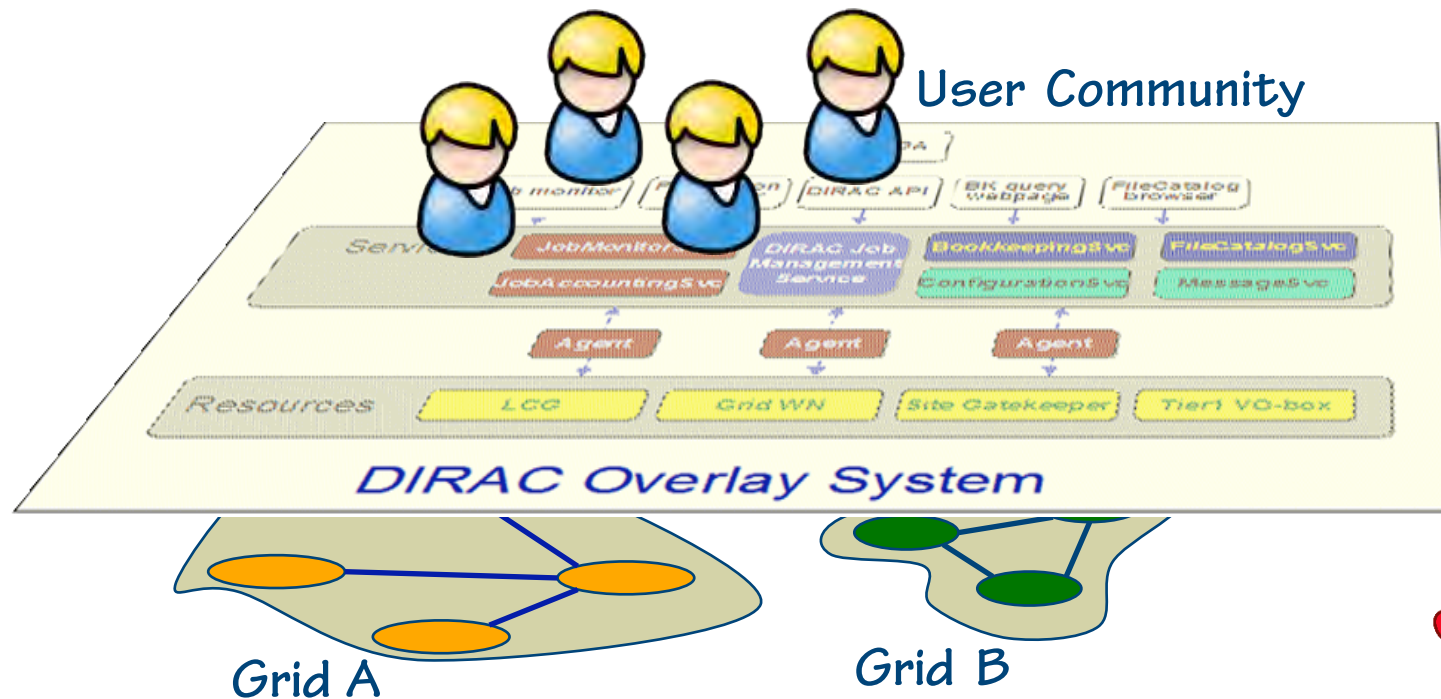
# Belle II Computing Model

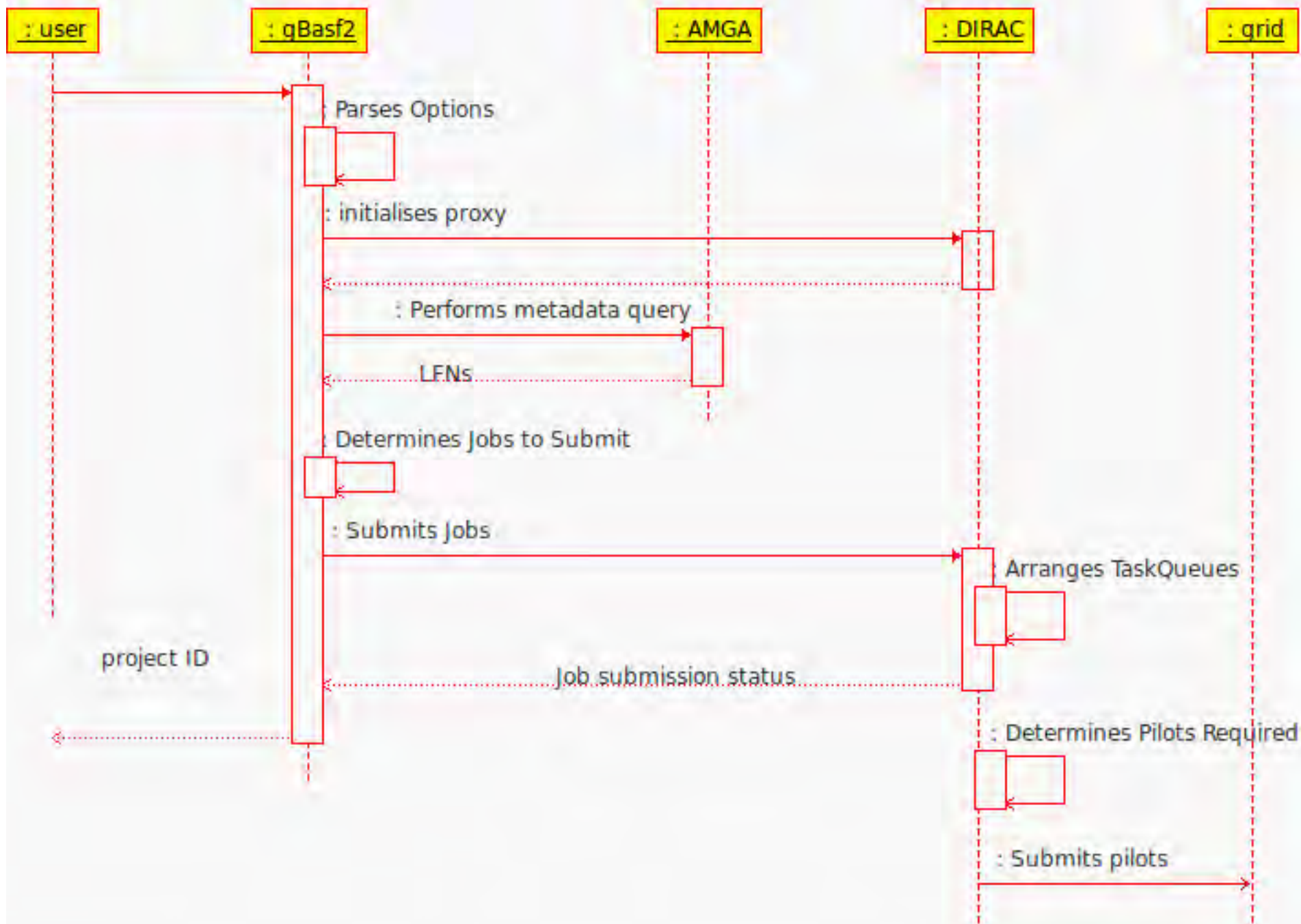


◆ DIRAC (developed by LHCb)

Distributed Infrastructure with Remote Agent Control

- Pilot jobs
- Modular structure that enabled it possible to submit to different backends.
- It provides many features that we would have to develop





# User Interface : gbasf2

python steering file, same as for ofline basf2 jobs,  
but with additional parameters for the grid job

## Basf2 Steering File options

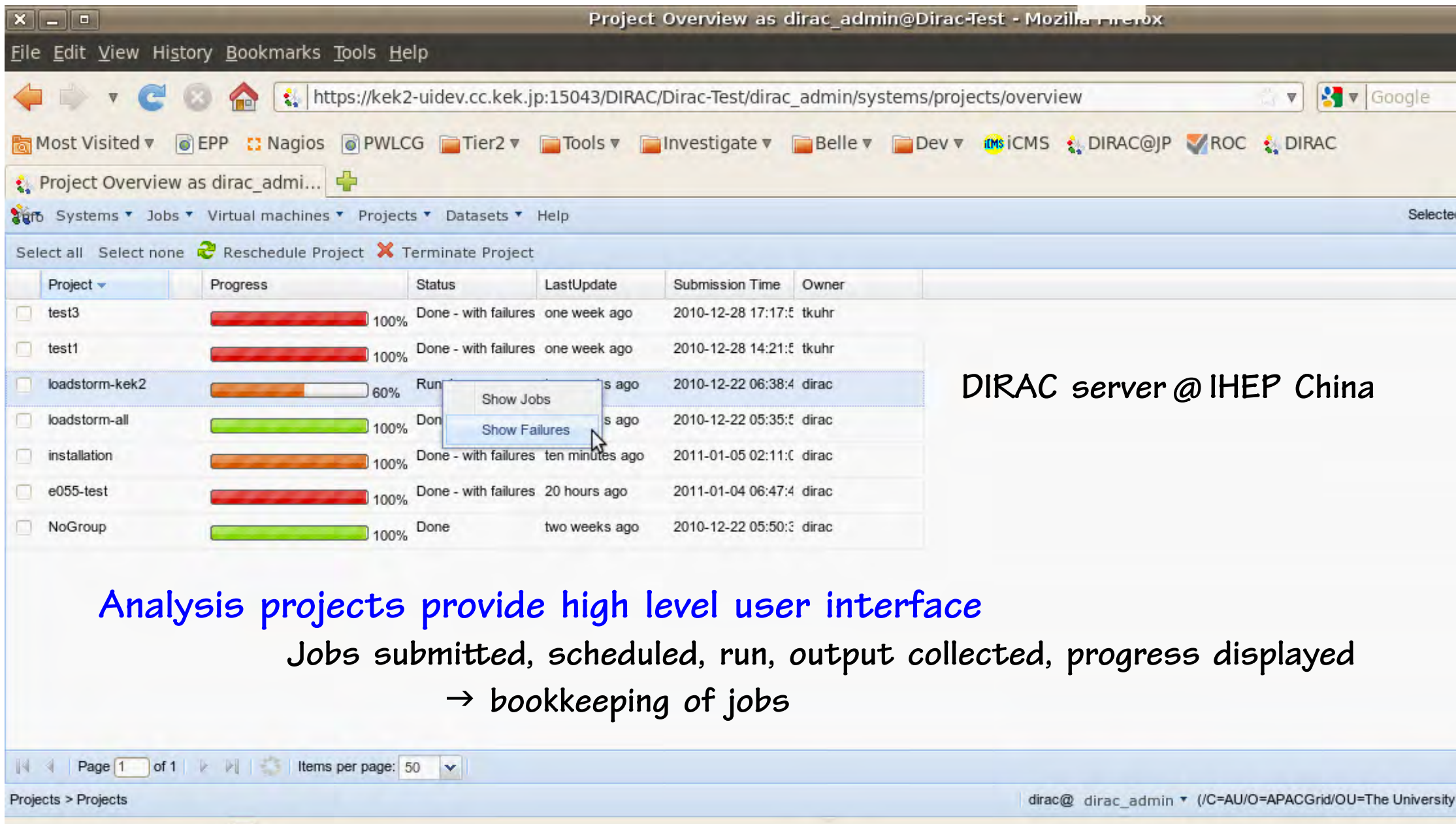
The default configuration option for gBasf2 is to set a number of variables in your normal basf2 steering file:

```
#####  
# gBaf2 configuration      #  
#####  
#Name for project  
project='e055-test'  
# (optional) Job priority [0-10]  
priority='1'  
#Experiments (comma separated list)  
experiments='13,57'  
#Metadata query  
query='id > 10 and id < 15'  
#Type of Data ('data' or 'MC')  
type='data'  
#estimated Average Events per Minute (eg Mcprod = 40)  
evtpermin='45'  
# (optional) Files to be sent with the job  
inputsandboxfiles = 'file1.txt,file2.txt'  
# (optional) max events - the maximum number of events to use  
maxevents = '100000'
```

basf2 : Belle II  
Analysis  
Software  
Framework

You can then invoke gBasf2 using the steering file and it will do the rest:

```
./gbasf2.py -s steering_file.py
```



