# Belle-II Distributed Computing
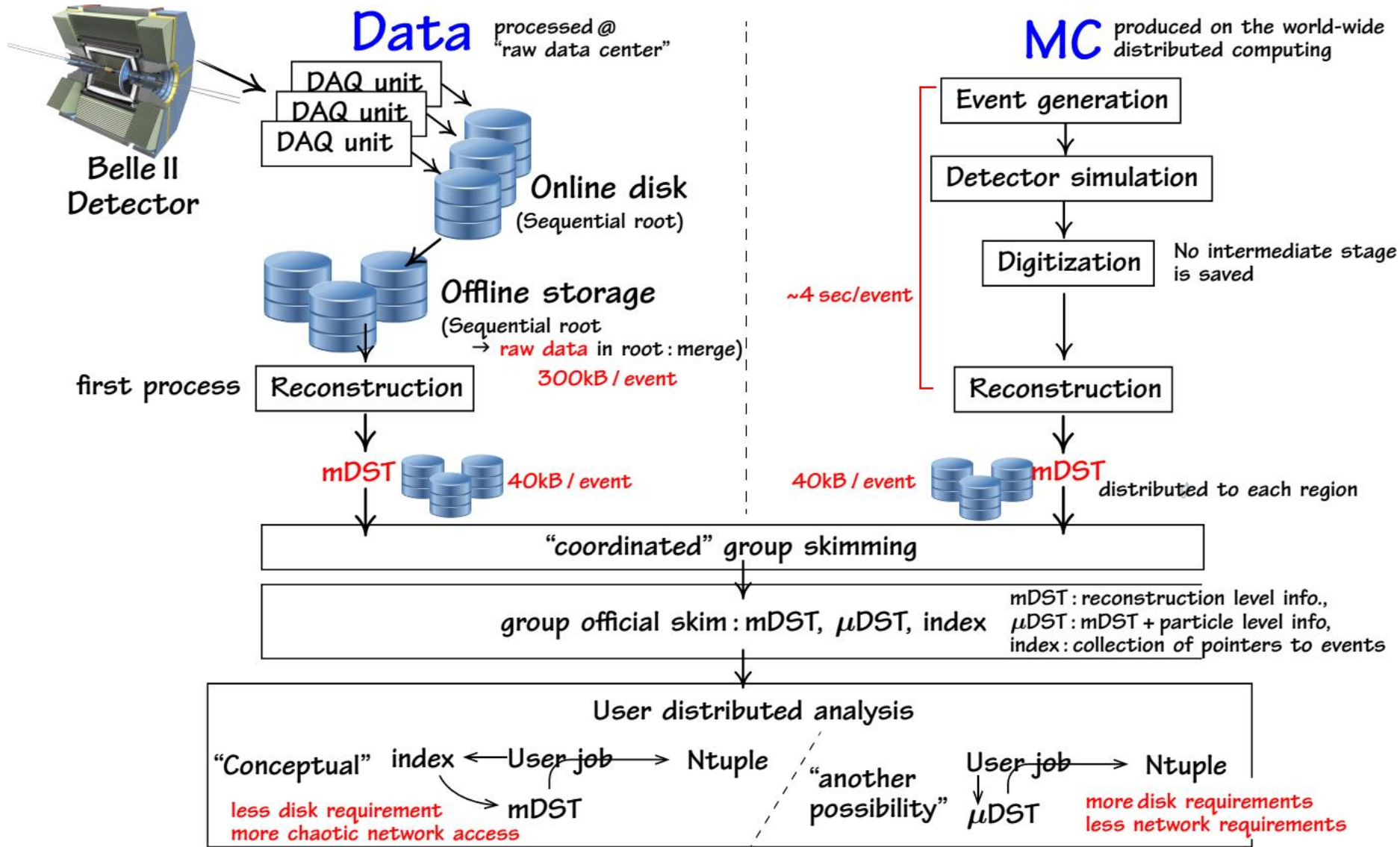
Dr. Silvio Pardi (INFN-Napoli)

JENNIFER Consortium General Meeting

10-12 June 2015

Data processed @ "raw data center"

Belle II Detector

DAQ unit
DAQ unit
DAQ unit

Online disk (Sequential root)

Offline storage (Sequential root → raw data in root : merge) 300kB / event

first process → Reconstruction

mDST 40kB / event

MC produced on the world-wide distributed computing

Event generation
Detector simulation
Digitization — No intermediate stage is saved
~4 sec/event
Reconstruction

40kB / event mDST distributed to each region

"coordinated" group skimming

group official skim : mDST, μDST, index

mDST : reconstruction level info.,
μDST : mDST + particle level info,
index : collection of pointers to events

User distributed analysis

"Conceptual" index ← User job → Ntuple
mDST
less disk requirement
more chaotic network access

"another possibility" User job → Ntuple
μDST
more disk requirements
less network requirements

