# Network issues on FR cloud

Eric Lançon (CEA-Saclay/Irfu)

# Network Usage

- Data distribution
- MC production
- Analysis
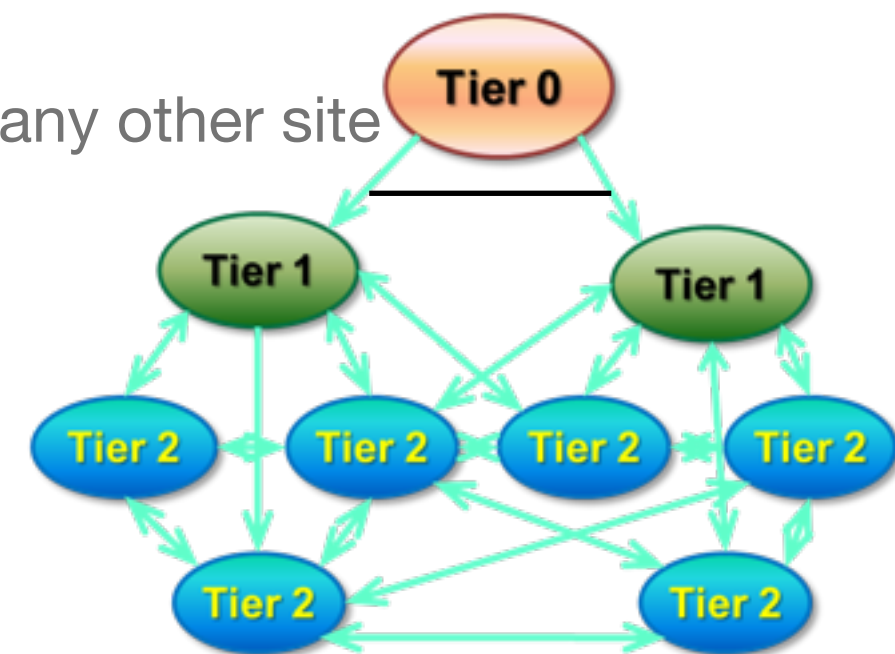- Distributed storage

# Network used for

- **Data distribution**, 2 components :

  - Pre-placed data (a la MONARC)

  - Dynamic data distribution (popular data to available sites)

- **MC production**

  - Within a given cloud

  - Across clouds

- **Analysis**

  - Retrieving results

  - Small data sets

- **Distributed storage**

  - To optimize resources

  - Simplify data management

## The 'old' computing model is dying

irfu

cea

saclay

# The ATLAS Data Model has changed

**USE THE NETWORK**

- Moved away from the historical model

- 4 recurring themes:

  - **Flat(ter) hierarchy**: Any site can replicate data from any other site

  - **Multi Cloud Production**

    - Need to replicate output files to remote Tier-1

  - **Dynamic data caching**: Analysis sites receive datasets from any other site "on demand" based on usage pattern

    - Possibly in combination with pre-placement of data sets by centrally managed replication of datasets

  - **Remote data access**: local jobs accessing data stored at remote sites

- **ATLAS is now heavily relying on multi-domain networks and needs decent e2e network monitoring**

irfu

cea

saclay

# ATLAS sites and connectivity

- ATLAS computing model has (will continue to) changed

  - More experience

  - More tools and monitoring

- New category of sites : Direct T2s (**T2Ds**)

  - Primary hosts for datasets (**analysis**) and for group analysis

  - Get and send data from different clouds
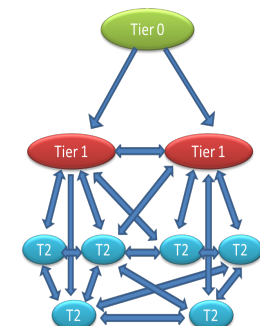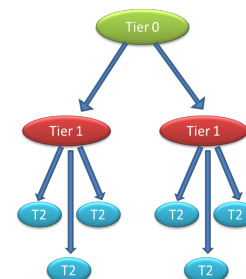
  - Participate in cross cloud production

## T2D: revising the criteria

New criteria - under evaluation

- All transfers from the candidate T2D to **9**/12 T1s for big files ('L') must be above 5 MB/s during the last week and during 3 out of the **5** last weeks.

- All transfers from **9**/12 T1s to the candidate T2D for big files must be above 5 MB/s during the last week and during 3 out of the **5** last weeks

http://gnegri.web.cern.ch/gnegri/T2D/t2dStats.html

**FR-cloud T2Ds** : BEIJING, GRIF-LAL, GRIF-LPNHE, IN2P3-CPPM, IN2P3-LAPP, IN2P3-LPC, IN2P3-LPSC

irfu
cea
saclay

# Network performance monitoring

- **Networking accounting :**

  - **Organized** (FTS) file transfers : http://dashb-atlas-data.cern.ch/ddm2/, not for direct transfers by users (dq2-get)

- **ATLAS 'sonar' :**

  - Calibrated file transfers by ATLAS Data Distribution system, from **storage to storage** : http://bourricot.cern.ch/dq2/ftsmon/

  - > 1 GB file transfers used to monitor and validate T2Ds

- **perfSONAR (PS) :**

  - **Network performance** (throughput, latency) : http://perfsonar.racf.bnl.gov:8080/exda/

  - Located as close as possible to storage at site and with similar connection hardware

irfu

cea

saclay

# T0 exports over a year

**Over 1GB/s**

Better LHC efficiency and higher trigger rate

**16,336 TB to T1s**



**1.2 GB/s**

LHC winter stop

**2,186 TB to CCIN2P3**

**150 MB/s**

**78% to T1s**

irfu

cea

saclay

# T1→T1

**23,966 TB**



**1GB/s**

*Popularity based*

*Pre-placement*

*Group data*

*Cross-cloud MC production*

*User requests*

irfu

cea

saclay

# T1→T2         T2→T1

## 54,491 TB         16,487 TB



**Reconstruction in T2s**         **User subscriptions!**
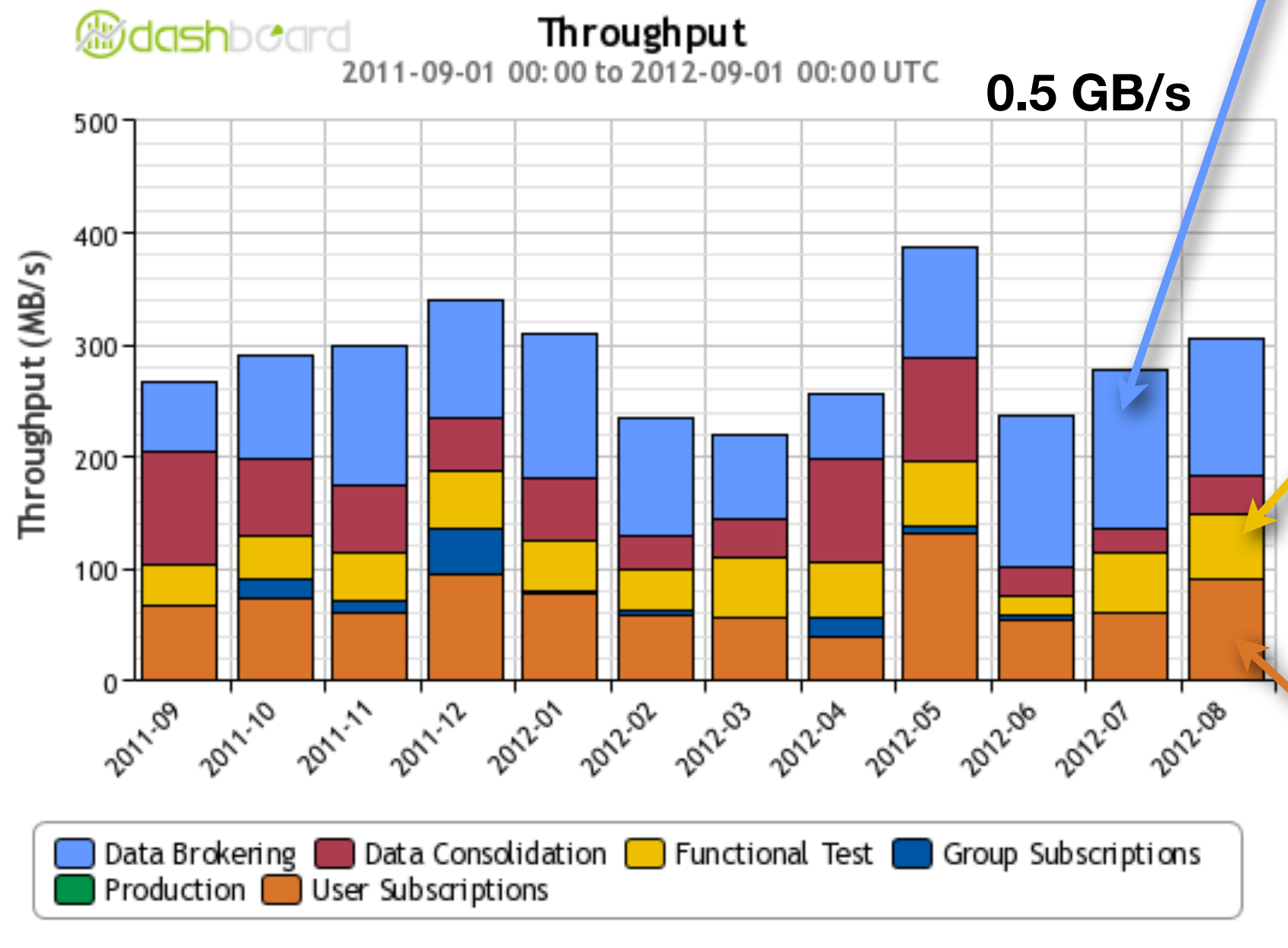
**Dynamic data placement** > **pre-defined**

*Not in original computing model*   9

**T2→T2**

**Dynamic data placement > pre-defined**

**9,050 TB**

Network mesh tests

User subscriptions
Group data at some T2s
+ Outputs of analysis

irfu
cea
saclay

# ALL together

## 131,473 TB



**Throughput**
2011-09-01 00:00 to 2012-09-01 00:00 UTC

5 GB/s

Legend: Data Brokering, Data Consolidation, Functional Test, Group Subscriptions, Production, T0 Export, User Subscriptions

**Data volume: T1s 60% of sources**



**Transfer Volume**
2011-09-01 00:00 to 2012-09-01 00:00 UTC

3 : 0 % (228.801 TB)
2 : 20 % (26854.312 TB)
0 : 18 % (23391.671 TB)
1 : 62 % (80997.954 TB)

**Activity: T2s 40% of sources**



**Transfer Successes**
2011-09-01 00:00 to 2012-09-01 00:00 UTC

3 : 1 % (2071182)
0 : 10 % (33460399)
2 : 41 % (143709377)
1 : 49 % (170600989)

irfu
cea
saclay

# ALL together

**Data volume: pre-placement ~2 times dynamic placement room for improvements**



**Activity: users ~30% of transfers**



irfu

cea

saclay

# Cross-cloud MC production

- The 'easy' part

- No need to be a T2D

- Only connection to remote T1 needed

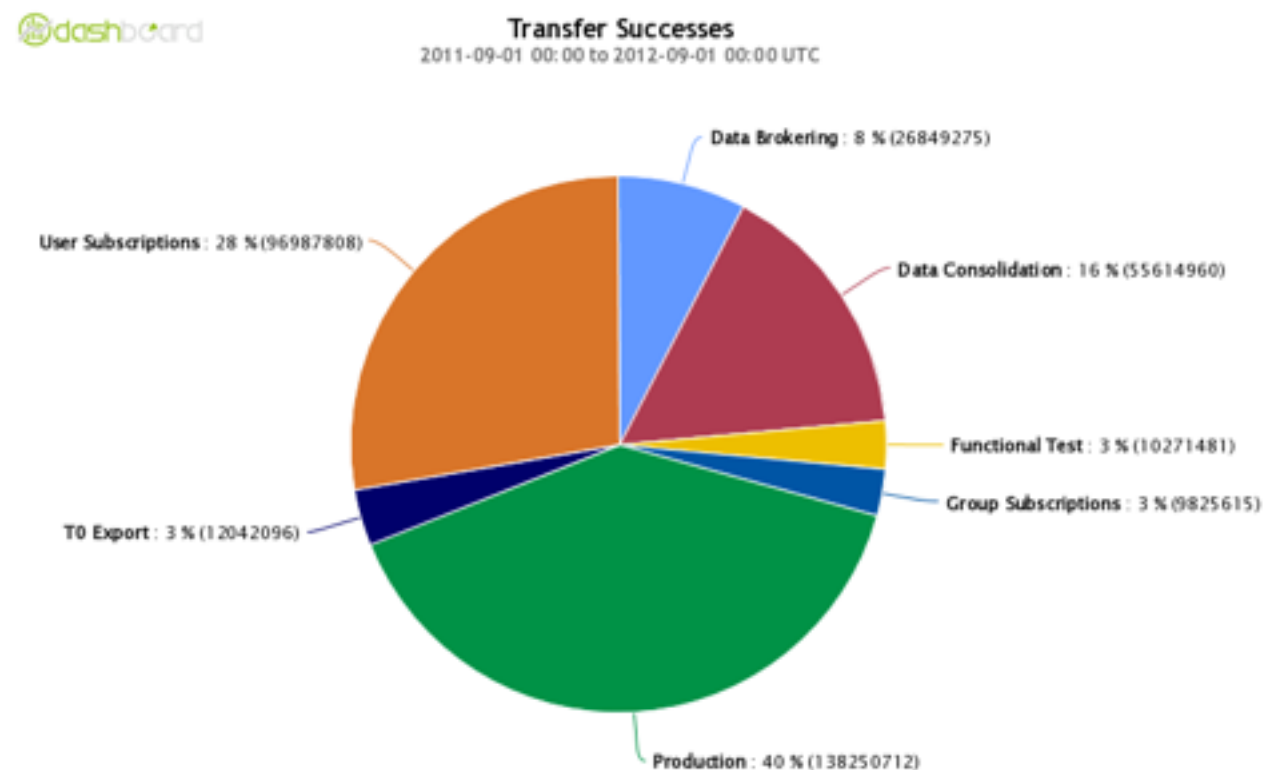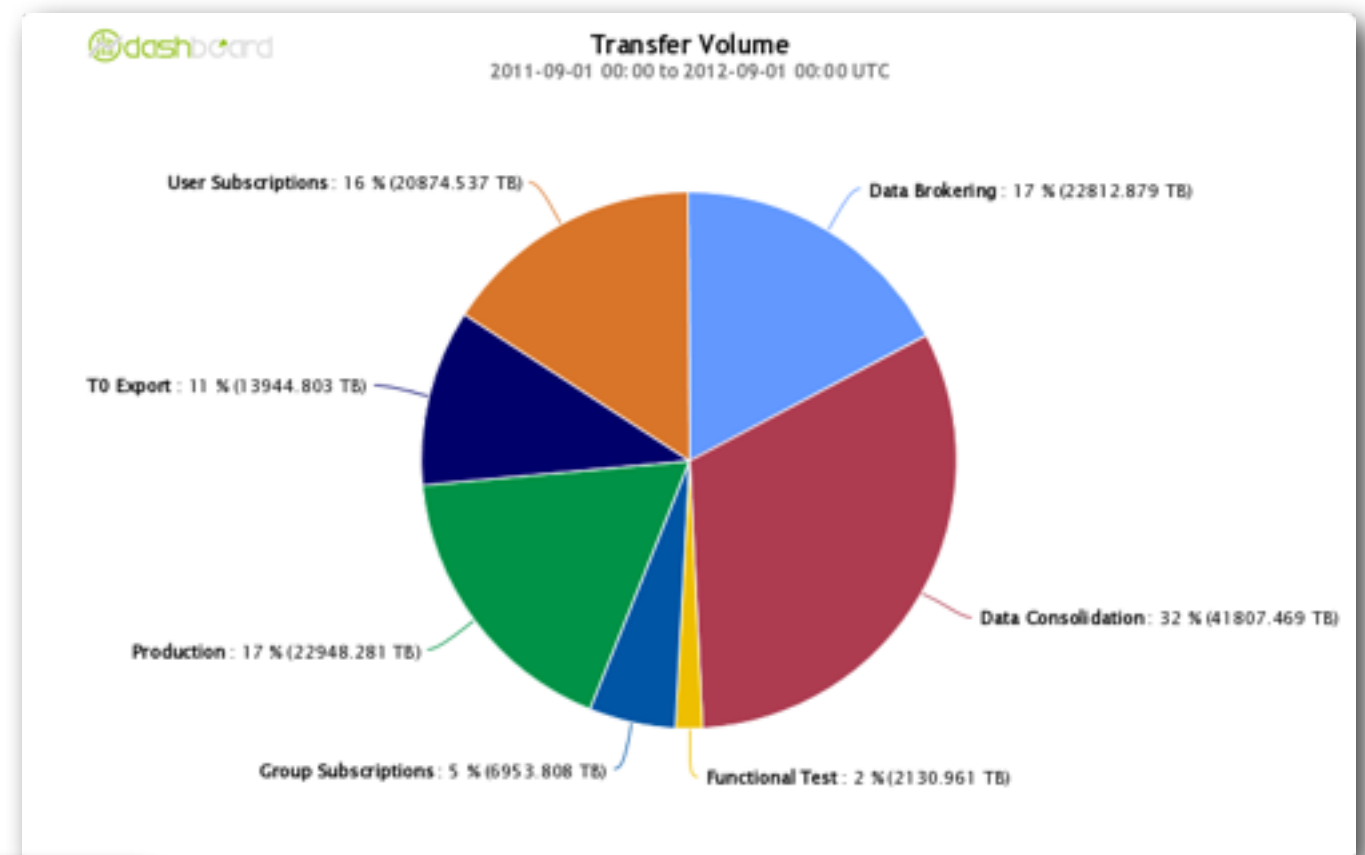- Example : NL cloud

  - 65 ! sites contributing



Completed jobs (Sum: 414,631)

JINR-LCG2 - 12.30%
SARA-MATRIX - 25.95%
50,994
107,616

SARA-MATRIX - 25.95% (107,616)    JINR-LCG2 - 12.30% (50,994)
RRC-KI - 7.51% (31,147)    RU-PROTVINO-IHEP - 7.00% (29,006)
NIKHEF-ELPROD - 6.62% (27,459)    IL-TAU-HEP - 4.25% (17,629)
WEIZMANN-LCG2 - 3.96% (16,424)    TECHNION-HEP - 3.75% (15,539)
PIC - 2.65% (11,003)    CSTCDIE - 2.34% (9,693)
GRIF-LAL - 1.84% (7,641)    DESY-HH - 1.47% (6,115)
SIGNET - 1.45% (5,994)    NDGF-T1 - 1.44% (5,950)
ARNES - 1.34% (5,574)    TR-10-ULAKBIM - 1.33% (5,523)
LRZ-LMU - 1.20% (4,992)    IN2P3-LAPP - 1.20% (4,959)
MPPMU - 1.03% (4,260)    plus 46 more

irfu

cea

saclay

# The French cloud

- The most 'exploded' cloud of ATLAS
- 4 Romanians sites at the far end of GEANT
- 2 sites in far east Beijing & Tokyo connected to CCIN2P3 via different paths

The "French" cloud

T1 : Lyon

T2s : 14 sites
- Annecy
- Clermont
- Grenoble
- Grif (3 sites)
- Lyon
- Marseille
- Beijing
- Romania x4
- Tokyo