Project Overview as dirac\_admin@Dirac-Test - Mozilla Firefox

File Edit View History Bookmarks Tools Help

https://kek2-uidev.cc.kek.jp:15043/DIRAC/Dirac-Test/dirac\_admin/systems/projects/overview

Most Visited EPP Nagios PWLCC Tier2 Tools Investigate Belle Dev iCMS DIRAC@JP ROC DIRAC

Project Overview as dirac\_admin@Dirac-Test

Systems Jobs Virtual machines Projects Datasets Help

Select all Select none Reschedule Project Terminate Project

Project	Progress	Status	LastUpdate	Submission Time	Owner
<input type="checkbox"/> test3	100%	Done - with failures	one week ago	2010-12-28 17:17:5	tkuhr
<input type="checkbox"/> test1	100%	Done - with failures	one week ago	2010-12-28 14:21:5	tkuhr
<input type="checkbox"/> loadstorm-kek2	60%	Run	s ago	2010-12-22 06:38:4	dirac
<input type="checkbox"/> loadstorm-all	100%	Don	s ago	2010-12-22 05:35:5	dirac
<input type="checkbox"/> installation	100%	Done - with failures	ten minutes ago	2011-01-05 02:11:0	dirac
<input type="checkbox"/> e055-test	100%	Done - with failures	20 hours ago	2011-01-04 06:47:4	dirac
<input type="checkbox"/> NoGroup	100%	Done	two weeks ago	2010-12-22 05:50:3	dirac

DIRAC server @ IHEP China

Analysis projects provide high level user interface

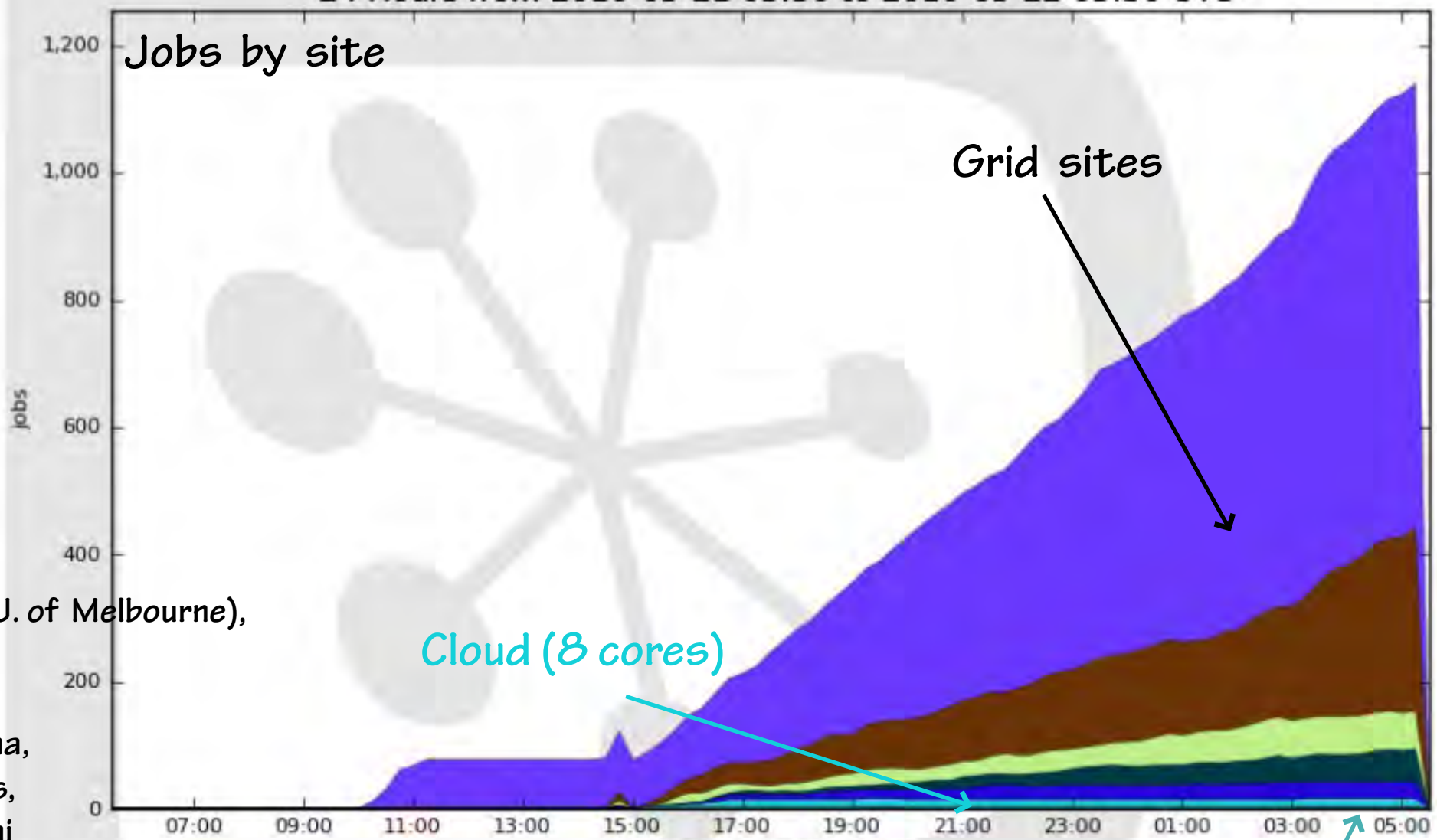
Jobs submitted, scheduled, run, output collected, progress displayed

→ bookkeeping of jobs

Page 1 of 1 Items per page: 50

Projects > Projects dirac@ dirac\_admin (/C=AU/O=APACGrid/OU=The University

24 Hours from 2010-05-21 03:30 to 2010-05-22 03:30 UTC

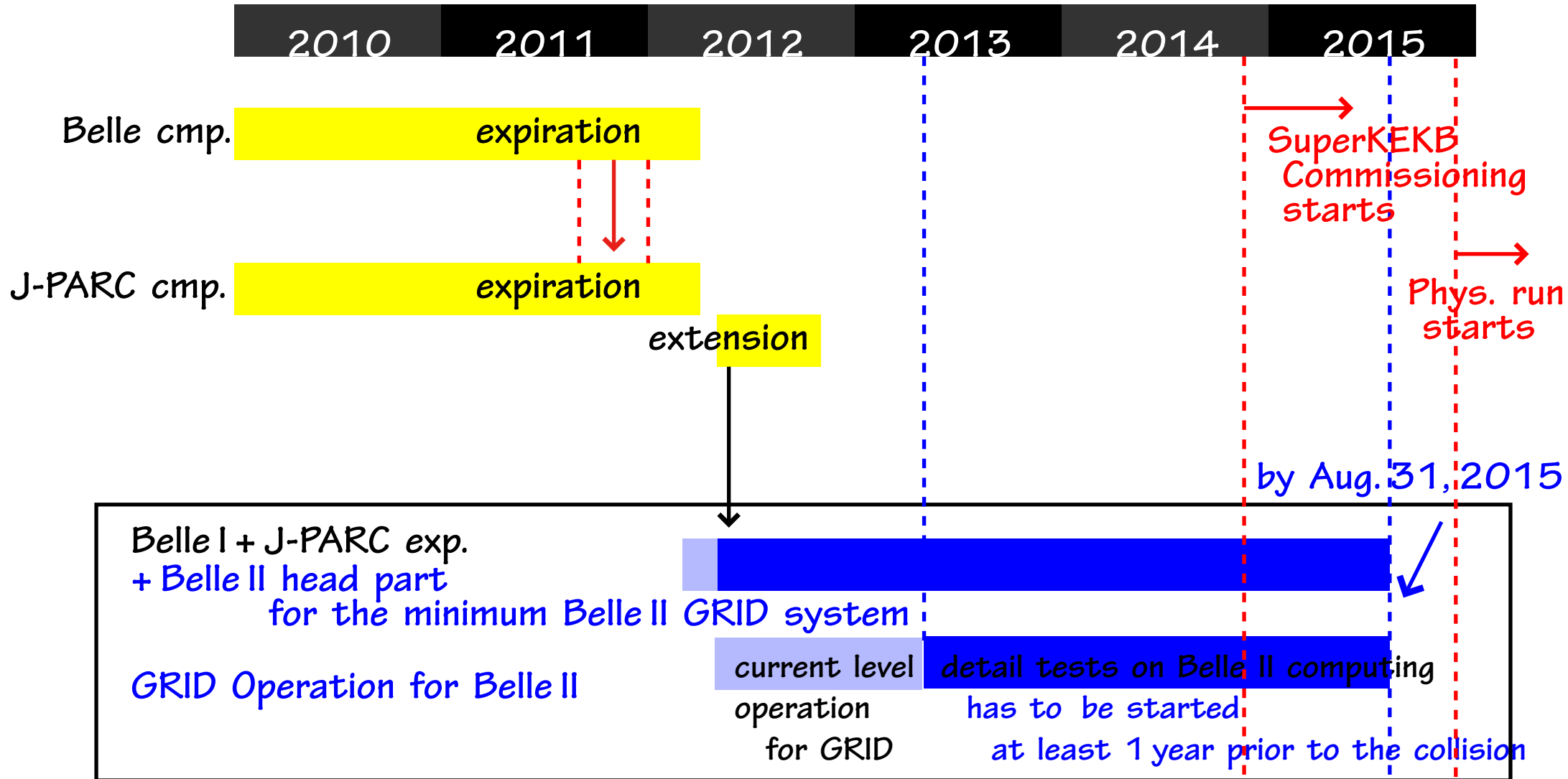


T.Fifield (U. of Melbourne),  
and  
M.Sevior  
A.Carmona,  
A.Casajus,  
R.Graciani  
& The DIRAC team

LCG.KEK2.jp	699.0	LCG.CYFRONET.pl	52.0	DIRAC.Barcelona.es	4.0
LCG.GRIDKA.de	289.0	LCG.CESNET.cz	27.0	ANY	1.0
LCG.IJS.si	60.0	DIRAC.Amazon.us	10.0		

Local resources (8 cores)