# The BELLE II Collaboration



23 countries/region
99 institutes
634 colleagues
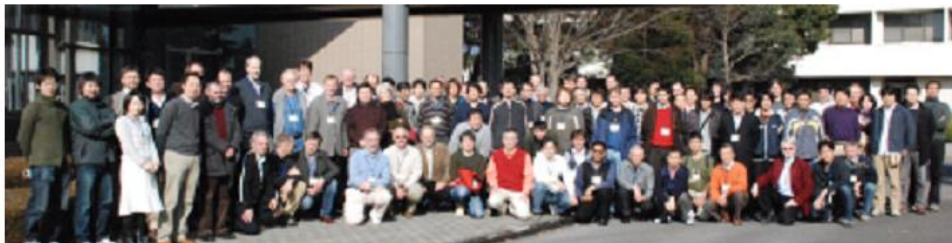
c.f.
ATLAS, 38 countries, 177 institutes, ~3000 members
CMS: 42 countries, 182 institutes, 4300 members
ALICE: 36 countries, 131 institutes, 1200 members
LHCb: 16 countries, 67 institues, 1060 members

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

# BELLE II Computing model

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

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

On going activities
- Resource Estimation
- Define strategy for analysis and data distribution
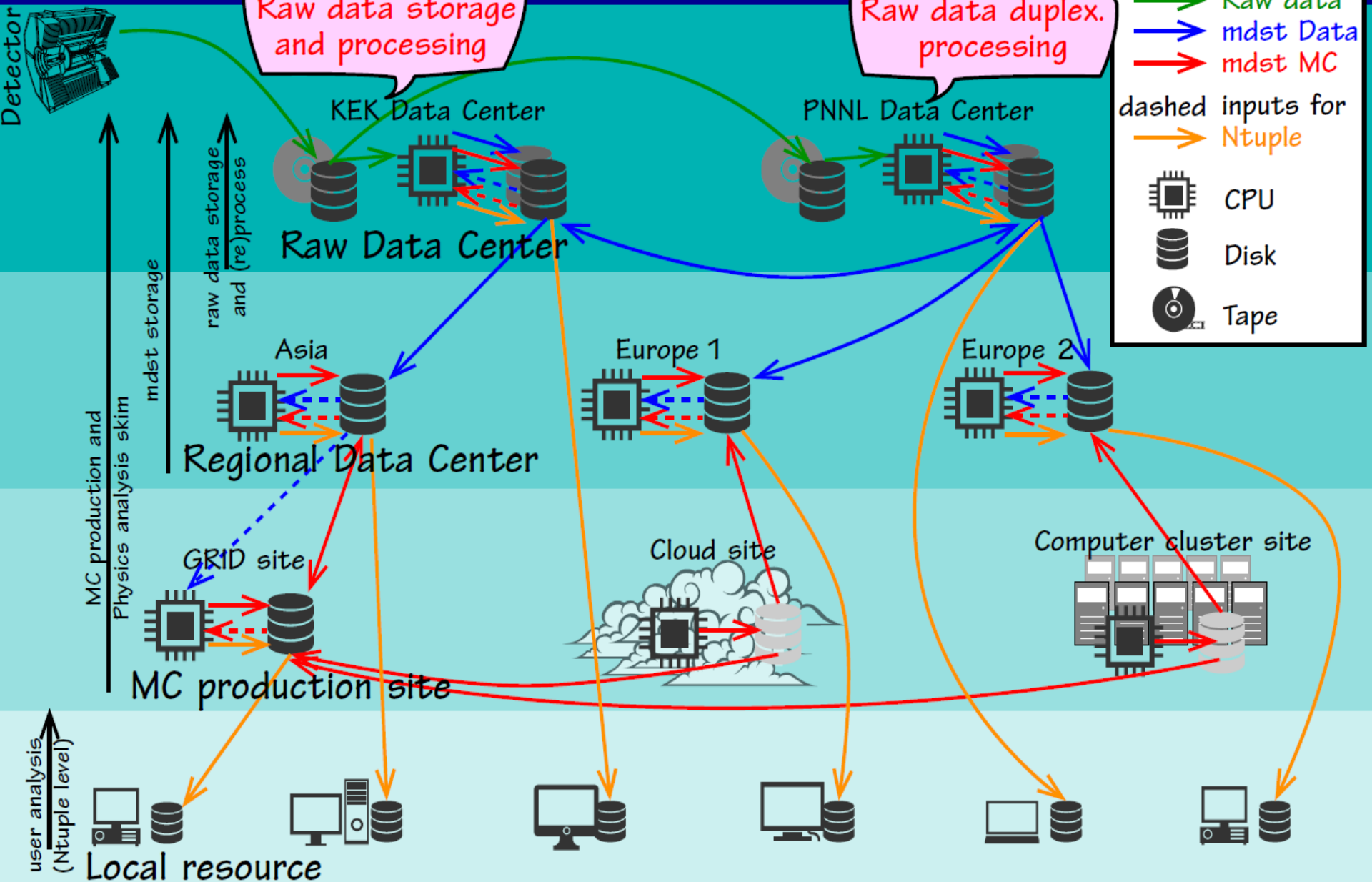- Individuating technologies

# Site Classification

The BELLE II Computing Sites are classified as follow:

- **Raw Data Center:** Who store the RAW Data and made data processing and/or data reprocessing.