# Data exchanges for FR-cloud
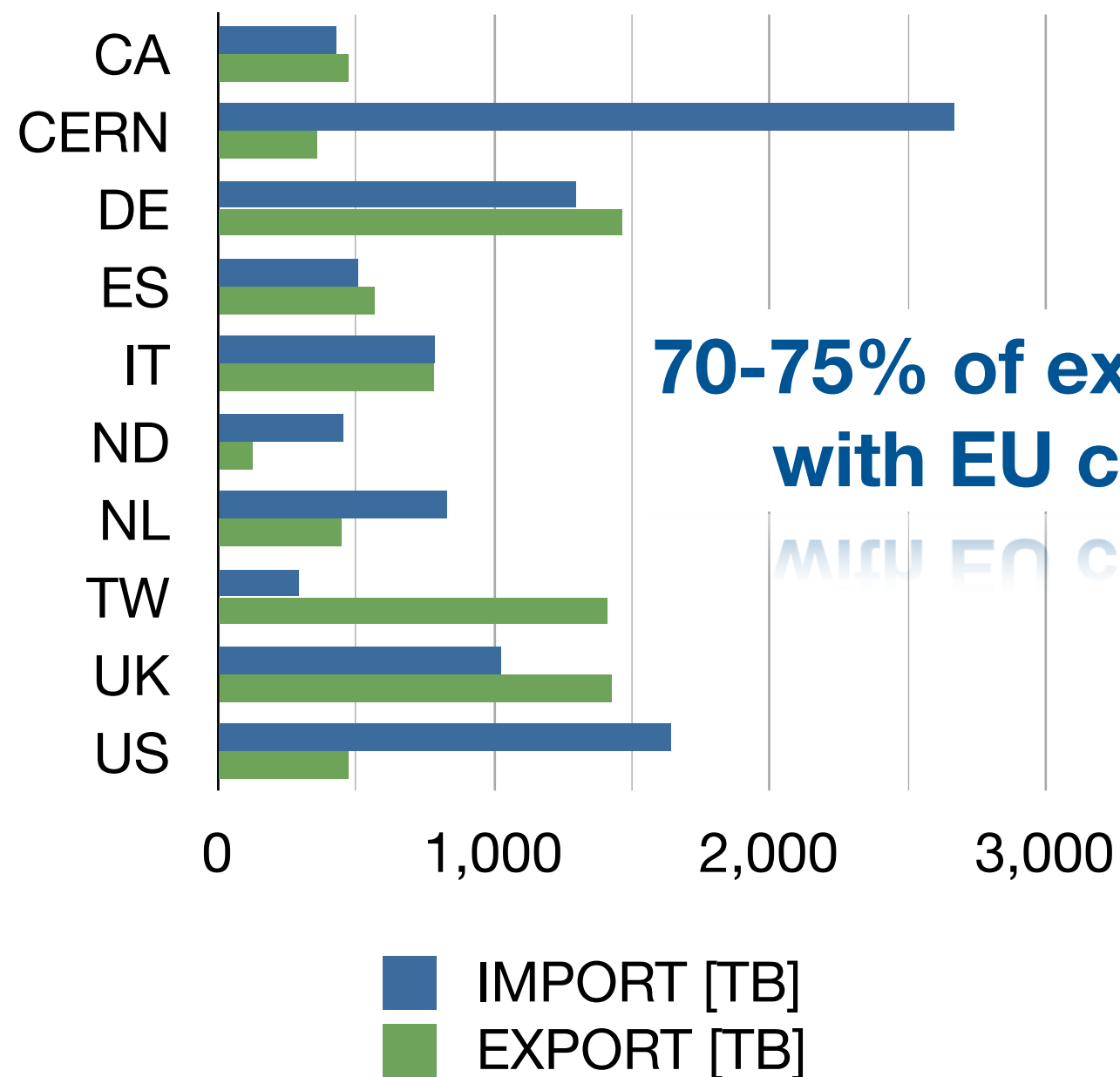
**Imports : 9,9 PB**

**French cloud exchanges**

**70-75% of exchanges with EU clouds**

**Exports : 7,5 PB**

IMPORT [TB]
EXPORT [TB]

[Sep. 2011 - Sep. - 2012]

Rencontre LCG-France, SUBATECH Nantes, septembre 2012

irfu
cea
saclay

cross-cloud production

# Data volume transferred to French T2s



**Transfer Volume**
2011-09-01 00:00 to 2012-09-01 00:00 UTC

IN2P3-LPSC : 11 % (798.031 TB)

GRIF-IRFU : 18 % (1384.526 TB)

IN2P3-LPC : 15 % (1149.589 TB)

GRIF-LAL : 13 % (995.591 TB)

IN2P3-LAPP : 16 % (1238.394 TB)

GRIF-LPNHE : 18 % (1374.564 TB)

IN2P3-CPPM : 9 % (659.162 TB)

## not proportional to number physicist nor CPUs

irfu

cea

saclay

# Connectivity within French cloud (ATLAS sonar)



**T2s→T1**

**T1→T2s**

5 MB/s

Beijing

# Origin of data transferred to 2 GRIF sites

LAL : T2D
connected to LHCONE

IRFU : ~~T2D connected to LHCONE~~



66 OTHERS : 3 % (25.286 TB)
TRIUMF-LCG2 : 3 % (33.168 TB)
TOKYO-LCG2 : 1 % (10.143 TB)
TAIWAN-LCG2 : 3 % (28.233 TB)
SARA-MATRIX : 6 % (62.392 TB)
RAL-LCG2 : 6 % (61.071 TB)
PIC : 4 % (43.856 TB)
NIKHEF-ELPROD : 3 % (26.205 TB)
NDGF-T1 : 4 % (41.840 TB)
INFN-T1 : 8 % (79.898 TB)
IN2P3-LPSC : 0 % (3.839 TB)
IN2P3-LPC : 1 % (8.419 TB)
IN2P3-LAPP : 1 % (7.347 TB)
IN2P3-CPPM : 0 % (4.512 TB)
BEIJING-LCG2 : 0 % (3.051 TB)
BNL-OSG2 : 12 % (120.938 TB)
CERN-PROD : 8 % (79.269 TB)
FZK-LCG2 : 8 % (80.758
GRIF-IRFU : 1 % (9.352 TB
GRIF-LPNHE : 1 % (10.025
IN2P3-CC : 26 % (255.988 TB)

**From 'everywhere'**



70 OTHERS : 1 % (18.964 TB)
WUPPERTALPROD : 0 % (1.617 TB)
TOKYO-LCG2 : 1 % (9.782 TB)
RO-07-NIPNE : 0 % (5.686 TB)
RO-02-NIPNE : 0 % (5.323 TB)
NDGF-T1 : 0 % (2.717 TB)
INFN-T1 : 1 % (9.197 TB)
IN2P3-LPSC : 2 % (20.934 TB)
IN2P3-LPC : 1 % (19.671 TB)
IN2P3-LAPP : 2 % (26.311 TB)
IN2P3-CPPM : 1 % (13.764 TB)
BEIJING-LCG2 : 1 % (7.873 TB)
BNL-OSG2 : 1 % (7.192 TB)
CERN-PROD : 1 % (18.742 TB)
DESY-HH : 0 % (2.017 TB)
DESY-ZN : 0 % (2.755 TB)
FZK-LCG2 : 1 % (14.503 TB)
GRIF-IRFU : 1 % (9.229 TB)
GRIF-LAL : 2 % (23.478 TB)
GRIF-LPNHE : 4 % (59.609 TB)
IN2P3-CC : 80 % (1105.160 TB)

**80% from CCIN2P3**

irfu

cea

saclay

# User & group analysis

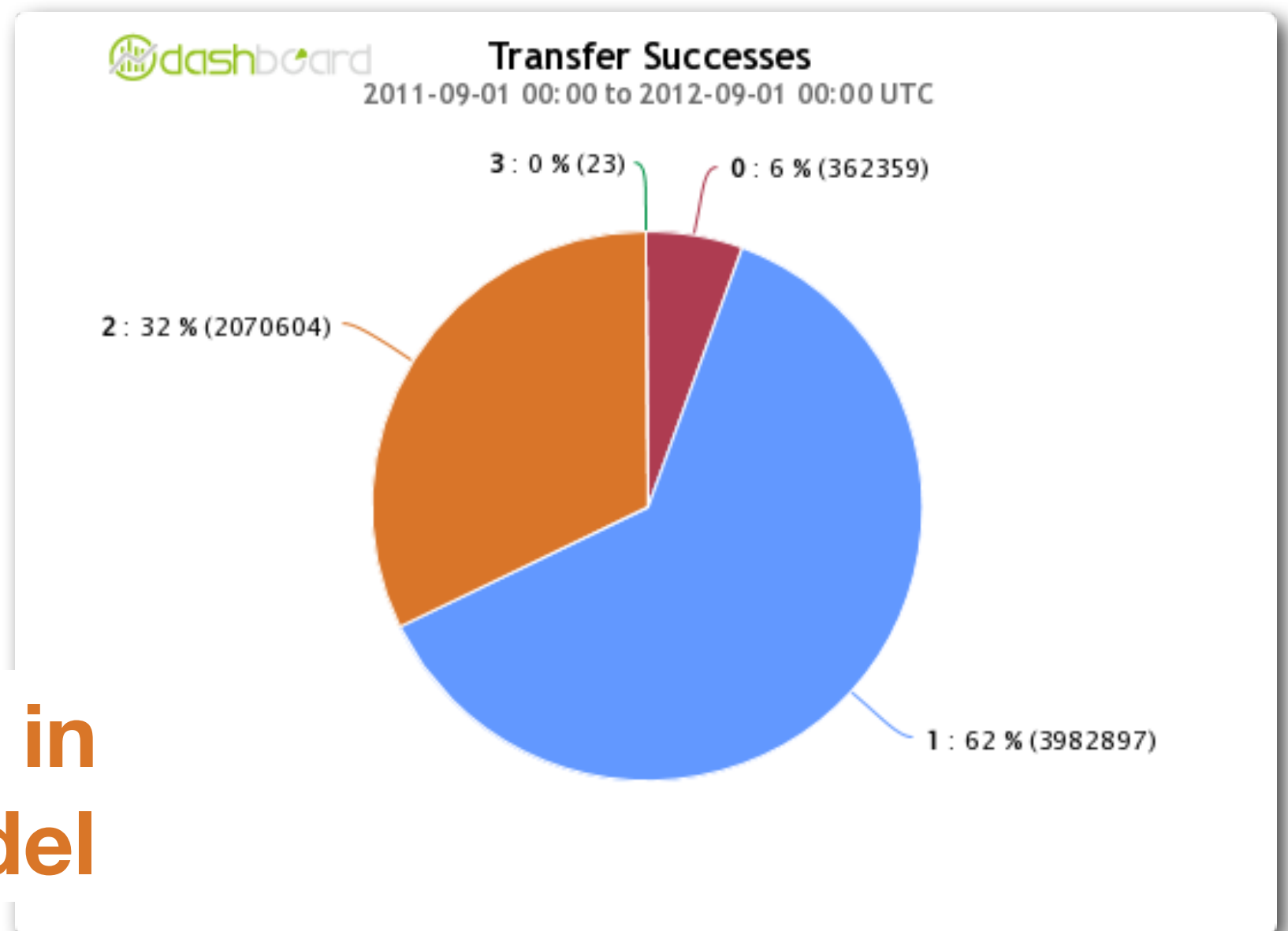- Most of users run final analysis at their local site : delivery of data or analysis job outputs to users **very sensitive** (the last transferred file determines the efficiency)

- 2 ways to get (reduced format) data

  - The majority : Let PanDA decide where to run the jobs ; where data are (in most of the cases). Outputs stored on SCRATCHDISK (buffer area) at remote sites have to be shipped back to user local site

  - Get limited volume of data at local site to run locally analysis

- Both imply data transfers from remote site to local site

  - Direct transfers for T2Ds

  - Through T1 for other T2s

irfu

cea
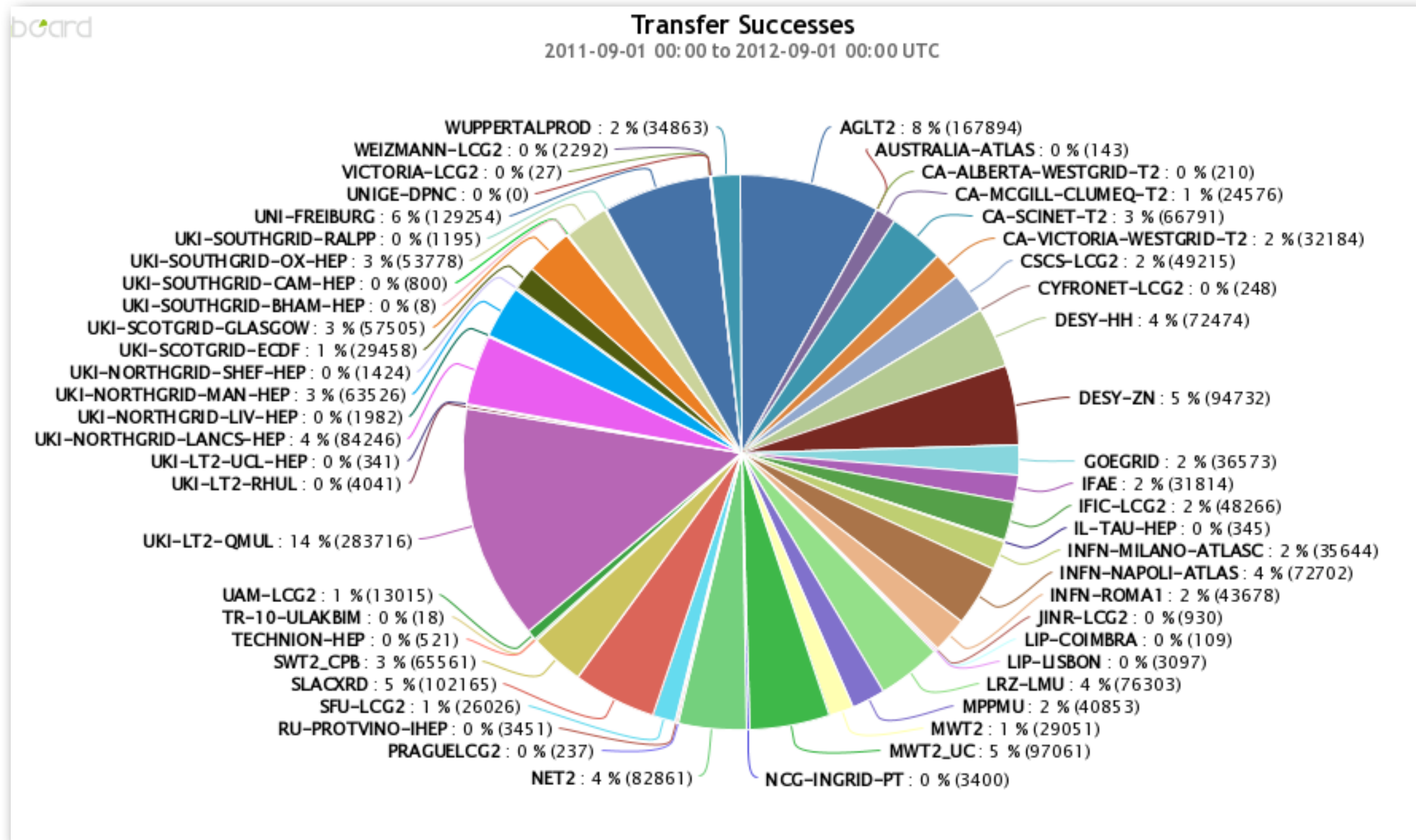
saclay

# User + Group transfers to FR-cloud sites

**⅓** of data transfers (not data volume) used for final analysis from **T2 sites** **outside** FR-cloud
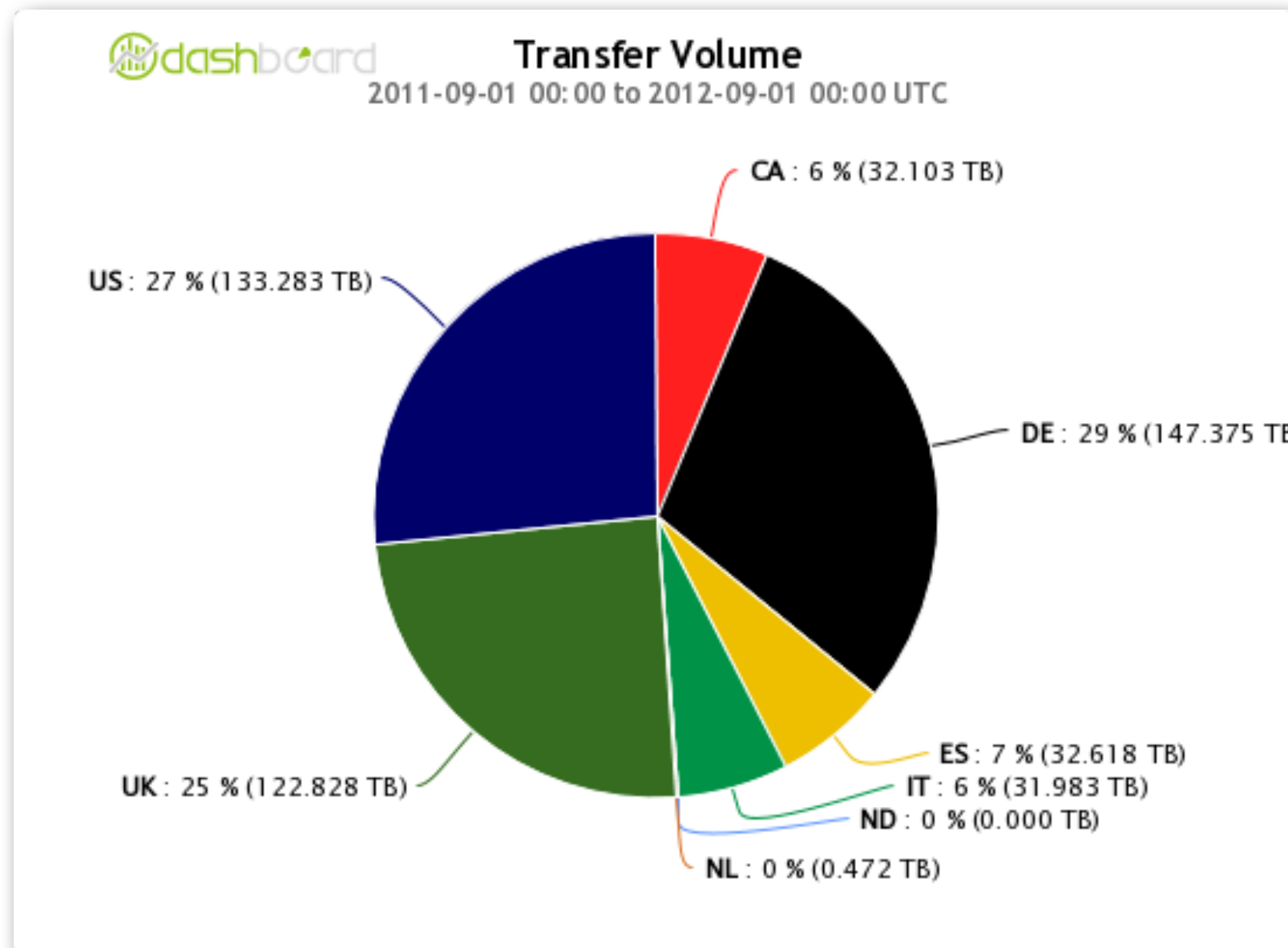
1,828 TB



**dashboard**

**Transfer Successes**
2011-09-01 00:00 to 2012-09-01 00:00 UTC

3 : 0 % (23)
0 : 6 % (362359)
2 : 32 % (2070604)
1 : 62 % (3982897)