# Schedule

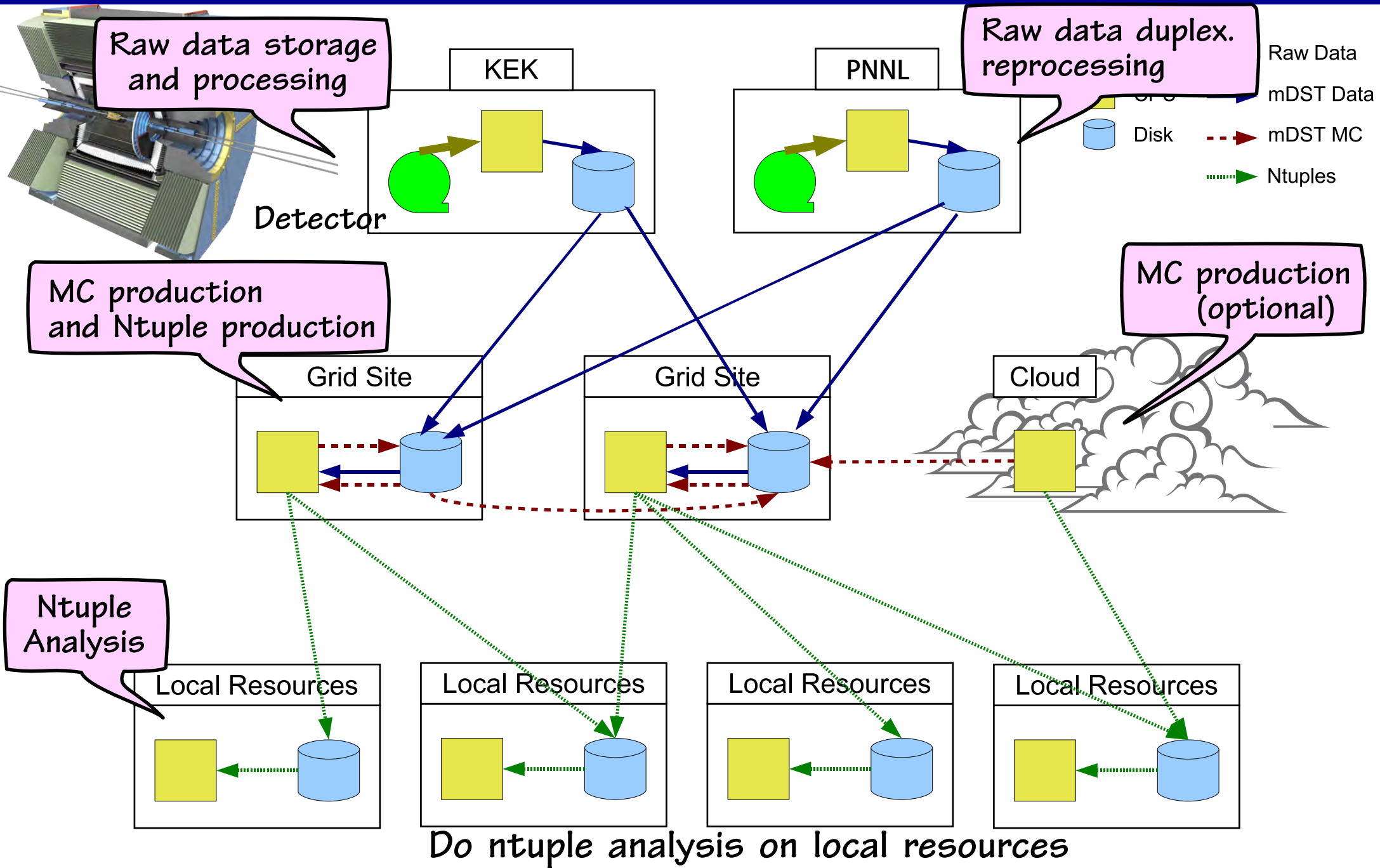




# Specification

	Belle current system	Final specification
Storage	Dell, SONY Petasite	DDN SFA 10000, IBM TS3500
	Disk/Tape	1.5PB/3PB
	# of tape drivers	~10
CPU	# of Cores	~4000
	and memories	80 workgroup servers 3840 (computing servers) (4GB) 3/5 used for Belle
	internal Disk drive	300 (4GB)+240(8GB) 3540 (computing servers) 3084 (4GB) 456 (8GB)
	OS	SL5
Band Width	transfer rate (single connection)	200MB/s disk: par 1 process par 1 thread staging disk (HSM): par 1 process par 1 thread tape (HSM): par 1 process par 1 thread: 150MB/s
	total throughput	50GB/s disk storage ctrl - CN staging disk ctrl (HSM) - CN

# Belle II Computing Model

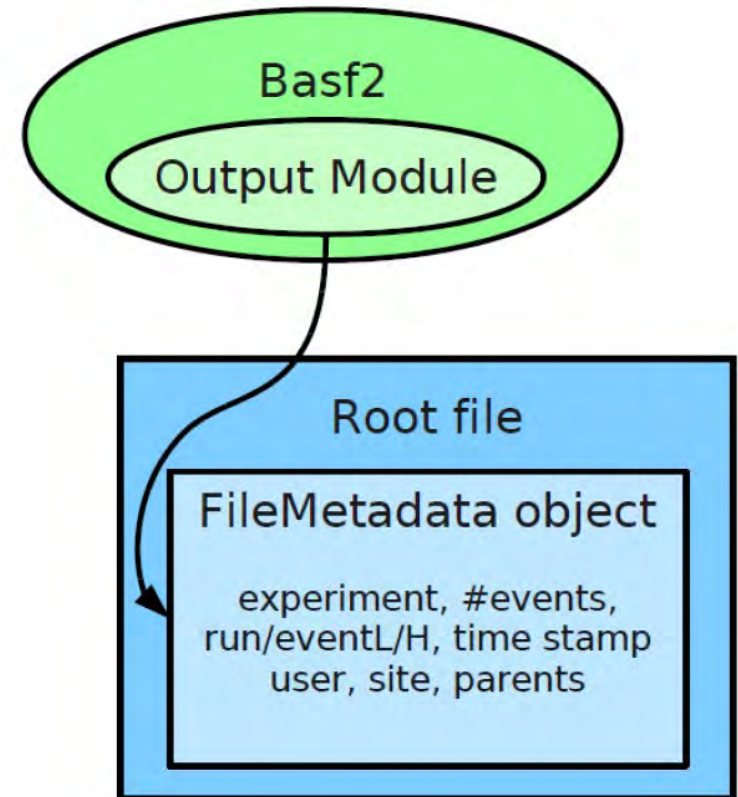


# Data Registration Tool

Metadata has to be extracted from the output files

Tool for Belle 1 data implemented

```
hep2.kisti.re.kr - PuTTY
LFN : //f3/belle/bdata/exp13.charm-00/evtgen-charm-00-all-e000013r001627-b200308
07_1600.mdst
Starting : 1316418750
data Type[real(1)/MC(2): 2
MC type : 1003
stream : 0
Exp No: 13
Run No: 1627
File No : f1627
ID No : 1627
EXP date : 1050501
EXP time : 162516
parentID :
Belle Library version : b20030807_1600
Belle Detector version : 1
End :1316418754
Run status : good
hostname : hep2.kisti.re.kr
Create Date : Mon Sep 19 16:52:34 KST 2011
Min evt : 1
Max evt : 8294
```

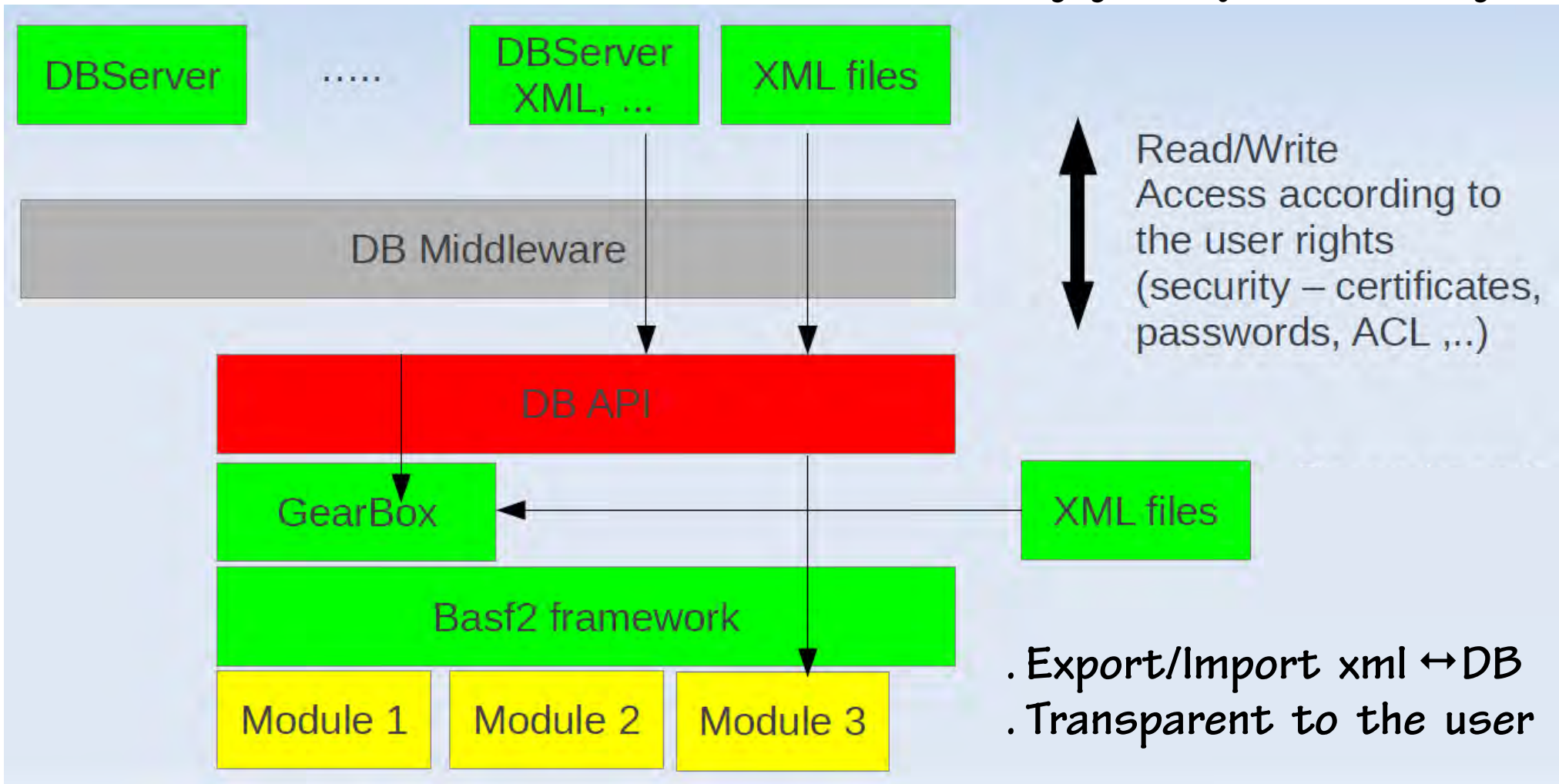


Metadata automatically stored in Belle II files

# Database

DB type: *Logger DB for online environment monitoring*  
*Configuration DB for data needed to start a run*  
*Condition DB for all data needed offline*