- **Regional Data Center**: Large data center that stores mDST and participates at the Monte Carlo production

- **MC Production site**: Data Center that produces and stores Monte Carlo simulations, that included:
  - Grid Site
  - Cloud Site
  - Computing Cluster Site

# RAW Data Distribution

We plan to have two full copy of RAW Data

RAW data are produced at KEK, replicated and stored at PNNL(USA) for the first 3 years.

Starting from the 4$^{th}$ year of operation they will be distributed in others RAW Data Centers. The current hypothesis is:

- PNNL(30%)
- Italy(20%)
- Germany(20%)
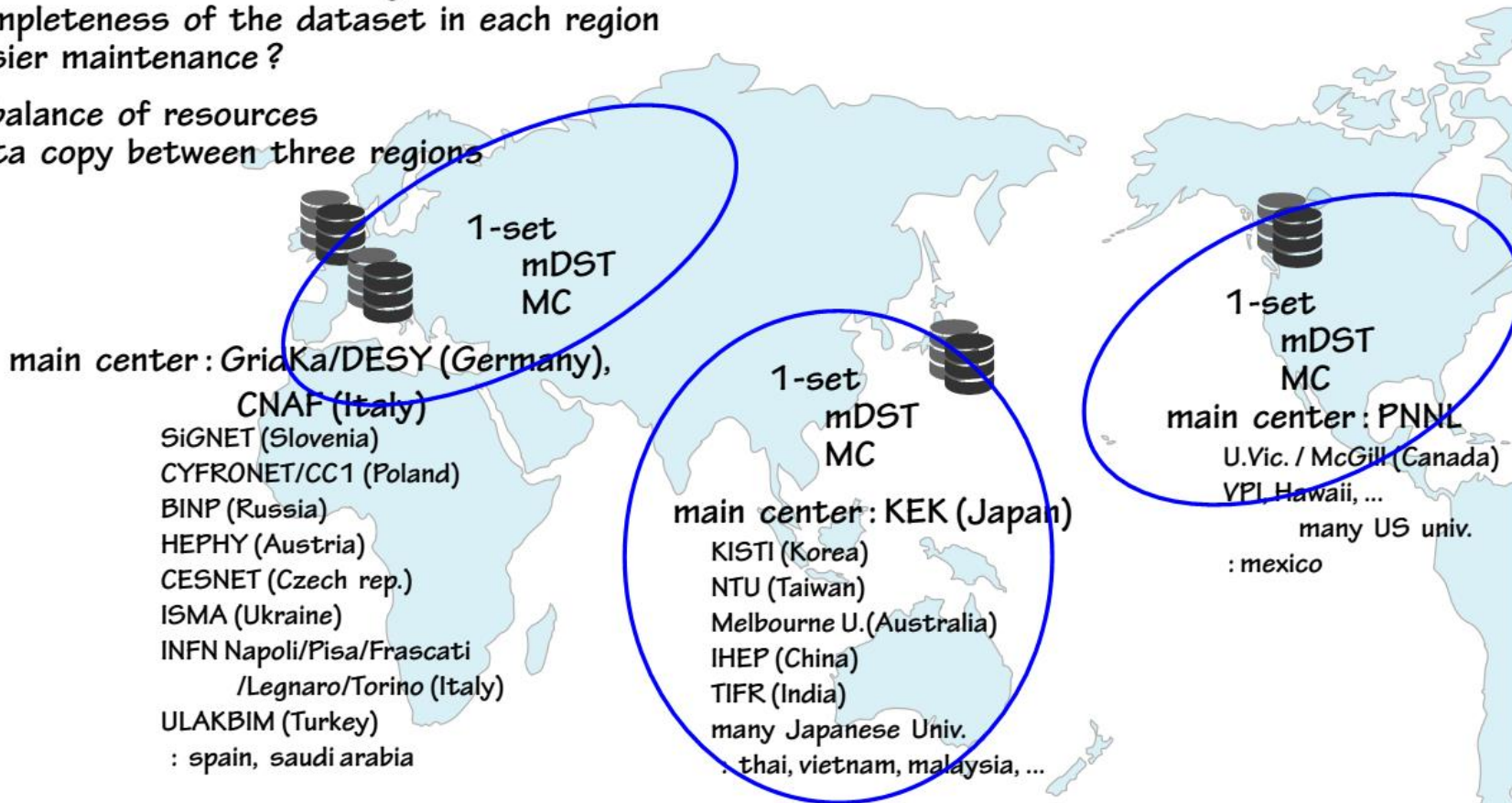- Canada(10%)
- Korea(10%)
- India(10%)