**Not allowed in original model**

irfu

cea

saclay

# Data come from 52 T2 sites



**Transfer Successes**
2011-09-01 00:00 to 2012-09-01 00:00 UTC

WUPPERTALPROD : 2 % (34863)
WEIZMANN-LCG2 : 0 % (2292)
VICTORIA-LCG2 : 0 % (27)
UNIGE-DPNC : 0 % (0)
UNI-FREIBURG : 6 % (129254)
UKI-SOUTHGRID-RALPP : 0 % (1195)
UKI-SOUTHGRID-OX-HEP : 3 % (53778)
UKI-SOUTHGRID-CAM-HEP : 0 % (800)
UKI-SOUTHGRID-BHAM-HEP : 0 % (8)
UKI-SCOTGRID-GLASGOW : 3 % (57505)
UKI-SCOTGRID-ECDF : 1 % (29458)
UKI-NORTHGRID-SHEF-HEP : 0 % (1424)
UKI-NORTHGRID-MAN-HEP : 3 % (63526)
UKI-NORTHGRID-LIV-HEP : 0 % (1982)
UKI-NORTHGRID-LANCS-HEP : 4 % (84246)
UKI-LT2-UCL-HEP : 0 % (341)
UKI-LT2-RHUL : 0 % (4041)

UKI-LT2-QMUL : 14 % (283716)

UAM-LCG2 : 1 % (13015)
TR-10-ULAKBIM : 0 % (18)
TECHNION-HEP : 0 % (521)
SWT2_CPB : 3 % (65561)
SLACXRD : 5 % (102165)
SFU-LCG2 : 1 % (26026)
RU-PROTVINO-IHEP : 0 % (3451)
PRAGUELCG2 : 0 % (237)
NET2 : 4 % (82861)

AGLT2 : 8 % (167894)
AUSTRALIA-ATLAS : 0 % (143)
CA-ALBERTA-WESTGRID-T2 : 0 % (210)
CA-MCGILL-CLUMEQ-T2 : 1 % (24576)
CA-SCINET-T2 : 3 % (66791)
CA-VICTORIA-WESTGRID-T2 : 2 % (32184)
CSCS-LCG2 : 2 % (49215)
CYFRONET-LCG2 : 0 % (248)
DESY-HH : 4 % (72474)

DESY-ZN : 5 % (94732)

GOEGRID : 2 % (36573)
IFAE : 2 % (31814)
IFIC-LCG2 : 2 % (48266)
IL-TAU-HEP : 0 % (345)
INFN-MILANO-ATLASC : 2 % (35644)
INFN-NAPOLI-ATLAS : 4 % (72702)
INFN-ROMA1 : 2 % (43678)
JINR-LCG2 : 0 % (930)
LIP-COIMBRA : 0 % (109)
LIP-LISBON : 0 % (3097)
LRZ-LMU : 4 % (76303)
MPPMU : 2 % (40853)
MWT2 : 1 % (29051)
MWT2_UC : 5 % (97061)
NCG-INGRID-PT : 0 % (3400)

irfu
cea
saclay

# ⅓ from non EU T2s



1/4 from UK (not on LHCONE)

Issues with distant T2s

# Beijing
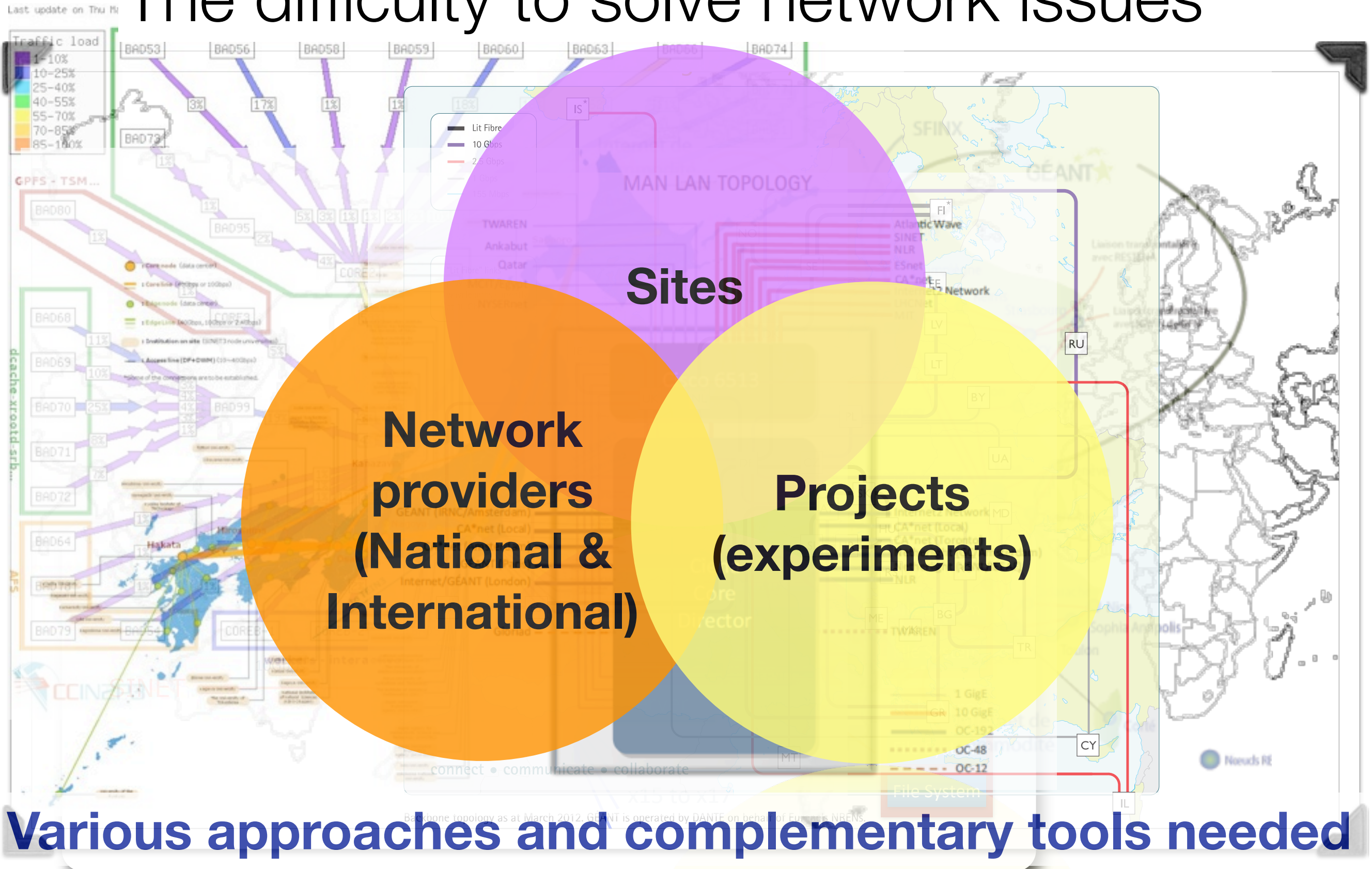
- **Beijing** :

  - Connected to Europe via GEANT/TEIN3 (except CERN : GLORIAD/KREONET)

  - RTT ~190 ms

- **Tokyo** :

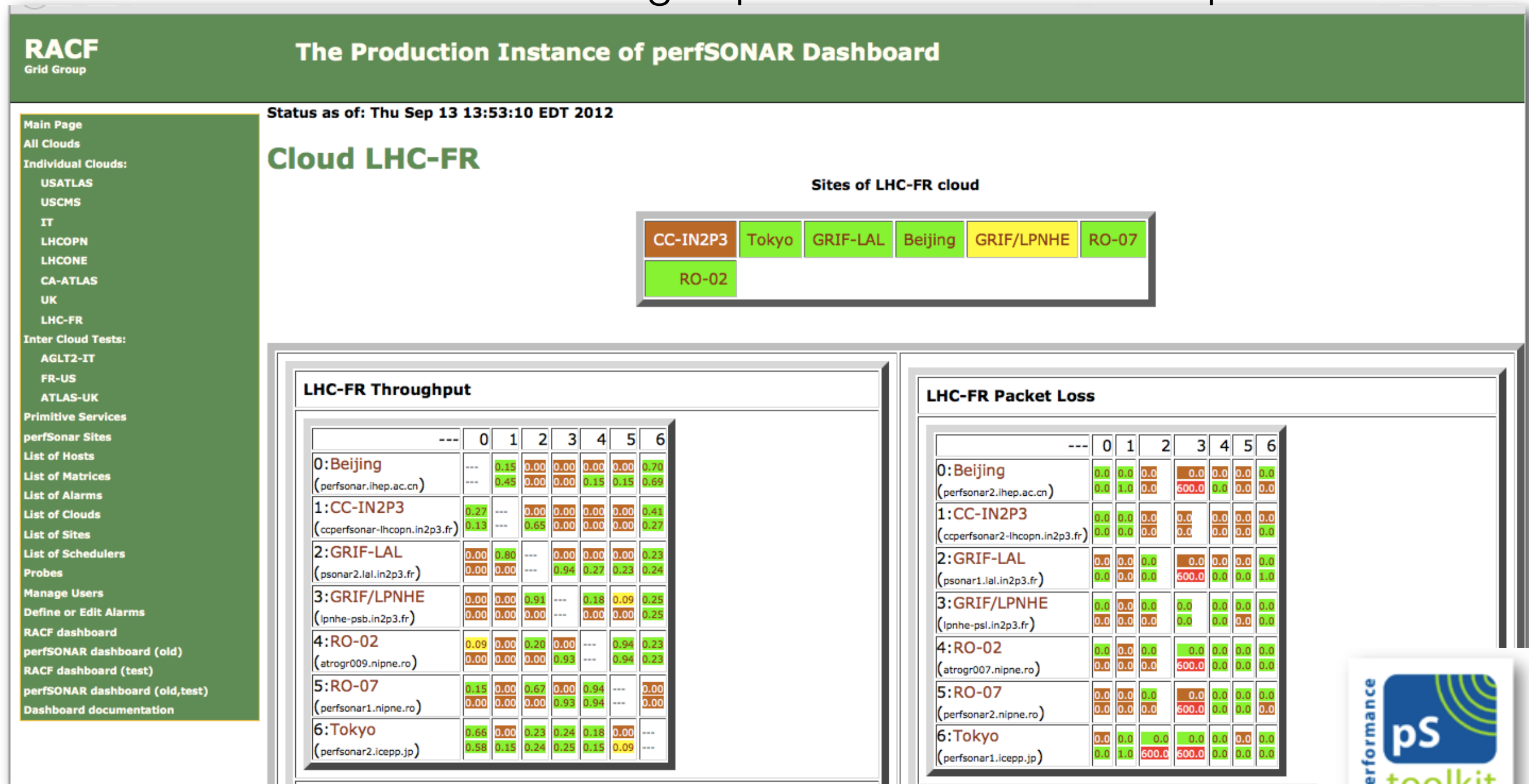  - Connected to Europe via GEANT/MANLAN/SINET4

  - RTT ~ 300 ms

- Several network operators on the path (Nationals, GEANT, …)

irfu

saclay

# The difficulty to solve network issues



**Sites**

**Network providers (National & International)**

**Projects (experiments)**

**Various approaches and complementary tools needed**

irfu

cea

saclay

# perfSonar dashboard of FR-cloud

## Being expanded as sites install perfSonar

# ATLAS transfers to Beijing since beg. 2011

**Beijing→ CCIN2P3**
**CCIN2P3 → Beijing**

*Performances changed over last year*
- *Asymmetry in transfer rate : why?*
- *Asymmetry reversed*

**Each 'event' explained sometime after some delay...**



irfu

cea

saclay

# ATLAS transfers to Tokyo since beg. 2011

**<span style="color:red">Tokyo→ CCIN2P3</span>**
**<span style="color:magenta">CCIN2P3 → Tokyo</span>**

*Performances changed over last year*
- *Asymmetry in transfer rate : why?*
- *Asymmetry reversed*

**Each 'event' explained sometime after some delay...**



irfu

cea

saclay

# Beijing from/to EU T1s

# Tokyo from/to EU T1s



**Each T1 is different**

**EU →Tokyo better for most T1s**

saclay

# Beijing - T1s asymmetry



**Throuhput [Mb/s]**

perfSonar (July)

- ○ T1->Beijing
- ○ Beijing->T1

400
300
200
100
0

IN2P3  KIT  PIC  CNAF  RAL  ND

FTSmon
(last 5 weeks)
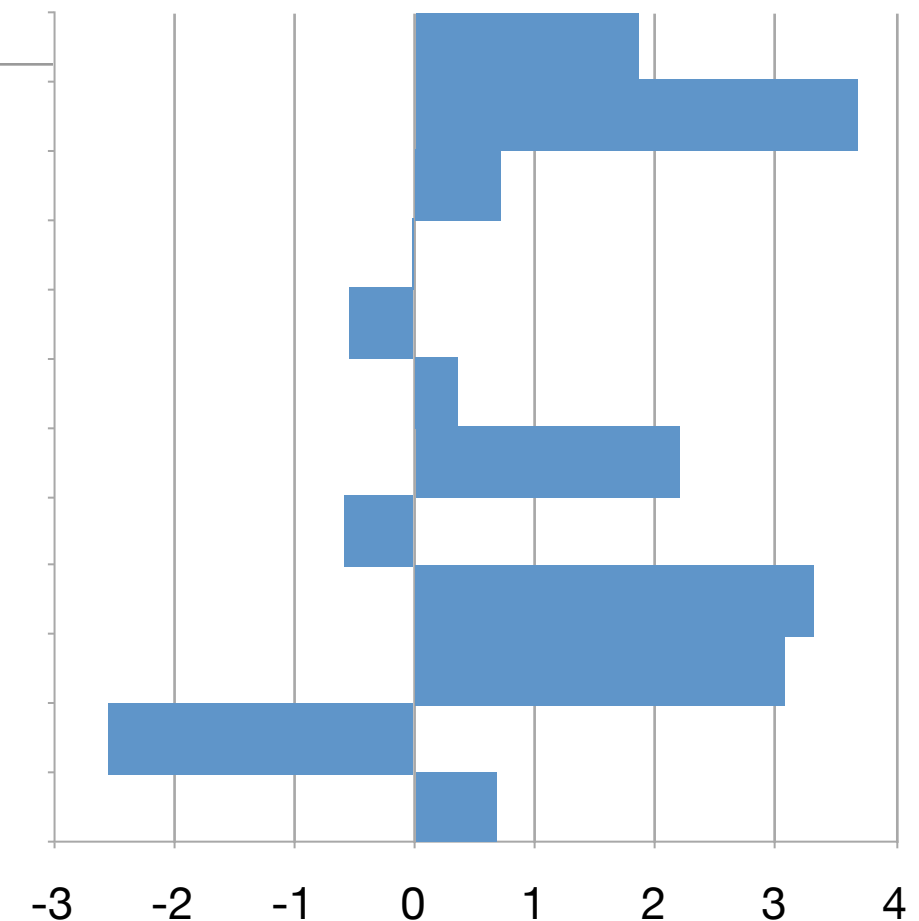
asymetry [MB/s]

TRIUMF-LCG2
INFN-T1
TAIWAN-LCG2
FZK-LCG2
PIC
BNL-OSG2
IN2P3-CC
NDGF-T1
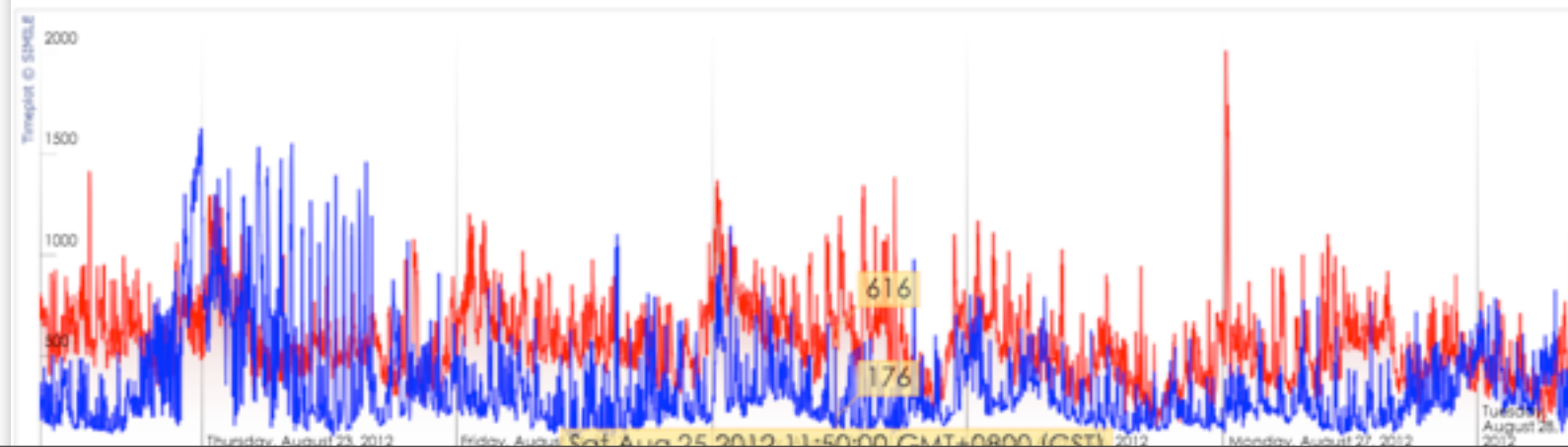NIKHEF-ELPROD
SARA-MATRIX
CERN-PROD
RAL-LCG2

-3  -2  -1  0  1  2  3  4

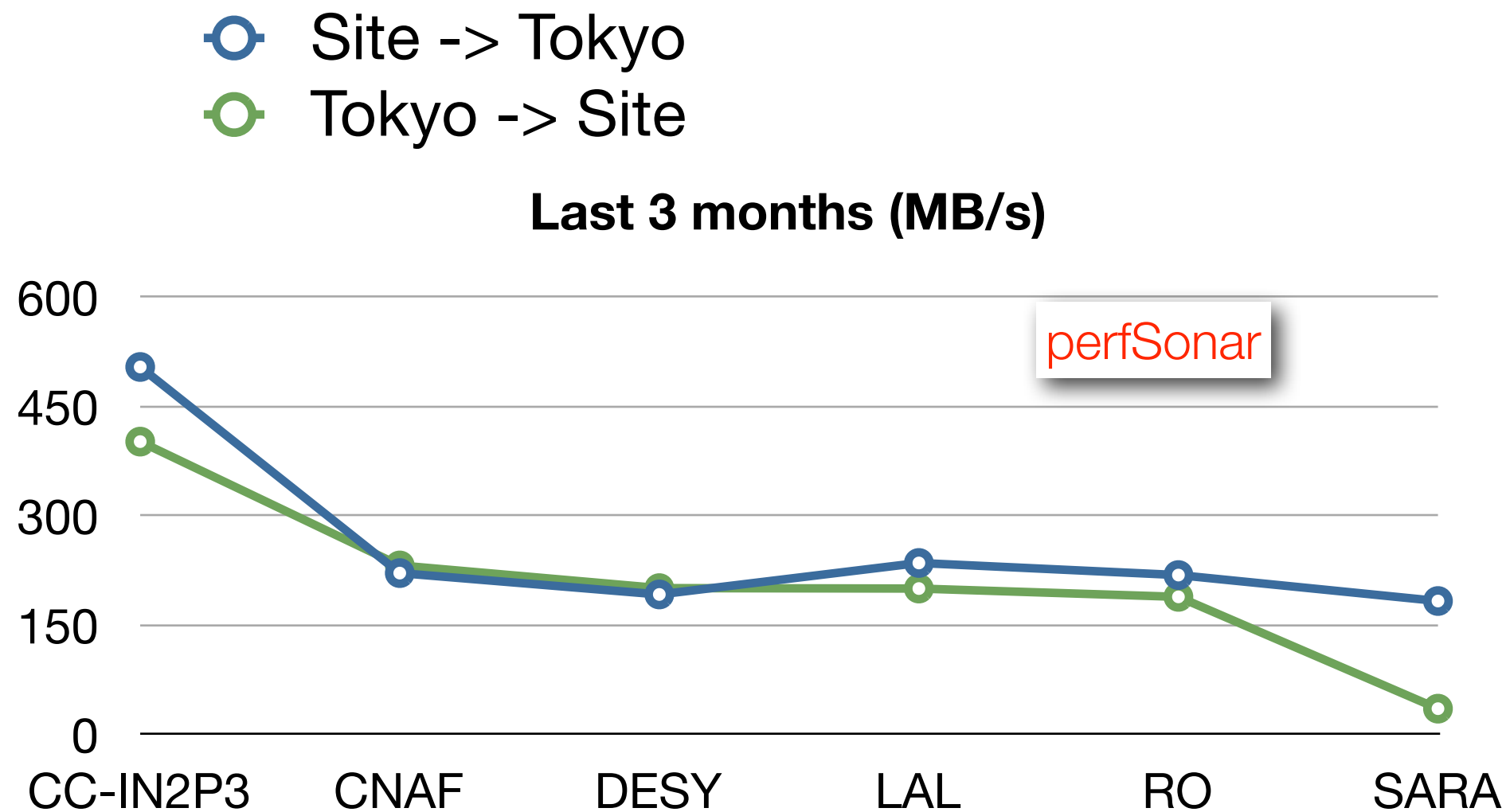## Throughput Monitor for Orient+ (Beijing-London)

This is the Throughput Monitor for the Link between Beijing and London (Orient+),
Thanks for the help from CERNet and Orient NOC,they gave IHEP the privilege to read the Router(CNGI-6IX) interface information.
We are reading the router interface information every 2 minutes.