*e.g. geometry, calibration, alignment*



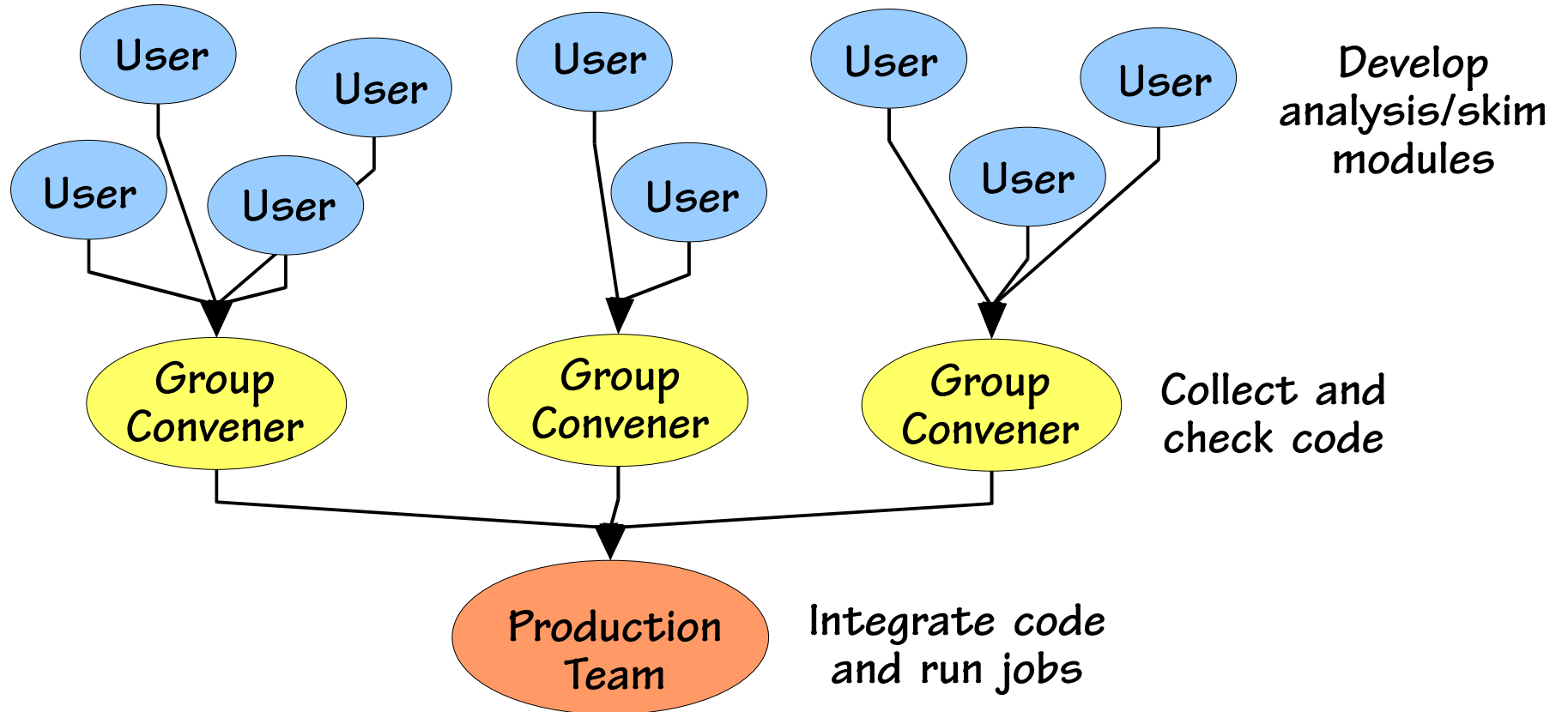
Offline software (basf2) installation is needed  
to run Belle II jobs on grid sites

Current solution:      Installation jobs, validation jobs  
Next step:              Investigate solutions provided by DIRAC  
Mid term:                Use Cern VM-FS

Needs support of sites,  
already installed in Ljubljana and Melbourne

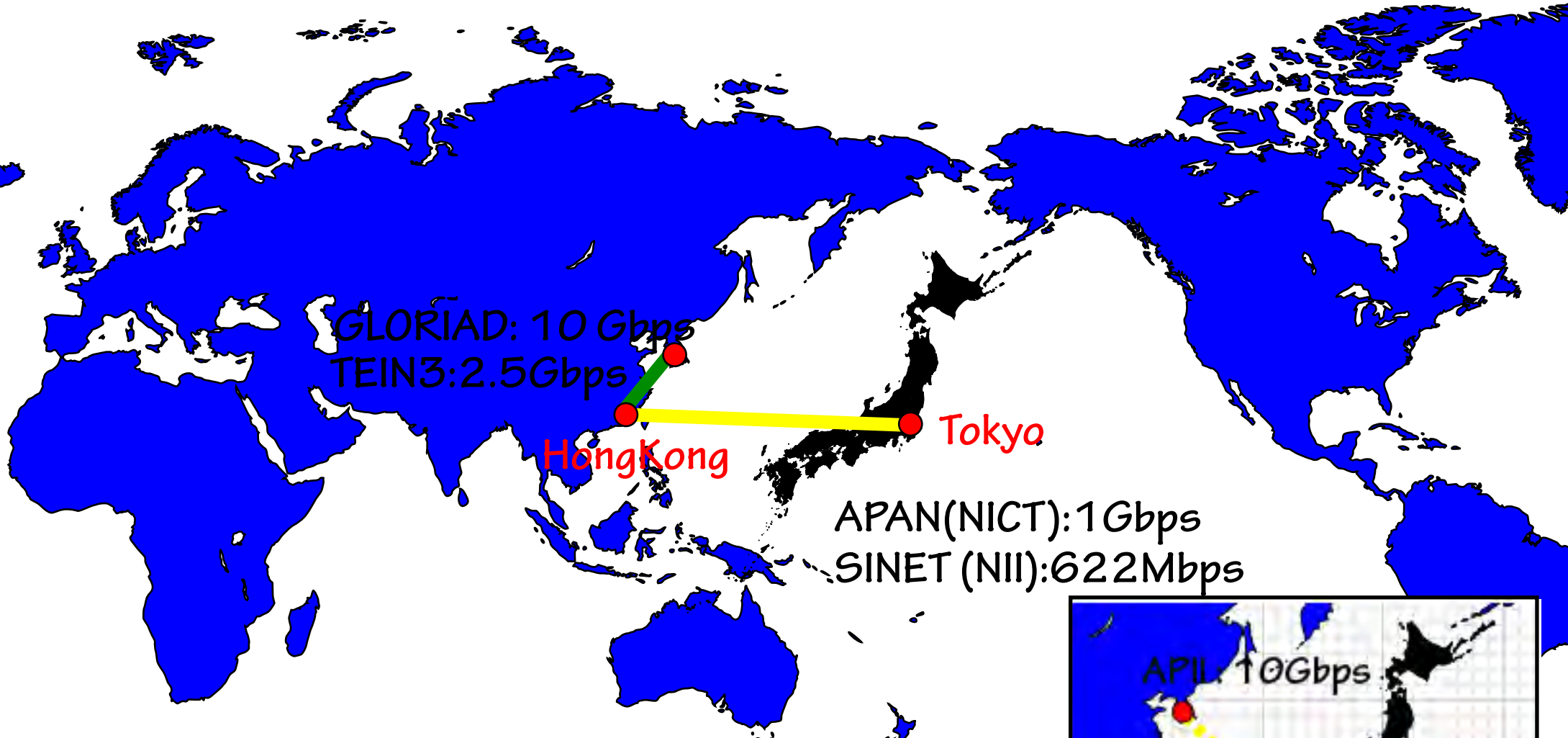
Server at CERN ?

# Organized Analysis



Problem: inefficient resource usage by many users  
complexity of grid environment for many users

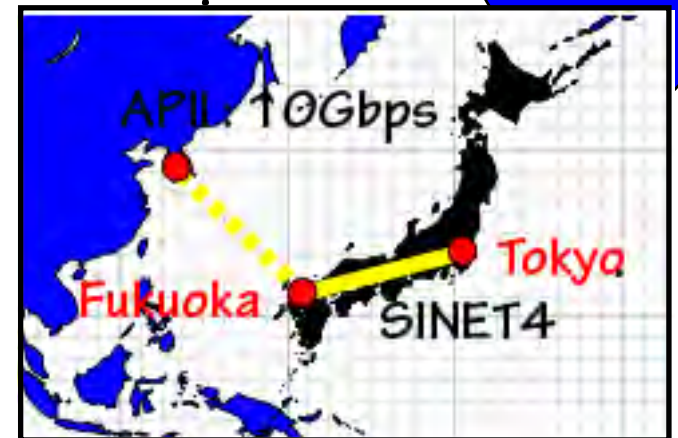
same as for official generic MC and signal MC

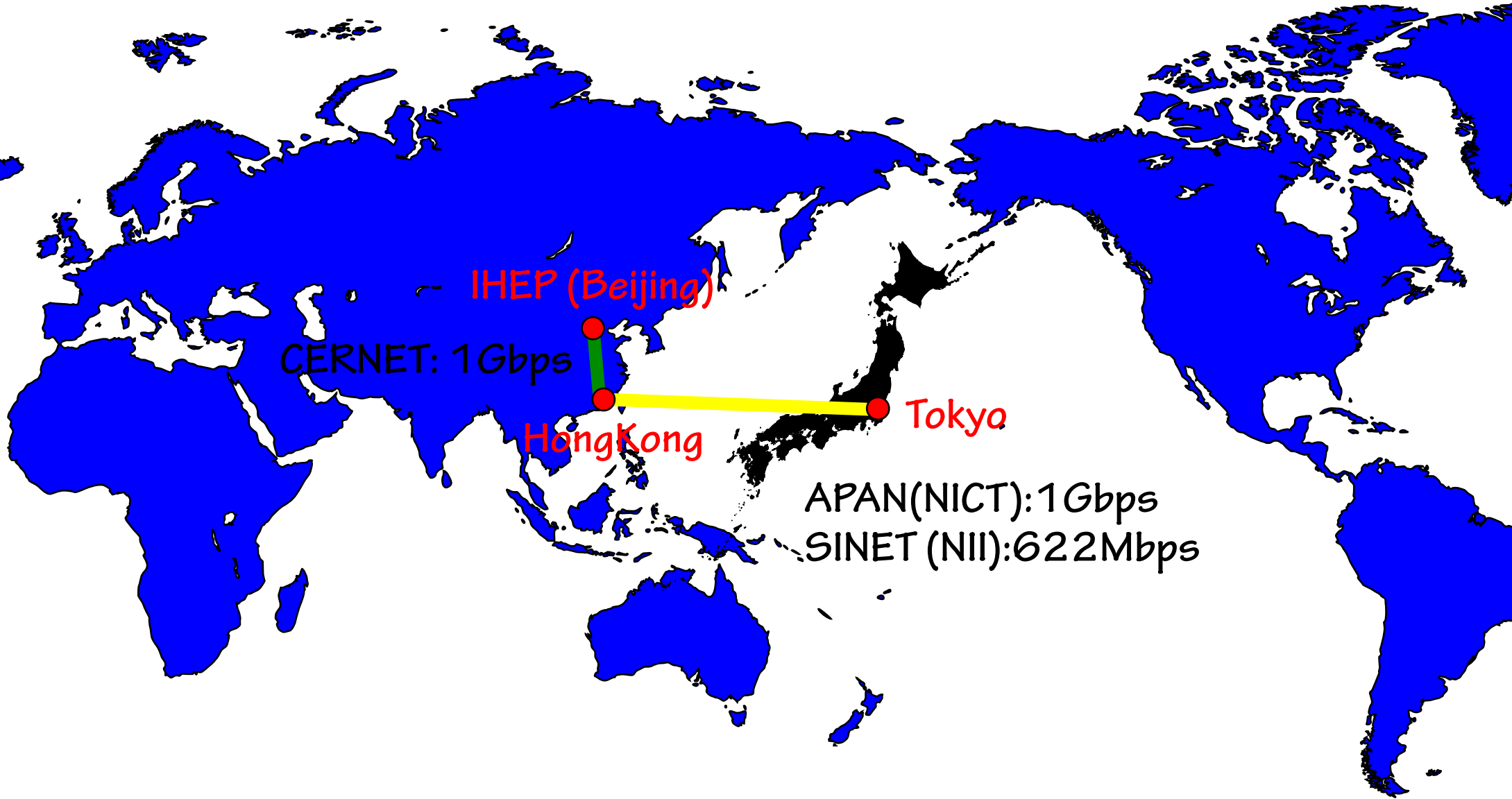


inside Korea (KOREN x KREONet)