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?
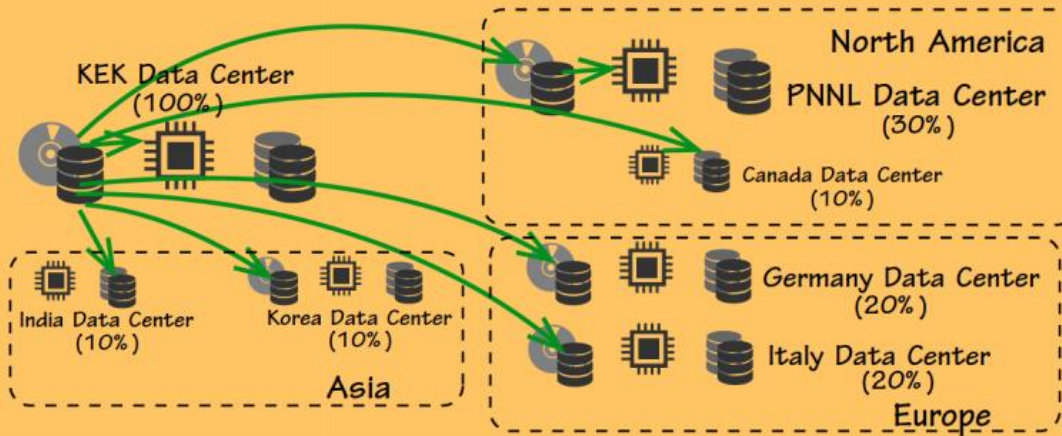
unbalance of resources
data copy between three regions

1-set
mDST
MC

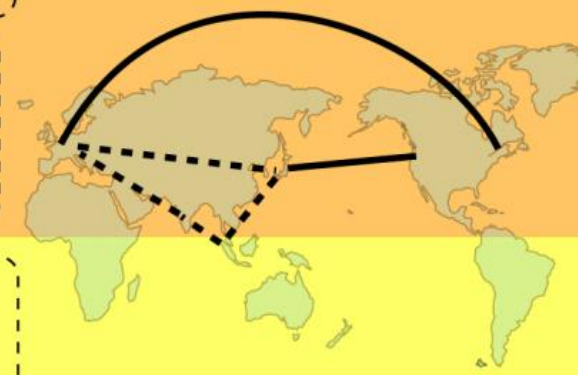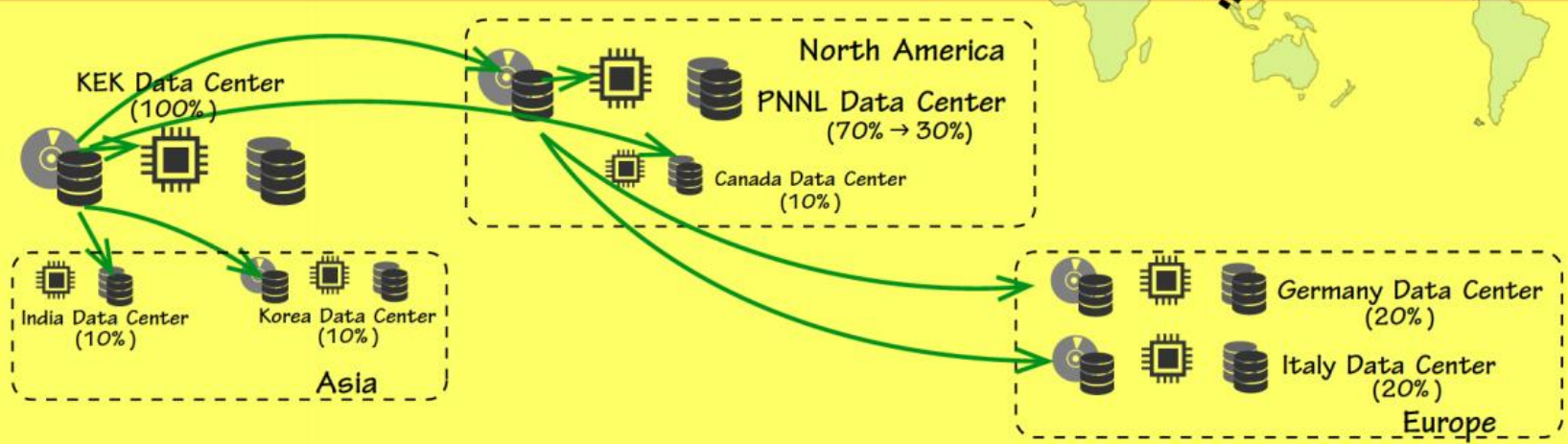main center: GridKa/DESY (Germany),
CNAF (Italy)
SiGNET (Slovenia)
CYFRONET/CC1 (Poland)
BINP (Russia)
HEPHY (Austria)
CESNET (Czech rep.)
ISMA (Ukraine)
INFN Napoli/Pisa/Frascati
/Legnaro/Torino (Italy)
ULAKBIM (Turkey)
: spain, saudi arabia

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

INBAND

KEK ■ PNNL ■ Korea ■ Germany ■ Italy ■ Slovenia ■ Canada ■ Australia

OUTBAND

PNNL ■ KEK ■ Germany ■ Italy ■ Korea ■ Canada ■ India

# GÉANT global connectivity

GÉANT and partner networks enabling user collaboration across the globe

September 2014

Legend:
- GÉANT Coverage
- RedCLARA Network
- EUMEDCONNECT3 Network
- TEIN Network
- AfricaConnect – UbuntuNet Alliance
- CAREN Network
- SILK-Afghanistan
- Other R&E Networks

