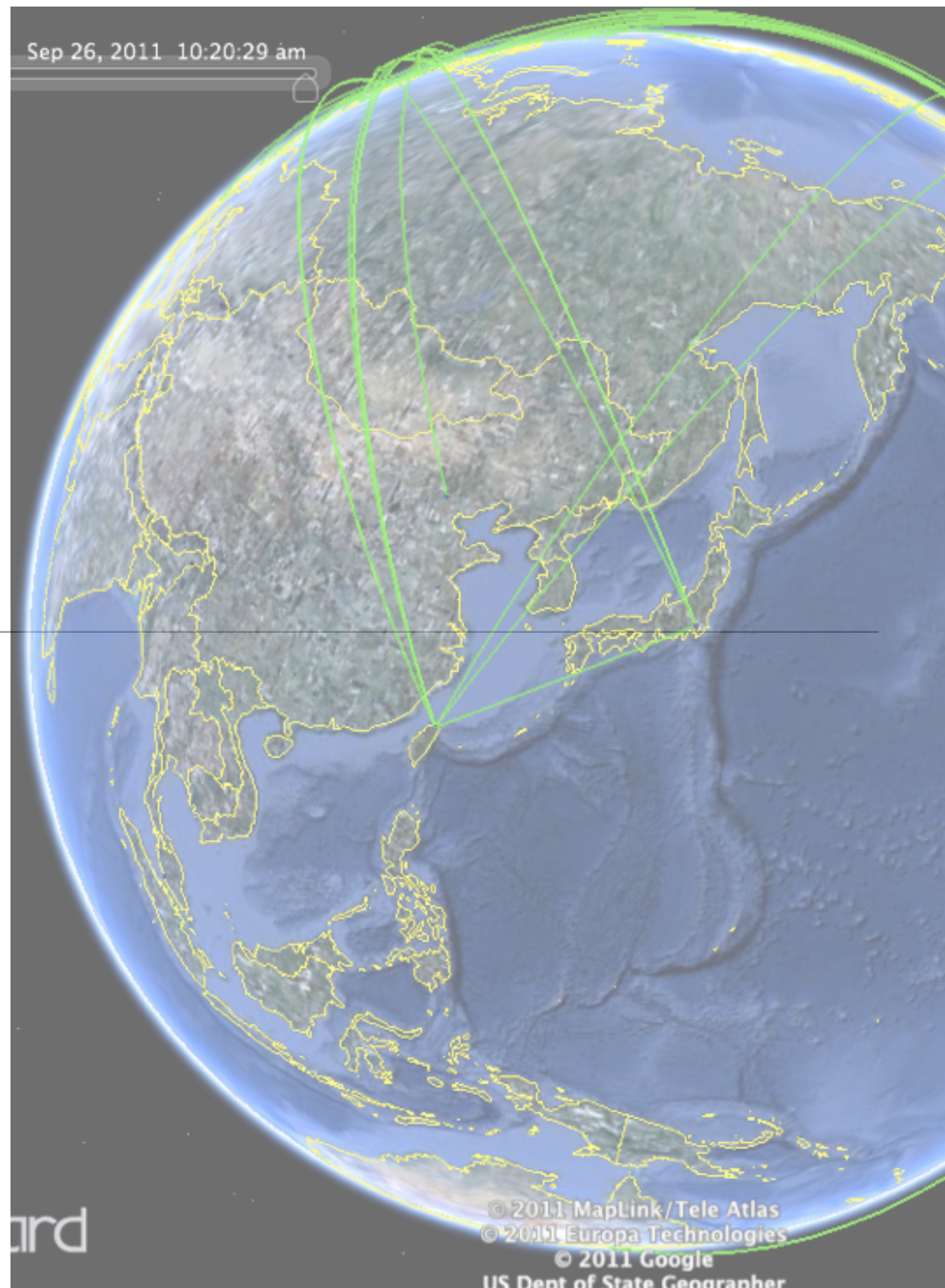


ATLAS activities on FR cloud

欢迎

歓迎

bun venit
bienvenue
Welcome



Some items about FR-cloud

- Who we are
- What we did over the last year
- Operation issues
 - Performances
 - Common procedures / communication (Squad/Luc)
 - Network (Monday)

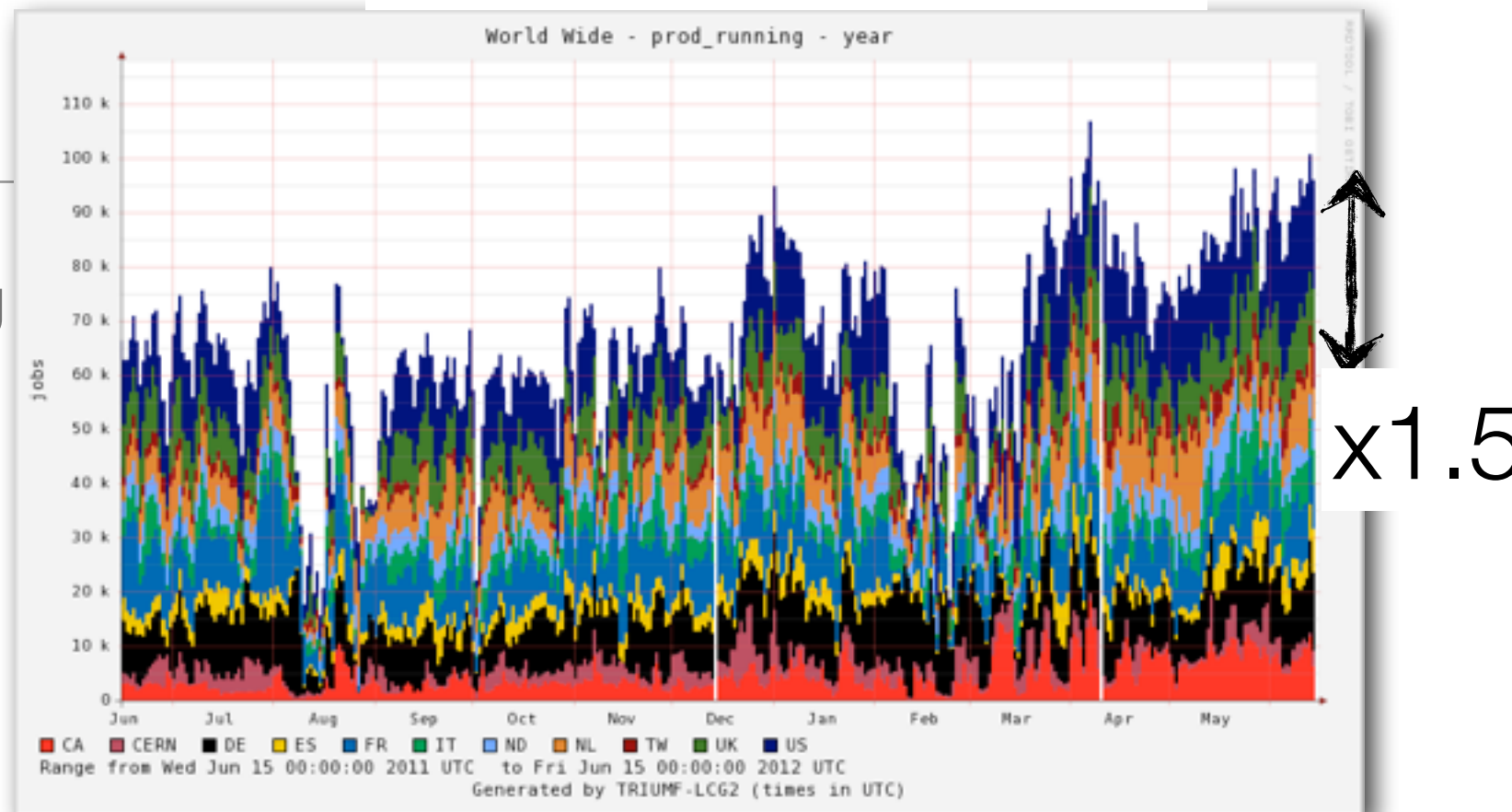
Items for discussions (random order) during the workshop

- ATLAS related site issues
- Network monitoring and tests
- T2Ds
- Monitoring
- CVMFS
- Squid/FronTier
- multi-core
- Xrootd federation
- cloud computing

ATLAS activities [May 2011 - May 2012]

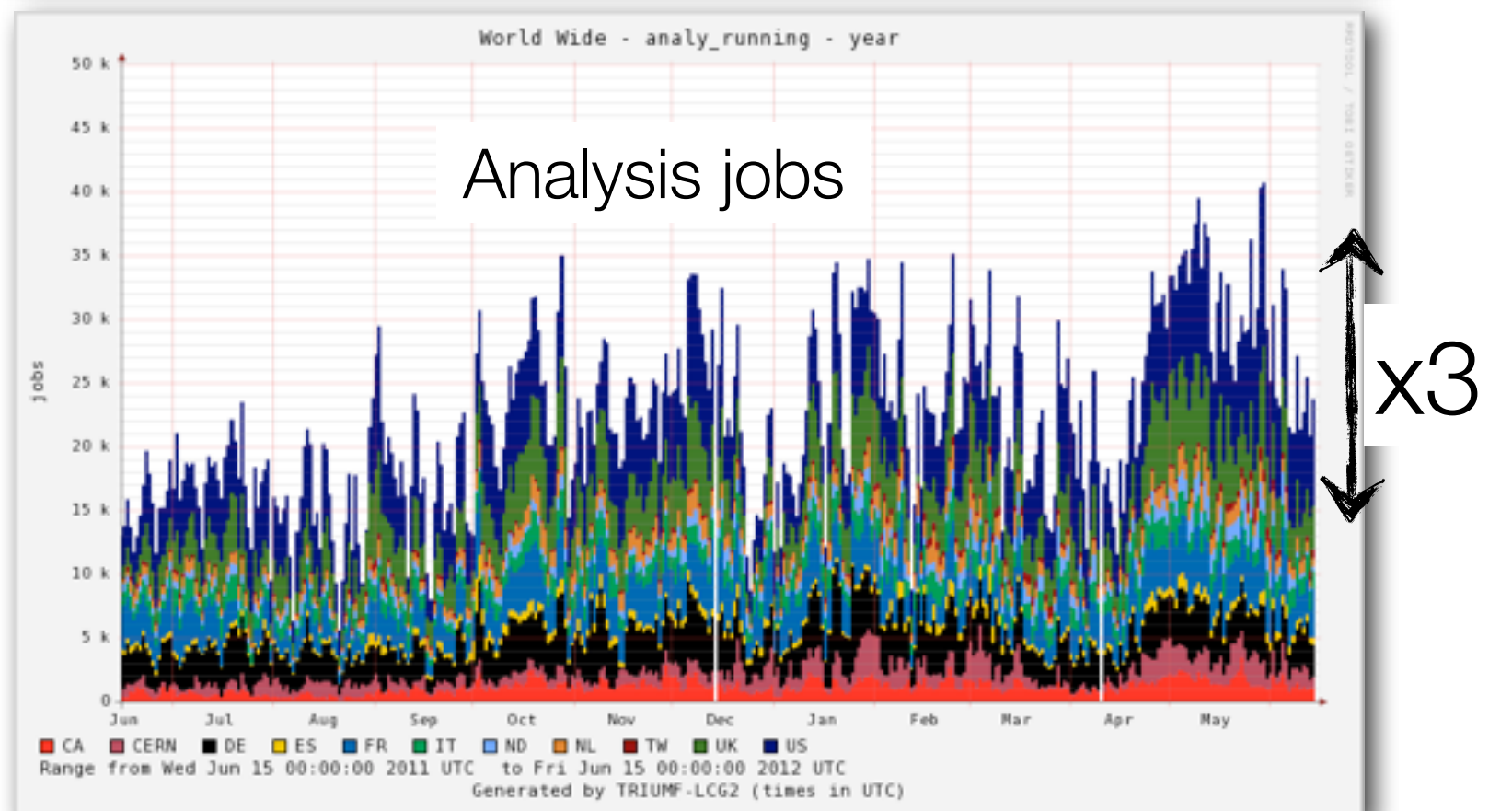
- Heavy data & MC processing

Data & MC processing jobs

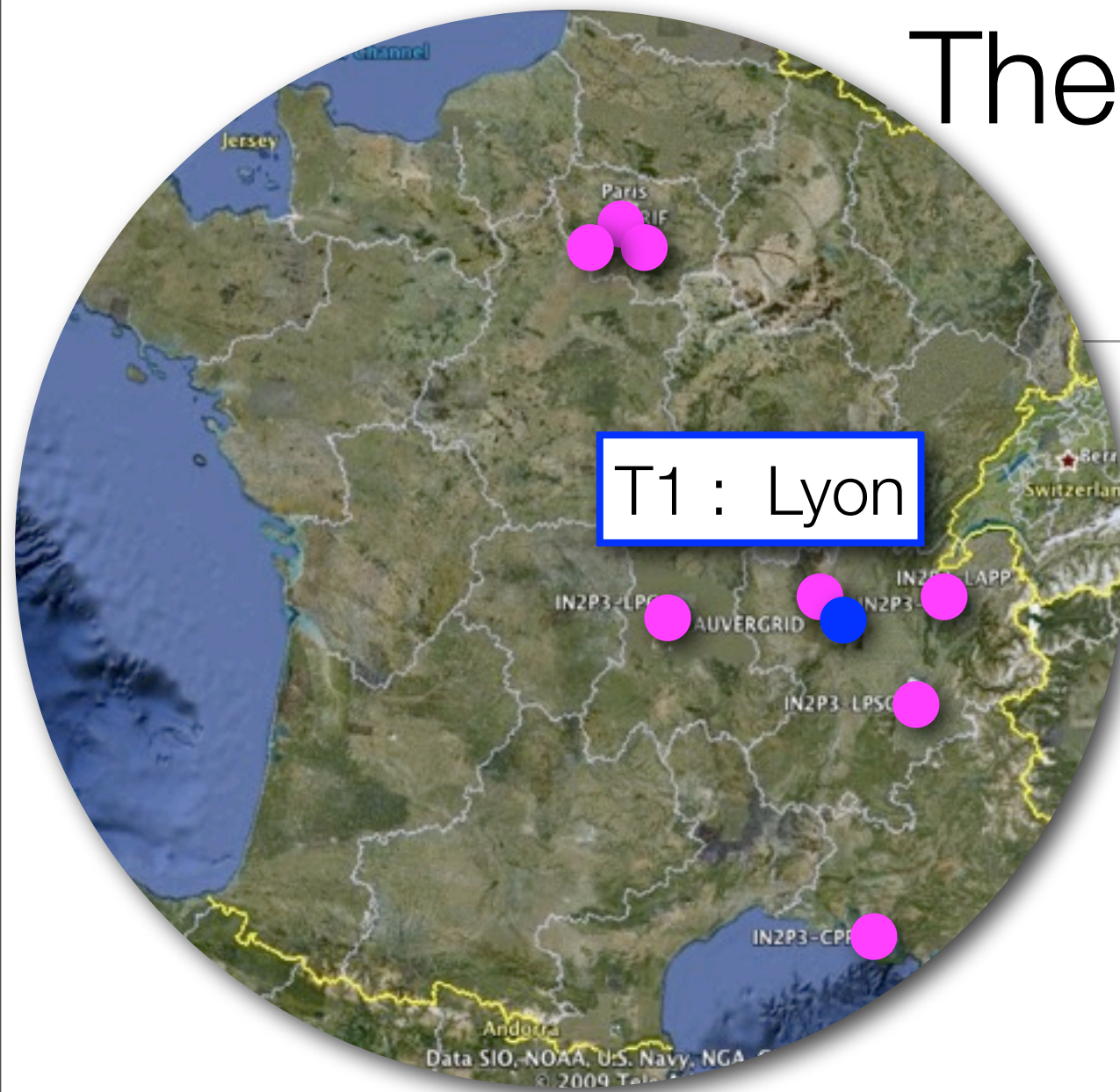


- Increasing load of analysis

Analysis jobs



The “French” cloud



T1 : Lyon

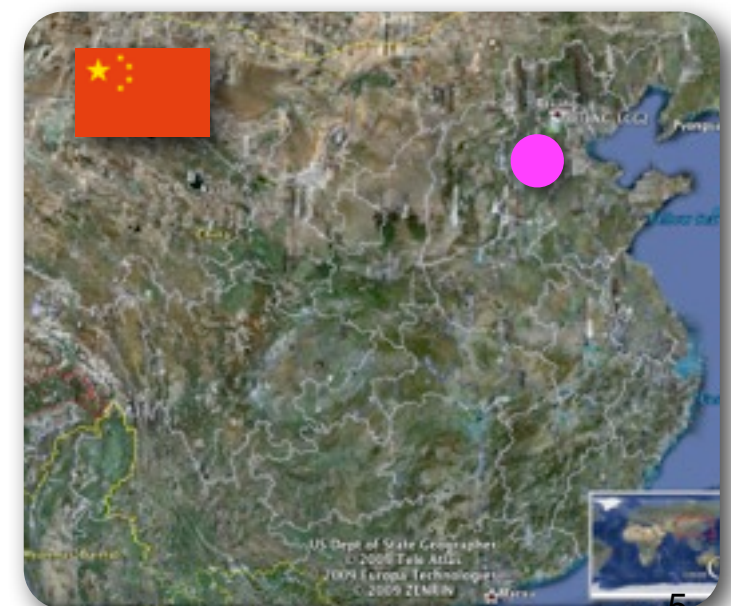
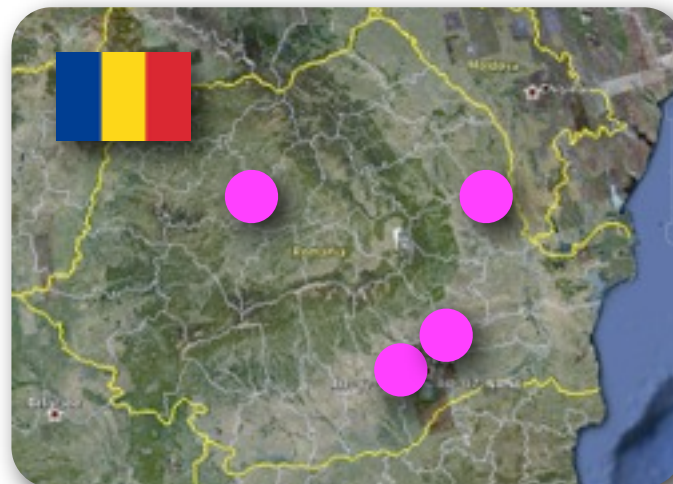
T2s : 14 sites