London -> Beijing vs. Beijing - > London
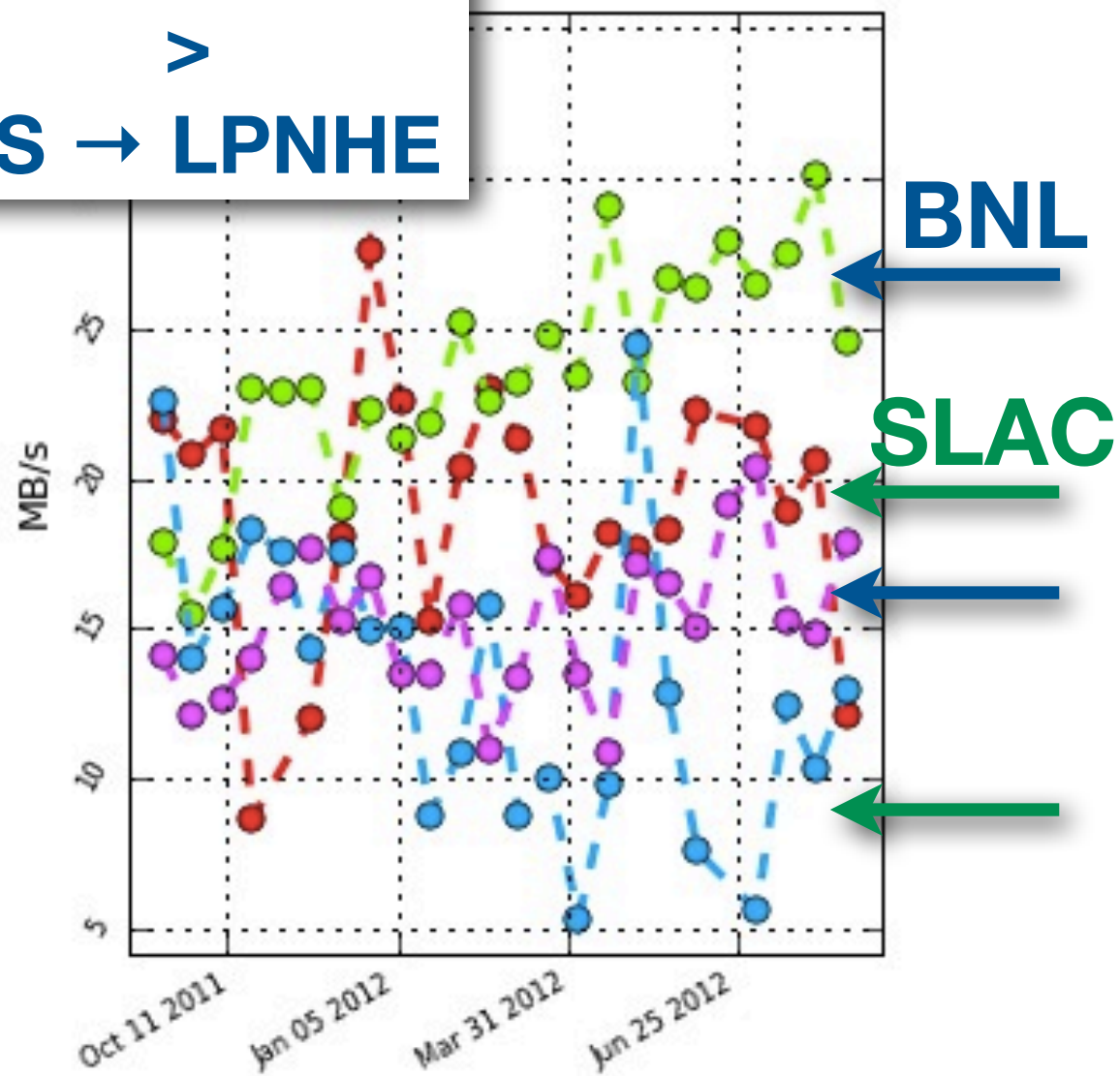
**Beijing -> T1s
>
T1s -> Beijing**

33

# Tokyo ↔ EU as seen by perfSonar

Site -> Tokyo
Tokyo -> Site

**Last 3 months (MB/s)**

perfSonar

irfu
cea
saclay

# US (LHCONE) ↔ GRIF (LHCONE)

irfu

cea

saclay

# DISTRIBUTED STORAGE / REMOTE ACCESS

- Better used of storage resources (disk prices!)
- Simplification of data management
- Eventually remote access (with caching at both ends); direct reading or file copy
- Bandwidth and stability needed

On going projects
- Storage : dCache, dpm,...
- Protocol : Xrootd, HTTP/WebDAV

# Existing distributed dCache systems : 2 examples

## MWT2 (Chicago)

## NORDUGRID

# Xrootd federation project in US    REMOTE ACCESS

## R&D Activity to Production

BNL Tier 1
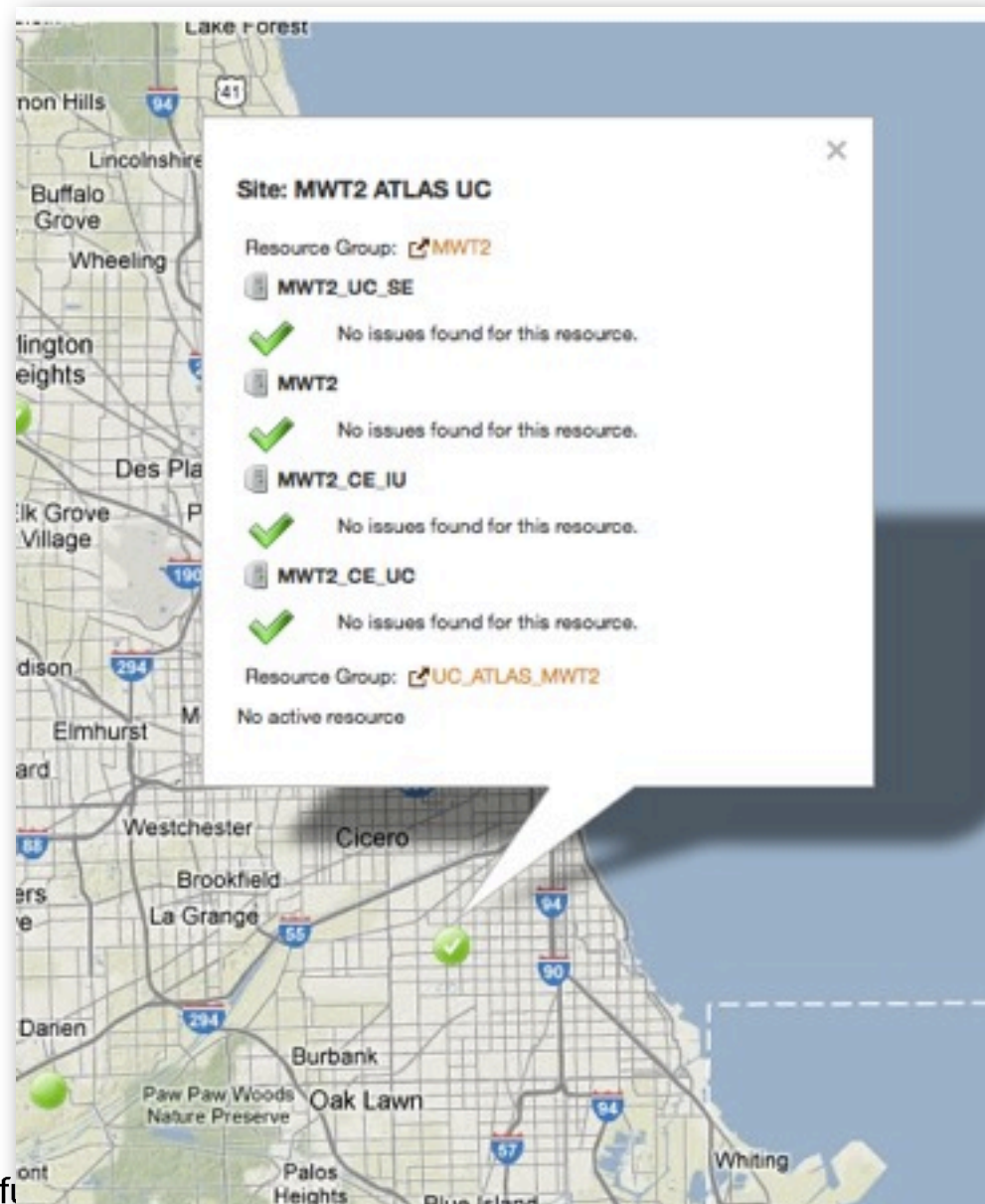AGLT2 (Tier 2)
MWT2 (Tier 2)
SWT2 (Tier 2)
SLAC (Tier 2)
ANL (Tier 3)
BNL (Tier 3)
Chicago (Tier 3)
Duke(Tier 3)
OU (Tier 3)
SLAC (Tier 3)
UTA (Tier 3)
NET (Tier 2)

- 2011 R&D project FAX (Federating ATLAS data stores using Xrootd) was deployed over US Tier 1, Tier 2s and some Tier3s

- Feasibility testing monitoring, site integrations

- In June 2012 extended effort to European sites as an ATLAS-wide project



2

## EU federation tests

### Four levels of redirection: site-cloud-zone-global



Start locally - expand search as needed

15

irfu

cea

saclay

# Topology validation

- Launch jobs to every site, test reading of site-specific files at every other site

- Parse client logs to infer resulting redirection
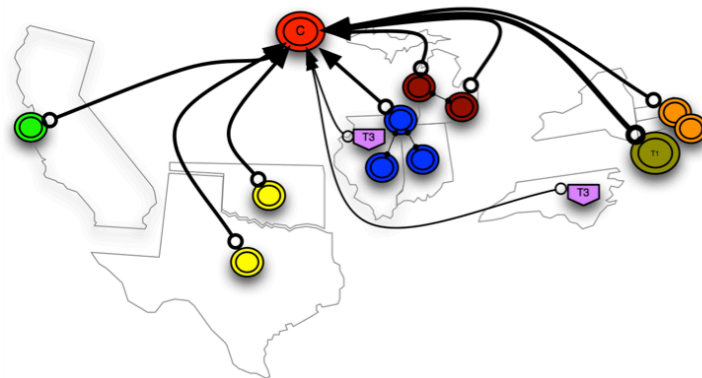


result from Ilija Vukotic

US regional

UK regional

DE regional

US-central regional

EU regional

http://ivukotic.web.cern.ch/ivukotic/FAX/index.asp

XRD redirector  federating site  16

## WAN Read Tests (basis for "cost matrix")



read time for 100% default cache

CERN to MWT2 (140 ms rtt)  277.93

MWT2 federation internal (<5 ms)  123.85

SLAC to MWT2 (50 ms)  222.83

100% default cache

Time to read file

Spread grows with network latency

Overall WAN performance adequate for a class of use cases

Date & Time

from Ilija Vukotic

17

irfu

cea

saclay

# HTTP/WebDAV

Storage Federations using standard web protocols

dCache.org

- Project with CERN DM under the umbrella of EMI but not limited to the EMI funding period.

- Definition of TEG:

  - "Collection of disparate storage resources managed by co-operating but independent administrative domains transparently accessible via a common name space"

- We do it with standard HTTP/WebDAV

Atlas WS 2012, CERN| dCache.org| 10 Sep 2012 | 11

## Use http urls for input files

- DDM enabled web redirection service
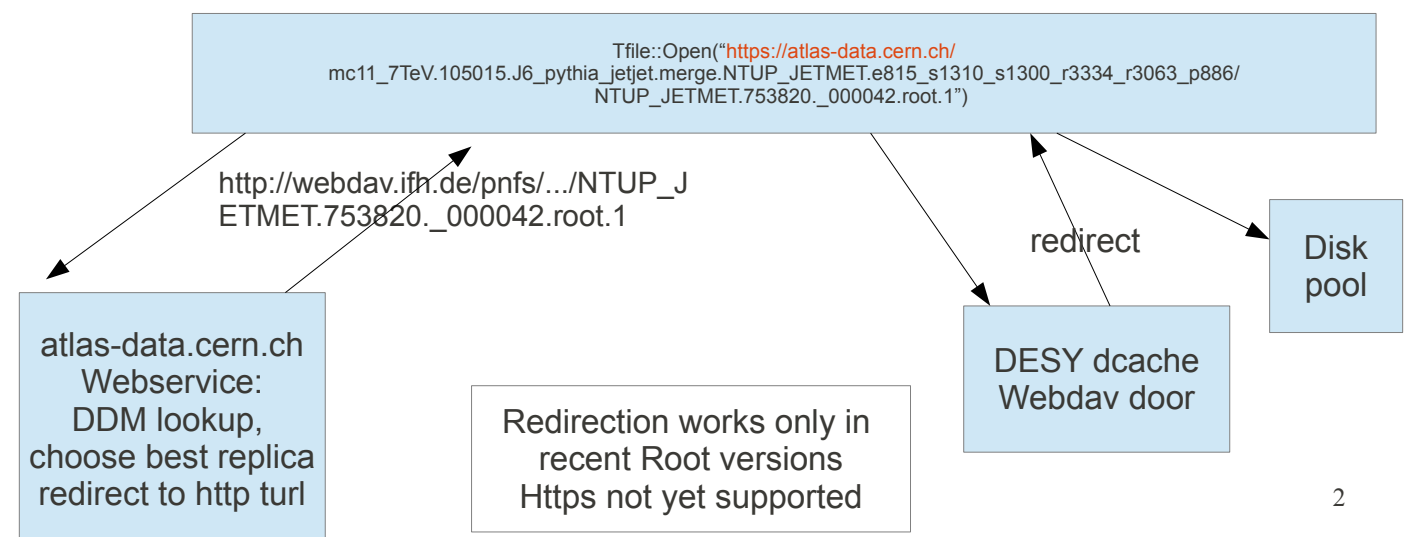
  - generic url including dataset and lfn

  - redirects to http turl in dcache/dpm storage

Tfile::Open("https://atlas-data.cern.ch/
mc11_7TeV.105015.J6_pythia_jetjet.merge.NTUP_JETMET.e815_s1310_s1300_r3334_r3063_p886/
NTUP_JETMET.753820._000042.root.1")

http://webdav.ifh.de/pnfs/.../NTUP_J
ETMET.753820._000042.root.1

redirect

Disk pool

atlas-data.cern.ch
Webservice:
DDM lookup,
choose best replica
redirect to http turl