KOREN : universities

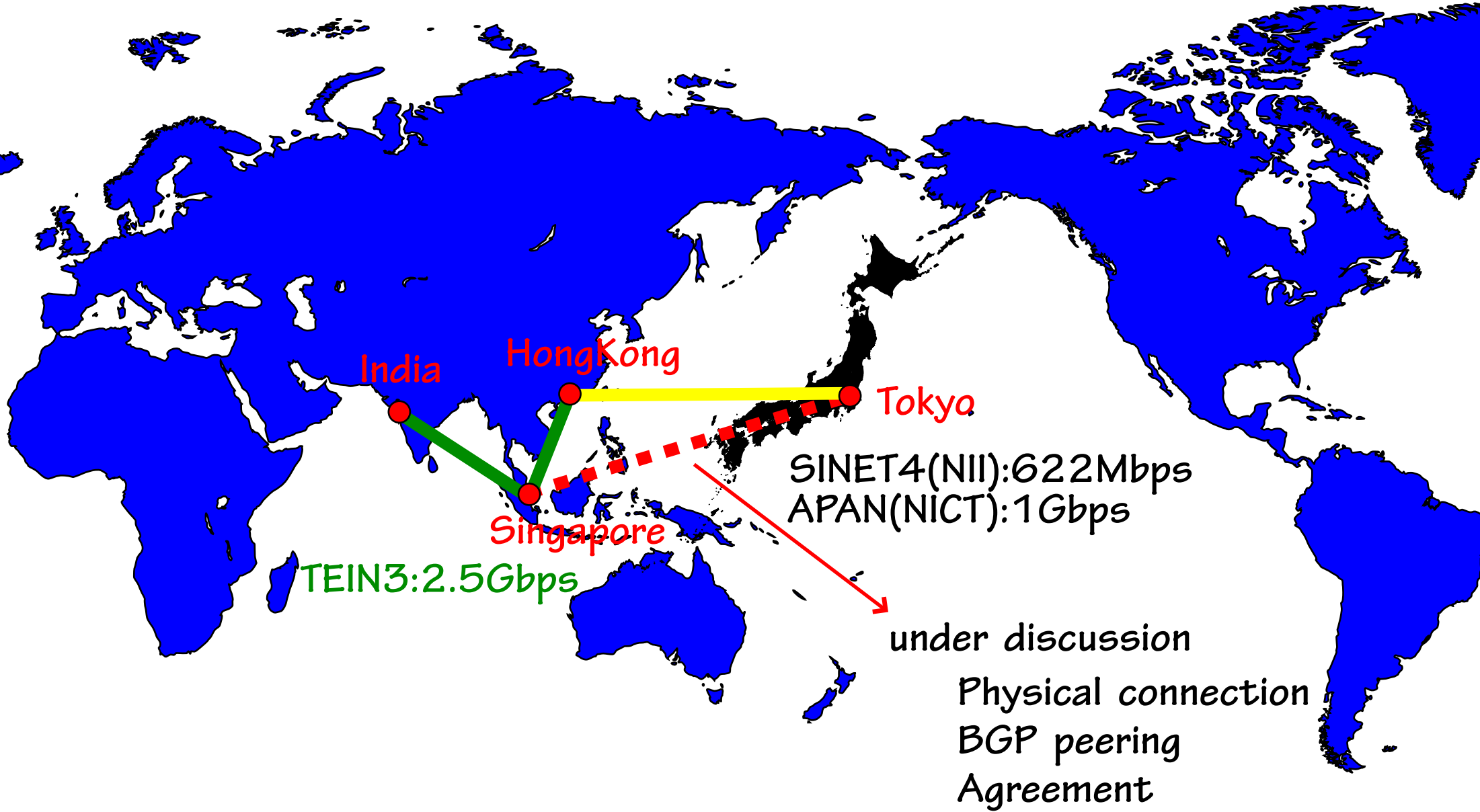
KREONet : KISTI

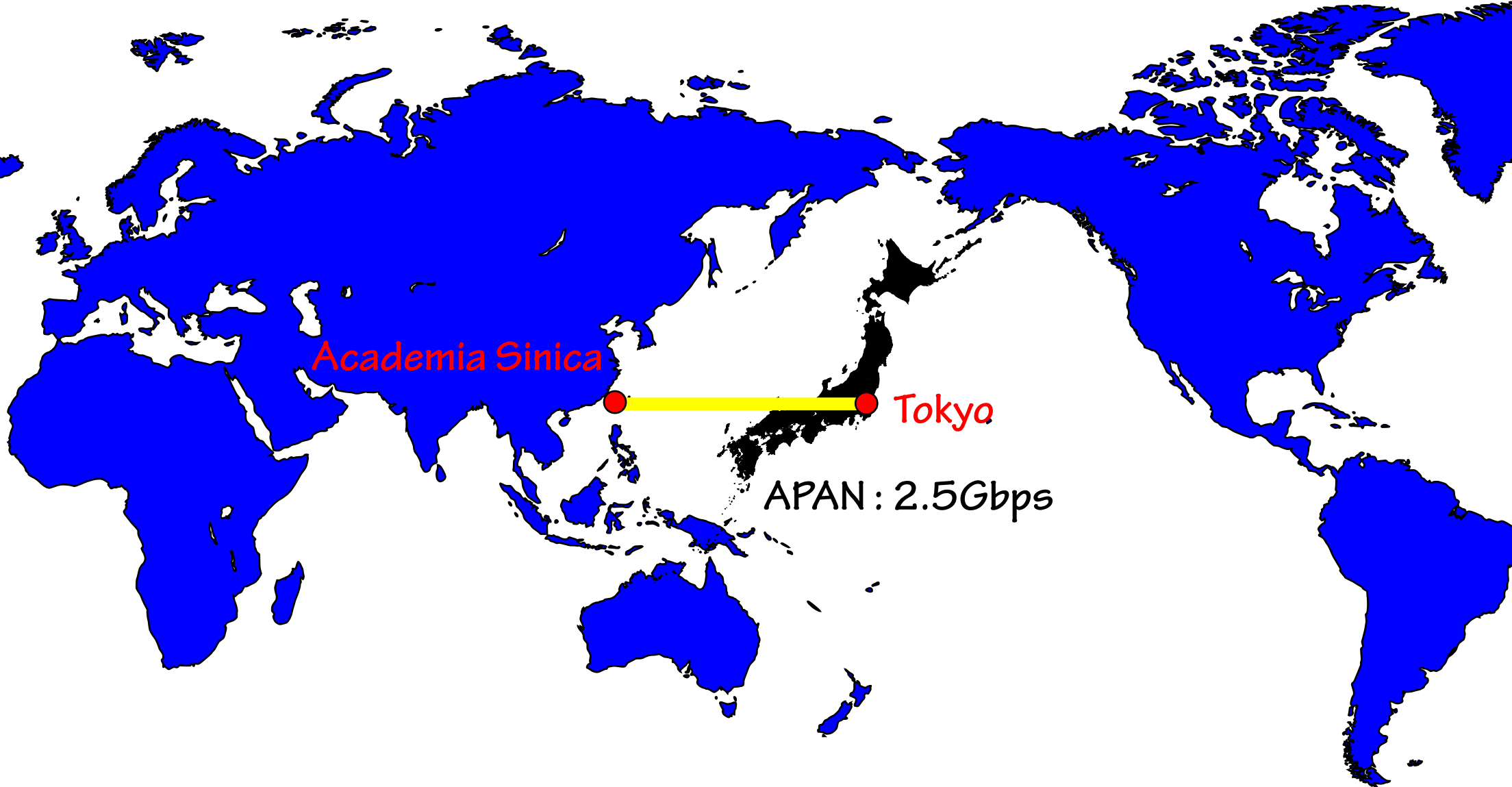




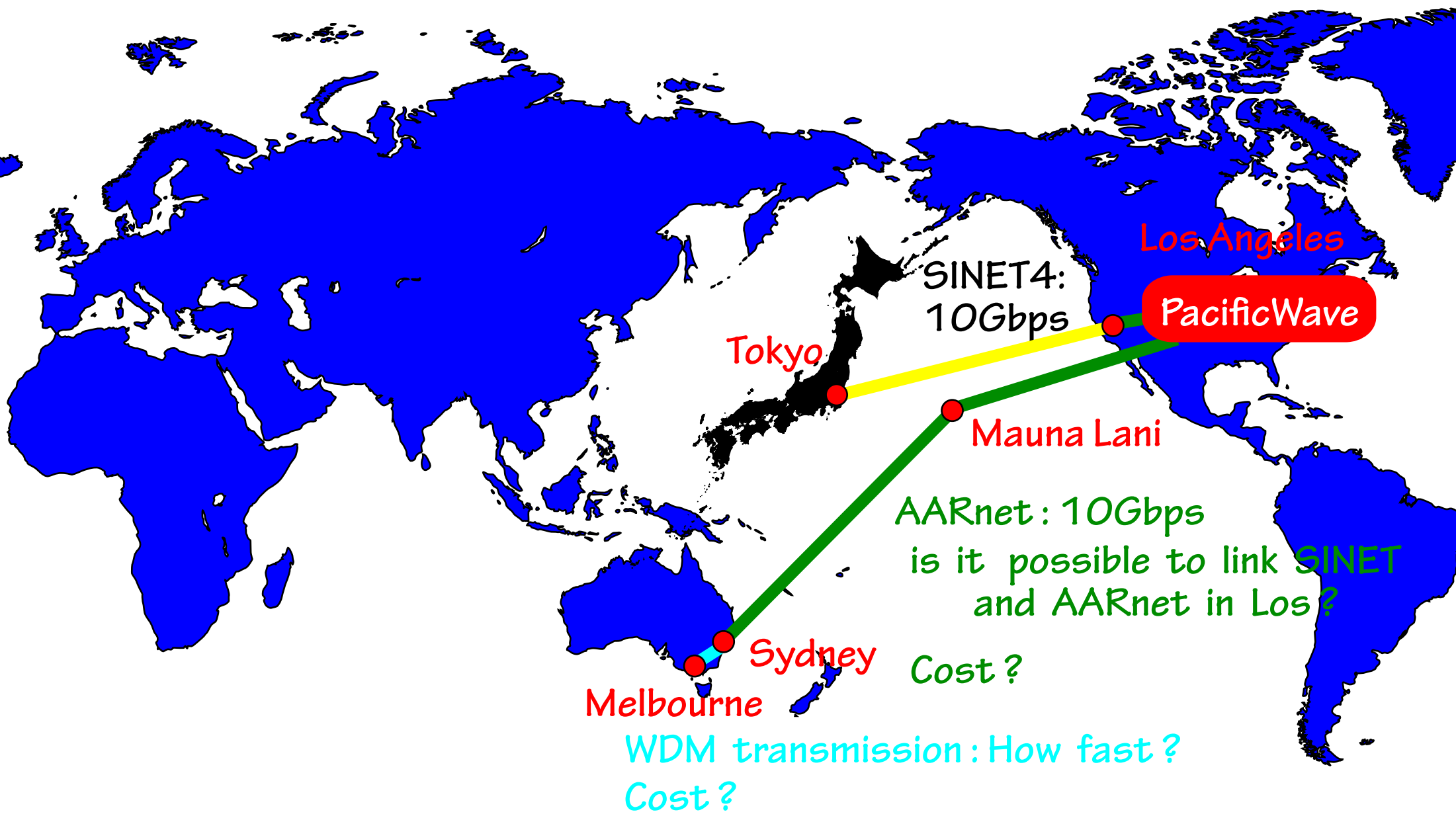


# Network Connection to India

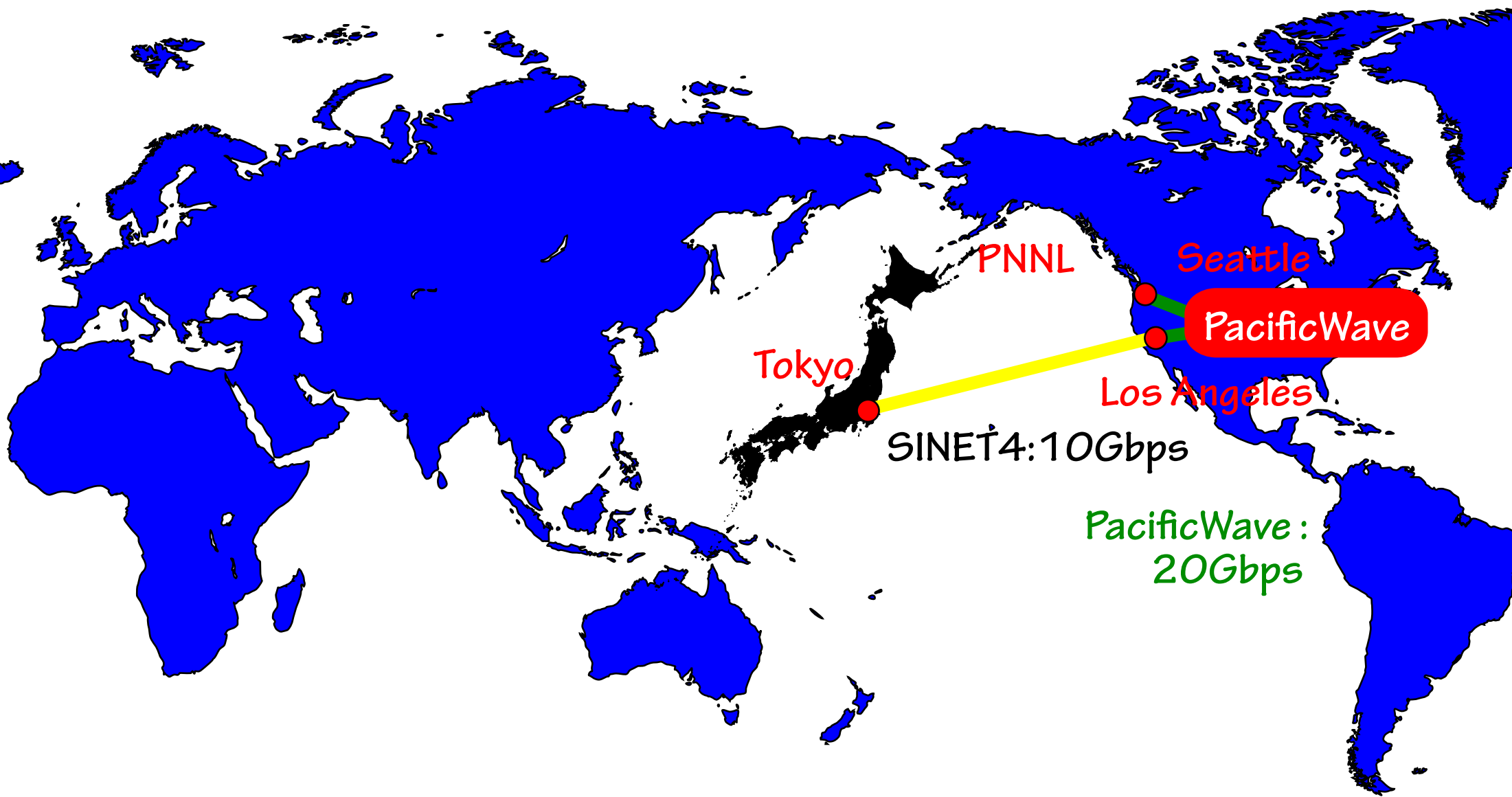




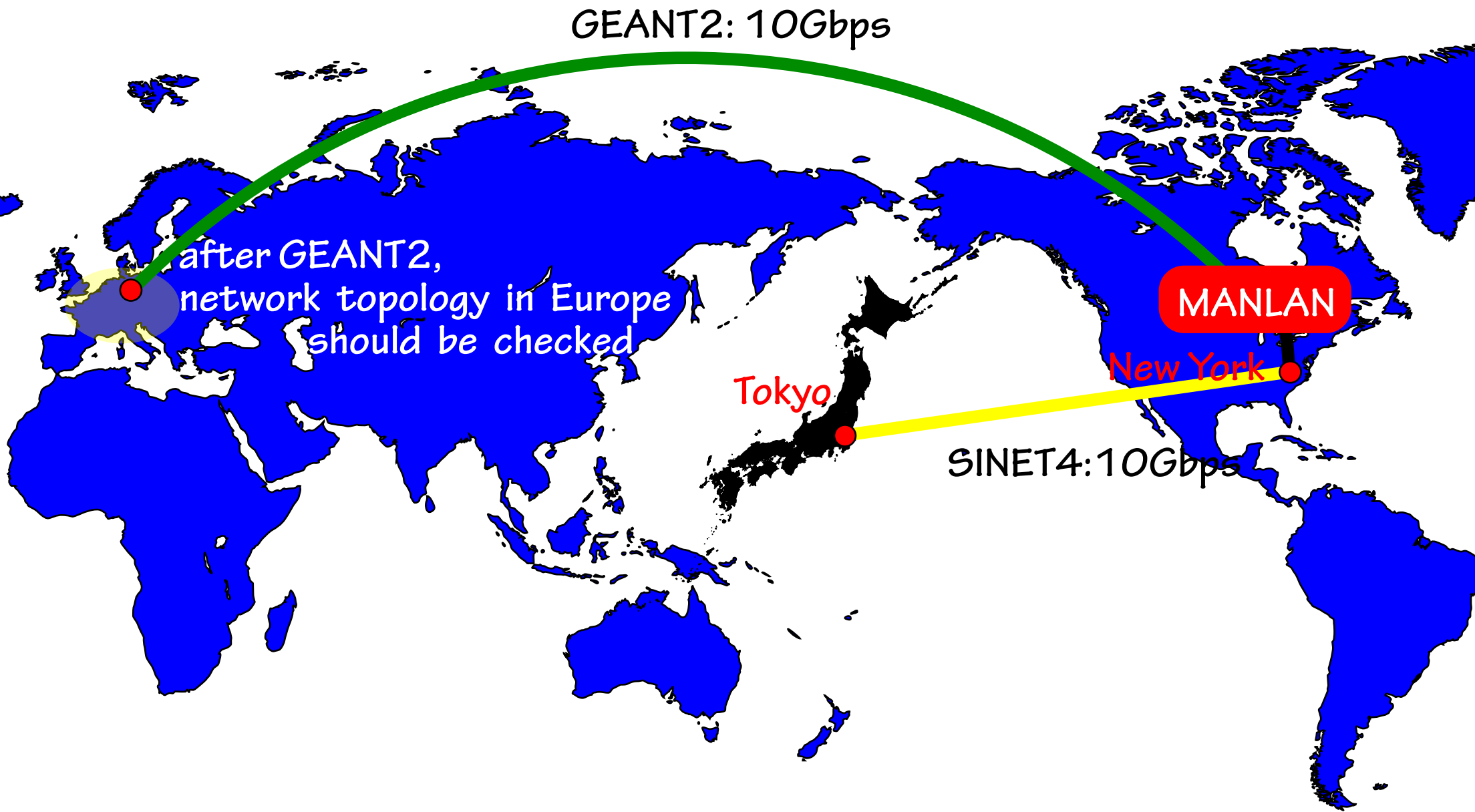