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

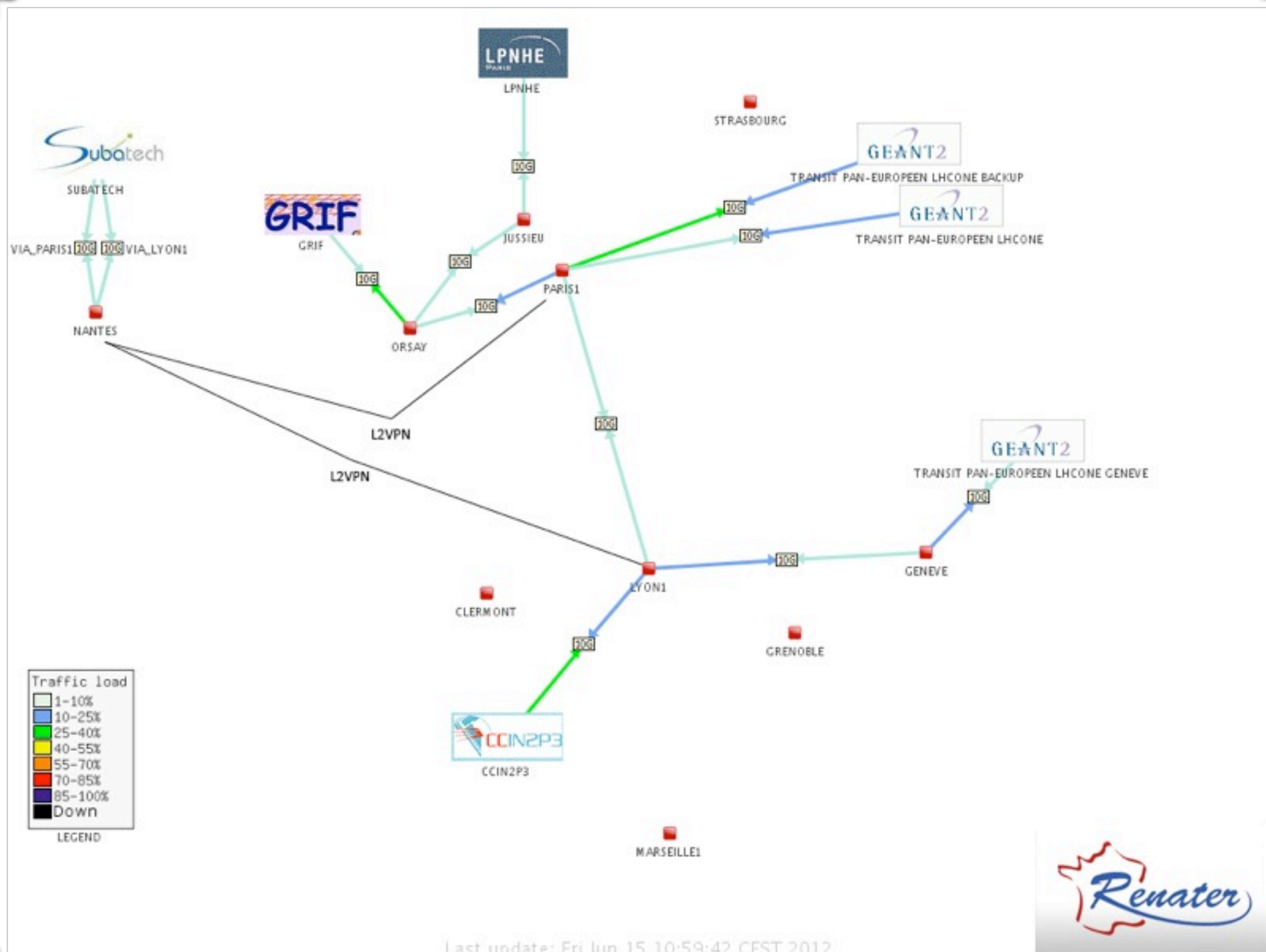


T3 :

- 2 in 2010 (LPSC, CPPM)
- 0 in 2011, T3s became T2s



Most of FR T2s are connected to LHCONE



Cloud main features

- 4 countries
- Several time zone
- Large RTT spread
 - Beijing ~190 ms
 - Tokyo ~290 ms
 - Romanian sites ~55ms
 - Grif ~6 ms

communication

Network

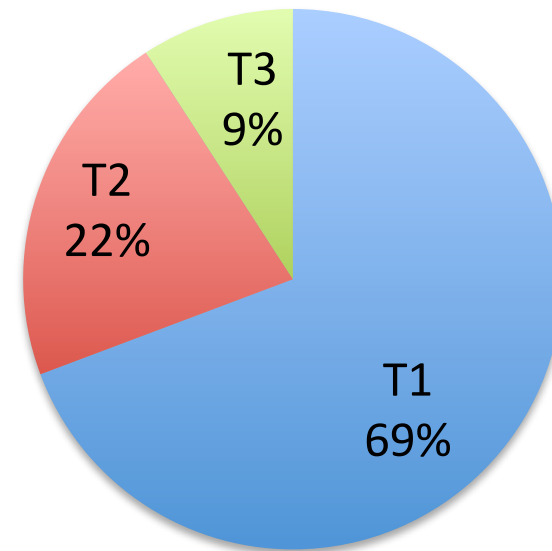
Cloud organization

- **CAF** (ATLAS France Computing)
 - One (at least) representative per French laboratory
 - ‘Experts’ from software & distributed computing
- ATLAS T1 team (Ghita, Manoulis)
- Squad-FR
 - Team : 4 people, rotating ~each week
- S. Crepe (LPSC), Ch. Beau (LPNHE), E. Leguiriec (CPPM), L. Poggioli (LAL)
- Regular CAF meetings :
 - \approx monthly meeting at Lyon
 - Review cloud activities
 - Part of the meeting dedicated to T1
 - Part dedicated to T2s (with T2 representatives)
- Monthly LCG-France technical meetings

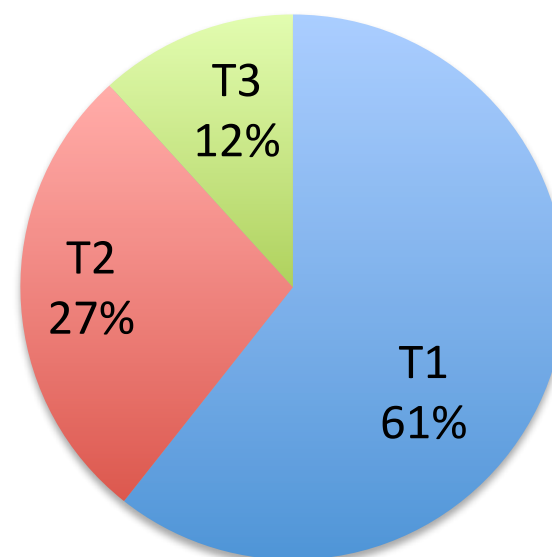
ATLAS resources at Lyon

- ATLAS : **45%** of LCG-France budget
- 2012 : LCG-France budget cut again...
- **T1** (reprocessing, MC processing):
- **T2** (MC production, analysis)
- **T3** (grid and non grid resources for French users) : 50/50

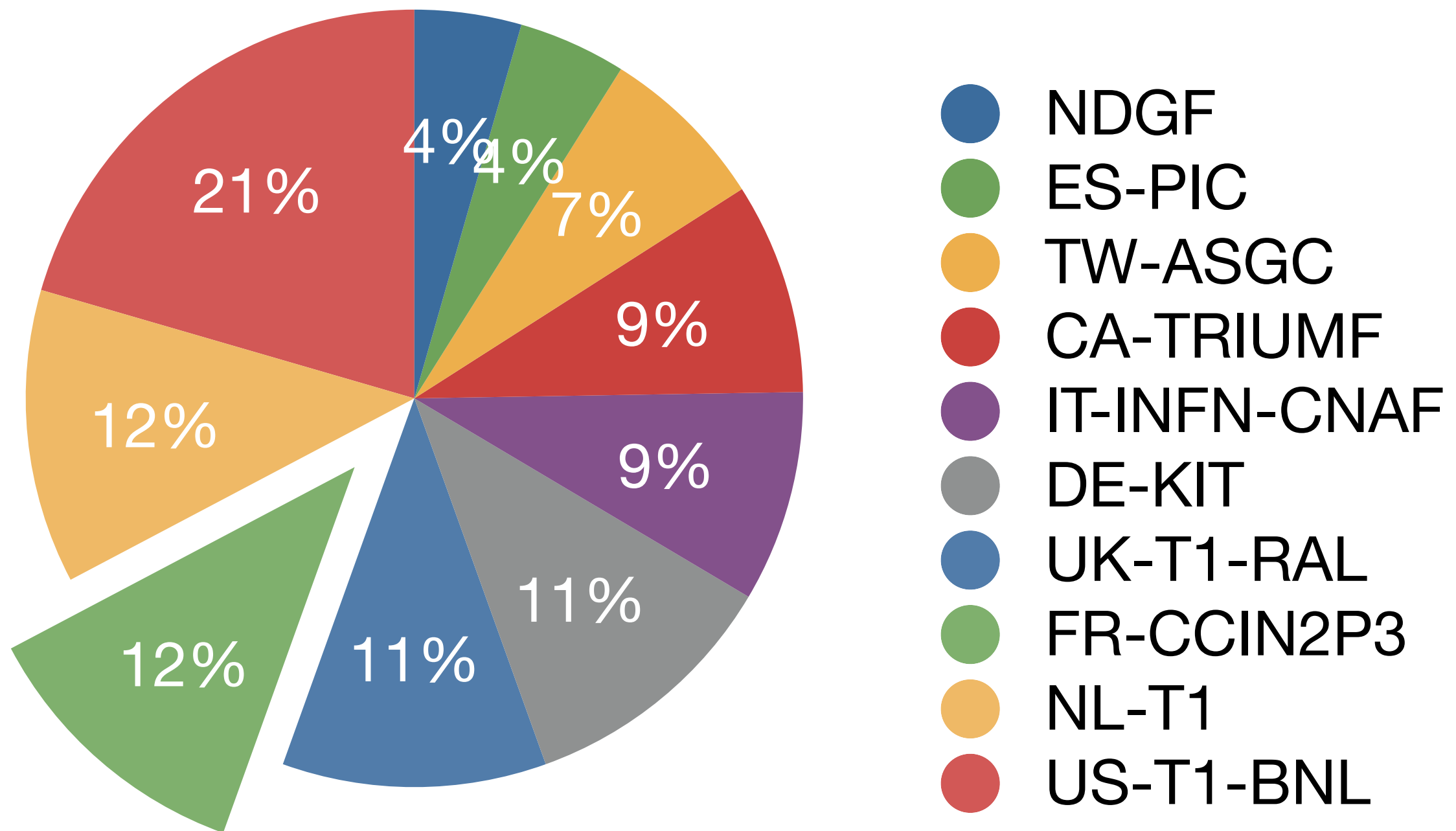
CPU share



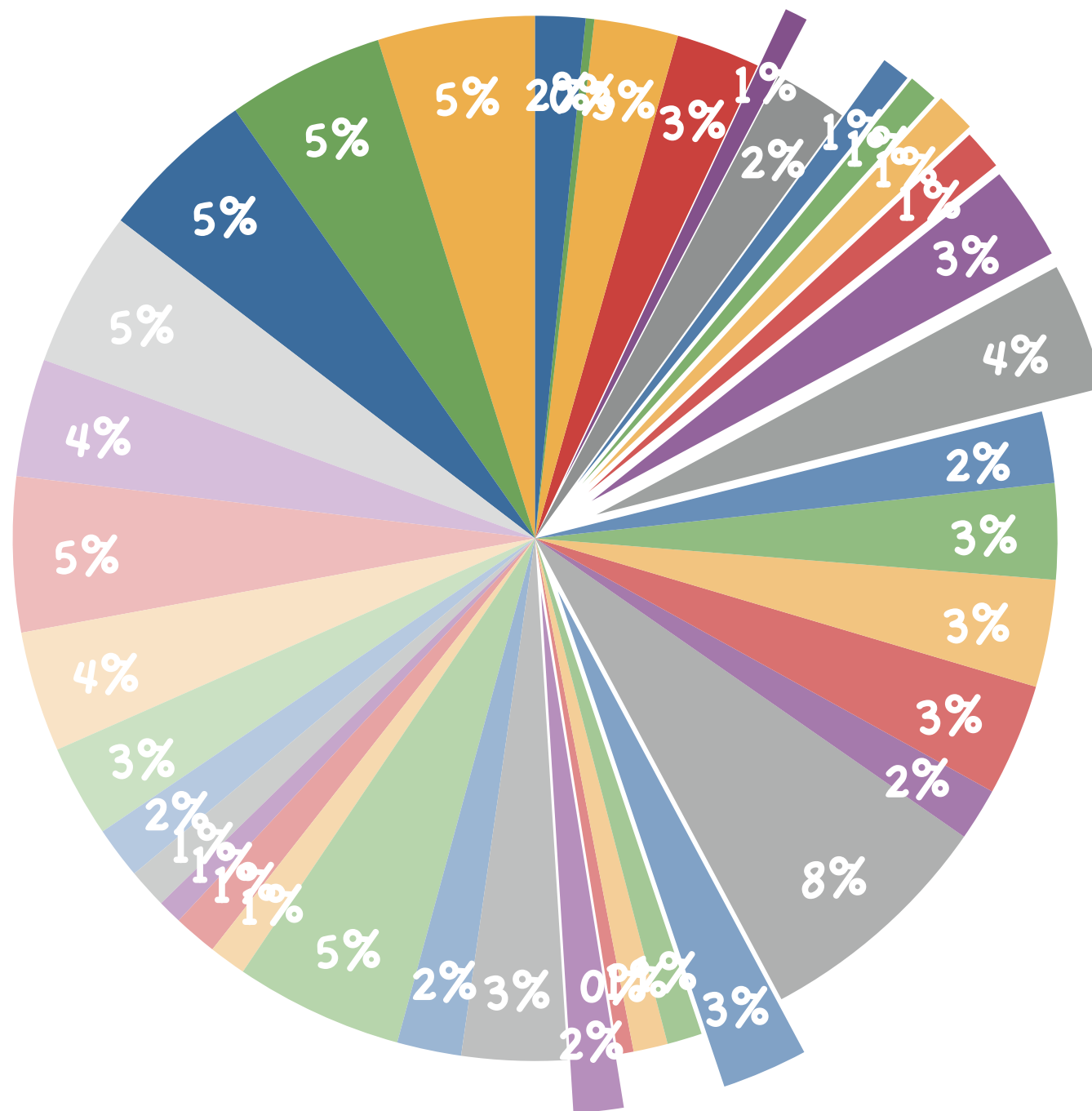
Disk share



T1 2012 disk pledges



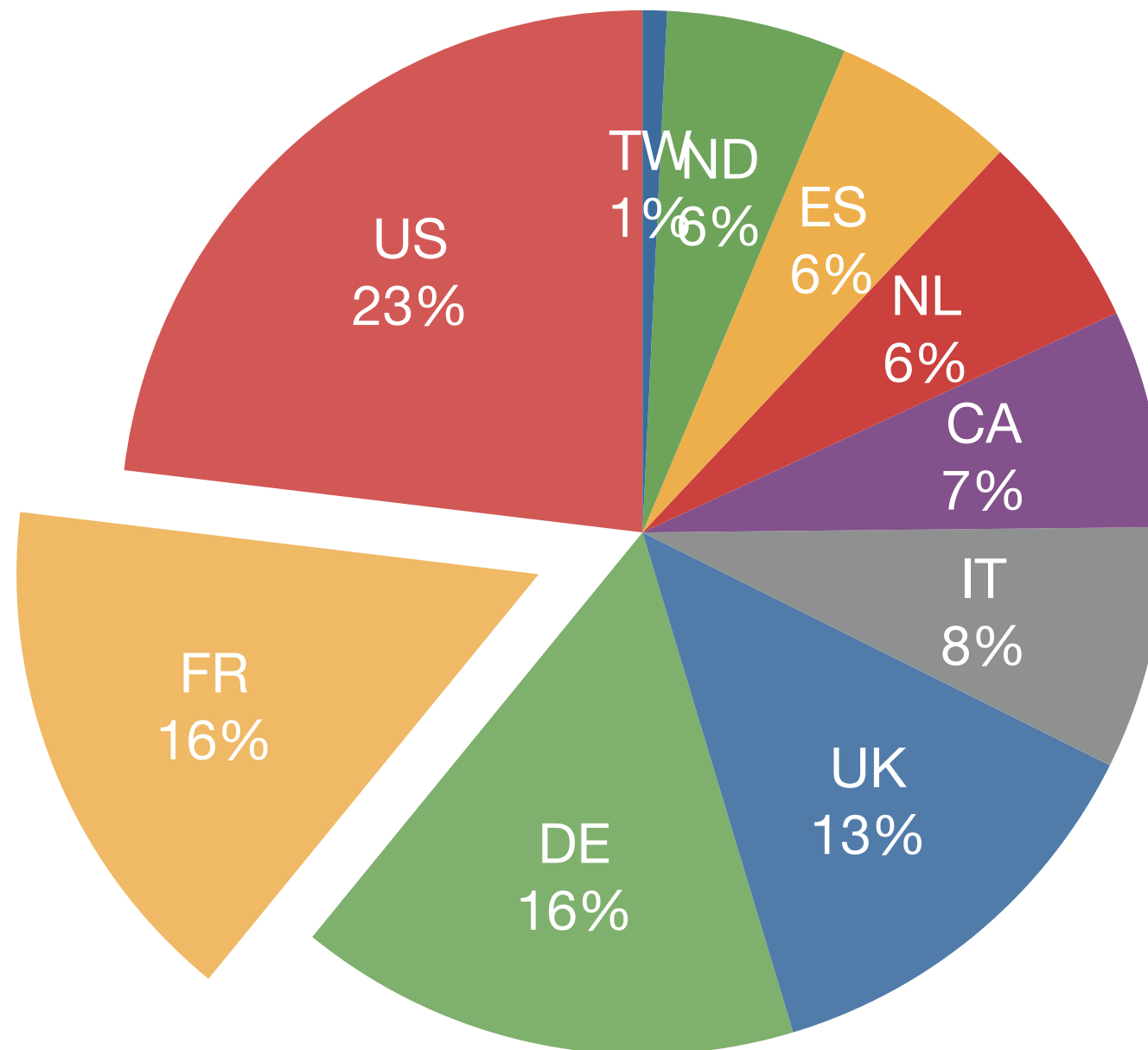
2012 T2 disk pledges



16% of ATLAS

- University of Melbourne Australia
- Austrian Tier-2 Federation Austria
- Canada-East Federation Canada
- Canada-West Federation Canada
- IHEP, Beijing China
- FZU AS, Prague Czech Republic
- CPPM, Marseille France
- LPSC Grenoble France
- LPC, Clermont-Ferrand France
- LAPP, Annecy France
- CC-IN2P3 AF France
- GRIF, Paris France
- ATLAS Federation, HH/Goe Germany
- ATLAS Federation, Munich Germany
- ATLAS Federation DESY Germany
- ATLAS Federation FR/W Germany
- IL-HEP Tier-2 Federation Israel
- INFN T2 Federation Italy
- ICEPP, Tokyo Japan
- UNINETT SIGMA Tier-2 Norway
- Polish Tier-2 Federation Poland
- LIP Tier-2 Federation Portugal
- Romanian Tier-2 Federation Romania
- Russian Data-Intensive GRID Russian Federation
- SiGNET Slovenia
- ATLAS Federation Spain
- SNIC Tier-2 Sweden
- CHIPP Switzerland
- Taiwan Analysis Facility Federation Taipei
- Turkish Tier-2 Federation Turkey
- SouthGrid UK
- ScotGrid UK
- London Tier 2 UK
- NorthGrid UK
- Northeast ATLAS T2 USA
- Great Lakes ATLAS T2 USA
- Midwest ATLAS T2 USA
- SLAC ATLAS T2 USA
- Southwest ATLAS T2 USA