Redirection works only in recent Root versions
Https not yet supported

DESY dcache
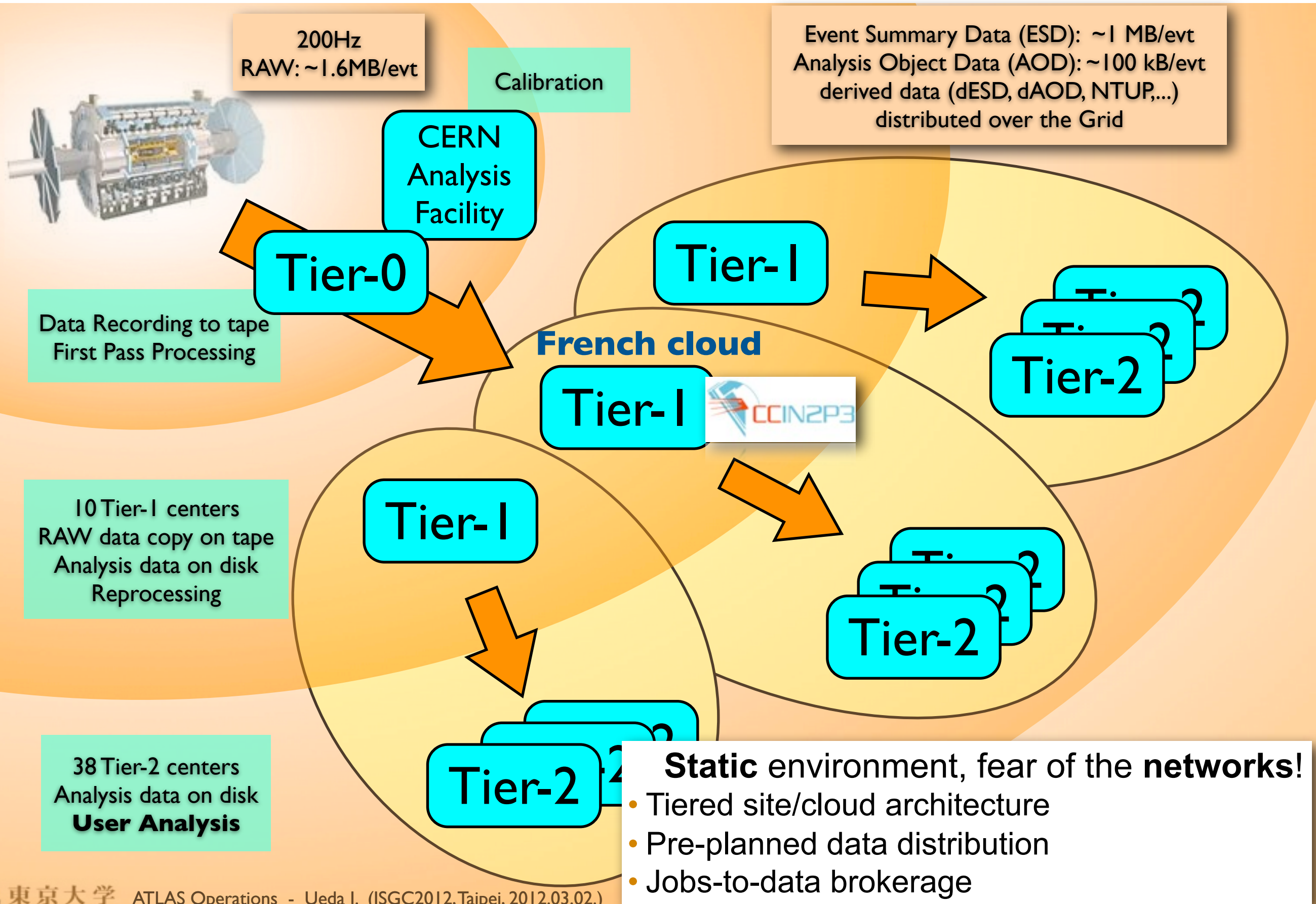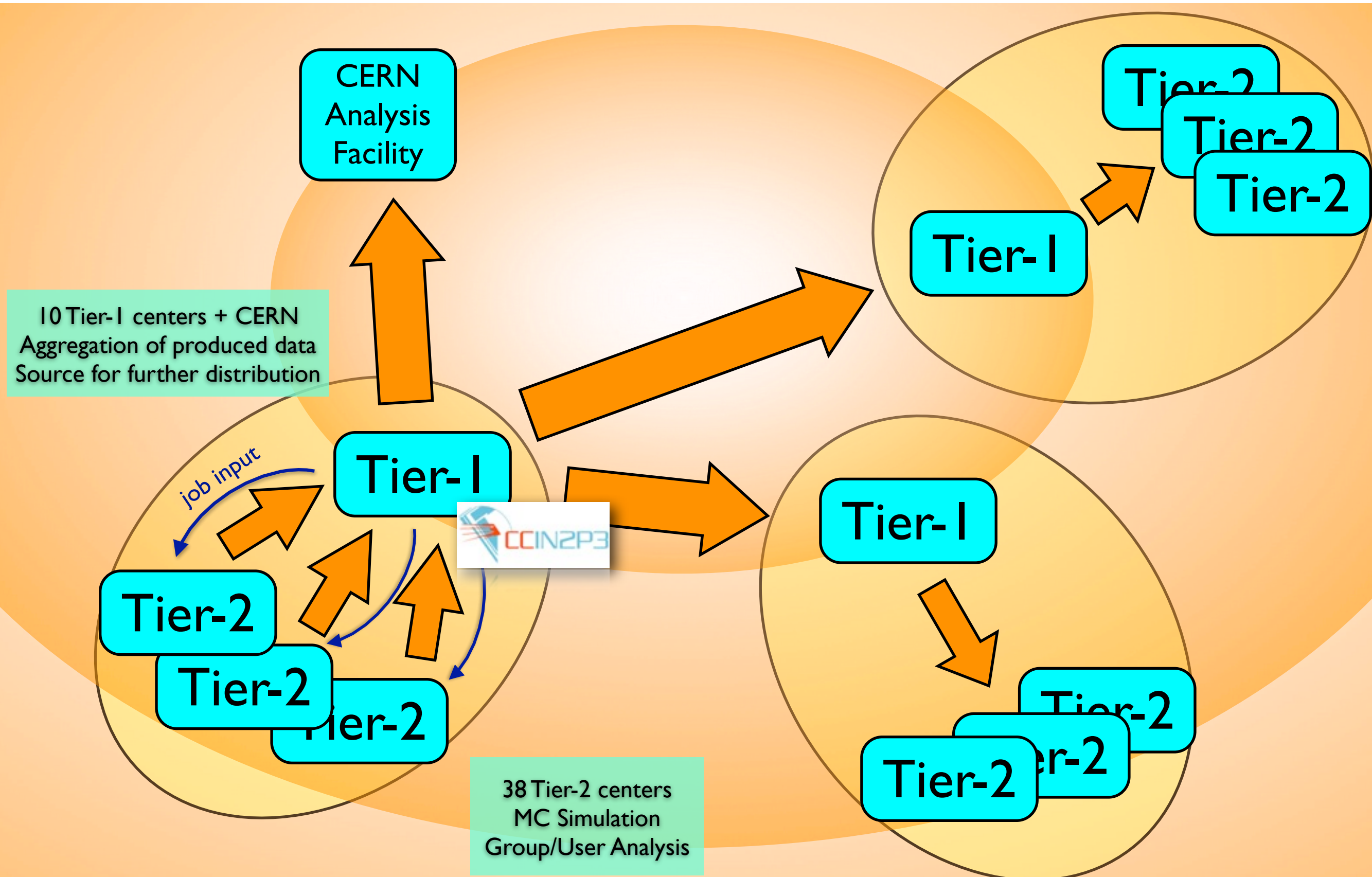Webdav door

2

irfu

cea

saclay

# Summary

- Thanks to the high performances of networks

- ATLAS computing model has changed significantly: simplification of data and workflow management

- Would have been impossible to handle current data volume (LHC performing beyond expectations) and LHC running extension up to spring 2013 with initial model

- More efficient use of storage resources (reduce replica counts; direct sharing of replicas across sites)

- Ongoing projects (distributed storage, remote access) will further change the landscape

irfu

cea

saclay

BACKUP

# ATLAS Computing Model: T0 Data Flow

# ATLAS Computing Model:  MC Data Flow



CERN
Analysis
Facility

Tier-2

Tier-2

Tier-2

Tier-1

10 Tier-1 centers + CERN
Aggregation of produced data
Source for further distribution

job input

Tier-1

CCIN2P3

Tier-1

Tier-2

Tier-2

Tier-2

Tier-2

Tier-2

Tier-2

38 Tier-2 centers
MC Simulation
Group/User Analysis

# Data Processing Model Revised



CERN Analysis Facility

10 Tier-1 centers + CERN
Aggregation of produced data
Source for further distribution

Tier-1

Tier-2

Tier-2

Tier-2-D

Tier-1

Tier-1

Tier-2

Tier-2

Tier-2

Tier-2-D

Tier-2

Tier-2-D

job input

Tier-2-Direct : 'Super' Tier-2s
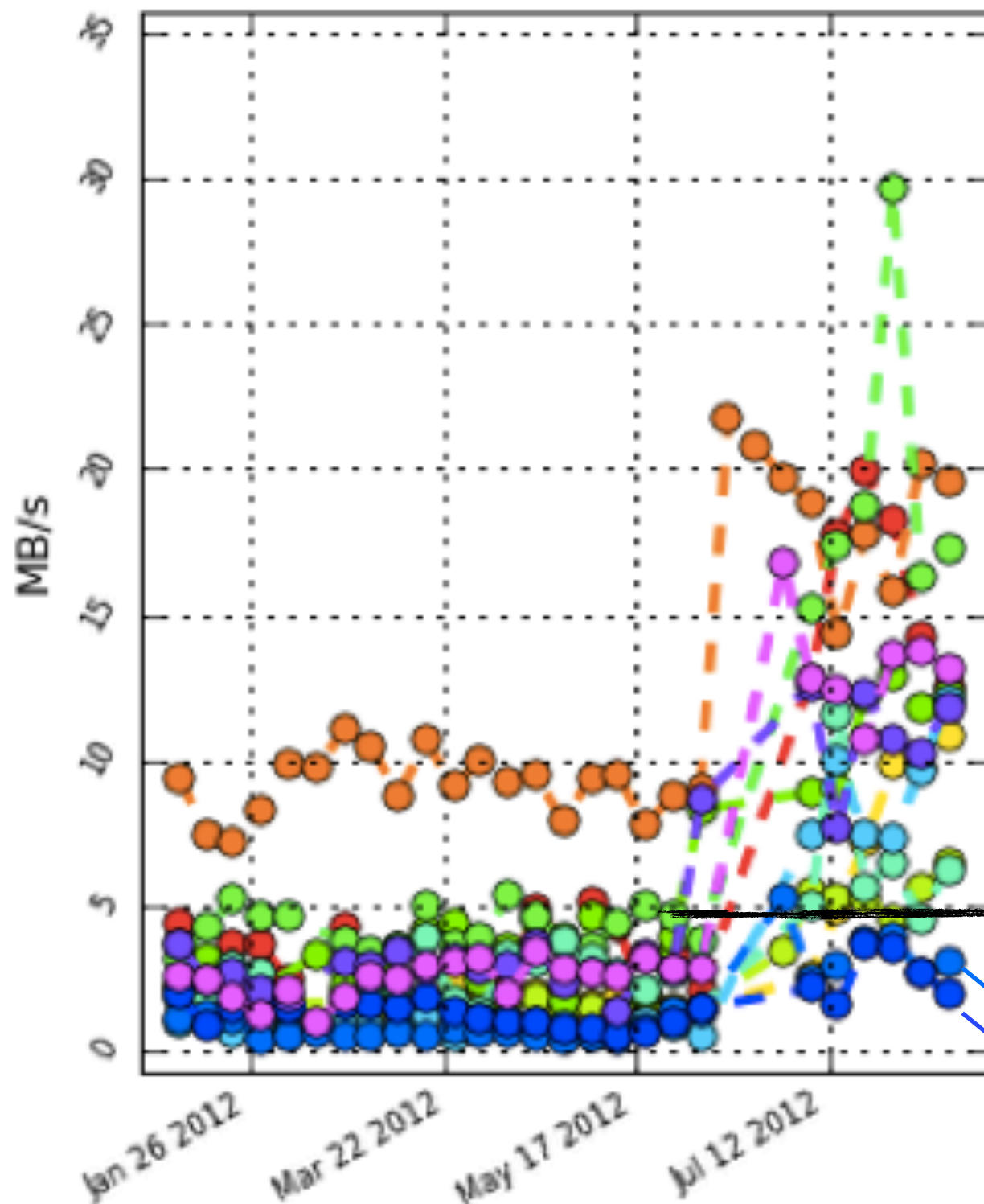
# T1s -> IN2P3-CPPM



June 25th
connected to LHCONE
@ 10 Gb/s

FZK-LCG2 - IN2P3-CPPM  (1360 files)
IN2P3-CC - IN2P3-CPPM  (106227 files)
RAL-LCG2 - IN2P3-CPPM  (2017 files)
BNL-OSG2 - IN2P3-CPPM  (4383 files)
SARA-MATRIX - IN2P3-CPPM  (2031 files)
INFN-T1 - IN2P3-CPPM  (1898 files)
PIC - IN2P3-CPPM  (541 files)
NDGF-T1 - IN2P3-CPPM  (1571 files)
TAIWAN-LCG2 - IN2P3-CPPM  (275 files)
TRIUMF-LCG2 - IN2P3-CPPM  (186 files)
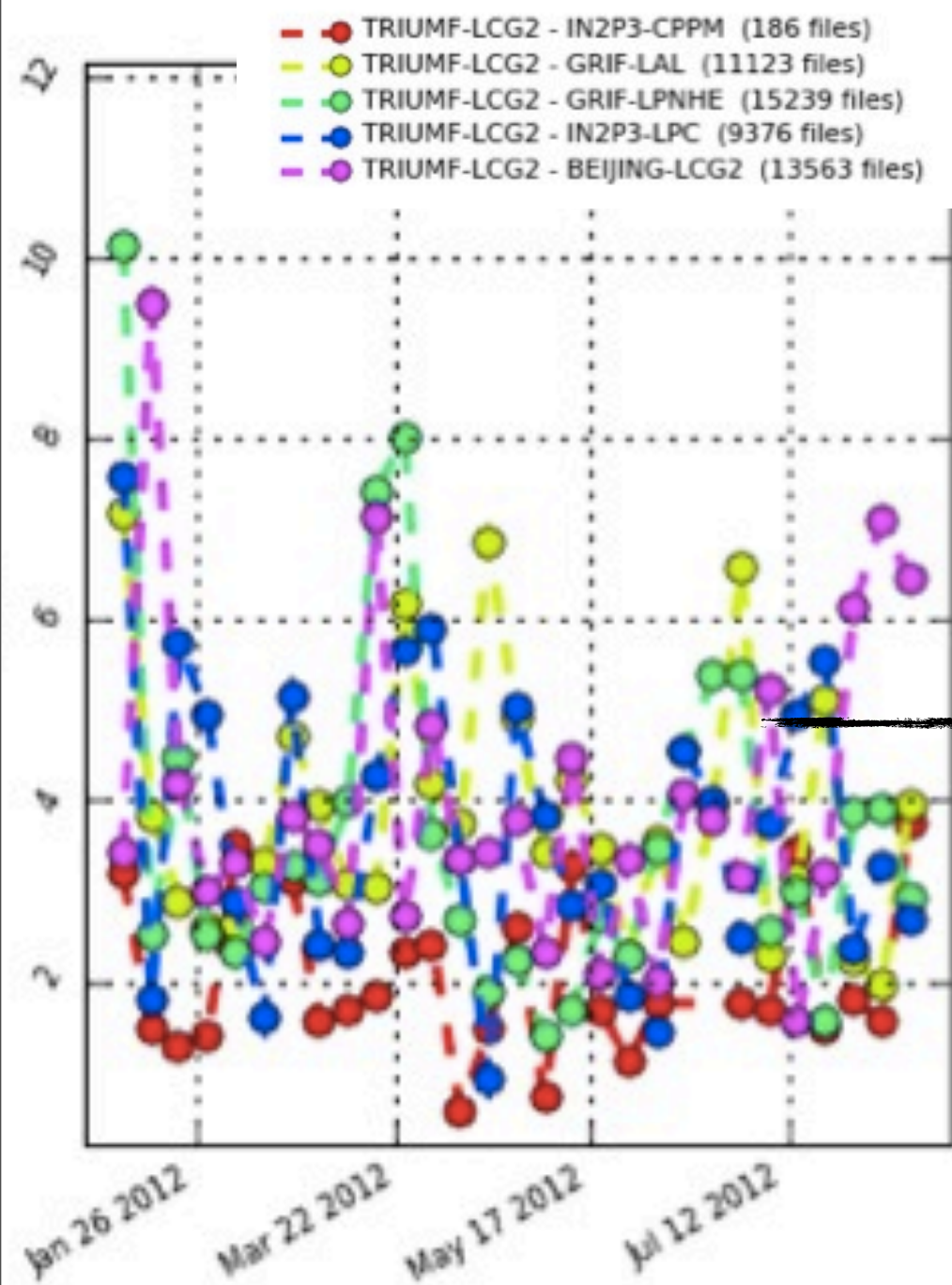NIKHEF-ELPROD - IN2P3-CPPM  (292 files)
CERN-PROD - IN2P3-CPPM  (3122 files)

5 MB/s

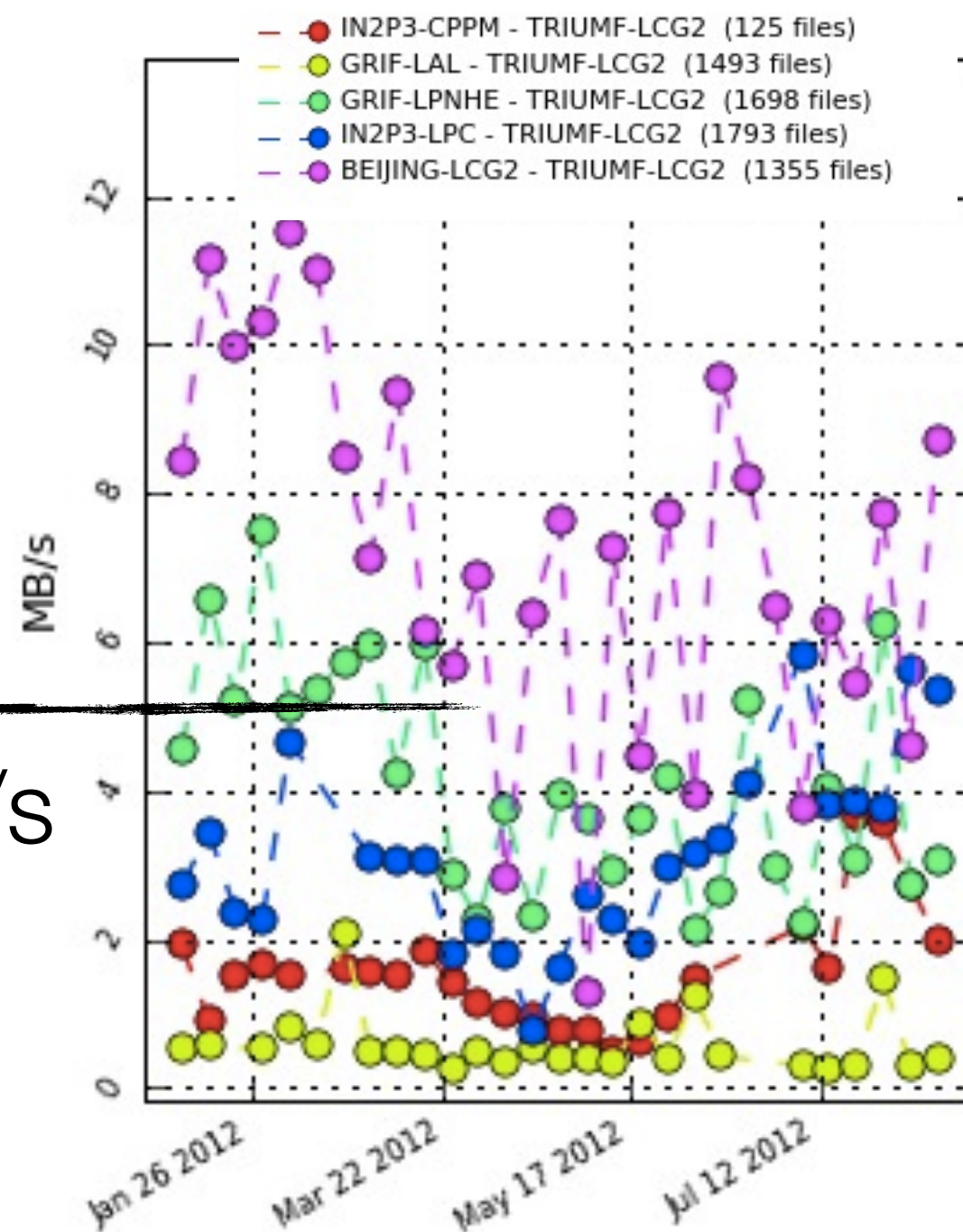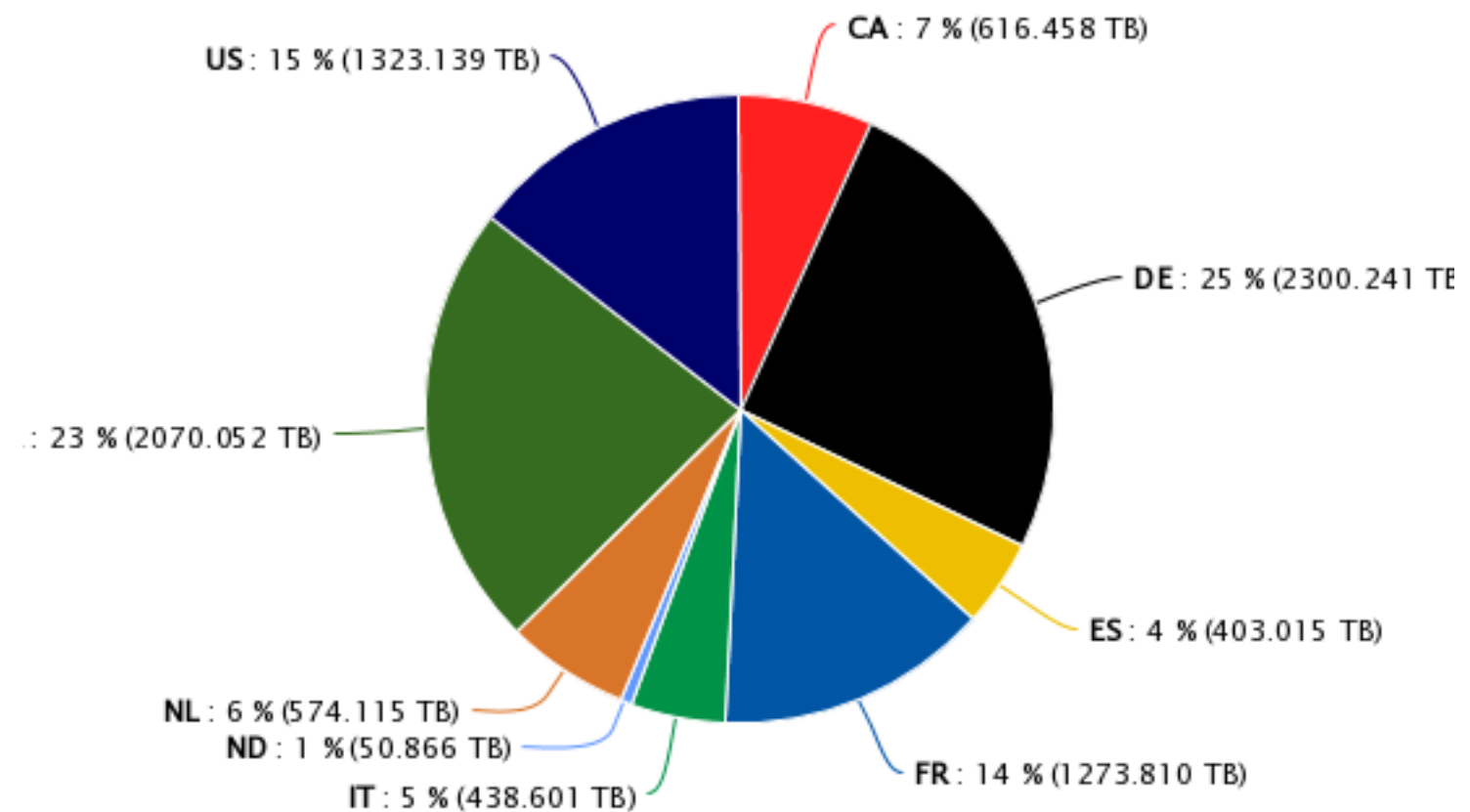TRIUMF