# Network Connection to Melbourne



# Network Connection to USA (PNNL)



# Network Connection to Europe

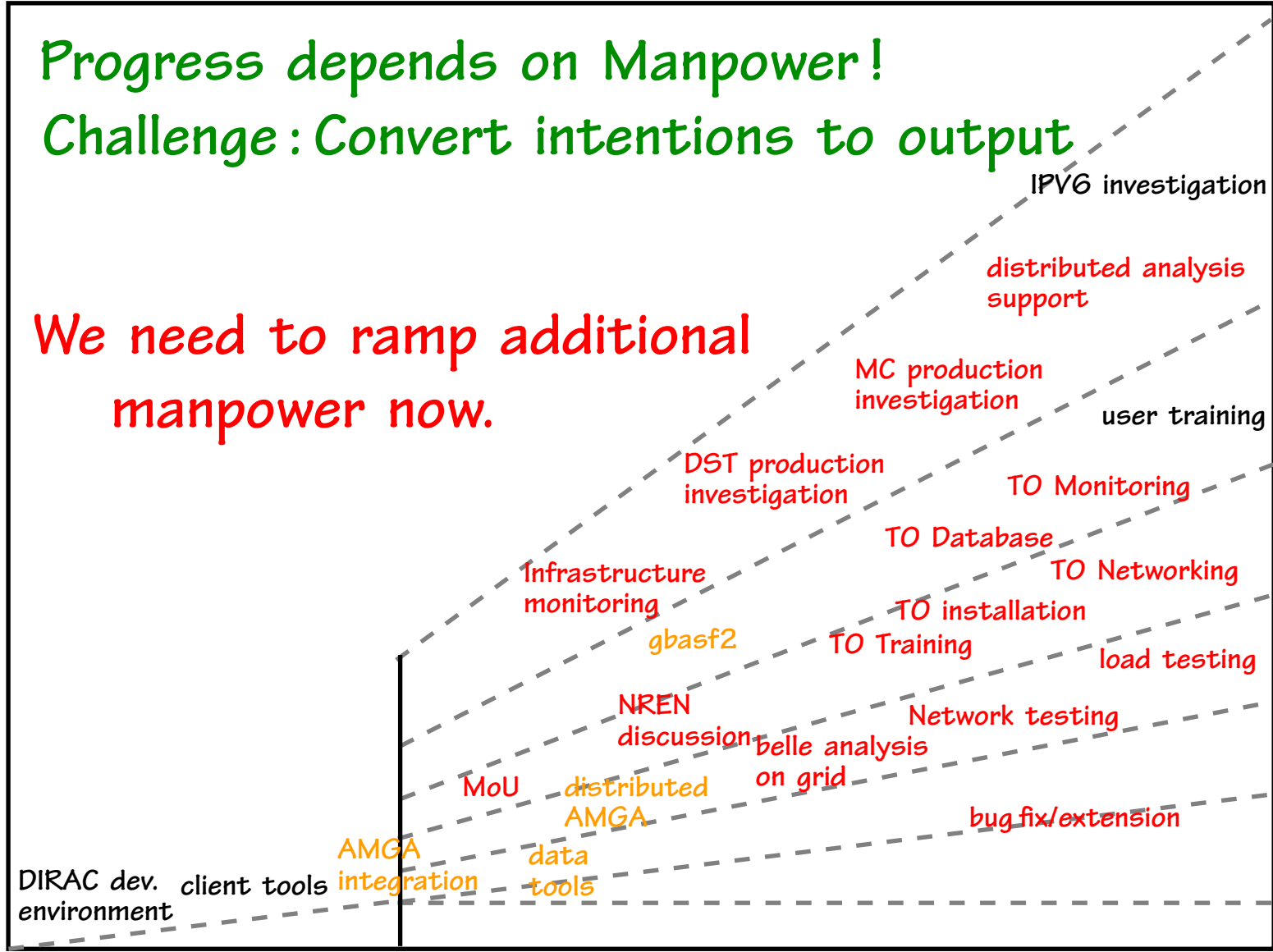


# Manpower on Computing

Progress depends on Manpower!  
 Challenge: Convert intentions to output

We need to ramp additional manpower now.