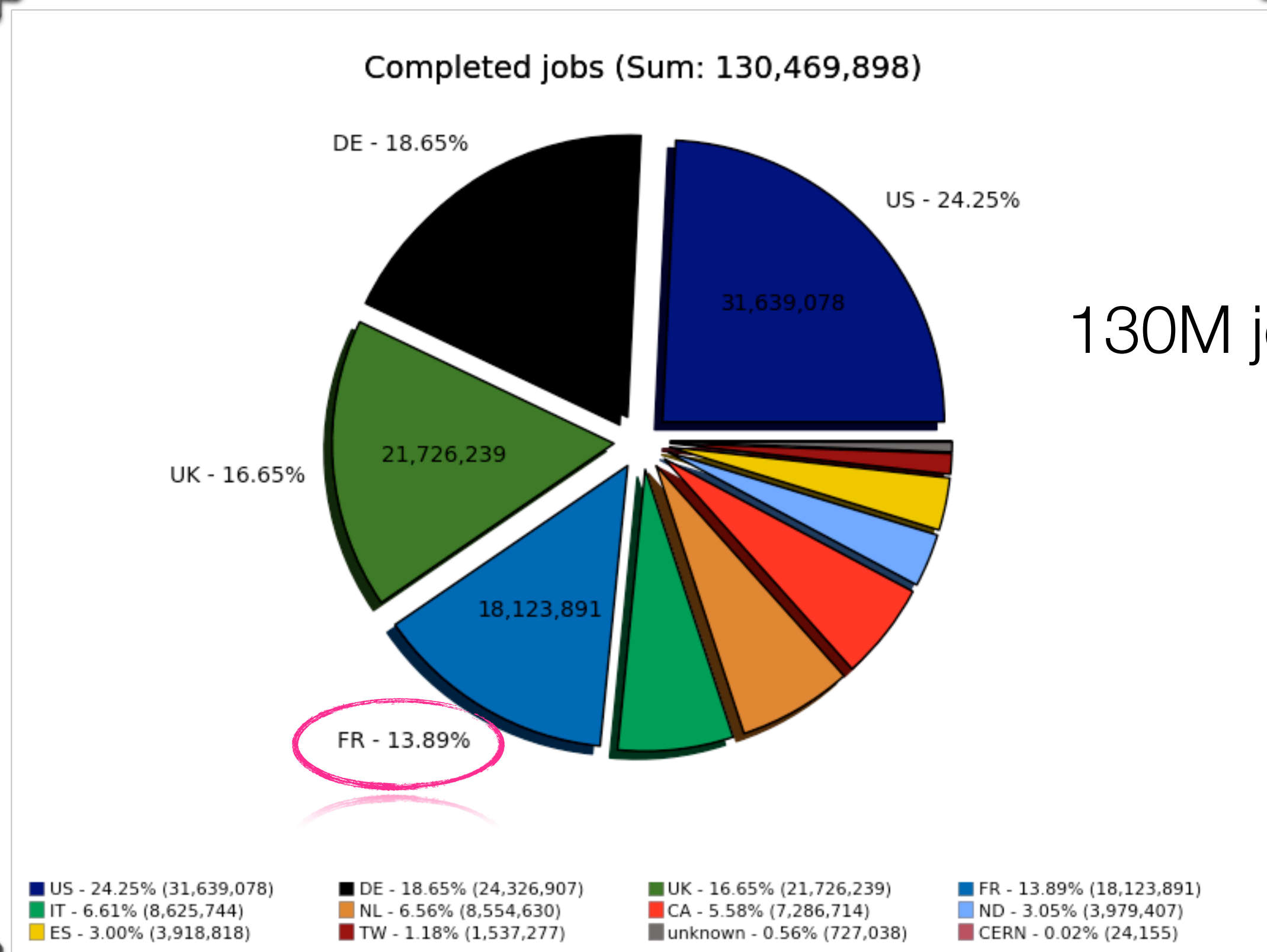
T2 2012 disk pledge by cloud



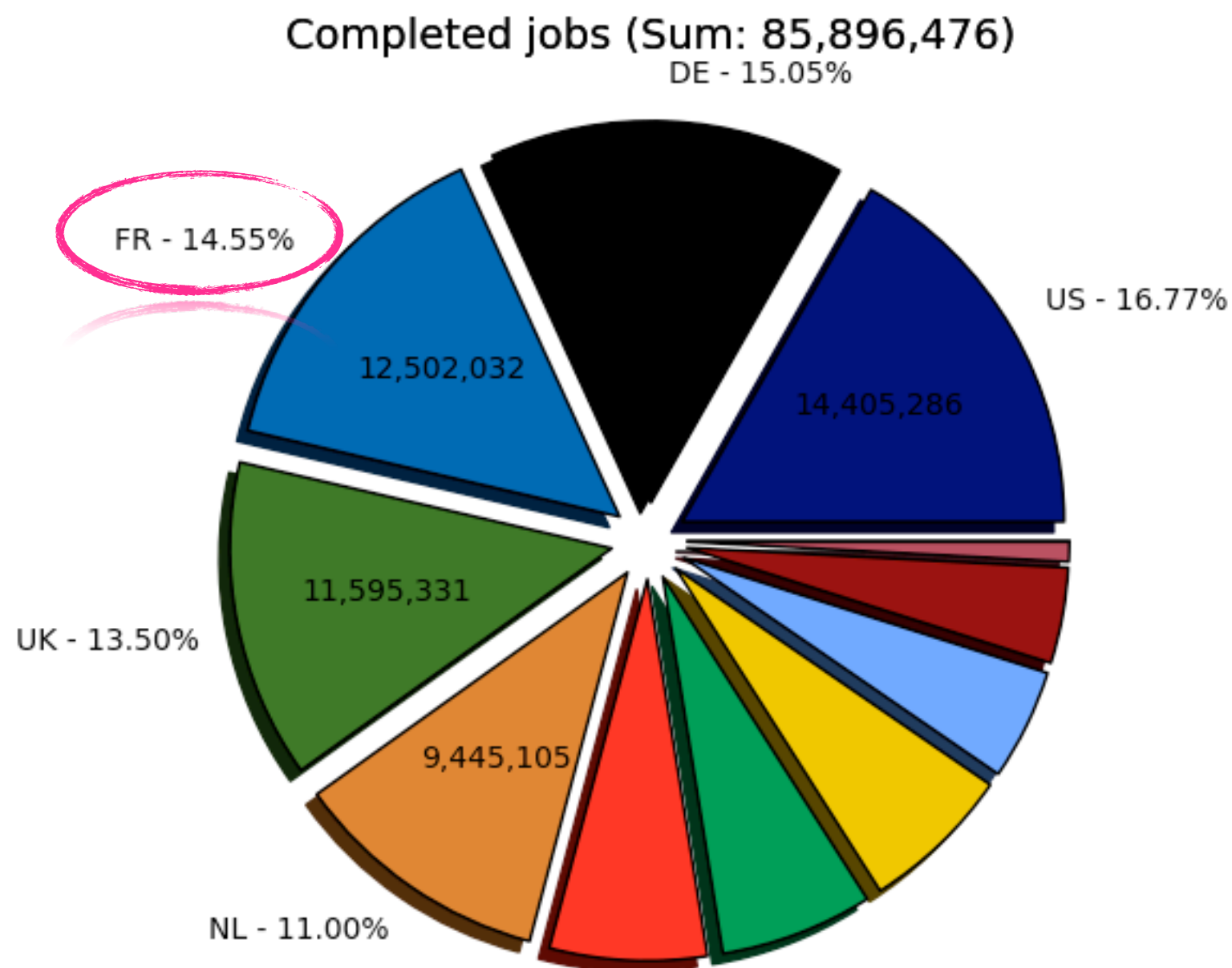
FR-cloud share :
12 - 16 % of ATLAS activities

Performances over last year
[May, 2011 - May,2012]

Analysis jobs [May, 2011 - May, 2012]



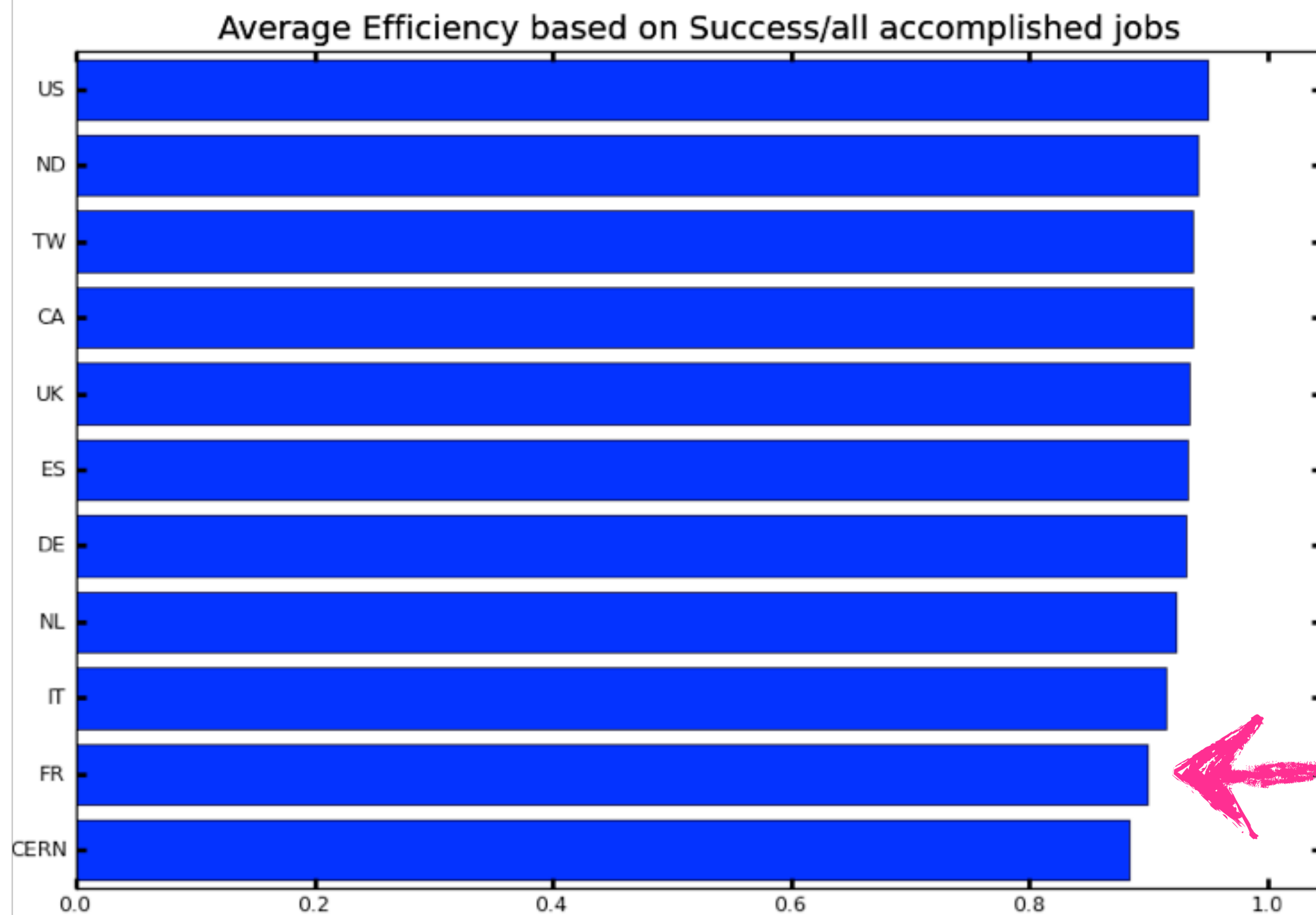
Production jobs [May, 2011 - May, 2012]



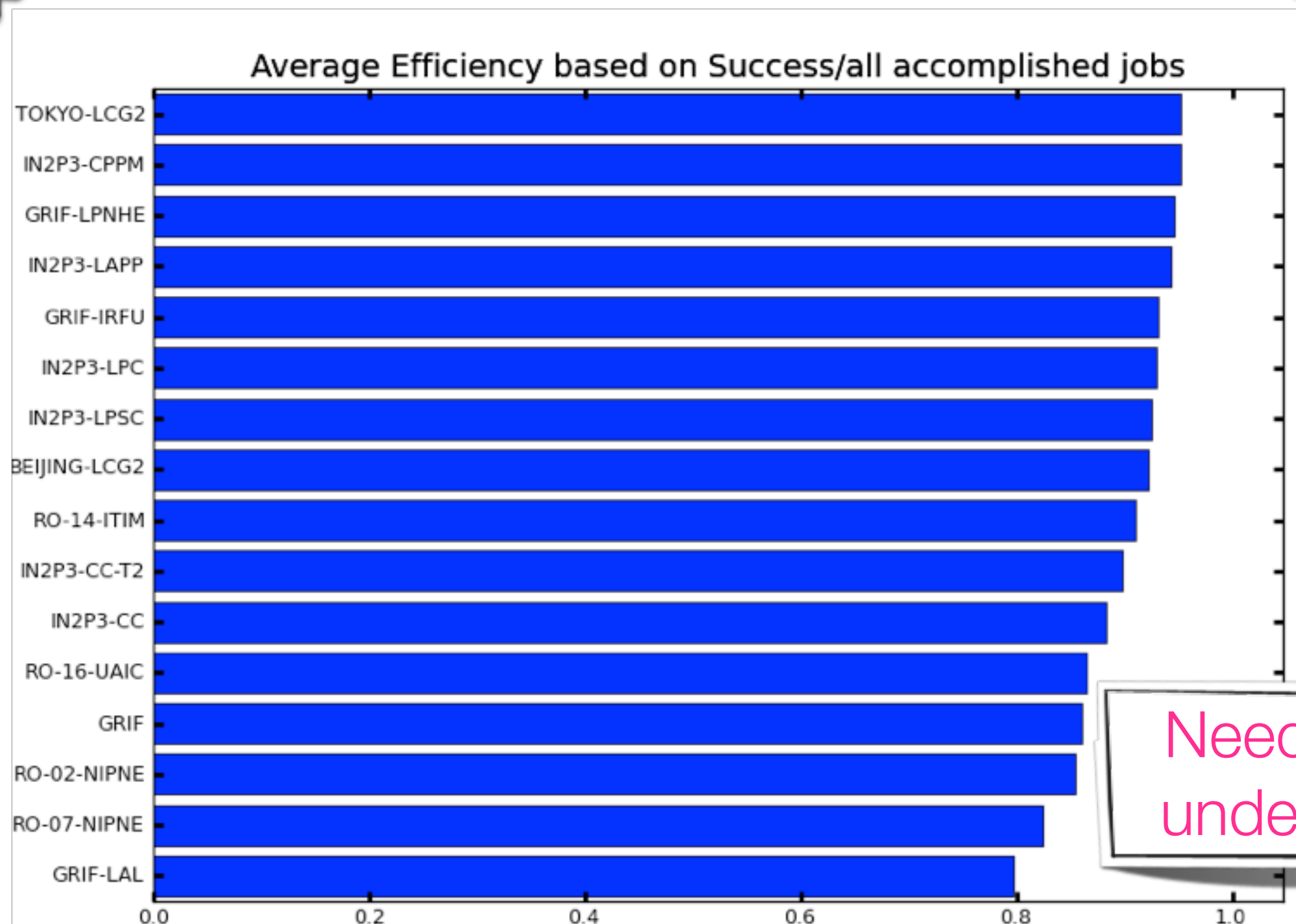
86M jobs

US - 16.77% (14,405,286)	DE - 15.05% (12,926,429)	FR - 14.55% (12,502,032)	UK - 13.50% (11,595,332)
NL - 11.00% (9,445,105)	CA - 6.65% (5,713,111)	IT - 6.47% (5,558,499)	ES - 6.42% (5,510,819)
ND - 4.71% (4,042,701)	TW - 4.06% (3,483,851)	CERN - 0.83% (713,311)	unknown - 0.00% (0.00)

Production jobs [May, 2011 - May, 2012]

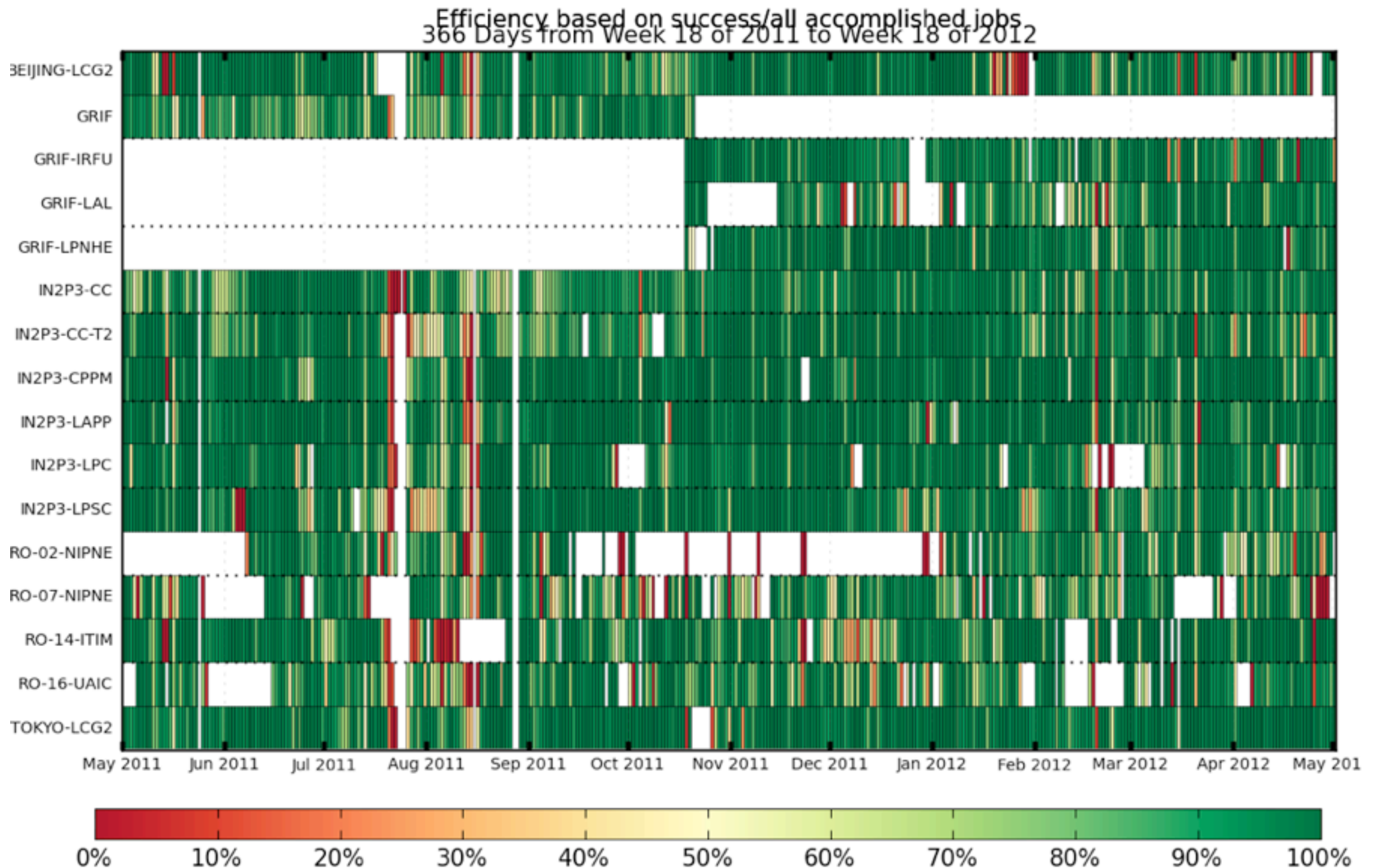


Production jobs [May, 2011 - May, 2012]