Dark Shading: Connected to regional network
Light Shading: Eligible to connect to regional network

Speed legend:
- 30 Gbps
- 10 Gbps
- 10 Gbps backup
- >1 Gbps<10 Gbps
- ≤1 Gbps

**North America**
100Gb Paris-NY (30Gb on ACE)
340Gb from ESNet
200Gb from ANA-200
*(Not yet on the map)*

**Asia**
10Gb to SINET4 (via N.A.)
2.5Gb to TEIN (Mumbai)
10Gbit to ORIENT+/TEIN
622+155Mb to CAREN
*Future deployments:*
10Gb to TEIN (Singapore)
2x10Gb to SINET5 (direct)

**Africa**
2x10Gb to UbuntuNet Alliance
2x622Mb to EumedCONNECT3

**Latin America**
5Gb to RedCLARA

**GÉANT connects 65 countries outside of Europe, reaching all continents through international partners**

▶ 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

Jason Zurawski, Joe Metzger (Mar. 23, 2015)
http://www.es.net/assets/pubs_presos/20150323-OSG-ATLASCMS-Zurawski-v2.pdf

# DATA Management - Storage System

*Belle II*

SRM based storage systems. The most common technologies in the Sites are:

- DPM
- dCache
- STORM
- Bestman2

In evaluation webDav and xrootd for direct access.

▶ Storage Element Accounting System:

  ■ Provide health of each registered Storage Element (SE)

  ■ Overall storage availability at each SE

  ■ Group and user level accounting

▶ Data Transfer System:

  ■ Provides tools to transfer data between sites

  ■ Provide tools for users to retrieve their samples

  ■ Network Monitoring

▶ Data Integrity

  ■ Insures that the physical data is pristine and consistent with the File-Catalogs (FC)

# DATA Transfer System(DTS)

*Belle II*

- ► FTS3 Server:
  - Previously used for Data Challenges and latest Monte Carlo campaign
  - Requires continuous studies of the FTS3 transfers to tune the channels as needed (FTS3 optimization not perfect yet)

- ► DIRAC Integration:
  - The BelleDIRAC test server is using v6r12
  - Belle II FTS3 DIRAC Agent was developed to automate the Data Challenges
  - Transformation DIRAC/FTS3 data transfer working

- ► Networking:
  - Belle II perfSONAR mesh is now deployed providing automated network monitoring
  - Developing DIRAC agent to access perfSONAR results using REST API

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

GRID Middlewares

European Middleware Initiative

ARC

Open Science Grid

**D**istributed
**I**nfrastructure with
**R**emote
**A**gent
**C**ontrol
(originally developed for LHCb)

KEK

DIRAC THE INTERWARE

A part of DIRAC

VMDIRAC

SSH tunnel or DIRAC SiteDirector

batch system

Clusters w/o middleware
GE, TORQUE, LSF,...
Direct submission

BINP, NSU, many universities in Japan

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

Cracow

cc1
Cracow Cloud One

cloud site | cloud site | cloud site

CREAM CE

Dynamic Torque

cloud site | cloud site | cloud site

SiteDirector

HTCondor Cloud Scheduler

cloud site | cloud site | cloud site

SLURM SiteDirector

HTCondor VM Manager

cloud site

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

Melbourne

UVic

PNNL

HPC

Academic clouds

Commercial clouds, Amazon EC2, etc

DIRAC main servers @ KEK

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

**Napoli (Italy)**

VOMS @ KEK

AMGA          +          LFC   : has been working well
recent improvement

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

FTS3 : getting integrated

CernVM
File system          cvmfs is used for software distribution

cvmfs is used for software installation for most of sites

10-15% of current MC production is cloud-based

Clouds at Belle-II member sites

Opportunistic (private and commercial) clouds

- 3rd MC Campaign: ~ April 1 – May 15, 2014

- Simulation and reconstruction, with background mixing → mdst data

- 2x previous CPU#: 11k concurrent jobs; > 80 kHS max

- ~30 sites contributing

- 4.2G events produced

→ Very successful; also updated analysis and grid software

→ To obtain useful data for physics studies **new extensive MC production started this week**

**Normalized CPU usage by Site**

30 Days from 2014-04-05 to 2014-05-05



Max: 81,947, Min: 223, Average: 40,033, Current: 223

| | | | | | |
|---|---|---|---|---|---|
| LCG.DESY.de | 22.3% | LCG.CNAF.it | 4.3% | LCG.KISTI.kr | 1.1% |
| LCG.KIT.de | 9.3% | LCG.CYFRONET.pl | 3.9% | LCG.McGill.ca | 1.0% |
| LCG.UA-ISMA.ua | 6.3% | LCG.SIGNET.si | 3.6% | LCG.Legnaro.it | 0.9% |
| DIRAC.BINP.ru | 6.0% | DIRAC.UVic.ca | 3.3% | LCG.TORINO.it | 0.6% |
| LCG.KEK2.jp | 5.7% | LCG.Melbourne.au | 3.3% | SSH.KMI.jp | 0.4% |
| LCG.KMI.jp | 5.0% | LCG.CESNET.cz | 3.1% | DIRAC.Niigata.jp | 0.4% |
| DIRAC.PNNL.us | 5.0% | LCG.ULAKBIM.tr | 2.0% | OSG.Nebraska.us | 0.3% |
| LCG.PISA.it | 4.9% | LCG.Frascati.it | 1.5% | OSG.FNAL.us | 0.2% |
| LCG.Napoli.it | 4.3% | DIRAC.KrakowCloud.pl | 1.4% | ... plus 6 more | |