6 FTE  
 5 FTE  
 4 FTE  
 3 FTE  
 2 FTE  
 1 FTE



2010

2011

2012

2013

# Belle II GRID sites

GRID  
middleware

gLite

OSG

- ◆ Australia : LHC Tier2/3, Belle VO, Cloud system
- ◆ Austria : LHC Tier2
- ◆ China (IHEP) : LHC Tier2, DIRAC server
- ◆ Czech Republic : LHC Tier2, Belle VO
- ◆ Germany : LHC Tier1/2, Belle VO
- ◆ India : LHC Tier2, Belle II data center planned
- ◆ Japan (KEK) : Belle VO
- ◆ Korea (KISTI) : LHC Tier2, Belle VO
- ◆ Poland : LHC Tier2/3, Belle VO, Cloud system
- ◆ Russia : LHC Tier2
- ◆ Slovenia : LHC Tier2, Belle VO
- ◆ Taiwan : LHC Tier1/2
- ◆ USA : OSG @ PNNL planned, Belle VO @ OSG exists

2012, Jan. : KEK-PNNL meeting @ Richland

2012, Feb. : KEK-India institutes meeting @ Kolkata

2012, Mar. : Belle GRID site meeting @ Munich

# Plan

extension

expiration

by Aug. 31, 2015



sim. software



test MC prod. in KEK



test MC prod. in GRID (KEK+possible sites)



MC challenge in GRID (KEK+possible sites)



cosmic data ? (real)



Belle II computing federation ?



GRID Operation for Belle II

SuperKEKB Commissioning starts

Phys. run starts

current level operation for GRID

detail tests on Belle II computing has to be started at least 1 year prior to the collision



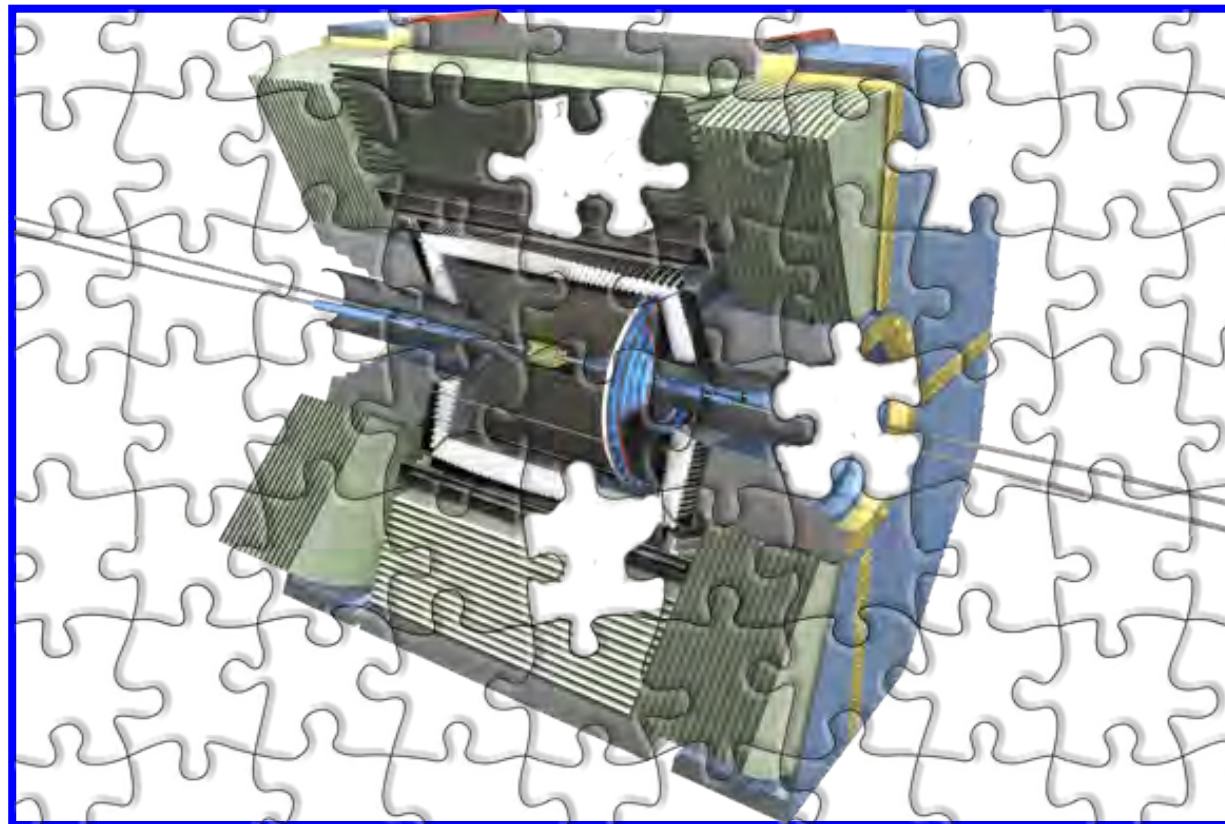
# Summary

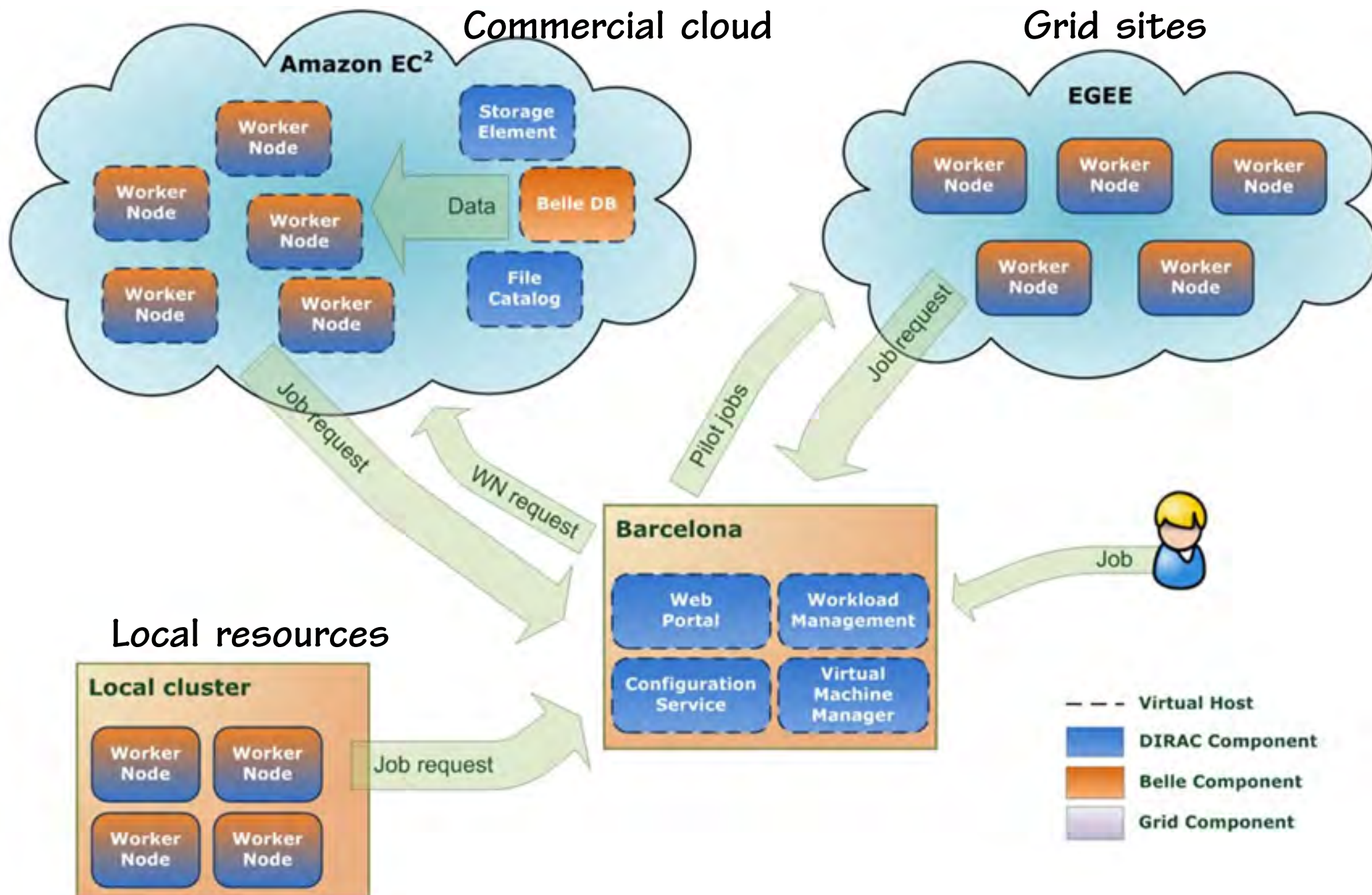
Deployment of Belle II grid site collaboration started