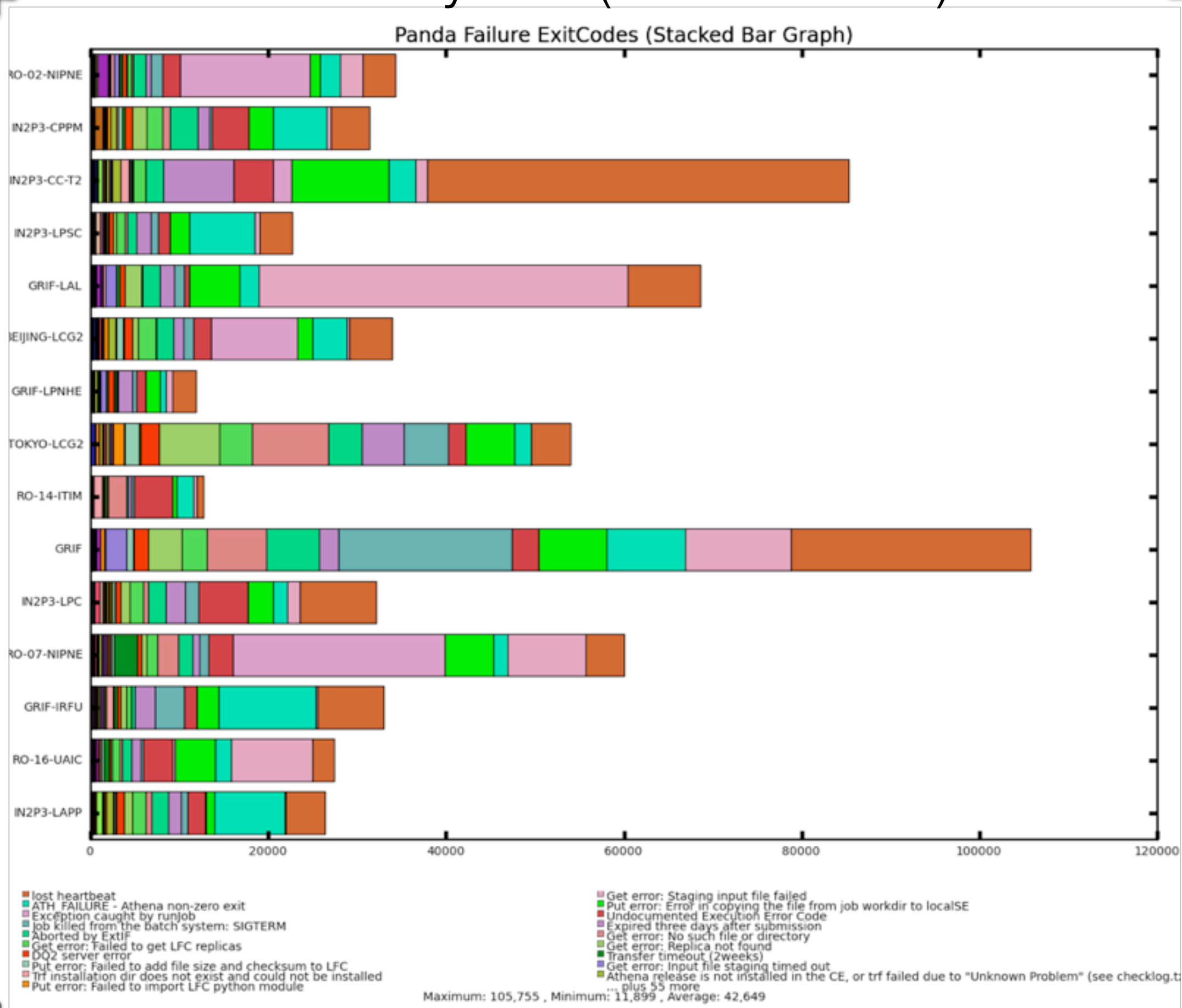


Need to be understood

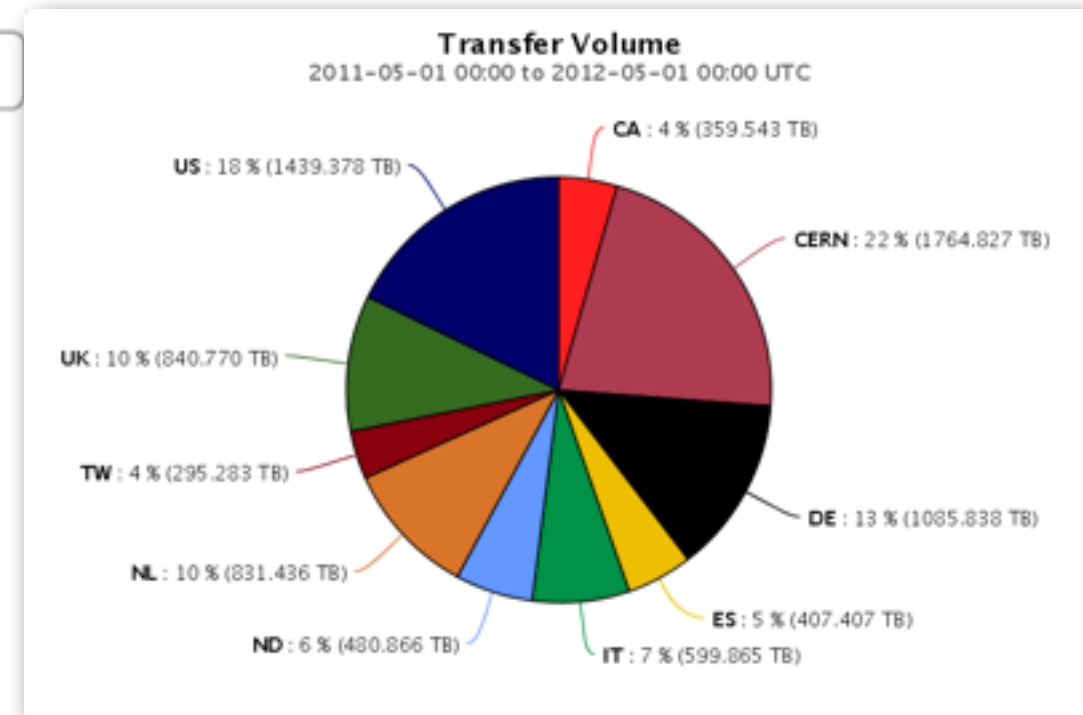
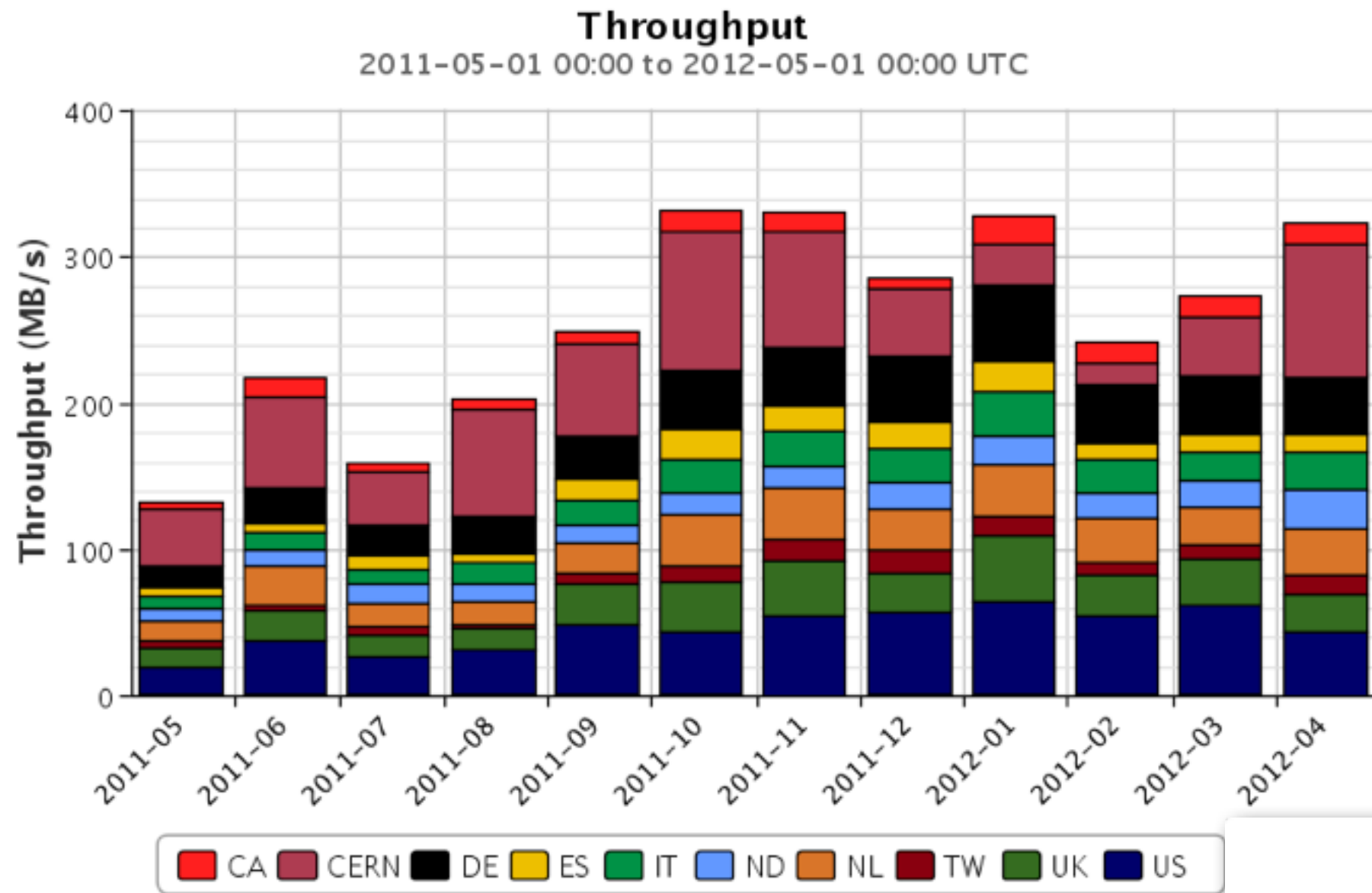
Production efficiency vs time by site



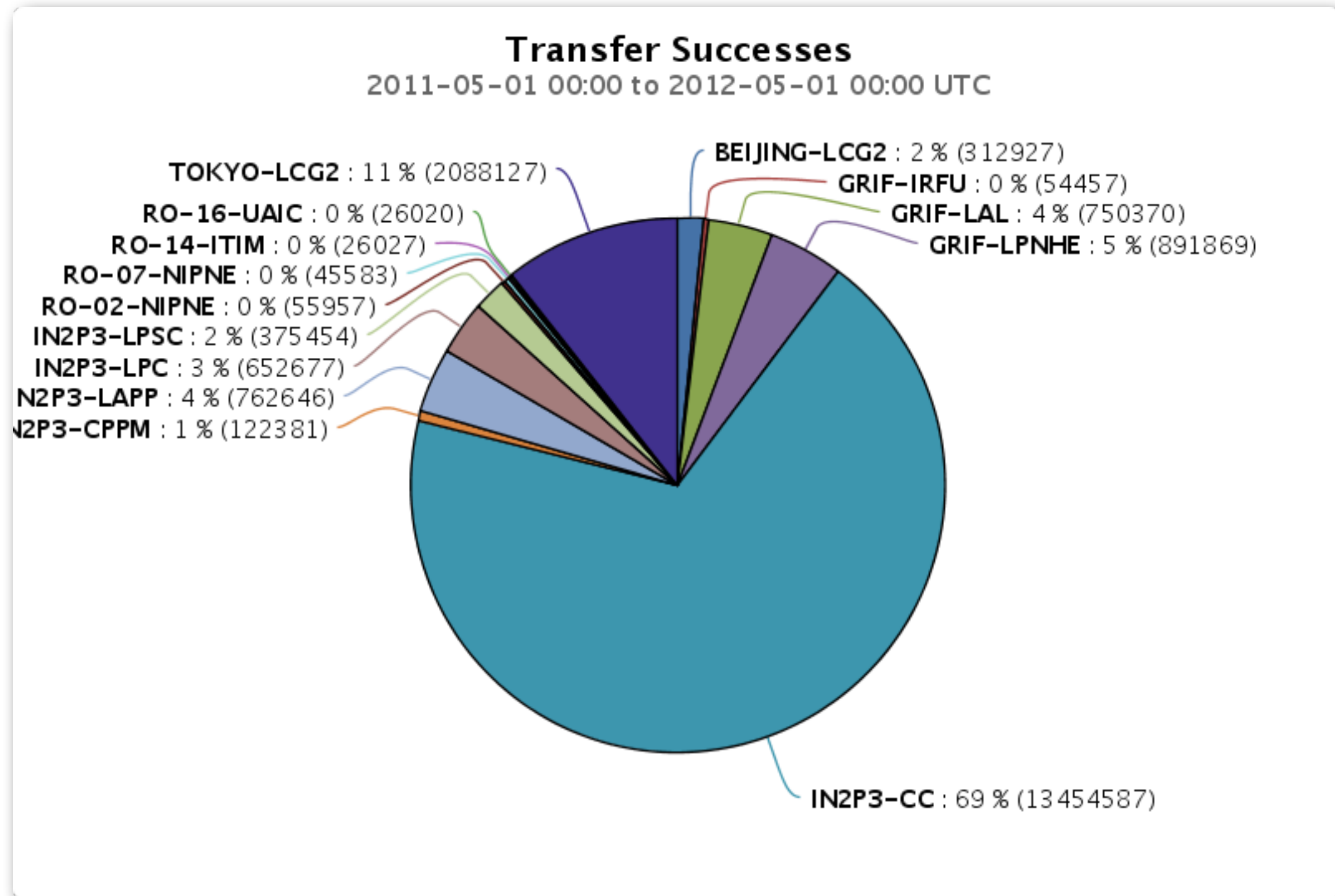
errors by site (T1 excluded)



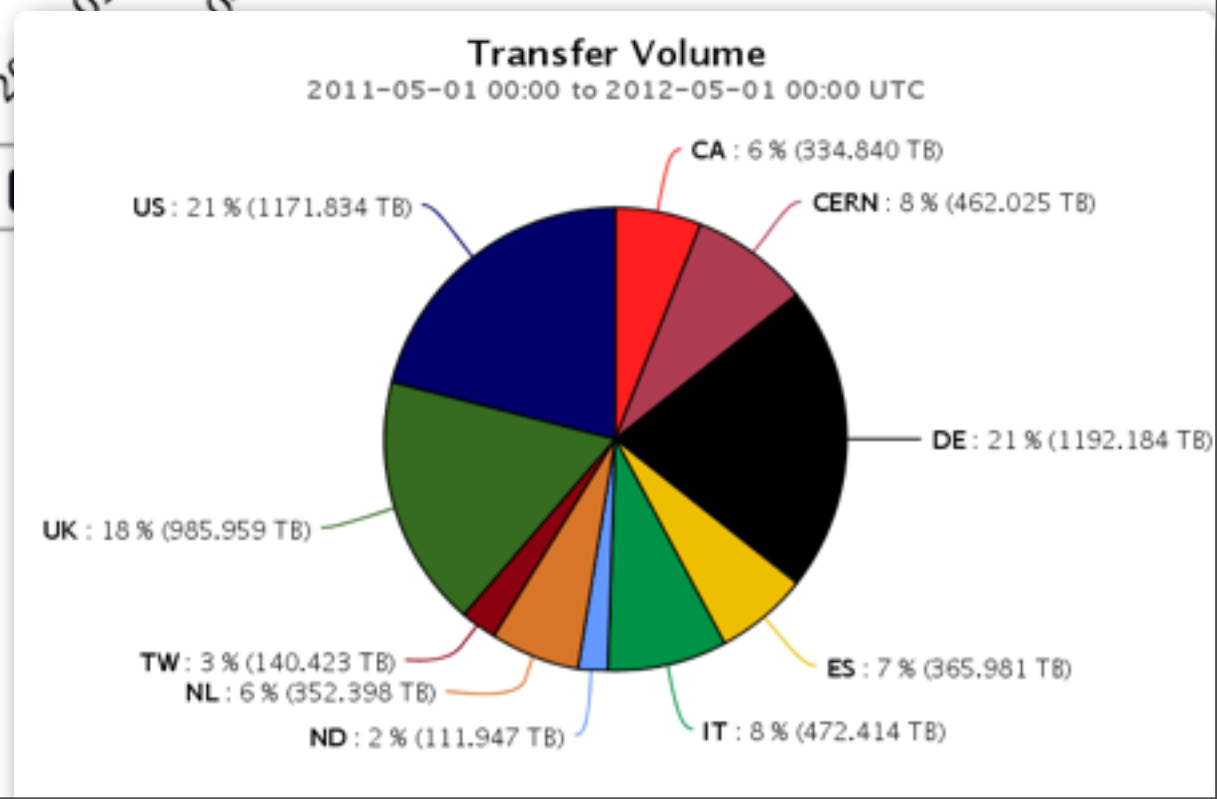
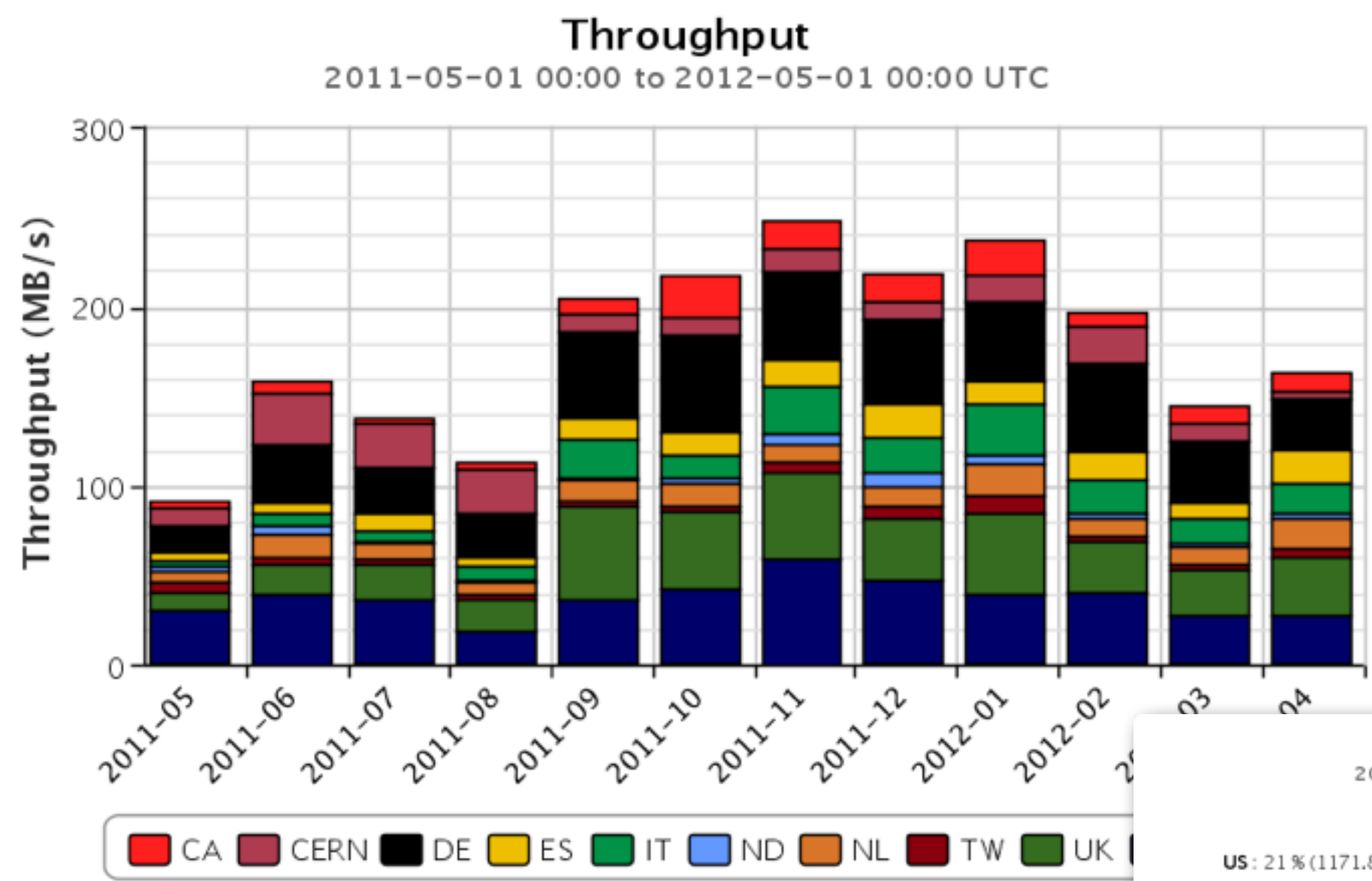
Transfer to FR cloud [May 2011 - May 2012] : **8.1 PB**