Generated on 2014-05-05 05:22:05 UTC

# MC Campaign 2015

Testing campaign done in April/May 2015 used to tune and validate the new components implemented in the distributed system.

In the next month a new massive MC production will start over the current computing infrastructure.

# **Conclusion**

- Belle II community is very active in developing the Distributed Computing infrastructure.

- There are several on-going activities but also some achieved results.

- Belle joined the LHCONE Network

- The current choices are based on open and flexible solutions avoiding technology lock-in and are mainly compliant with the other HEP communities in term of standard.

- We can take advantage from the CERN experiences for several tools, but other components must be developed ad hoc.

# REFERNCES

[1] Computing at the Belle-II experiment Authors: Hara Takanori
http://indico.cern.ch/event/304944/session/15/contribution/550


[2] Utilizing cloud computing resources for BelleII Authors: Sevior Martin HARA Takanori Sobie
Randy http://indico.cern.ch/event/304944/session/7/contribution/294

[3] "Belle II Networking Overview" MALACHI SCHRAM et.
https://indico.cern.ch/event/376098/contribution/12/material/slides/0.pdf

# BACKUP

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

# Methodology

Goal: Estimate the International network traffic that will be generated by the Belle 2 collaboration and then individuate the requirements.

Estimation of the In-Band and Out-Band peaks for each site by decupling the different data flows and by adopting a tolerance factor of a 50%.

More specifically 5 data flows are considered :

- RAW Data

- mDST from Data - after data taking

- mDST from Data - reconstruction process

- mDST-MC  from Monte Carlo production during data taking

- mDST-MC  form Monte Carlo production during Data reconstruction

# RAW DATA

**Involved sites**

- RAW data are produced at KEK and replicate at PNNL
- Starting from the 4$^{th}$ year of operation one hypothesis is the following distribution: PNNL(30%), Italy(20%), Germany(20%), Canada(10%), Korea(10%), and India(10%)

Event size= 300k,Month=8, Tollerance:50%

| | 2015 | 2015 | 2015 | 2015 | 2015 | 2015 | 2015 | 2015 | 2015 | 2015 |
|---|---|---|---|---|---|---|---|---|---|---|
| Event Rate | 0,00 | 0,00 | 1,76E+09 | 2,29E+09 | 2,12E+10 | 4,73E+10 | 7,70E+10 | 9,33E+10 | 9,33E+10 | 9,33E+10 |
| RAW Data (PB) | 0,00 | 0,00 | 0,00 | 0,62 | 5,78 | 12,90 | 21,00 | 25,46 | 25,46 | 25,46 |
| Network Req+t | 0,00 | 0,00 | 0,00 | 0,38 | 3,49 | 7,80 | 12,70 | 15,39 | 15,39 | 15,39 |

Two possible scenarios for data distribution since the 4$^{th}$ year of operation

- **Scenario 1: KEK send the RAW data directly to all the involved RAW-Data centers**

- **Scenario 2: KEK send the 80% of the RAW data to PNNL that store the 30% and distribute 10% to Canada, 20% to Germany and 20% to Italy.**

Max 15Gbps

Max of 8 Gbps (PNNL)

**RAW Data - KEK-OutBand**

| | | |
|---|---|---|
| 20,00 | | |
| 15,00 | | |
| 10,00 | | |
| 5,00 | | |
| 0,00 | | |

2015 2015 2015 2015 2015 2015 2015 2015 2015 2015

**RAW Data - IN-BAND**

Effect of the new RAW data distribution strategy

| | |
|---|---|
| 12,00 | PNNL |
| 10,00 | Korea |
| 8,00 | Germany |
| 6,00 | Italy |
| 4,00 | Canada |
| 2,00 | India |
| 0,00 | |

2015 2016 2017 2018 2019 2020 2021 2022 2023 2024

RAW Data - PNNL-OutBand

Max 12 Gbps (PNNL)

RAW Data - KEK-OutBand

RAW Data - IN-BAND

# mDST from RAW data

**mDST are produced during data taking**

- at KEK after data taking (60%)
- at PNNL after data taking (40%)
- at PNNL and :Italy, Germany, India, Korea and Canada after Y4

**and during data reprocessing**

- at PNNL (100%) for the first 3 years of operation
- at PNNL and :Italy, Germany, India, Korea and Canada after Y4