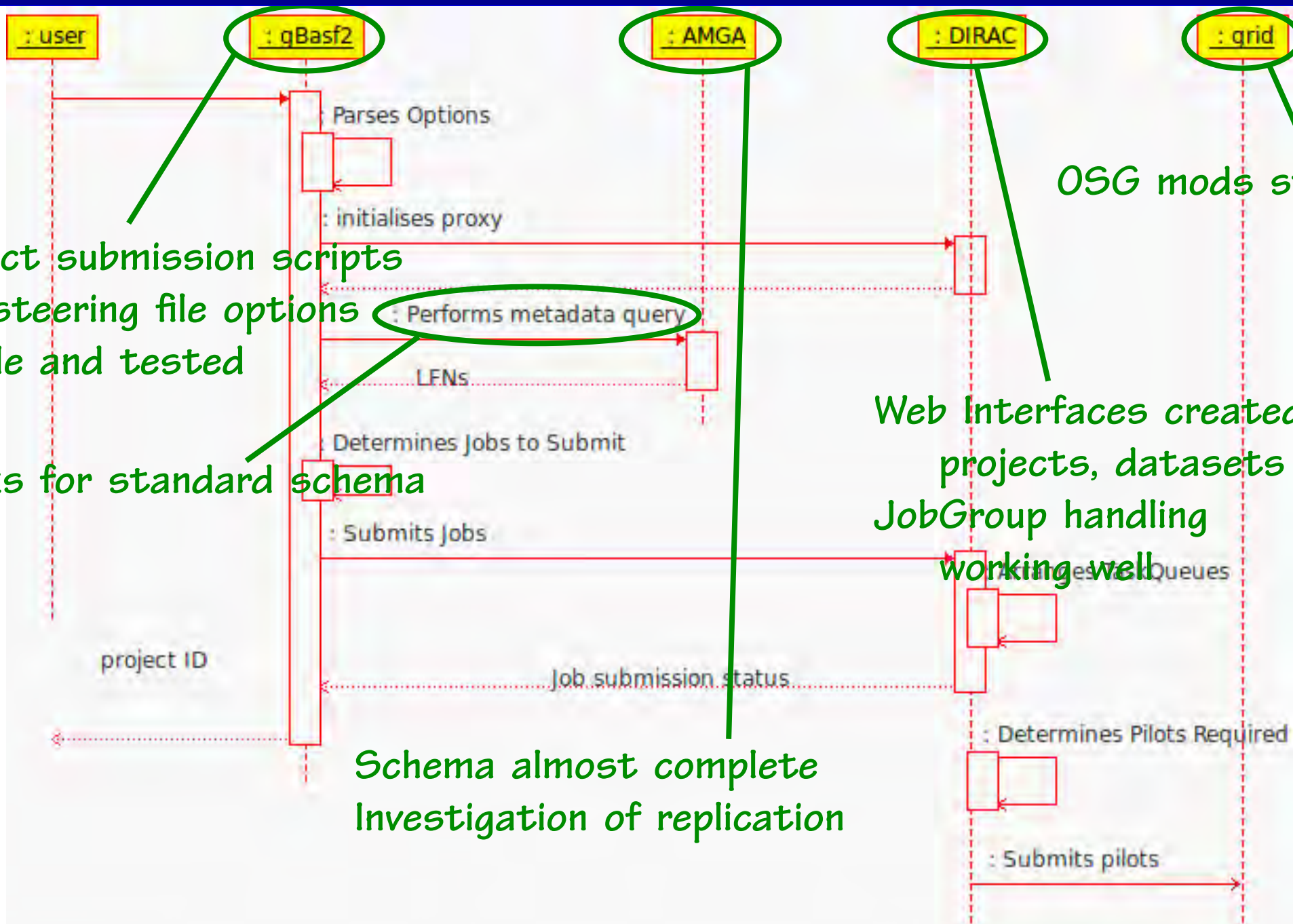
Many pieces of puzzle are getting ready

user interface, analysis projects, dataset tools, metadata catalog,  
data registration tool, software installation, database, New computer ...

but still many missing pieces







Project submission scripts  
and steering file options  
usable and tested

Works for standard schema

: Performs metadata query

Schema almost complete  
Investigation of replication

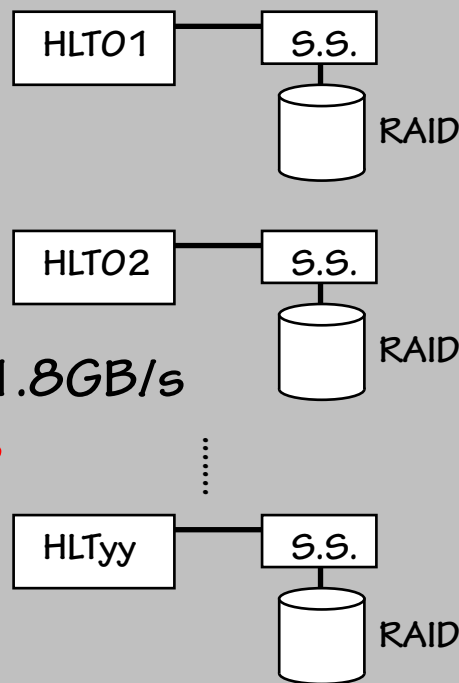
OSG mods stable

Web Interfaces created for  
projects, datasets  
JobGroup handling  
working well

## Plan in 2011

- Belle real Data registration to AMGA
- Test of large-scale Data Handling w/ Belle data @ KEK
- Prototype  $\rightarrow$  Close end-user support
- interface between DAQ spool disk  
and offline (GRID-accessible) storage : need manpower @ KEK

Online



$300\text{KB/event} + 6\text{kHz} = 1.8\text{GB/s}$

8 hours run = **51.84TB**

```
expXXXXrunYYYY.dst-00-HLT01
expXXXXrunYYYY.dst-01-HLT01
expXXXXrunYYYY.dst-02-HLT01
...
expXXXXrunYYYY.dst-zz-HLT01

expXXXXrunYYYY.dst-00-HLT02
expXXXXrunYYYY.dst-01-HLT02
expXXXXrunYYYY.dst-02-HLT02
...
expXXXXrunYYYY.dst-zz-HLT02

...

expXXXXrunYYYY.dst-00-HLTyy
expXXXXrunYYYY.dst-01-HLTyy
expXXXXrunYYYY.dst-02-HLTyy
...
expXXXXrunYYYY.dst-zz-HLTyy
```

Offline

Data copy scheme  
network bandwidth  
network path

idea of "Prompt-reco"  
check Itoh-san's talk

For rawdata process, MC (25%), analysis

safety factor  
(x 2) applied

Fiscal year	2015	2016	2017	2018	2019	2020
Tape [PB]	3	20	58	107	157	208
Disk [PB]	0.7	4.6	14	26	38	50
CPU [kHepSPEC]	7	42	100	132	145	155
WAN [Gbit/s]	0.6	3.8	8.8	11	11	12

5 months  
for DST

another 5 months  
for MC prod.

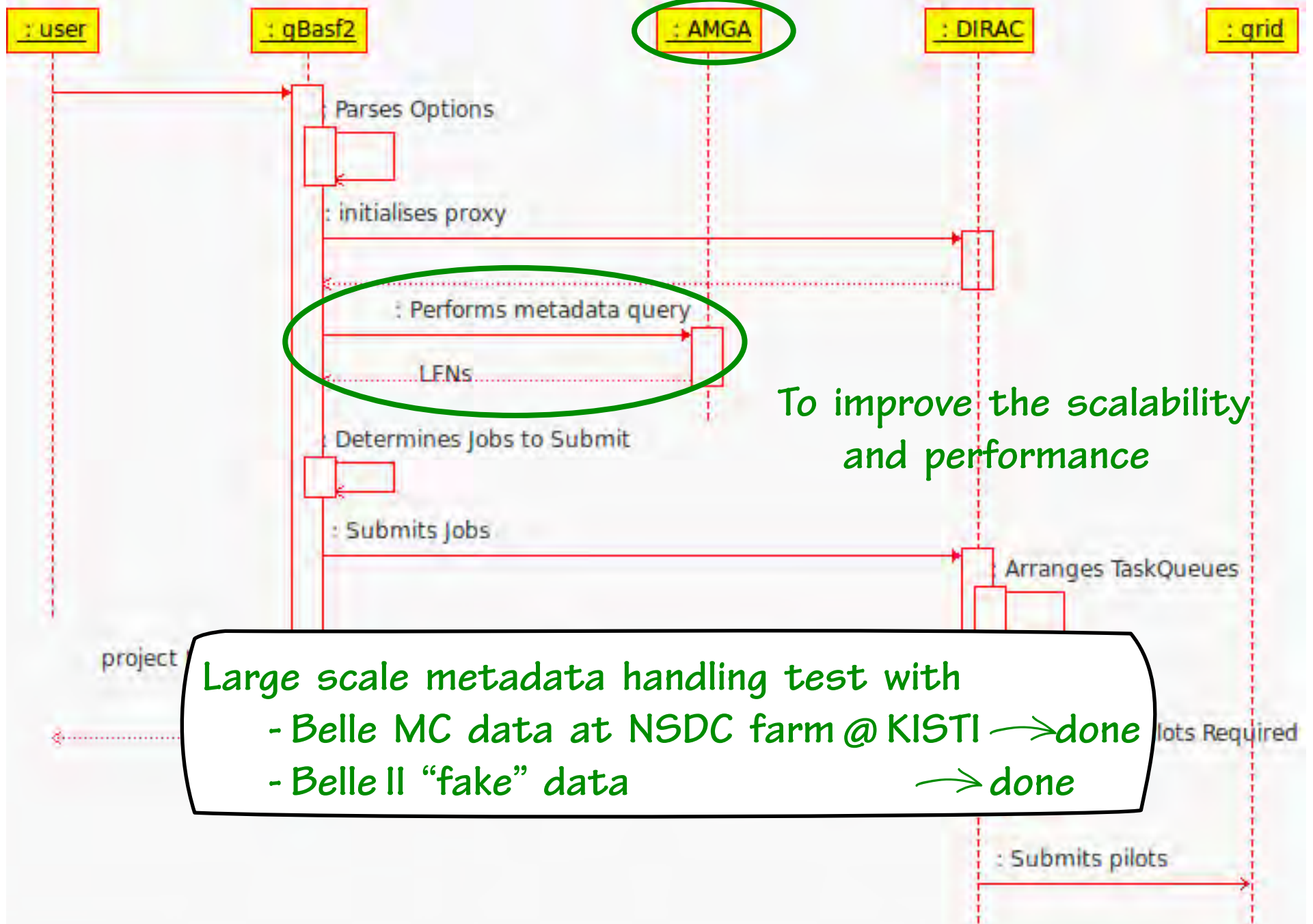
Belle  
 ~1.5PB disk  
 ~3.5PB tape  
 ~45 kHEPSpec (3/5 used)  
 (= ~4,000 cores (3GHz))

For MC (9.2%), analysis.  
and 1 copy of mDST datasets

Fiscal year	2015	2016	2017	2018	2019	2020
Disk [PB]	0.25	1.7	5	9.1	13	18
CPU [kHepSPEC]	2	12	29	38	42	45
WAN [Gbit/s]	0.084	0.48	1.1	1.4	1.5	1.5

in case that sites host a raw data copy

Fiscal year	2015	2016	2017	2018	2019	2020
Tape [PB]	3	20	58	107	157	208
CPU [kHepSPEC]	3	19	28	52	53	50
WAN [Gbit/s]	1.7	10	23	29	31	31



## via networks

running transfer

to Nagoya U.       $\sim 100\text{MB/s}$  is achieved (by tuning network parameters) by B-SE

all HadronBJ, tau data (copied before the earthquake)

1-stream MC ( $\sim 100\text{TB}$ ): copied (after the earthquake) by Hayasaka@Nagoya

to PNNL       $\sim 100\text{MB/s}$  is achieved (network parameter tuning + HPN-ssh)  
by S.Suzuki@KEKCRG, J.Schroeder, T.Carlson@PNNL

Y5S HadronBJ ( $\sim 10\text{TB}$ ): copied (after the earthquake) by G.Tatishvili@PNNL

Y4S HadronBJ data transfer: on-going

transfer test

to KISTI       $\sim 20\text{MB/s}$  (no network parameter tuning) by T.Khan@KISTI

to Karlsruhe       $\sim 3\text{MB/s}$  (network param. tuning + HPN-ssh, multi-stream)

$\sim 20\text{MB/s}$  (with 100-stream gridftp) by T.Kuhr@KIT

Assuming 2TB data with 10MB/s transfer speed, it takes  $\sim 2.5$  days

However, need to consider heavy load to the file servers for this read access  
as well as for the index-to-mdst file conversion (copy to HDs and ship to LC is another solution)