irfu

cea

saclay

46

# IN2P3-CPPM -> T1s



5 MB/s

TW
TRIUMF

# TRIUMF <-> FR T2Ds

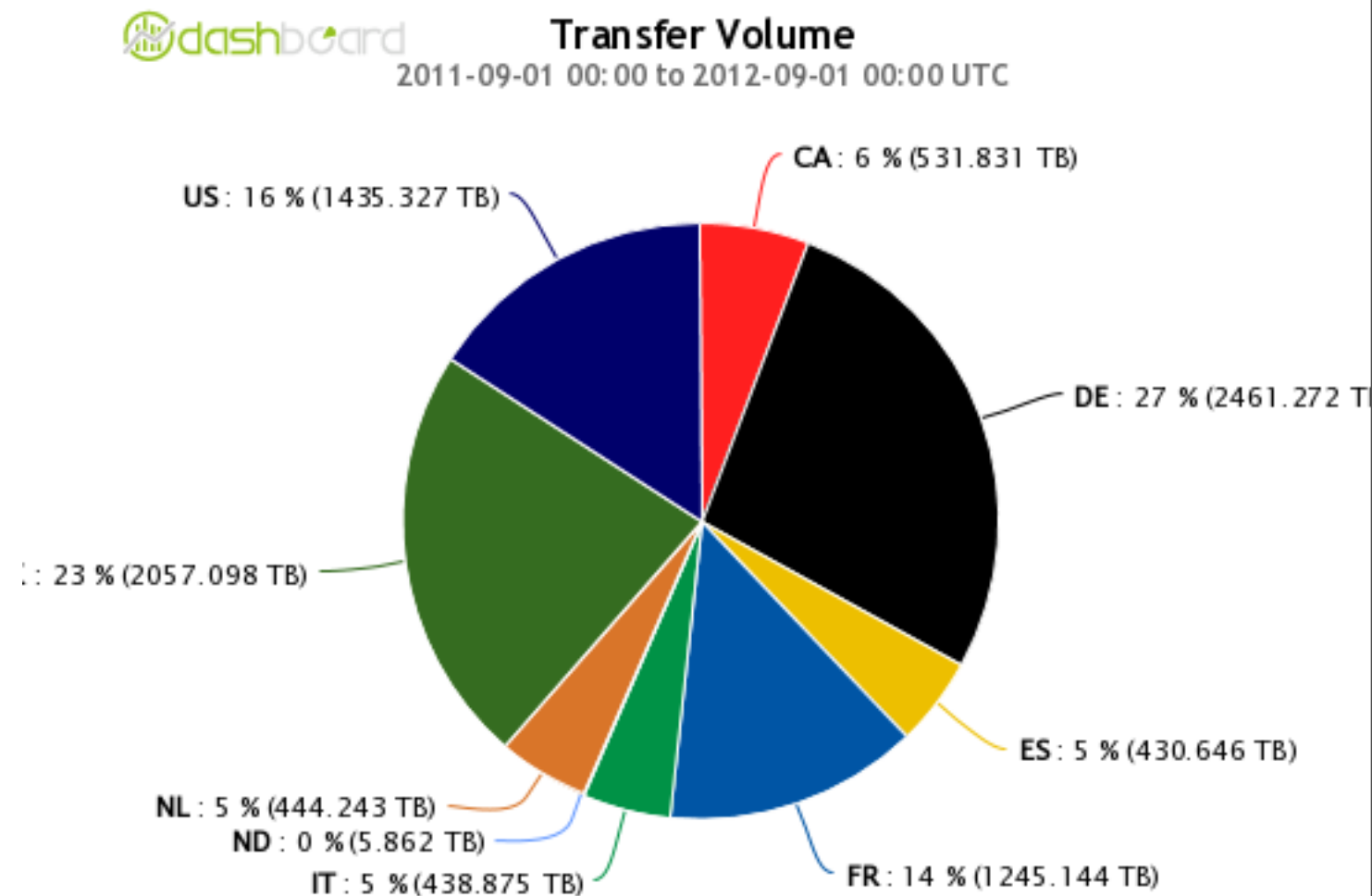**~None of T2Ds ever reaches the 5 MB/s canonical parameter value**



5 MB/s

irfu

cea

saclay

T2-T2 destination

Transfer Volume
2011-09-01 00:00 to 2012-09-01 00:00 UTC

CA : 7 % (616.458 TB)
US : 15 % (1323.139 TB)
DE : 25 % (2300.241 TB)
: 23 % (2070.052 TB)
ES : 4 % (403.015 TB)
NL : 6 % (574.115 TB)
ND : 1 % (50.866 TB)
IT : 5 % (438.601 TB)
FR : 14 % (1273.810 TB)

T2-T2 source

Transfer Volume
2011-09-01 00:00 to 2012-09-01 00:00 UTC

CA : 6 % (531.831 TB)
US : 16 % (1435.327 TB)
DE : 27 % (2461.272 TB)
: 23 % (2057.098 TB)
ES : 5 % (430.646 TB)
NL : 5 % (444.243 TB)
ND : 0 % (5.862 TB)
IT : 5 % (438.875 TB)
FR : 14 % (1245.144 TB)

irfu

cea

saclay

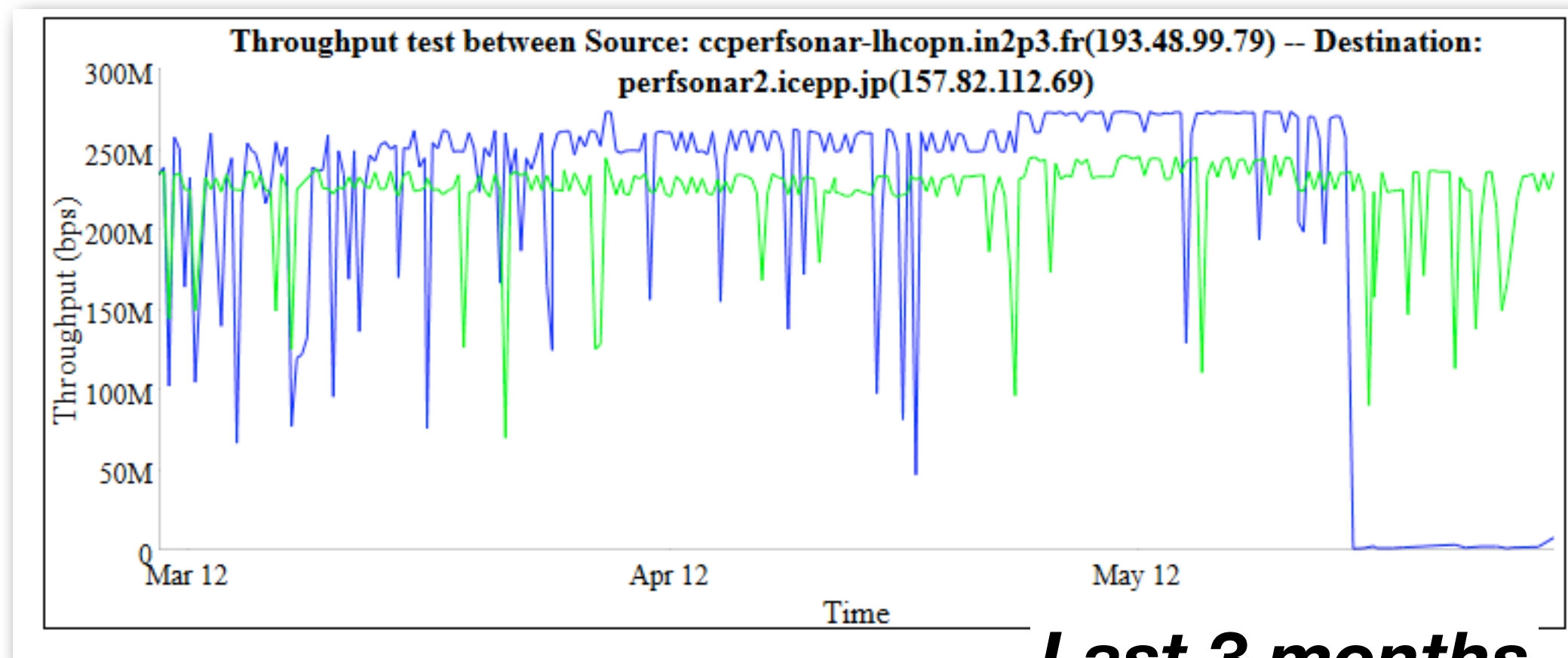Rencontre LCG-France, SUBATECH Nantes, septembre 2012

# Network throughput measured with perfSONAR

**CCIN2P3 → Tokyo**     **Tokyo → CCIN2P3**



**5%**

***Last 3 months***
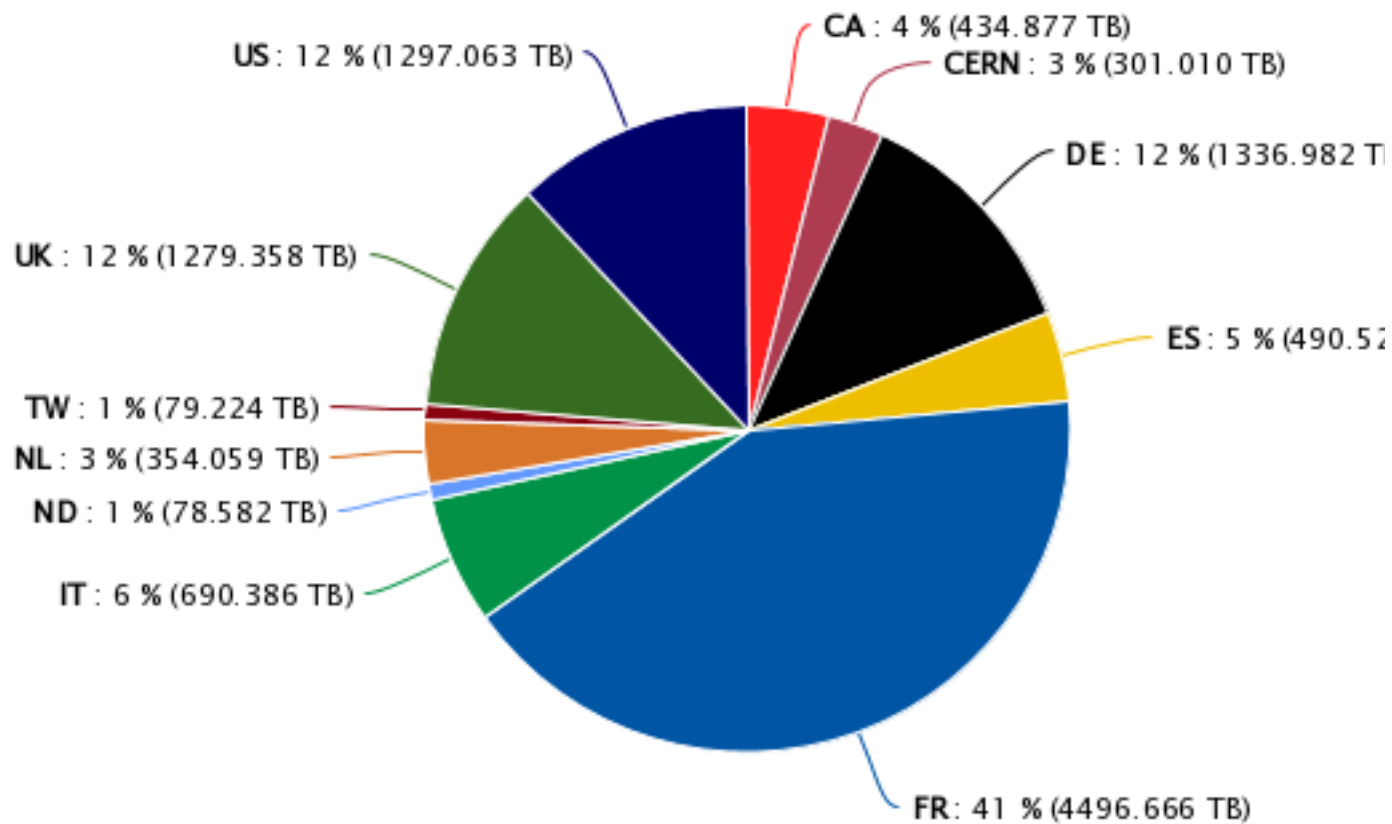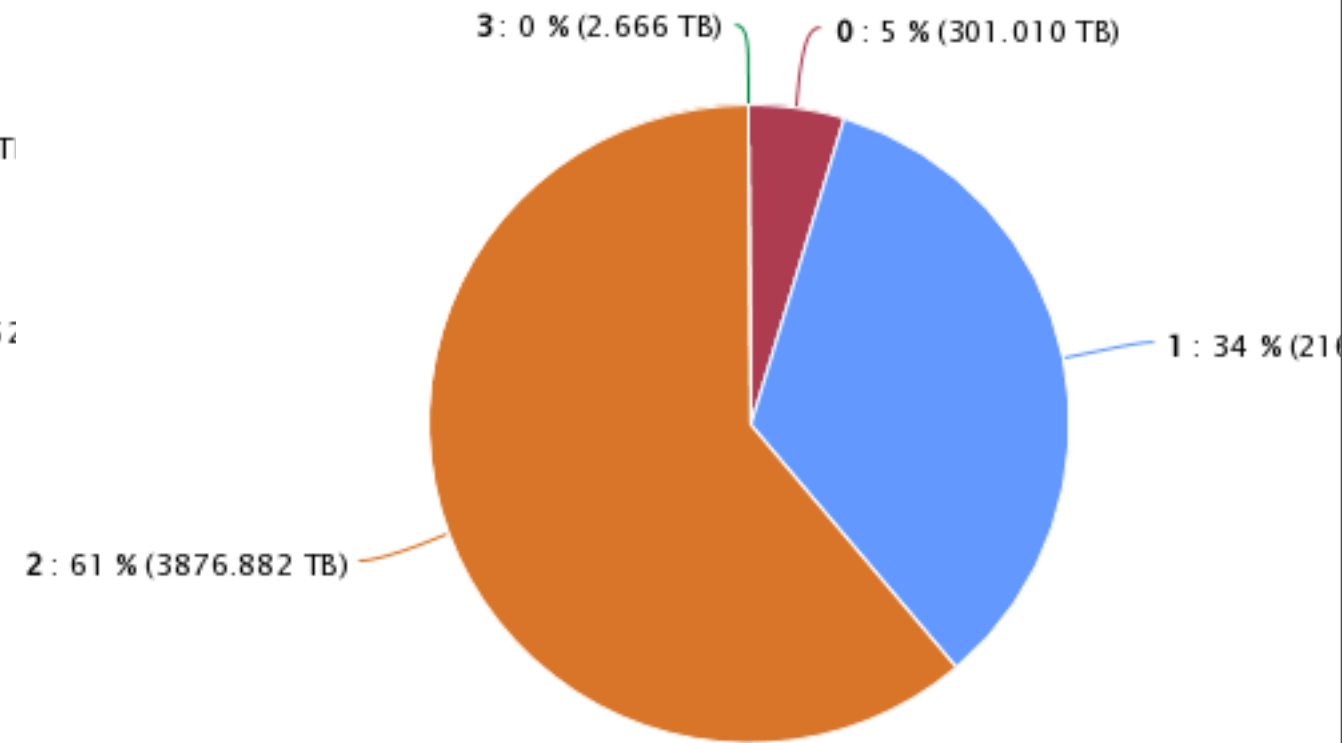
*No so stable*
*better by ~5% for CCIN2P3 →Tokyo*

irfu

cea

saclay

http://psps.perfsonar.net/

# CCIN2P3 Exports



**Transfer Volume**
2011-09-01 00:00 to 2012-09-01 00:00 UTC

- CA : 4 % (434.877 TB)
- CERN : 3 % (301.010 TB)
- US : 12 % (1297.063 TB)
- DE : 12 % (1336.982 T
- UK : 12 % (1279.358 TB)
- ES : 5 % (490.52
- TW : 1 % (79.224 TB)
- NL : 3 % (354.059 TB)
- ND : 1 % (78.582 TB)
- IT : 6 % (690.386 TB)
- FR : 41 % (4496.666 TB)



**Transfer Volume**
2011-09-01 00:00 to 2012-09-01 00:00 UTC

- 3 : 0 % (2.666 TB)
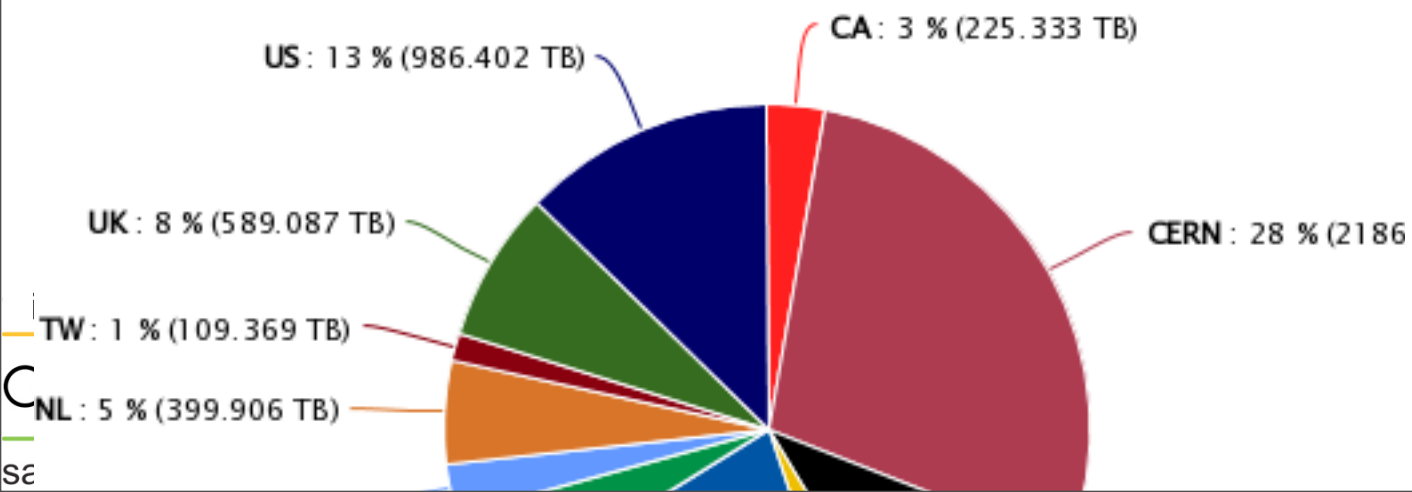- 0 : 5 % (301.010 TB)
- 1 : 34 % (216
- 2 : 61 % (3876.882 TB)