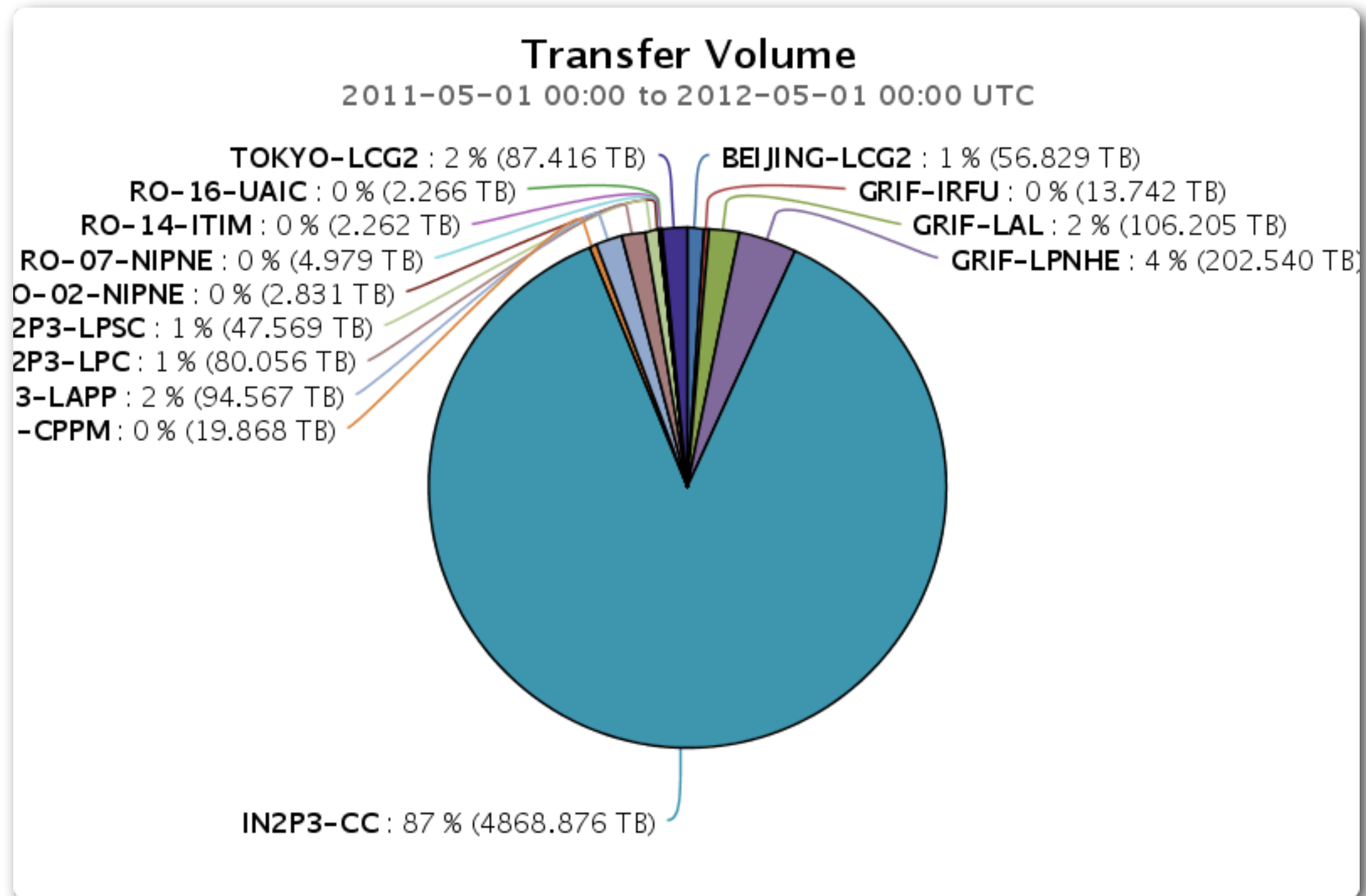
import destination : **70% to T1**



Export from FR-cloud [May, 2011 - May, 2012] : **5.6 PB**



export : source **87% from T1**



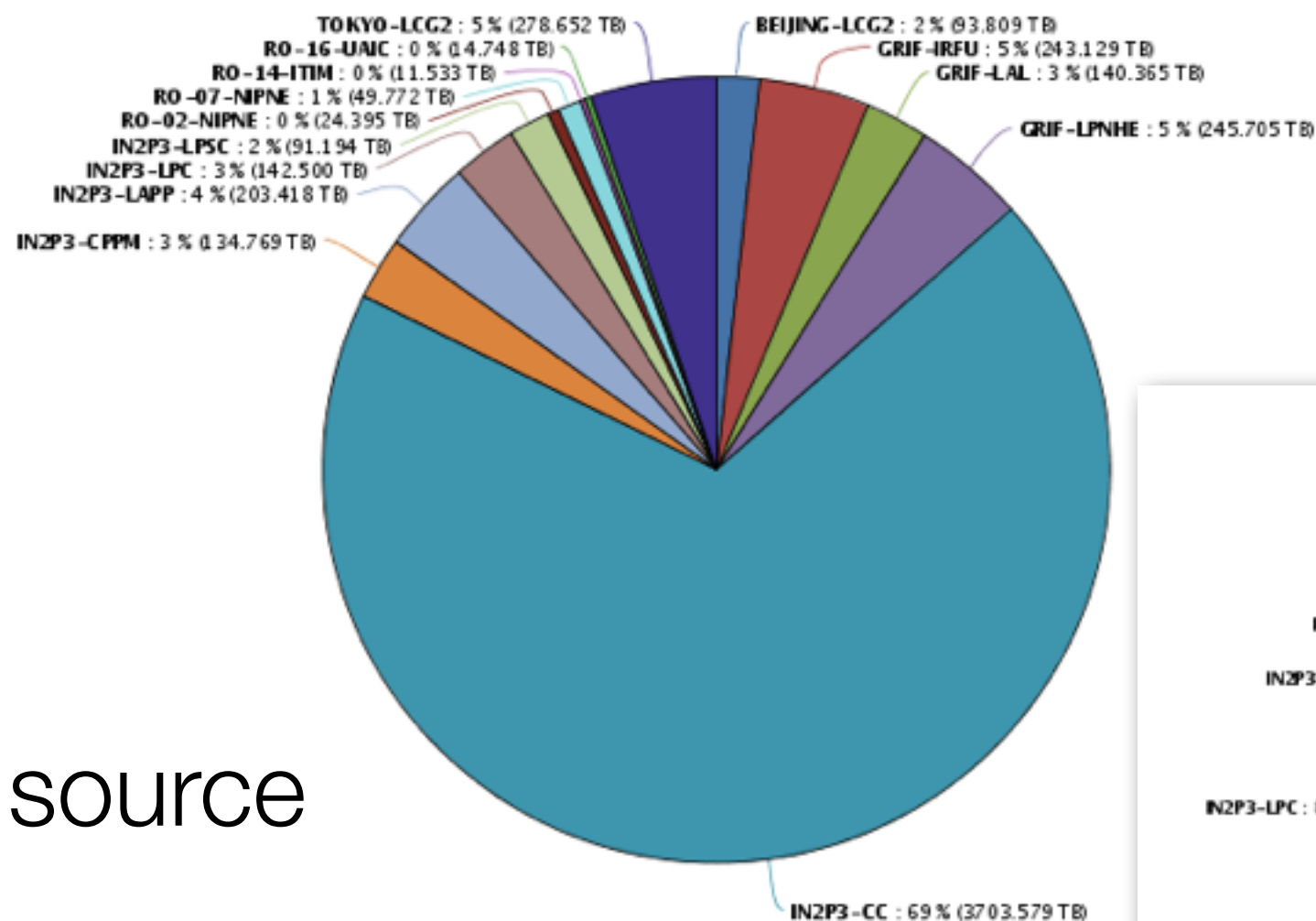
Intra cloud transfers [May,2011 - May,2012]

Not so uniform



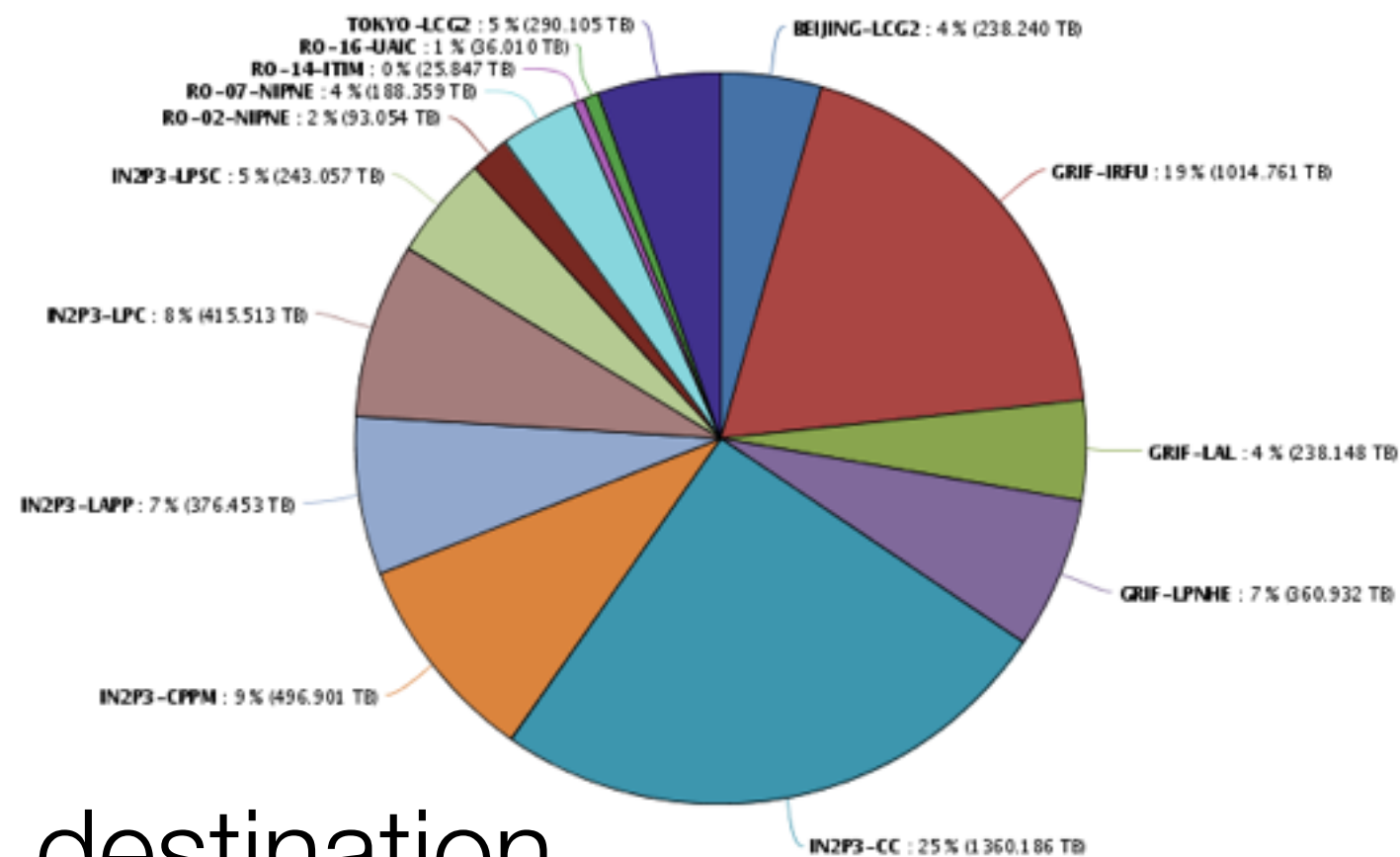
FR-cloud **intra** volume of data transfer : **5.4 PB**

Transfer Volume
2011-05-01 00:00 to 2012-05-01 00:00 UTC



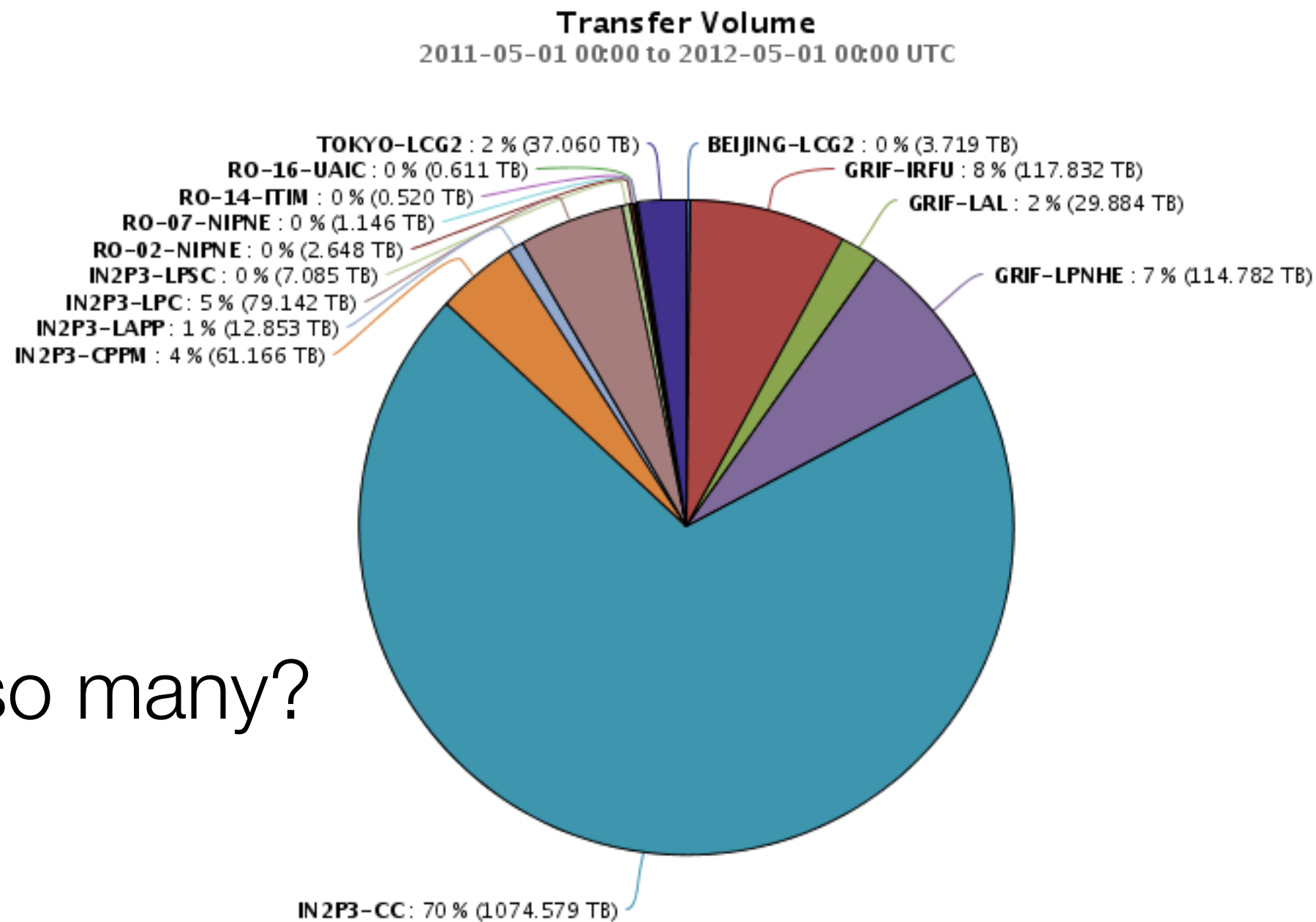
source

Transfer Volume
2011-05-01 00:00 to 2012-05-01 00:00 UTC



destination

user subscription destination on FR-cloud : **1.5 PB**

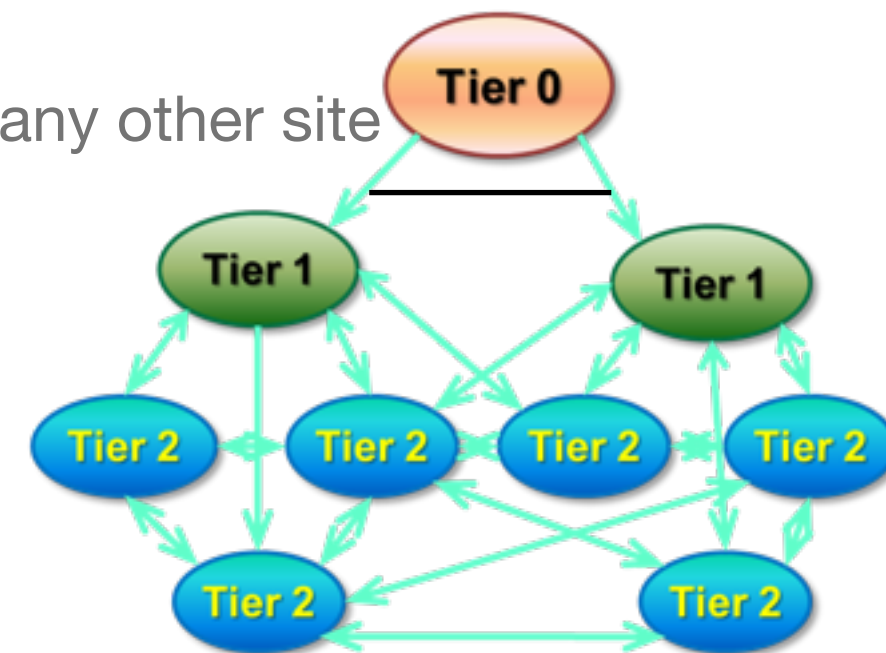


why so many?

correlated to 'local' activity at sites ?