**mDST after data-taking are distributed in that way:**

- KEK -> PNNL (60%)
- KEK -> Australia, Canada (20%)
- KEK -> Asian Regional Center (RC) (100%)
- PNNL -> European Distributed RC (100%)
- Many to May connection after y4

**mDST form data reprocessing are distributed in that way:**

- from PNNL to KEK, European Distributed RC, The Asian Regional Center
- Many to may connection after y4

INBAND

OUTBAND

KEK  PNNL  Korea  Germany  Italy  Slovenia  Canada  Australia

PNNL  KEK  Germany  Italy  Korea  Canada  India

# mDST-Monte Carlo

**Two flows of mDST from MC production**

- **mDST-MC produced during data taking**
- **mDST-MC produced during data reprocessing**

All the produced mDST-MC are distributed in that way:

Every MC Production Site sent the produced mDST-MC data to 2 other sites.

The single end-to-end connection are not defined yet.

Is a many-to-may network connection

# mDST-Monte Carlo

**mDST-MC - data taking:**

**Event size= 40k,Months=11, Tollerance:50%**

**mDST-MC - data reprocessing:**

**Event size=40k, Months=12, Tollerance:50%**

| | 2015 | 2016 | 2017 | 2018 | 2019 | 2020 | 2021 | 2022 | 2023 | 2024 |
|---|---|---|---|---|---|---|---|---|---|---|
| **Event Rate** | 0,00 | 0,00 | 2,34E+08 | 3,05E+08 | 2,82E+09 | 6,31E+09 | 1,03E+10 | 1,24E+10 | 1,24E+10 | 1,24E+10 |
| **MDST - Data (PB)** | 0,00 | 0,00 | 0,01 | 0,01 | 0,10 | 0,23 | 0,37 | 0,45 | 0,45 | 0,45 |
| **# Reprocessing** | 0 | 0 | 0 | 4 | 4 | 2 | 0,5 | 0,5 | 0,5 | 0,5 |
| **MC-Stream** | 0 | 0 | 0 | 20 | 20 | 10 | 5 | 5 | 5 | 5 |
| **MC-MDST(PB) Per Year** | 0,00 | 0,00 | 0,00 | 0,22 | 2,05 | 2,29 | 1,87 | 2,26 | 2,26 | 2,26 |
| **MC-MDST-TOT(PB)** | 0,00 | 0,00 | 0,00 | 0,22 | 2,28 | 4,57 | 6,44 | 8,70 | 10,96 | 13,23 |
| **MC-MDST-Reprocessing** | 0,00 | 0,00 | 0,00 | 0,89 | 9,10 | 9,14 | 3,22 | 4,35 | 5,48 | 6,61 |
| **MC-Ch(PB)** | 1 | 2 | 4 | 5,00 | 2,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 |
| **MC-MDST** | 0,00 | 0,00 | 0,00 | 1,11 | 11,16 | 11,43 | 5,09 | 6,61 | 7,74 | 8,88 |
| **BW MC Gbit/S** | 0,00 | 0,00 | 0,00 | 0,10 | 0,89 | 1,00 | 0,81 | 0,99 | 0,99 | 0,99 |
| **BW Repro Gbit/s** | 0,00 | 0,00 | 0,00 | 0,35 | 3,63 | 3,65 | 1,28 | 1,74 | 2,19 | 2,64 |
| **BW MC-Ch Gbit/s** | 1,62 | 1,94 | 1,94 | 1,62 | 0,55 | 0 | 0 | 0 | 0 | 0 |