## To Lyon



**Transfer Volume**
2011-09-01 00:00 to 2012-09-01 00:00 UTC

- CA : 3 % (225.333 TB)
- US : 13 % (986.402 TB)
- CERN : 28 % (2186
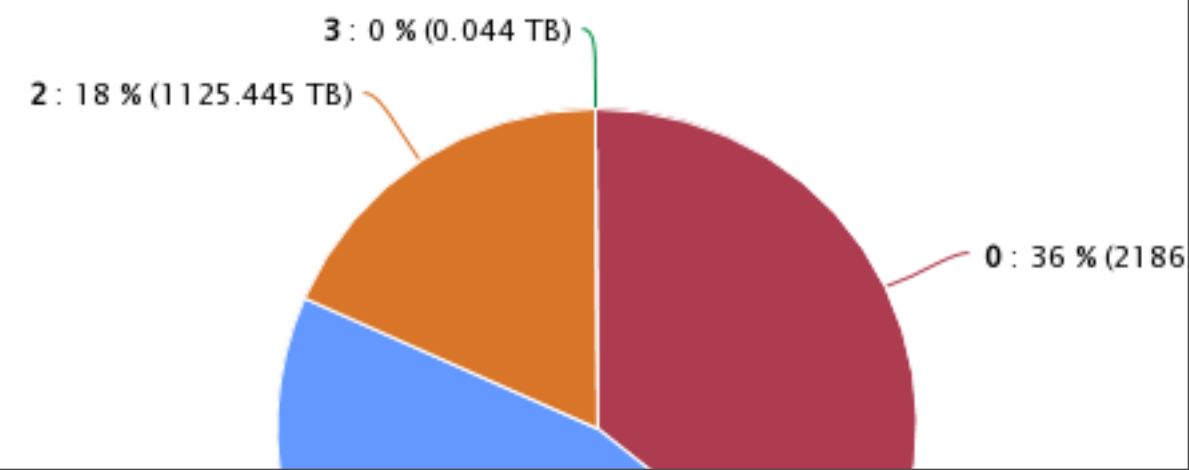- UK : 8 % (589.087 TB)
- TW : 1 % (109.369 TB)
- NL : 5 % (399.906 TB)



**Transfer Volume**
2011-09-01 00:00 to 2012-09-01 00:00 UTC

- 3 : 0 % (0.044 TB)
- 2 : 18 % (1125.445 TB)
- 0 : 36 % (2186

# T2s Exports



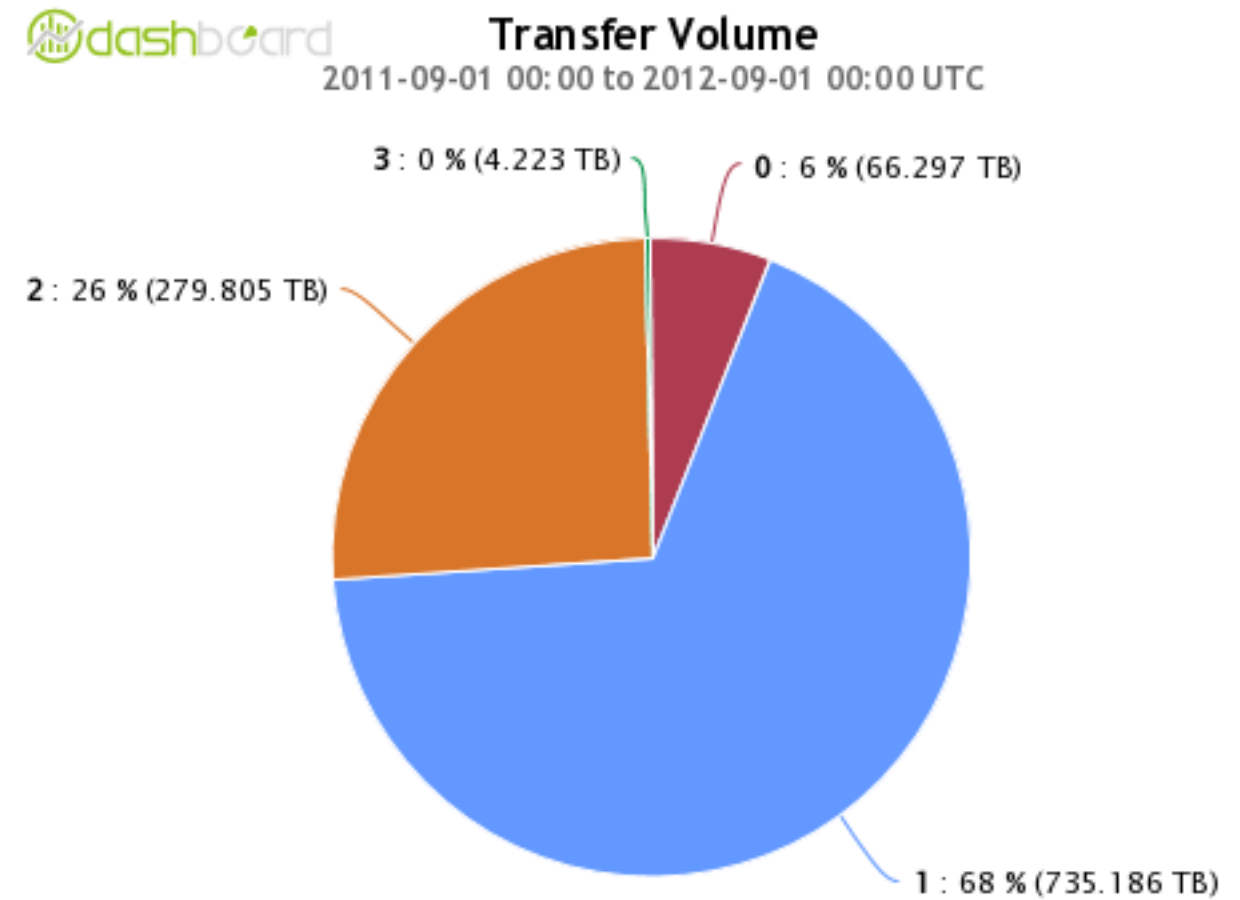**Transfer Volume**
2011-09-01 00:00 to 2012-09-01 00:00 UTC

US : 5 % (153.022 TB)
UK : 5 % (156.349 TB)
TW : 4 % (149.752 TB)
: 3 % (108.119 TB)
2 % (51.172 TB)
% (105.908 TB)
CA : 1 % (49.547 TB)
CERN : 2 % (66.297 TB)
DE : 5 % (155.196 TB)
ES : 3 % (90.148 TB)
FR : 68 % (2278.107 TB)



**Transfer Volume**
2011-09-01 00:00 to 2012-09-01 00:00 UTC

3 : 0 % (4.223 TB)
0 : 6 % (66.297 TB)
2 : 26 % (279.805 TB)
1 : 68 % (735.186 TB)

# T2s Import



**Transfer Volume**
2011-09-01 00:00 to 2012-09-01 00:00 UTC

US : 9 % (865.283 TB)
UK : 6 % (560.168 TB)
TW : 3 % (352.190 TB)
NL : 5 % (525.417 TB)
ND : 4 % (371.988 TB)
CA : 3 % (295.440 TB)
CERN : 6 % (562.607 TB)
DE : 6 % (597.169 TB)
ES : 3 % (296.676 TB)



**Transfer Volume**
2011-09-01 00:00 to 2012-09-01 00:00 UTC

3 : 0 % (0.613 TB)
2 : 6 % (308.471 TB)
0 : 11 % (562.607 TB)

# destination on FR cloud



**Transfer Volume**
2011-09-01 00:00 to 2012-09-01 00:00 UTC

- TOKYO-LCG2 : 10 % (1047.597 TB)
- SDU-LCG2 : 0 % (0.000 TB)
- RO-16-UAIC : 0 % (31.531 TB)
- RO-14-ITIM : 0 % (25.730 TB)
- RO-07-NIPNE : 4 % (404.366 TB)
- RO-02-NIPNE : 1 % (129.739 TB)
- IN2P3-LPSC : 8 % (798.031 TB)
- IN2P3-LPC : 11 % (1149.589 TB)
- IN2P3-LAPP : 12 % (1238.394 TB)
- IN2P3-CPPM : 7 % (659.162 TB)
- GRIF-LPNHE : 14 % (1374.564 TB)
- GRIF-LAL : 10 % (995.591 TB)
- GRIF-IRFU : 14 % (1384.526 TB)
- BEIJING-LCG2 : 8 % (837.680 TB)

saclay

# GEANT/TEIN3

GÉANT and sister networks enabling user collaboration across the globe

April 2011

Tokyo is far from CCIN2P3 : ~300 ms RTT (Round Trip Time)
Throughput ~ 1 / RTT

Data are transferred from site to site through a lot of
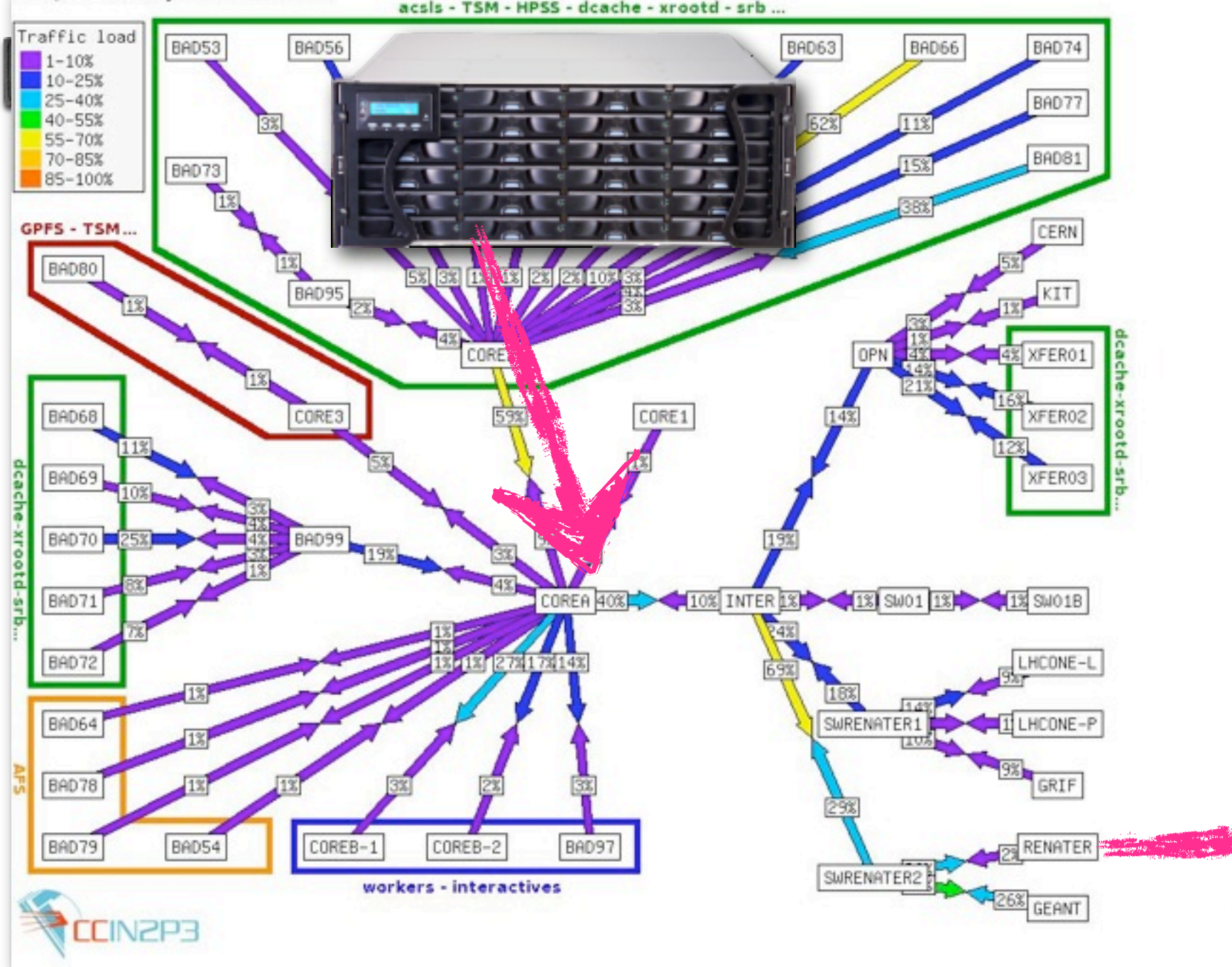networks (multi-hop) and software layers