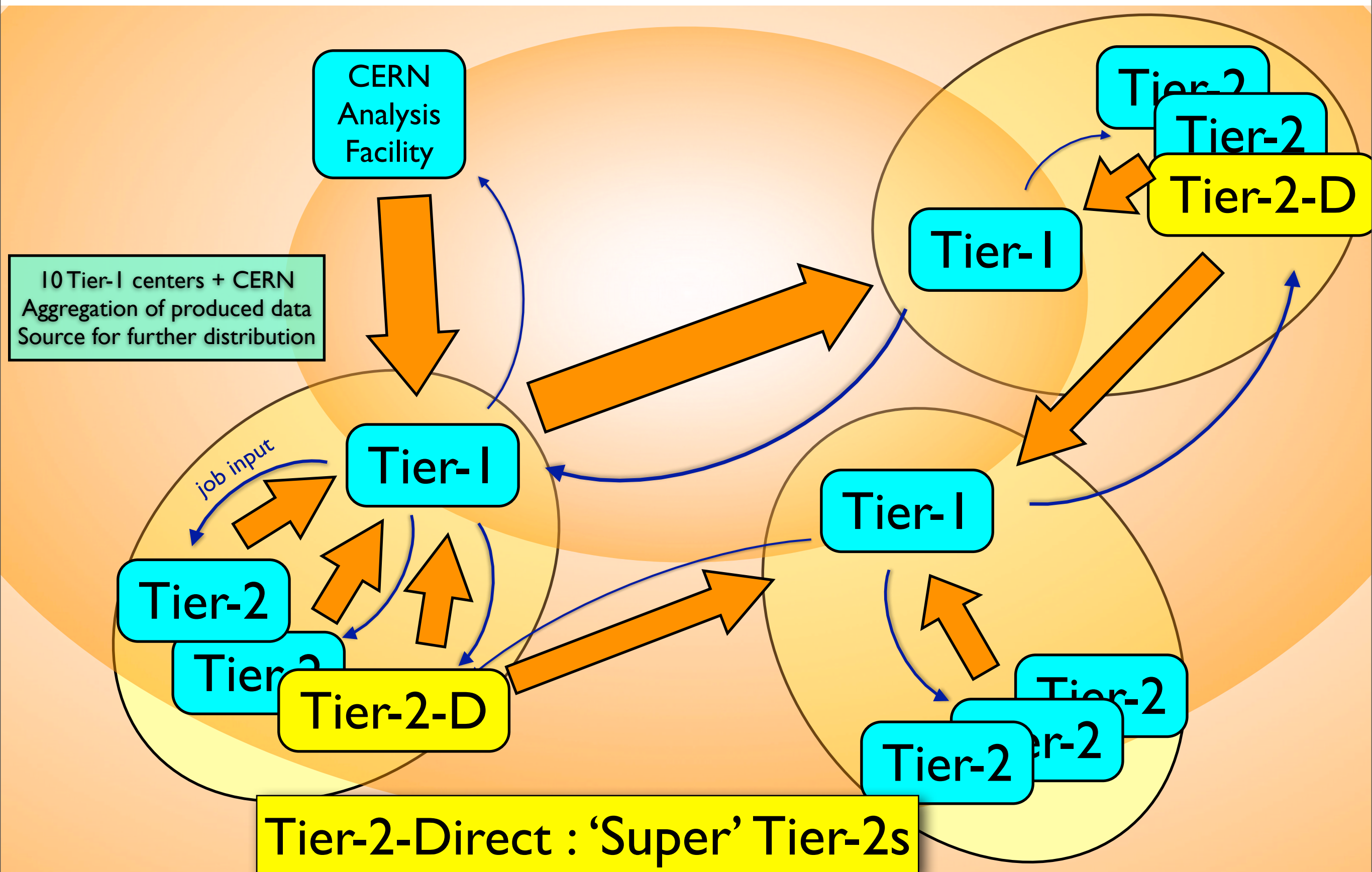
The ATLAS Data Model has changed

- 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**



Thank you Dan v.d. Steer

Data Processing Model Revised



ATLAS sites and connectivity

- ATLAS computing model is evolving
 - Experience of one year of data taking
 - Tools and monitoring are getting more mature
- New category of sites : Direct T2s (**T2Ds**)
 - primary hosts for datasets (analysis) and for group analysis
 - get and send data from different sites
- Connectivity is important for T2Ds (not only...)

Aim :
all FR-cloud T2s
should be T2D

validated T2Ds

- CA : CA-SCINET-T2, CA-VICTORIA-WESTGRID-T2, SFU-LCG2_DATADISK, CA-MCGILL
- DE : DESY-HH, DESY-ZN, MPPMU, LRZ-LMU, CSCS-LCG2, GOEGRID, UNI-FREIBURG,
- ES : IFIC-LCG2, IFAE, UAM-LCG2, LIP-LISBON, NCG-INGRID-PT
- FR : GRIF-LPNHE, GRIF-LAL, IN2P3-LAPP, IN2P3-LPC, IN2P3-LPSC, BEIJING-LCG2
- IT : INFN-NAPOLI-ATLAS, INFN-MILANO-ATLASC, INFN-ROMA1
- UK : UKI-LT2-QMUL, UKI-NORTHGRID-MAN-HEP, UKI-SCOTGRID-GLASGOW, UKI-NO
- US : AGLT2, MWT2_UC, NET2, SLACXRD, SWT2 CPB

Network performance monitoring

- Monitoring tools developed in 2011
 - ATLAS 'sonar' : 'calibrated' file transfers by the ATLAS Data Distribution system, from storage to storage
 - perfSONAR : network performance (throughput, latency, traceroute, ...)
- Non stable network performances : also side effects on storage system at CCIN2P3

more on Monday

- T2Ds: Tier2s directly connected to Tier1s of different clouds
 - “directly” in the DDM topology*ATLAS transfers from T2D to any T1*
- T2Ds should
 - demonstrate “good” connectivity from/to every T1/T2D
 - provide a certain level of commitment*Good connectivity*
- T2Ds could be de-commissioned if they degrade
 - Performances monitored*
- There is no “maximum number” of T2Ds
 - We should not create too many channels (FTS performance)

T2D: revising the criteria

Current criteria

- All transfers from the candidate T2D to 10/12 T1s for big files ('L') must be above 5 MB/s during the last week and during 3 out of the 4 last weeks.
- All transfers from 10/12 T1s to the candidate T2D for big files must be above 5 MB/s during the last week and during 3 out of the 4 last weeks

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>

- New T2D candidates
 - Sites/clouds should contact Central Operations if they want to be a candidate T2D
- Site will start receiving more Sonar Tests (large files from/to all T1s and T2Ds)
 - If the performance is acceptable, results are reported to the ADC Weekly Meeting and the site is declared a T2D
 - Need for FTS configuration of channels at T1s and this is done approx every 3 months
- So far, T2Ds candidates are monitored by squads
 - Cloud and sites themselves should take the necessary actions (monitoring, improving performance, reporting results to Central Operations)

perfSONAR(-PS) and ATLAS

Needed on every FR-cloud site

- Being deployed on T2D sites
- To measure latency and throughput (2 machines) between sites (matrix)
- Details in Shawn McKee's (BNL) talk at last GDB
- Matrix monitoring available



<http://psps.perfsonar.net/>

Summary for LHCONE

- ❄ Our specific goal in setting up perfSONAR-PS for LHCONE is to acquire before and after network measurements for the selected early adopter sites. This is **not** the long-term network monitoring setup for LHCONE...that is **TBD**
- ❄ Details of which sites and how sites should setup the perfSONAR-PS installations is documented on the Twiki at: <https://twiki.cern.ch/twiki/bin/view/LHCONE/SiteList>
- ❄ In the next few slides I will highlight some of the relevant details

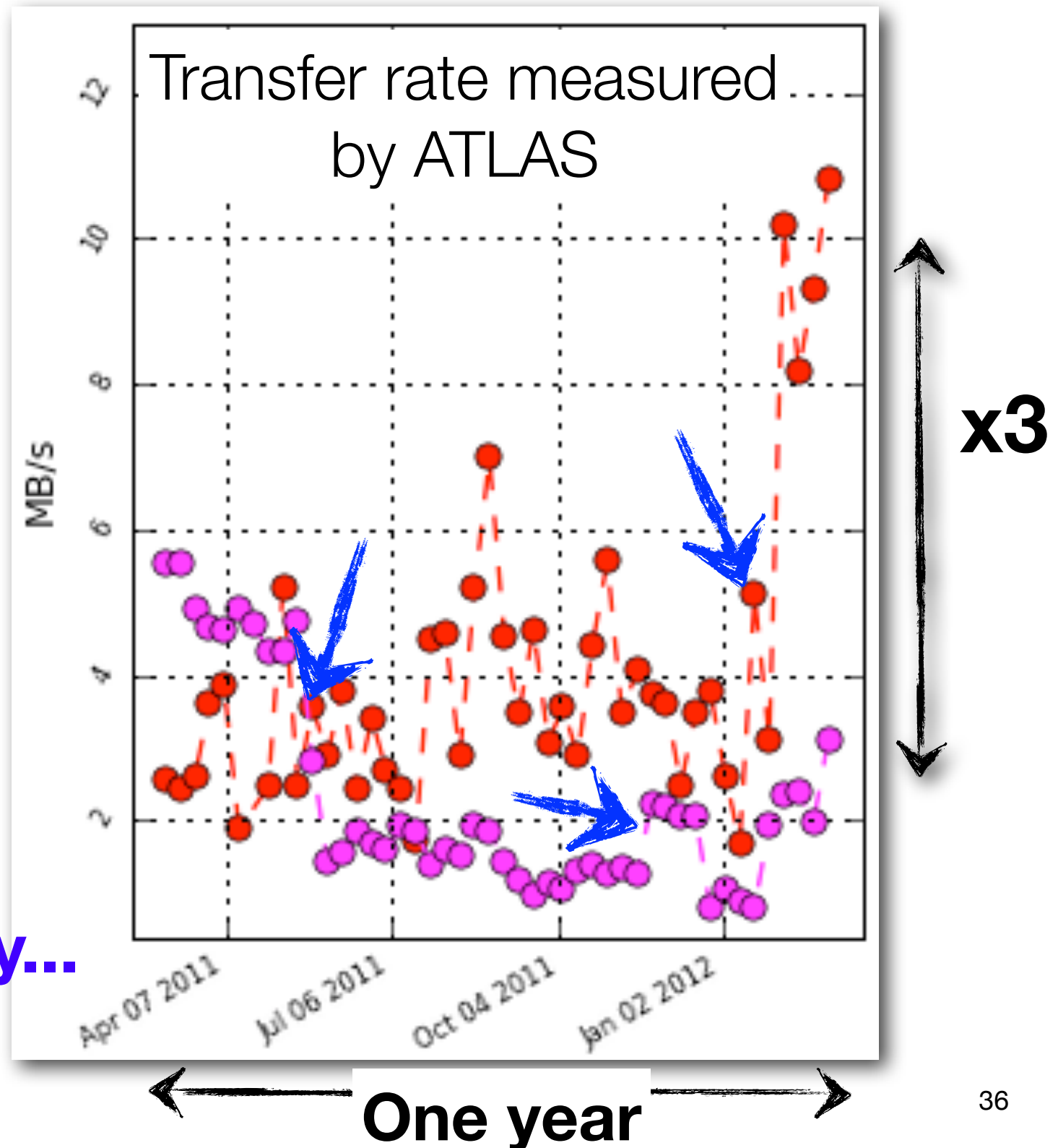
ATLAS transfers to/from Beijing over one year

Beijing → CCIN2P3

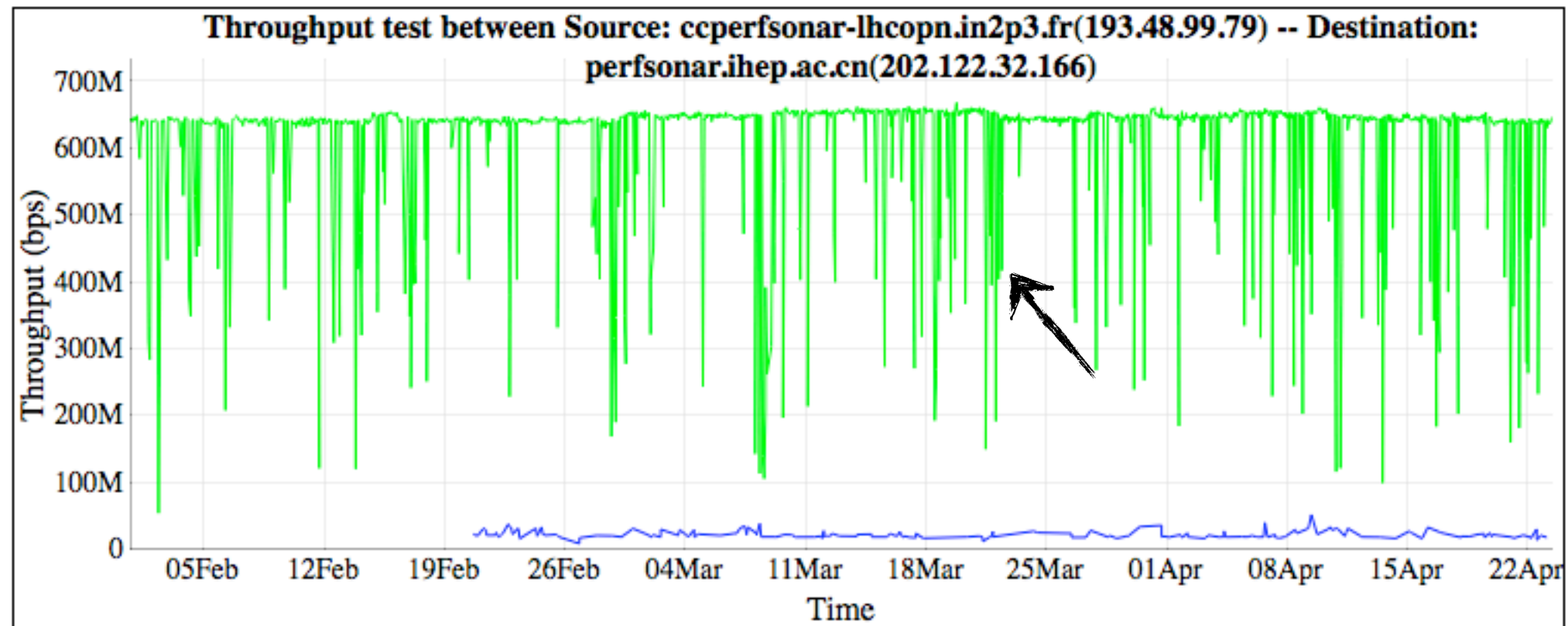
CCIN2P3 → Beijing

*asymmetry (why?) in
transfer rate
performance reversed
over last year*

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



Network throughput measured with perfsonar



<- 1 month

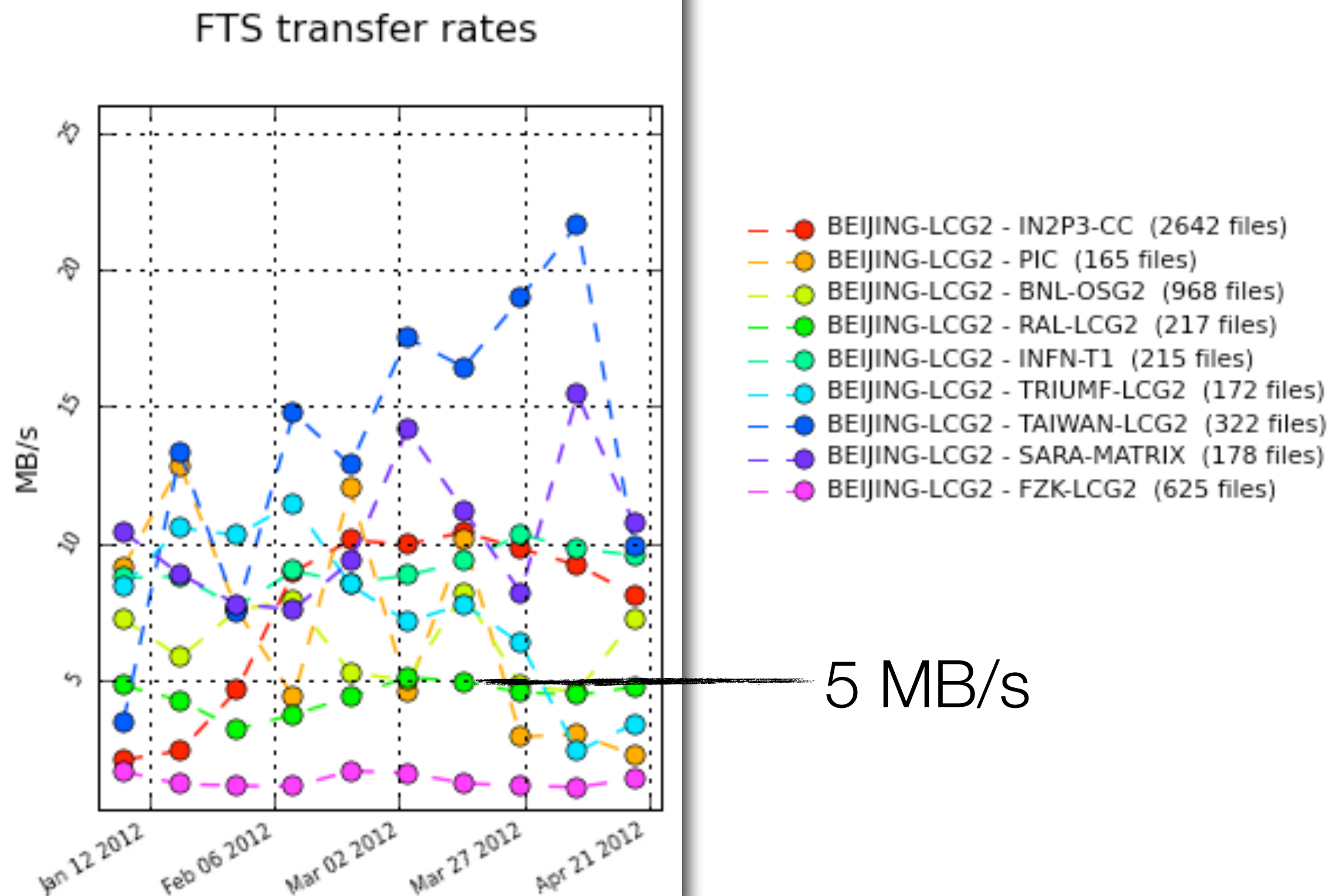
1 month ->

Timezone: GMT+0200 (CEST)