# mDST-Monte Carlo Resume

| | 2015 | 2016 | 2017 | 2018 | 2019 | 2020 | 2021 | 2022 | 2023 | 2024 | 2015 | 2016 | 2017 | 2018 | 2019 | 2020 | 2021 | 2022 | 2023 | 2024 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | In-Band | In-Band | In-Band | In-Band | In-Band | In-Band | In-Band | In-Band | In-Band | In-Band | Out-Band | Out-Band | Out-Band | Out-Band | Out-Band | Out-Band | Out-Band | Out-Band | Out-Band | Out-Band |
| KEK | 0,00 | 0,00 | 0,00 | 0,23 | 2,26 | 2,32 | 1,05 | 1,36 | 1,59 | 1,81 | 0,00 | 0,00 | 0,00 | 0,23 | 2,26 | 2,32 | 1,05 | 1,36 | 1,59 | 1,81 |
| PNNL | 0,00 | 0,00 | 0,00 | 0,14 | 1,36 | 1,39 | 0,63 | 0,82 | 0,95 | 1,09 | 0,00 | 0,00 | 0,00 | 0,14 | 1,36 | 1,39 | 0,63 | 0,82 | 0,95 | 1,09 |
| Korea | 0,00 | 0,00 | 0,00 | 0,05 | 0,45 | 0,46 | 0,21 | 0,27 | 0,32 | 0,36 | 0,00 | 0,00 | 0,00 | 0,05 | 0,45 | 0,46 | 0,21 | 0,27 | 0,32 | 0,36 |
| Germany | 0,00 | 0,00 | 0,00 | 0,12 | 1,18 | 1,21 | 0,55 | 0,71 | 0,82 | 0,94 | 0,00 | 0,00 | 0,00 | 0,12 | 1,18 | 1,21 | 0,55 | 0,71 | 0,82 | 0,94 |
| Italy | 0,00 | 0,00 | 0,00 | 0,11 | 1,09 | 1,11 | 0,50 | 0,65 | 0,76 | 0,87 | 0,00 | 0,00 | 0,00 | 0,11 | 1,09 | 1,11 | 0,50 | 0,65 | 0,76 | 0,87 |
| Slovenia | 0,00 | 0,00 | 0,00 | 0,02 | 0,18 | 0,19 | 0,08 | 0,11 | 0,13 | 0,14 | 0,00 | 0,00 | 0,00 | 0,02 | 0,18 | 0,19 | 0,08 | 0,11 | 0,13 | 0,14 |
| Australia | 0,00 | 0,00 | 0,00 | 0,03 | 0,27 | 0,28 | 0,13 | 0,16 | 0,19 | 0,22 | 0,00 | 0,00 | 0,00 | 0,03 | 0,27 | 0,28 | 0,13 | 0,16 | 0,19 | 0,22 |
| Canada | 0,00 | 0,00 | 0,00 | 0,03 | 0,27 | 0,28 | 0,13 | 0,16 | 0,19 | 0,22 | 0,00 | 0,00 | 0,00 | 0,03 | 0,27 | 0,28 | 0,13 | 0,16 | 0,19 | 0,22 |
| Austria | 0,00 | 0,00 | 0,00 | 0,02 | 0,18 | 0,19 | 0,08 | 0,11 | 0,13 | 0,14 | 0,00 | 0,00 | 0,00 | 0,02 | 0,18 | 0,19 | 0,08 | 0,11 | 0,13 | 0,14 |
| China | 0,00 | 0,00 | 0,00 | 0,03 | 0,27 | 0,28 | 0,13 | 0,16 | 0,19 | 0,22 | 0,00 | 0,00 | 0,00 | 0,03 | 0,27 | 0,28 | 0,13 | 0,16 | 0,19 | 0,22 |
| Czech Rep. | 0,00 | 0,00 | 0,00 | 0,01 | 0,09 | 0,09 | 0,04 | 0,05 | 0,06 | 0,07 | 0,00 | 0,00 | 0,00 | 0,01 | 0,09 | 0,09 | 0,04 | 0,05 | 0,06 | 0,07 |
| India | 0,00 | 0,00 | 0,00 | 0,03 | 0,27 | 0,28 | 0,13 | 0,16 | 0,19 | 0,22 | 0,00 | 0,00 | 0,00 | 0,03 | 0,27 | 0,28 | 0,13 | 0,16 | 0,19 | 0,22 |
| Malaysia* | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 |
| Mexico | 0,00 | 0,00 | 0,00 | 0,02 | 0,18 | 0,19 | 0,08 | 0,11 | 0,13 | 0,14 | 0,00 | 0,00 | 0,00 | 0,02 | 0,18 | 0,19 | 0,08 | 0,11 | 0,13 | 0,14 |
| Poland | 0,00 | 0,00 | 0,00 | 0,02 | 0,18 | 0,19 | 0,08 | 0,11 | 0,13 | 0,14 | 0,00 | 0,00 | 0,00 | 0,02 | 0,18 | 0,19 | 0,08 | 0,11 | 0,13 | 0,14 |
| Russia | 0,00 | 0,00 | 0,00 | 0,07 | 0,72 | 0,74 | 0,34 | 0,44 | 0,51 | 0,58 | 0,00 | 0,00 | 0,00 | 0,07 | 0,72 | 0,74 | 0,34 | 0,44 | 0,51 | 0,58 |
| Saudi Arabia | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 |
| Spain* | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 |
| Taiwan | 0,00 | 0,00 | 0,00 | 0,03 | 0,27 | 0,28 | 0,13 | 0,16 | 0,19 | 0,22 | 0,00 | 0,00 | 0,00 | 0,03 | 0,27 | 0,28 | 0,13 | 0,16 | 0,19 | 0,22 |
| Thailand* | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 |
| Turkey | 0,00 | 0,00 | 0,00 | 0,01 | 0,09 | 0,09 | 0,04 | 0,05 | 0,06 | 0,07 | 0,00 | 0,00 | 0,00 | 0,01 | 0,09 | 0,09 | 0,04 | 0,05 | 0,06 | 0,07 |
| Ukraine | 0,00 | 0,00 | 0,00 | 0,01 | 0,09 | 0,09 | 0,04 | 0,05 | 0,06 | 0,07 | 0,00 | 0,00 | 0,00 | 0,01 | 0,09 | 0,09 | 0,04 | 0,05 | 0,06 | 0,07 |
| Viet Nam | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 | 0,00 |

# Total Traffic Scenario 1



**Total-In-Band**

8Gbps PNNL
10Gbps Europe

**Total-Out-Band**

19Gbps-KEK

# Total Traffic Scenario 2