ICEPP 素粒子物理国際研究センター
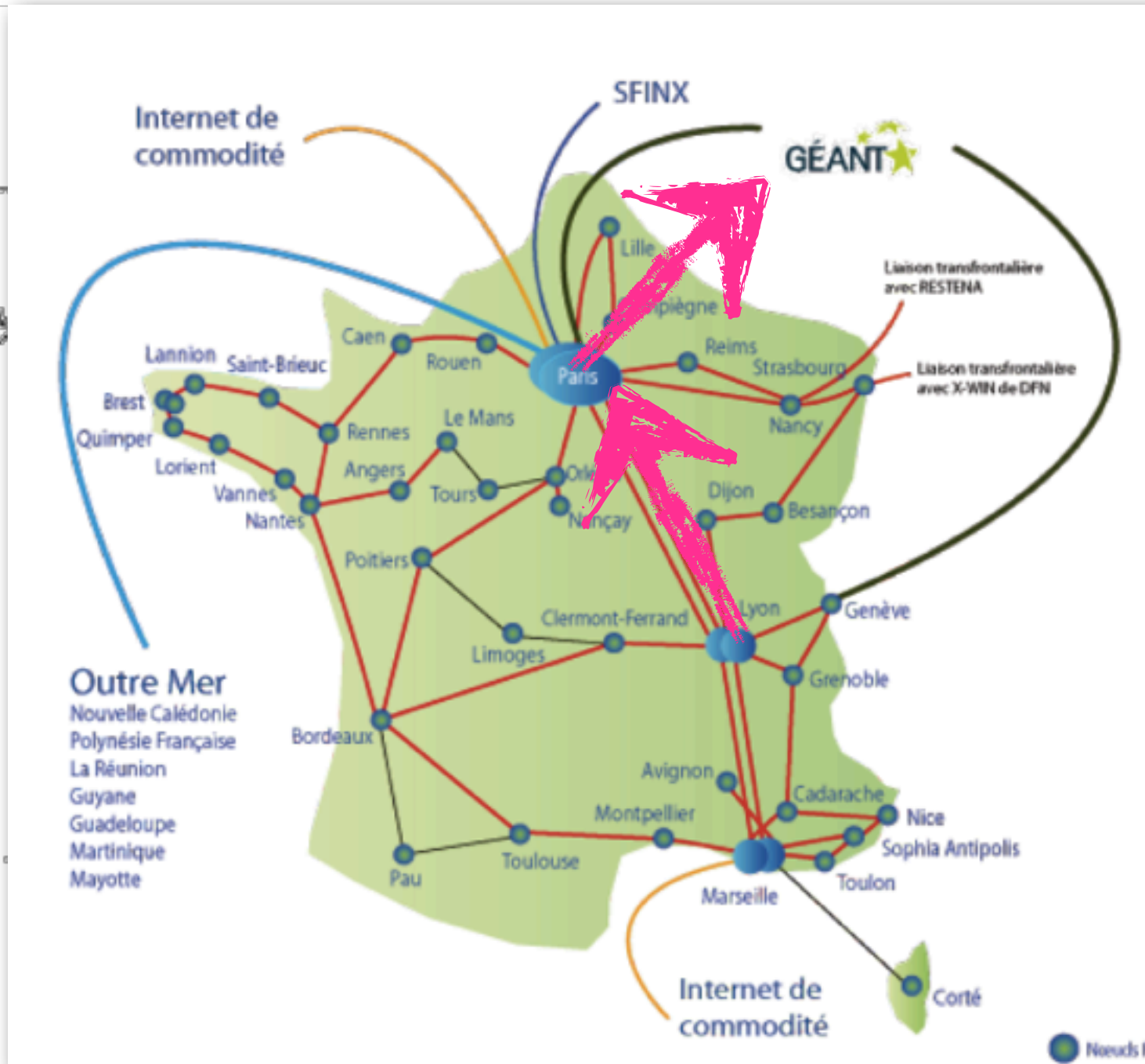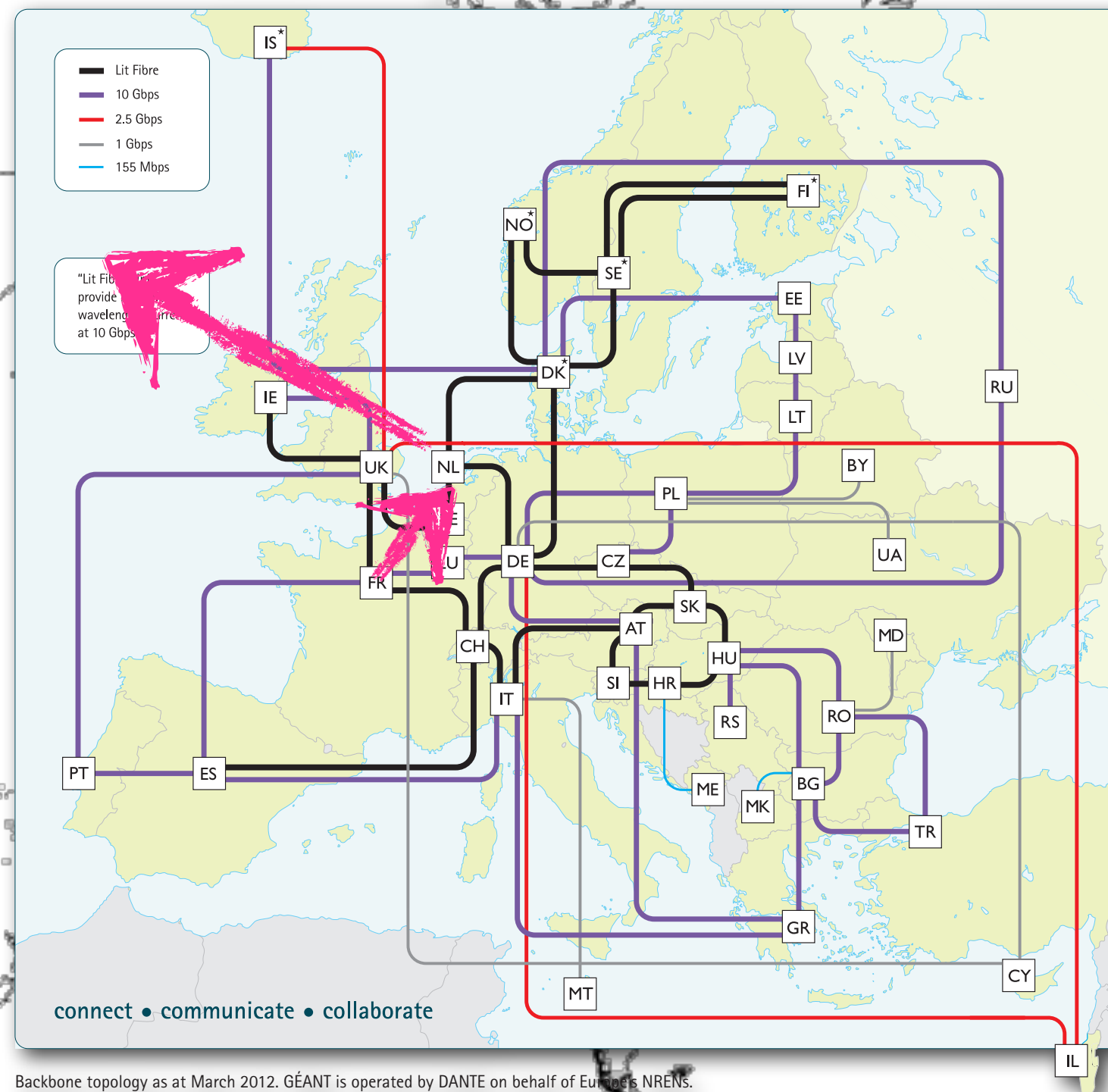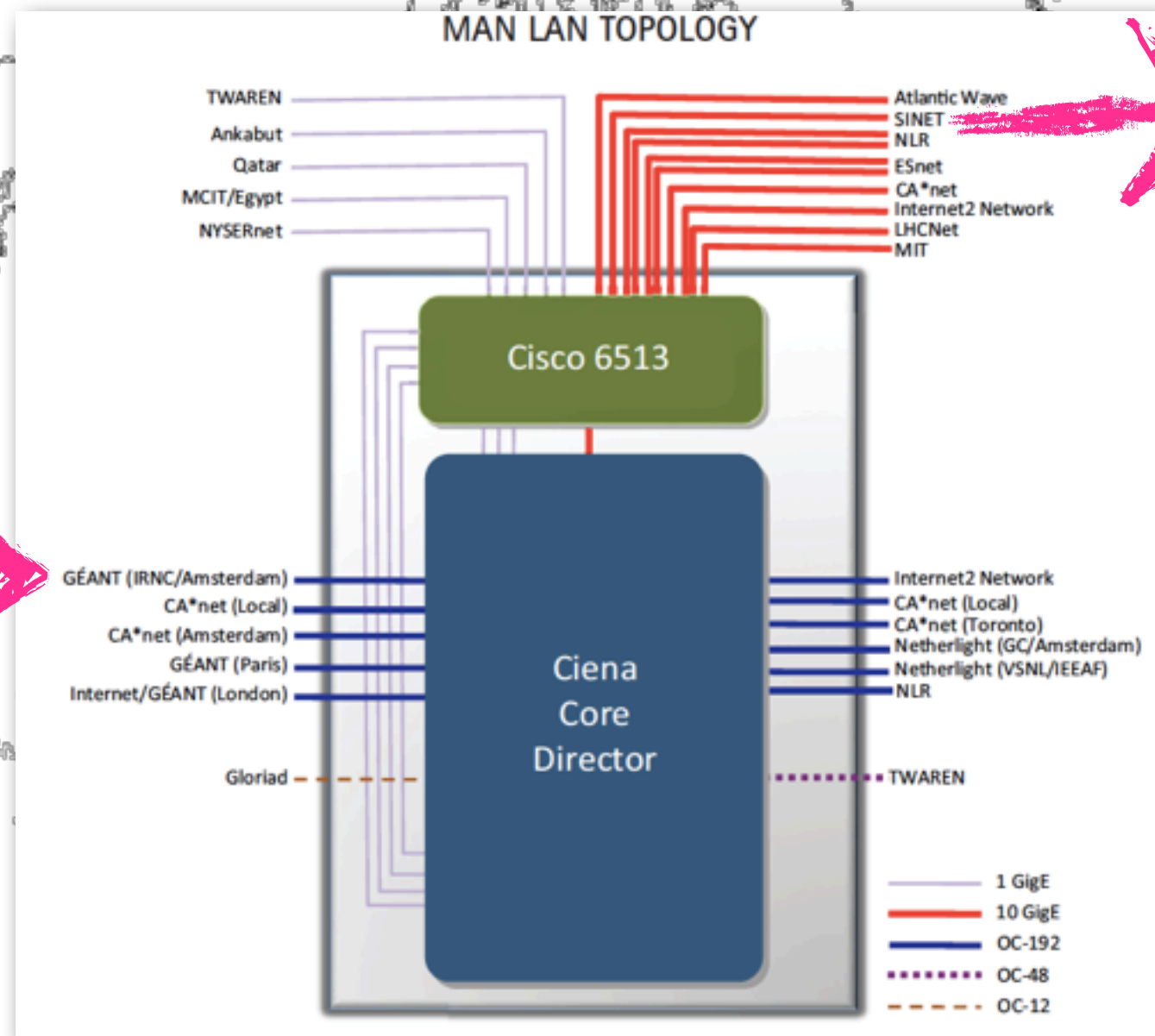International Center for Elementary Particle Physics
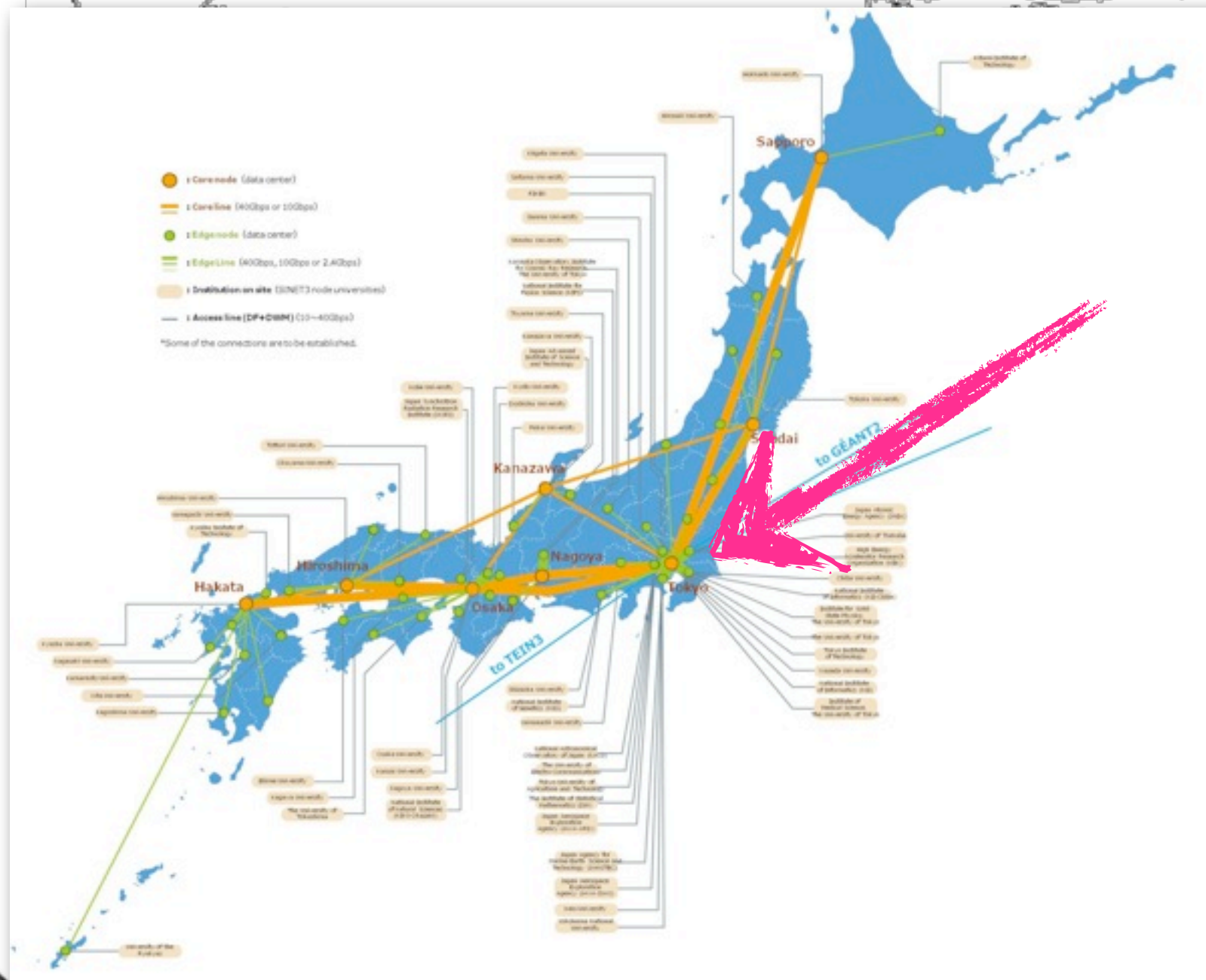
CCIN2P3

*ideal view*

*The reality* ➤

irfu

cea

saclay

irfu

cea

saclay

irfu

cea

saclay

ICEPP
素粒子物理国際研究センター
International Center for Elementary Particle Physics

TOKYO-LCG2

SINET — 10Gbps — Network Switch — 10GE — File Server — 8G-FC — Disk Array (File System, File System) / Disk Array (File System, File System) — 8G-FC

x15 to x17

irfu
cea
saclay

# T1s -> IN2P3-LPSC

*Heavy transfers from IN2P3-CC (T1) interference with FTS monitoring*



IN2P3-CC -> IN2P3-LPSC

**Throughput**
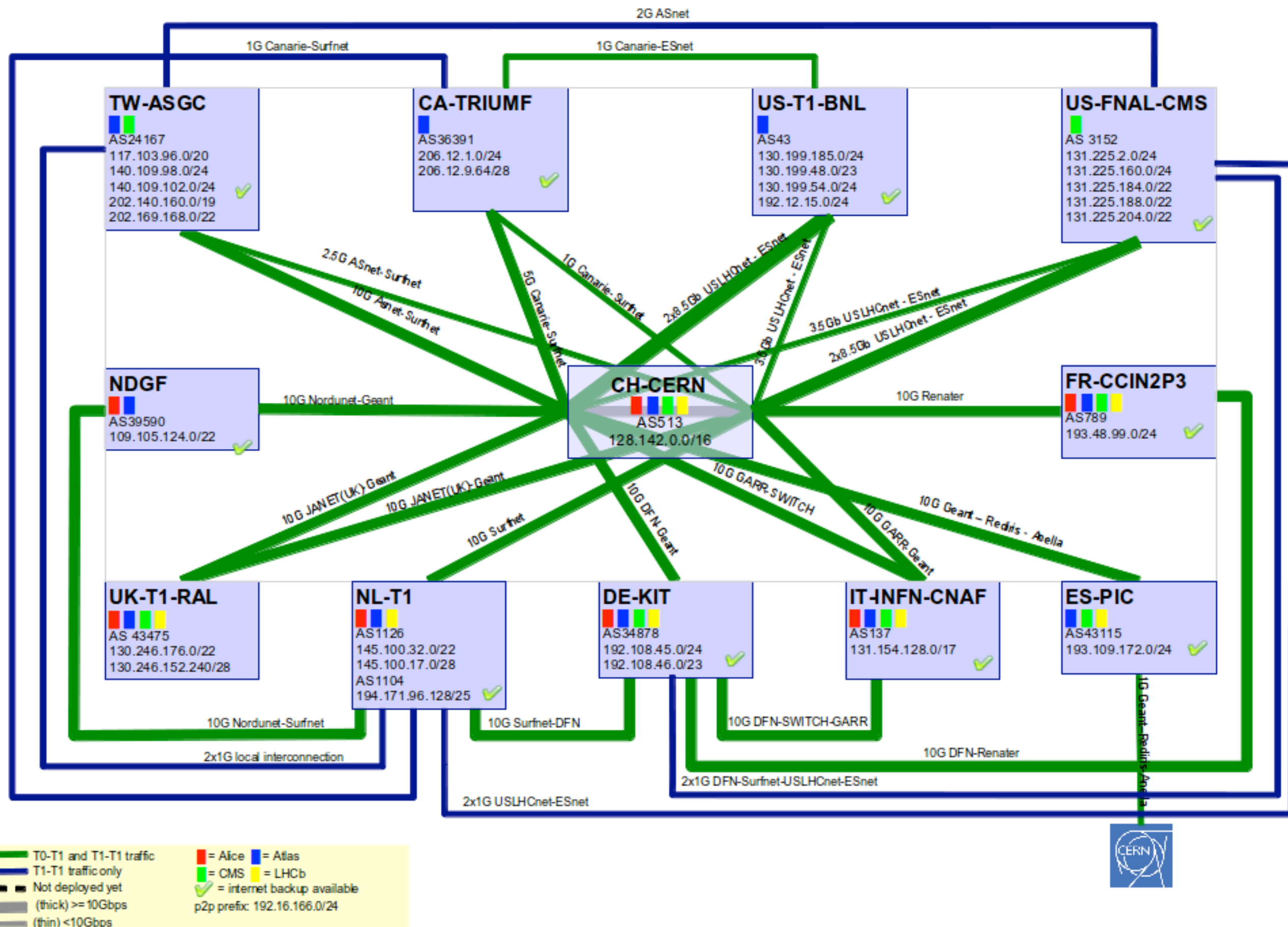2012-07-01 00:00 to 2012-08-28 00:00 UTC

Legend (left chart):
- CERN-PROD_DATADISK - IN2P3-LPSC (2060 files)
- BNL-OSG2_DATADISK - IN2P3-LPSC (4210 files)
- TRIUMF-LCG2_DATADISK - IN2P3-LPSC (2626 files)
- TAIWAN-LCG2_DATADISK - IN2P3-LPSC (919 files)
- SARA-MATRIX_DATADISK - IN2P3-LPSC (1989 files)
- NIKHEF-ELPROD_DATADISK - IN2P3-LPSC (706 files)
- FZK-LCG2_DATADISK - IN2P3-LPSC (2169 files)
- RAL-LCG2_DATADISK - IN2P3-LPSC (2255 files)
- IN2P3-CC_DATADISK - IN2P3-LPSC (21340 files)
- PIC_DATADISK - IN2P3-LPSC (1525 files)
- INFN-T1_DATADISK - IN2P3-LPSC (1588 files)
- NDGF-T1_DATADISK - IN2P3-LPSC (1786 files)

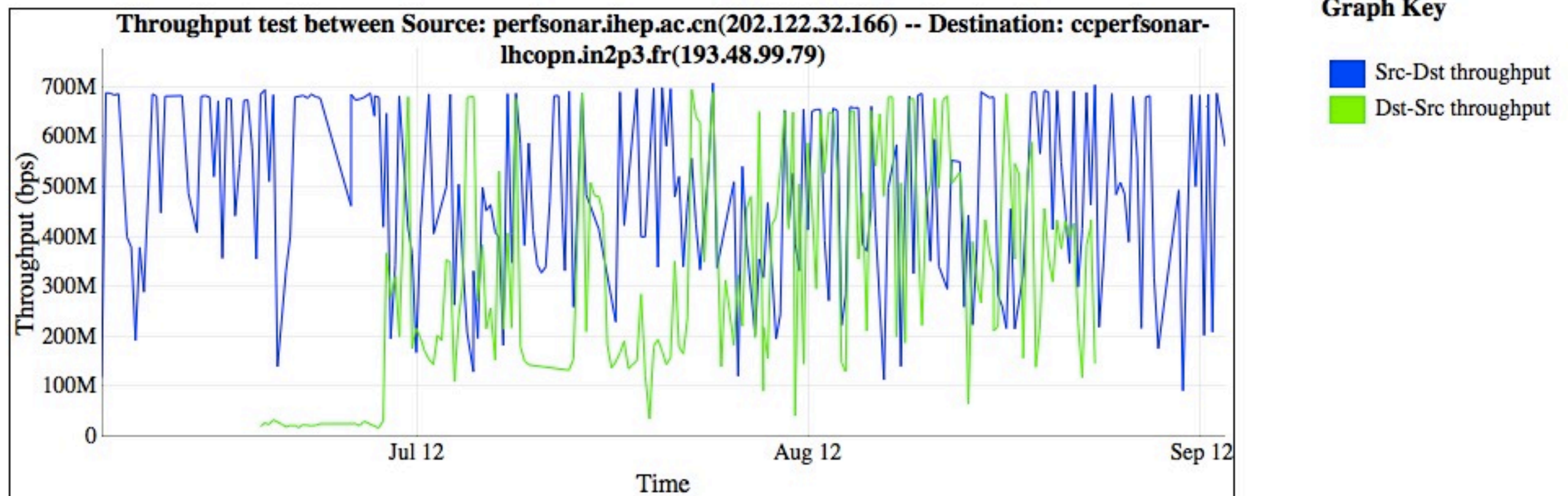Legend (right chart): Data Brokering, Data Consolidation, Functional Test, Production, User Subscriptions

irfu
cea
saclay

62

# IN2P3-CC ↔ Beijing as seen by perfSonar

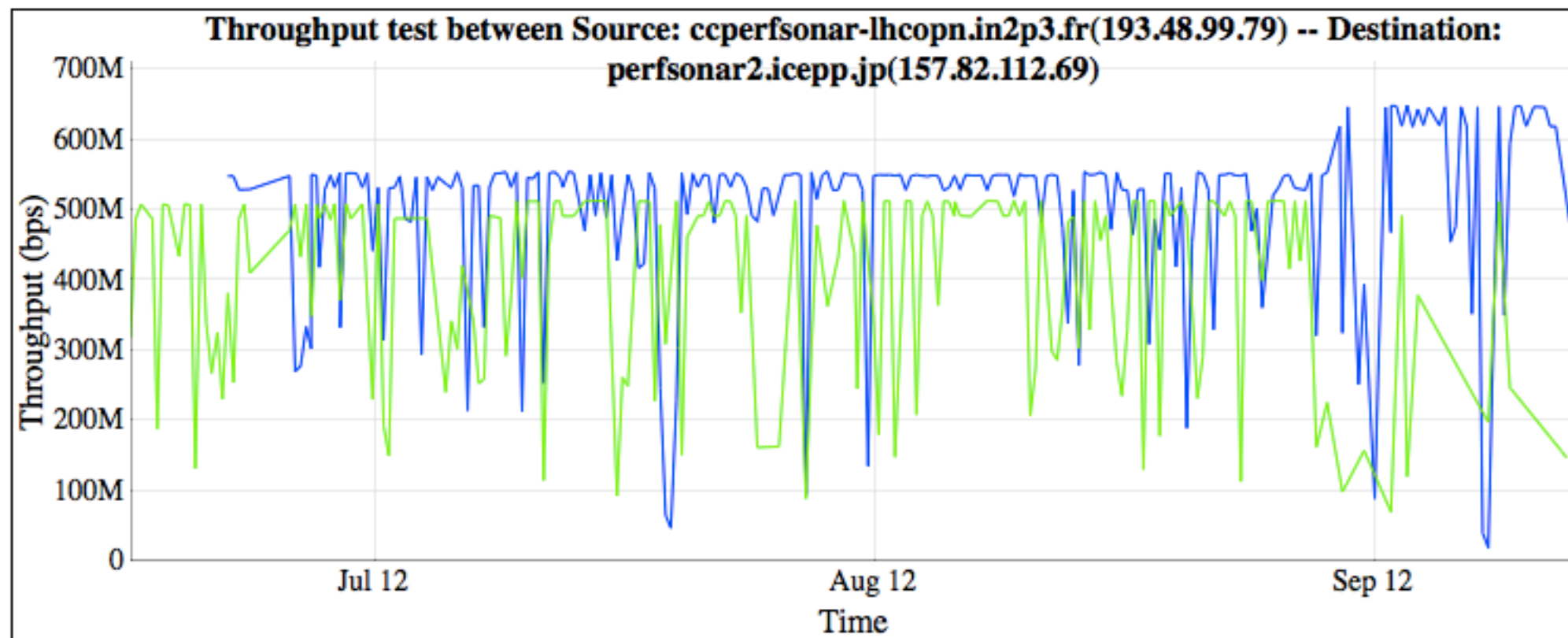Beijing -> IN2P3-CC    IN2P3-CC -> Beijing



- Link unstable
- Asymmetry

irfu

saclay

# IN2P3-CC ↔ Tokyo as seen by perfSonar

**IN2P3-CC-> Tokyo**    **Tokyo -> IN2P3-CC**