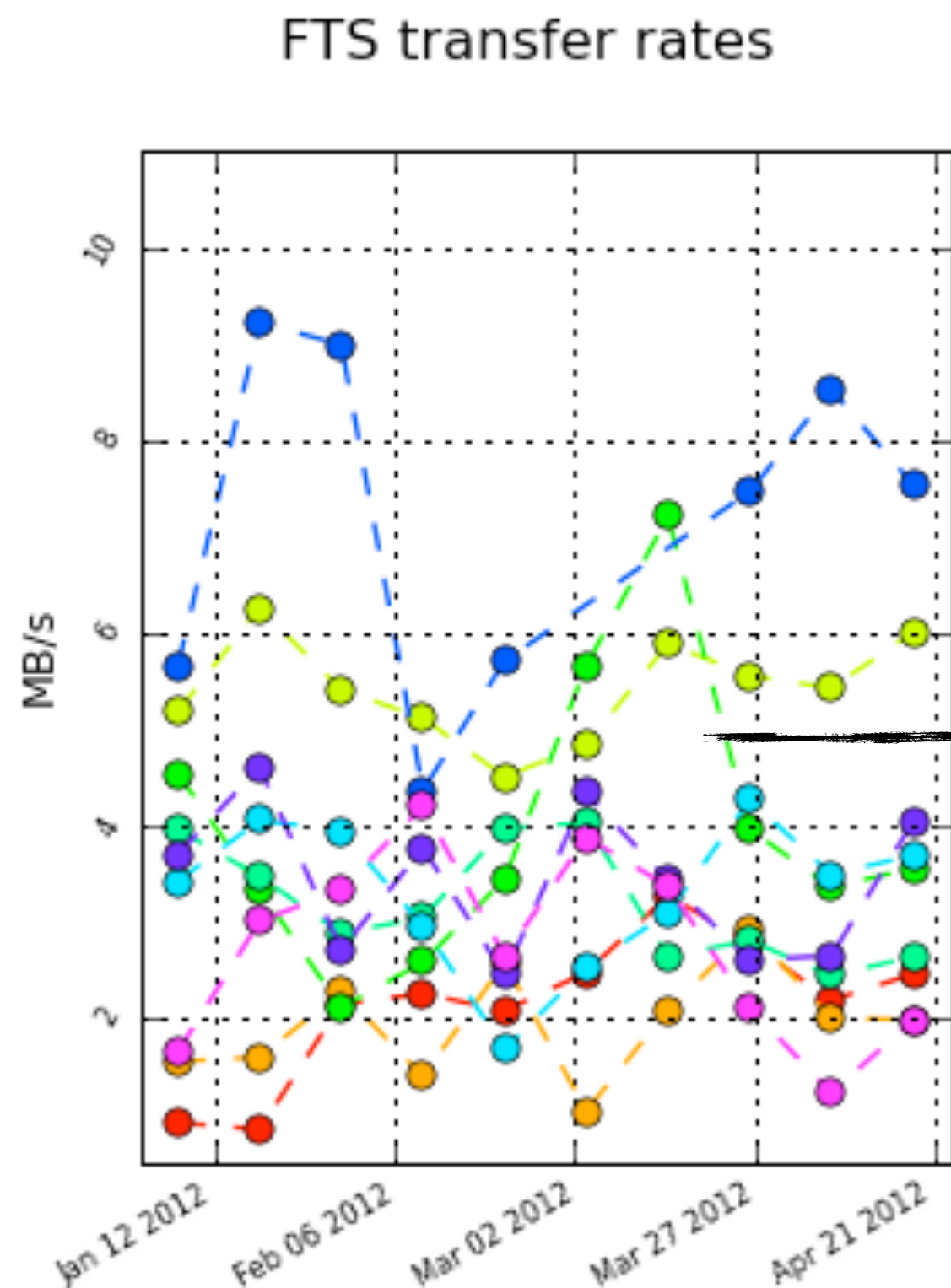
Beijing → CCIN2P3 **CCIN2P3 → Beijing**

Direction	Max throughput(bps)	Mean throughput(bps)	Min throughput(bps)
Src-Dst	52.76M	23.19M	9.49M
Dst-Src	669.33M	616.84M	55.11M

Beijing to T1s



T1s to Beijing



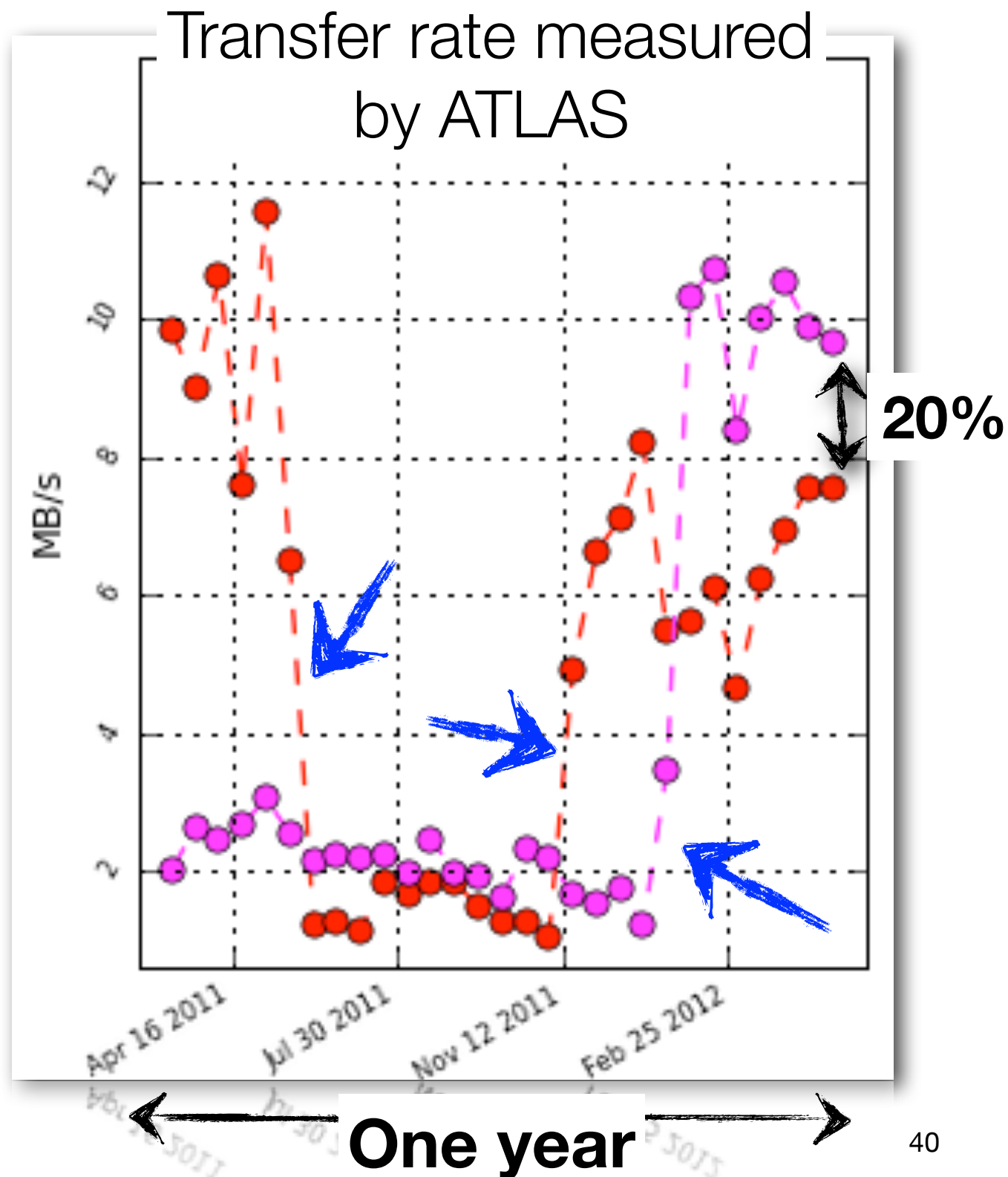
- IN2P3-CC - BEIJING-LCG2 (20899 files)
- PIC - BEIJING-LCG2 (1612 files)
- BNL-OSG2 - BEIJING-LCG2 (6152 files)
- RAL-LCG2 - BEIJING-LCG2 (4580 files)
- INFN-T1 - BEIJING-LCG2 (2619 files)
- TRIUMF-LCG2 - BEIJING-LCG2 (4510 files)
- TAIWAN-LCG2 - BEIJING-LCG2 (3193 files)
- SARA-MATRIX - BEIJING-LCG2 (2656 files)
- FZK-LCG2 - BEIJING-LCG2 (2790 files)

5 MB/s

each site is different

ATLAS transfers to/from Tokyo over one year

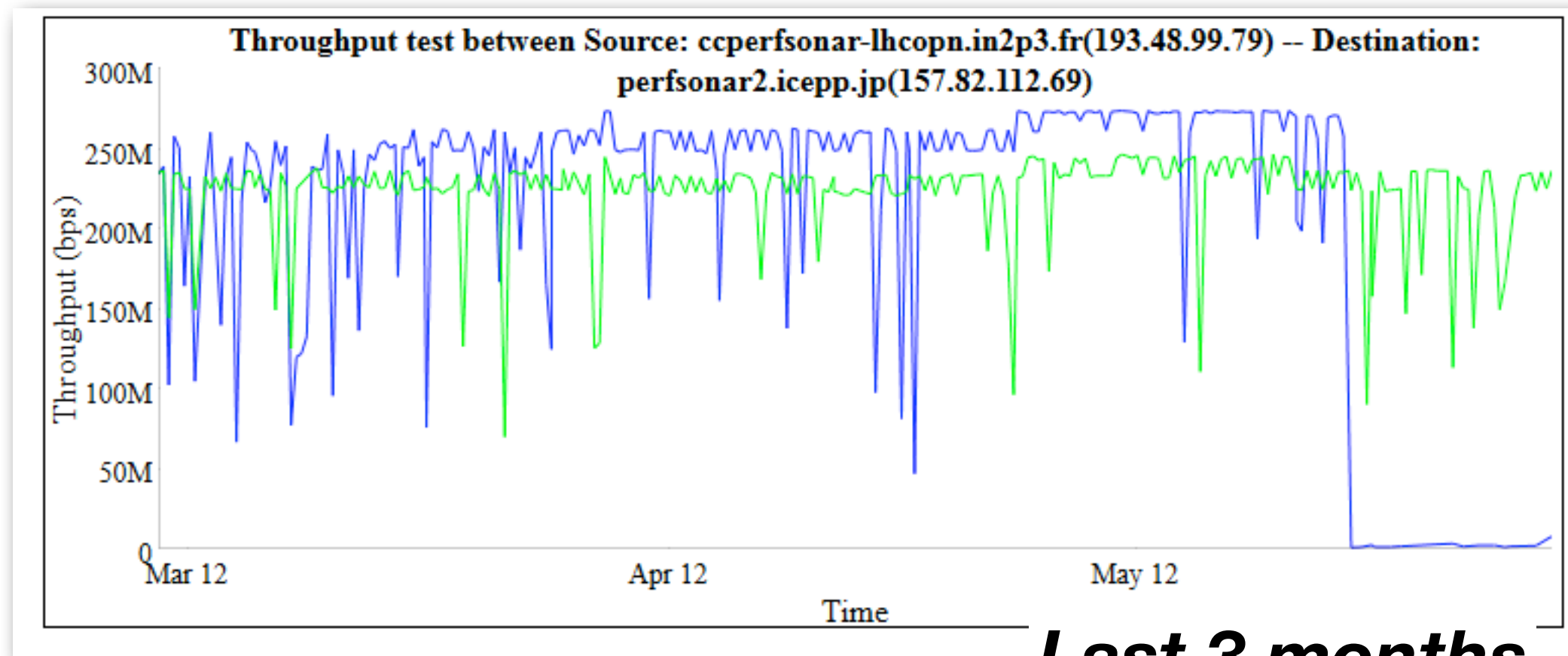
CCIN2P3 → Tokyo
Tokyo → CCIN2P3



Network throughput measured with perfSONAR

CCIN2P3 → Tokyo

Tokyo → CCIN2P3

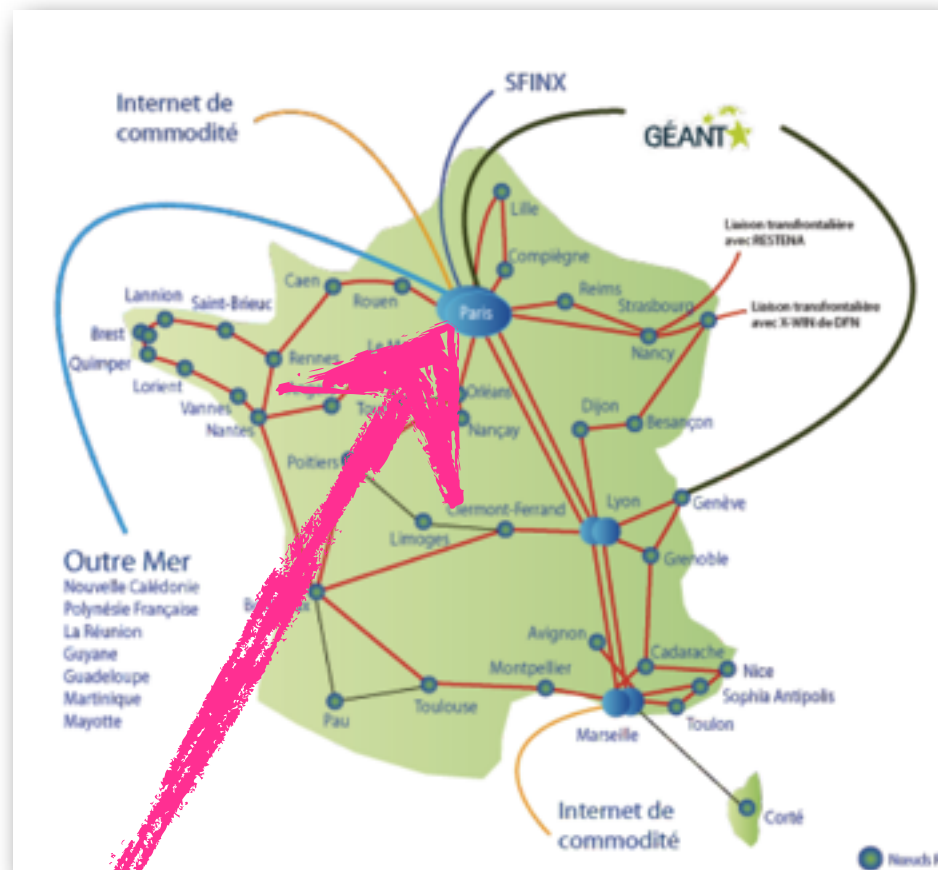


Last 3 months

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

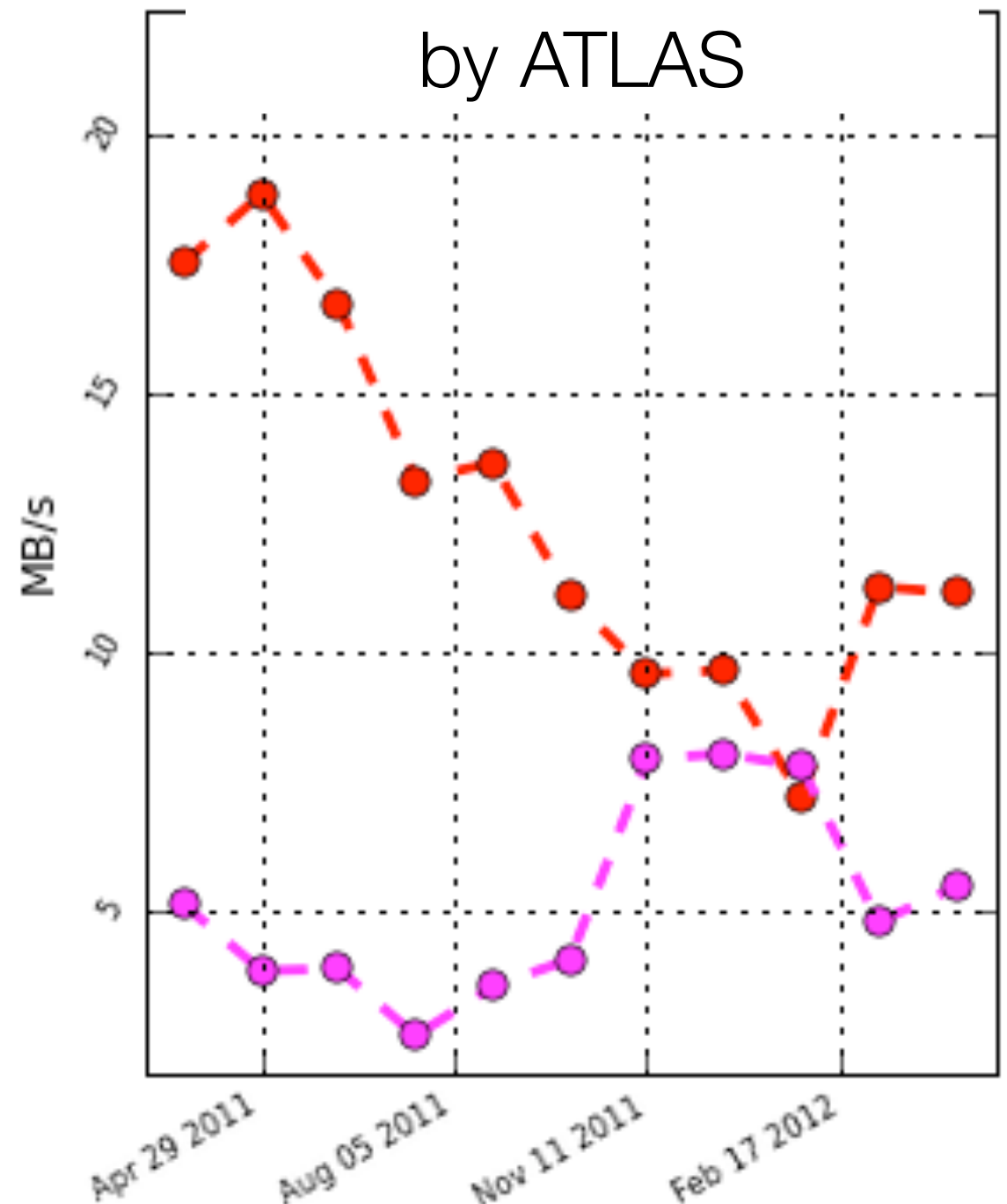
Different pattern for another French site

GRIF-LPNHE → Tokyo
Tokyo → GRIF-LPNHE

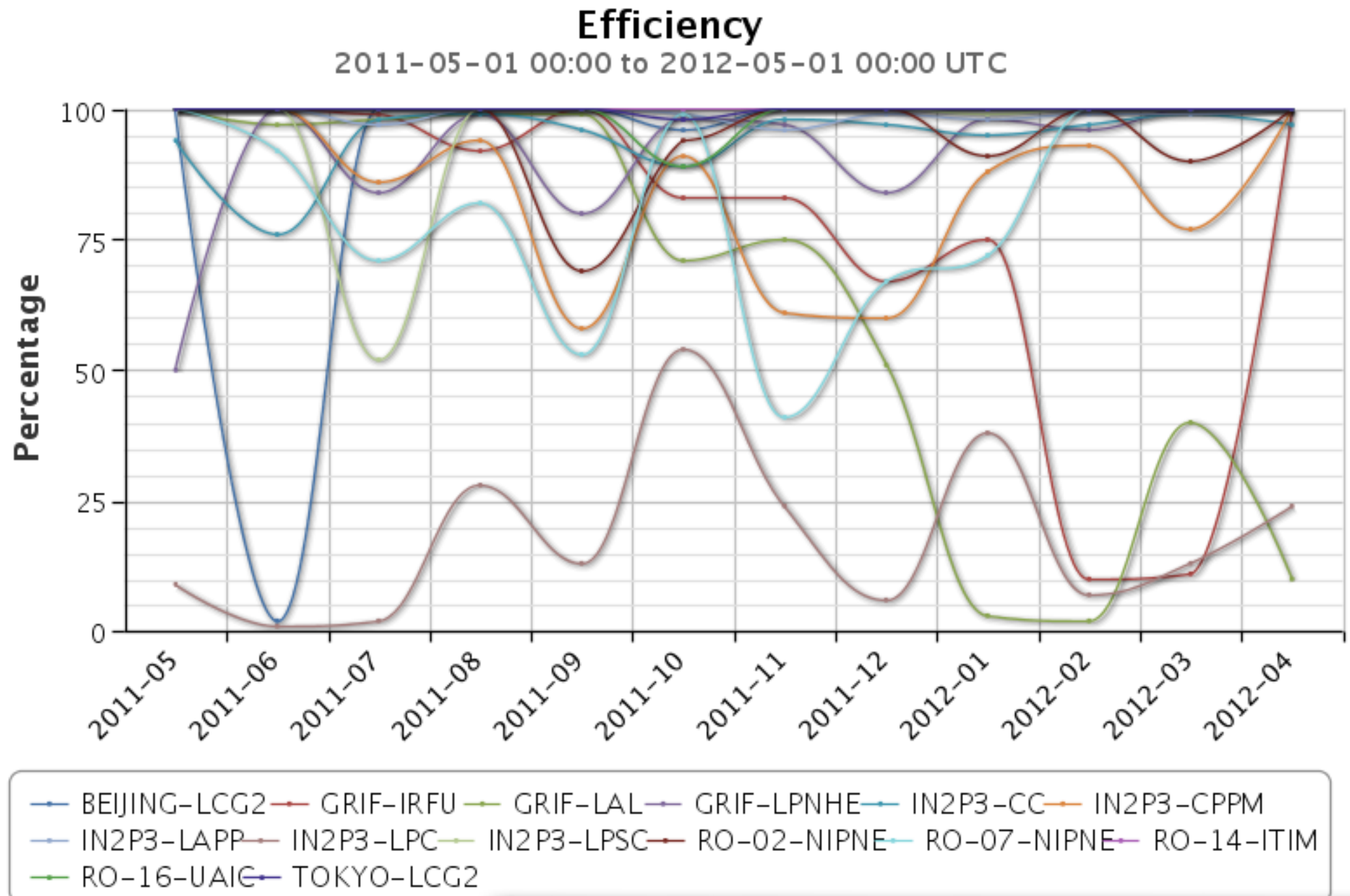


GRIF-LPNHE = Paris

Transfer rate measured
by ATLAS

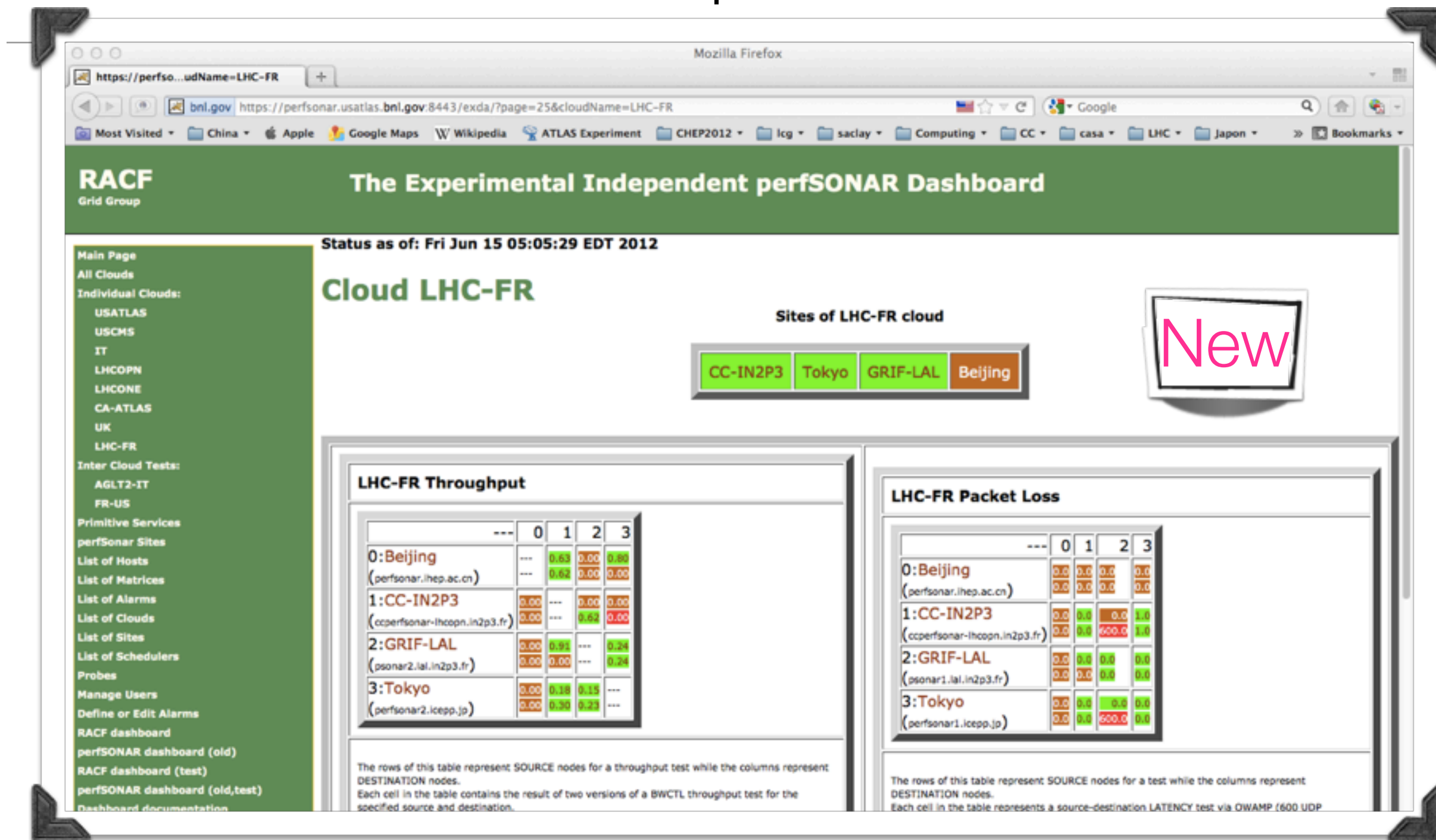


Efficiency of data transfer to Tokyo from FR-cloud sites



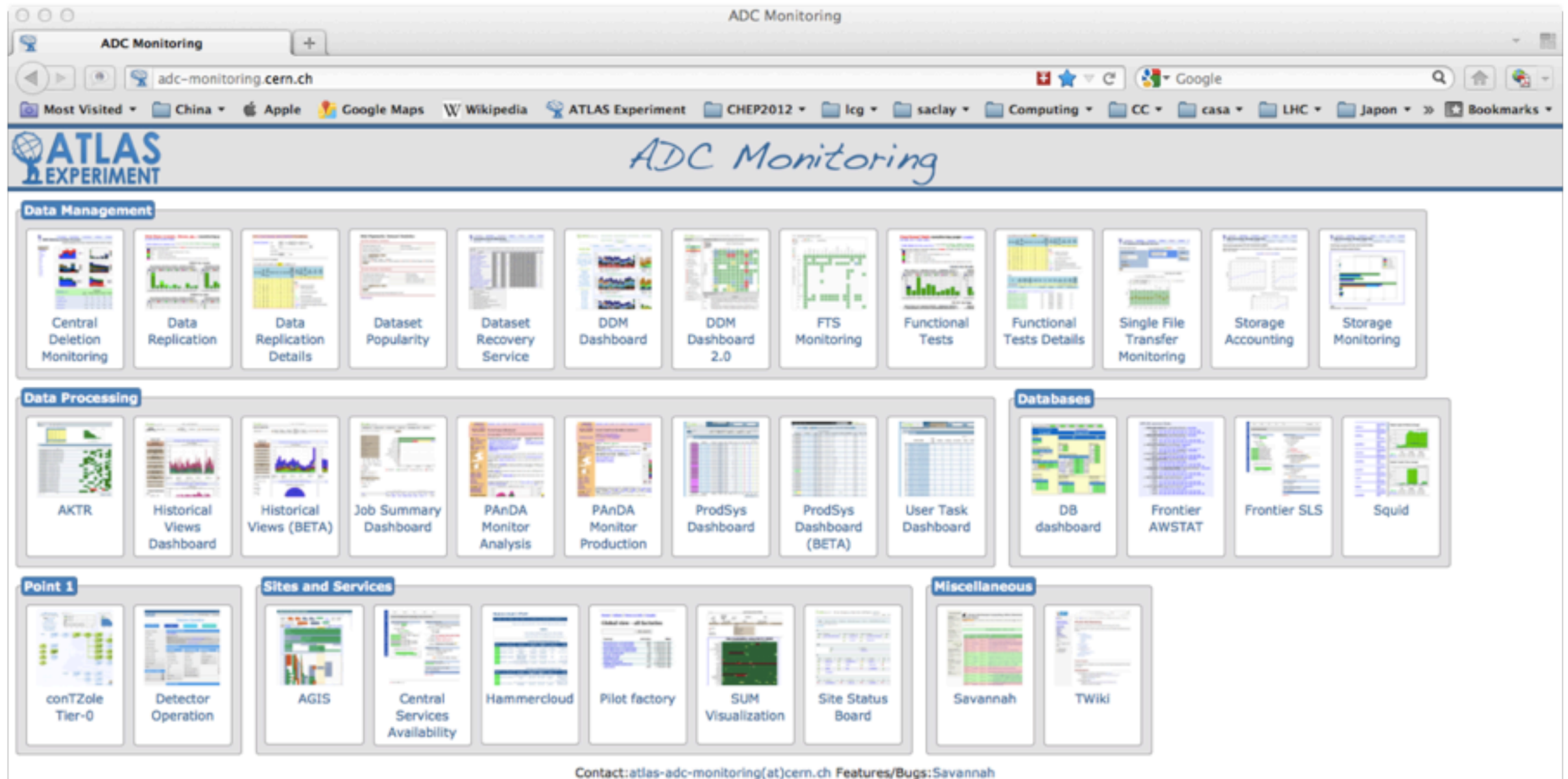
Each site is different, room for improvement everywhere...

FR-cloud view in ATLAS perfSonar dashboard



<https://perfsonar.usatlas.bnl.gov:8443/exda/?page=25&cloudName=LHC-FR>

ATLAS monitoring tools



Lot of new tools over last year

Too many entry points?

FR-Cloud synthesis

Site Status Board

dashb-atlas-ssb.cern.ch/dashboard/request.py/siteview#currentView=Cloud&find[pSC][0][sS]=&find[pSC][0][bR]=false&find[pSC][1][sS]=

Most Visited China Apple Google Maps Wikipedia ATLAS Experiment CHEP2012 lcg saclay Computing CC casa LHC Japon Bookmarks

Help Login Site Status for the ATLAS sites, v0.1.0_rc42

Index Expanded Table

Show 200 entries Copy Print Save views Cloud Search...

Site Name	Site Info		Downtime	DDM DT - Status	Panda Analysis status	Panda Production status	Panda Efficiency						SRM SAM 12 [%]	CE SAM 12	Functional Tests			SW releases - critical	GGUS	
	Tier	Cloud					Analy Activated Jobs	Analy Running Jobs	Analy Efficiency 12h [%]	Prod Activated Jobs	Prod Running Jobs	Prod Efficiency 12h [%]			Squid Functional Tests (SAM)	HC_AFT	HC_PFT		GGUS open	GGUS closed
BEIJING-LOG2	T2D	FR	ACTIVE	online	test	brokeroff	48	82	73	8	81	78	100	100	n/a	0	0	117/154	0-0-0-0-0	n/a
GRIF-IRFU	T2	FR	UNKNOWN	online	online	online	128	1003	86	8	176	100	n/a	n/a	n/a	100	no-test	170/170	0-0-0-0-0	0-0-0-0-0
GRIF-LAL	T2D	FR	UNKNOWN	online	online	online	1200	94	80	2075	566	90	n/a	n/a	n/a	100	100	169/170	0-0-0-0-0	0-0-0-0-0
GRIF-LPNHE	T2D	FR	UNKNOWN	online	online	online	126	326	86	303	326	88	n/a	n/a	n/a	100	100	170/170	0-0-0-0-0	0-0-0-0-0
IN2P3-CC	T1	FR	ACTIVE	online	online	online	21019	225	87	4742	4343	84	100	100	n/a	no-test	100	166/170	0-0-0-0-0	n/a
IN2P3-CC-T2	T2	FR	ACTIVE	online	online	online	28230	549	83	857	518	88	n/a	94	n/a	100	100	163/170	0-0-0-0-0	82418
IN2P3-CPPM	T2	FR	ACTIVE	online	online	online	0	48	100	3	463	85	100	100	n/a	100	100	166/170	0-0-0-0-0	82086
IN2P3-LAPP	T2D	FR	ACTIVE	online	online	online	36	122	86	808	282	87	100	100	n/a	100	100	167/170	0-0-0-0-0	0-0-0-0-0
IN2P3-LPC	T2D	FR	ACTIVE	online	online	online	800	208	84	138	258	80	100	89	n/a	100	100	169/170	82830	0-0-0-0-0
IN2P3-LPSC	T2D	FR	ACTIVE	online	online	online	405	146	86	304	286	82	100	100	n/a	100	100	168/170	0-0-0-0-0	0-0-0-0-0
RO-02-NIPNE	T2	FR	ACTIVE	online	online	online	122	79	85	26	131	87	100	94	n/a	100	100	129/154	0-0-0-0-0	n/a
RO-07-NIPNE	T2	FR	ACTIVE	online	online	online	133	113	87	840	818	88	100	100	n/a	100	100	170/170	0-0-0-0-0	n/a
RO-14-ITIM	T2	FR	ACTIVE	online	NoQueue	online	no data	no data	no data	165	208	100	100	100	n/a	no-test	100	120/154	0-0-0-0-0	82388
RO-16-UIAC	T2	FR	ACTIVE	online	NoQueue	online	no data	no data	no data	301	387	88	100	100	n/a	no-test	100	124/154	0-0-0-0-0	0-0-0-0-0
TOKYO-LOG2	T2	FR	ACTIVE	online	online	online	0	26	86	2041	2019	88	100	100	n/a	100	100	170/170	0-0-0-0-0	82366

Showing 1 to 15 of 15 entries DB query took 0.0653 s

First Previous 1 Next Last

CVMFS

Deployed on every site of FR-cloud?

- CVMFS (CERN Virtual Machine File System)
- Used by ATLAS (<https://twiki.cern.ch/twiki/bin/view/Atlas/CernVMFS>) for
 - software distribution
 - conditions data distribution
- Advantages
 - no need to worry about increasing size of NFS software area and bottleneck of shared area
 - increase of CPU efficiency
- Ongoing deployment on ATLAS sites ([http://dashb-atlas-ssb.cern.ch/dashboard/request.py/siteview#aaSorting\[0\]\[\]=1&aaSorting\[0\]\[\]=asc¤tView=cvmfs&highlight=false](http://dashb-atlas-ssb.cern.ch/dashboard/request.py/siteview#aaSorting[0][]=1&aaSorting[0][]=asc¤tView=cvmfs&highlight=false))

Xrootd federations



Use Case #2 Sharing storage amongst nearby Tier 2 sites



- In US - The 5 sites in ANALY_AGLT2 & ANALY_MWT2 are all within 7 ms RTT
- Midwest Tier 2 (MWT2) internal federation amongst 3 sub sites
- AGLT2 has storage federated between two sites
- ~ 4.4 PetaBytes of storage amongst sites
- Direct Read mode
- Most efficient for jobs that read part of data file (like Analysis jobs)



Applicable for some FR-cloud sites?

Multi-core processing

Panda Production Operations D...

+

panda.cern.ch:25880/server/pandamon/query?dash=prod

Google

Q

Most Visited

China

Apple

Google Maps

Wikipedia

ATLAS Experiment

CHEP2012

lcg

sacay

Computing

CC

casa

LHC

Japon

Bookmarks

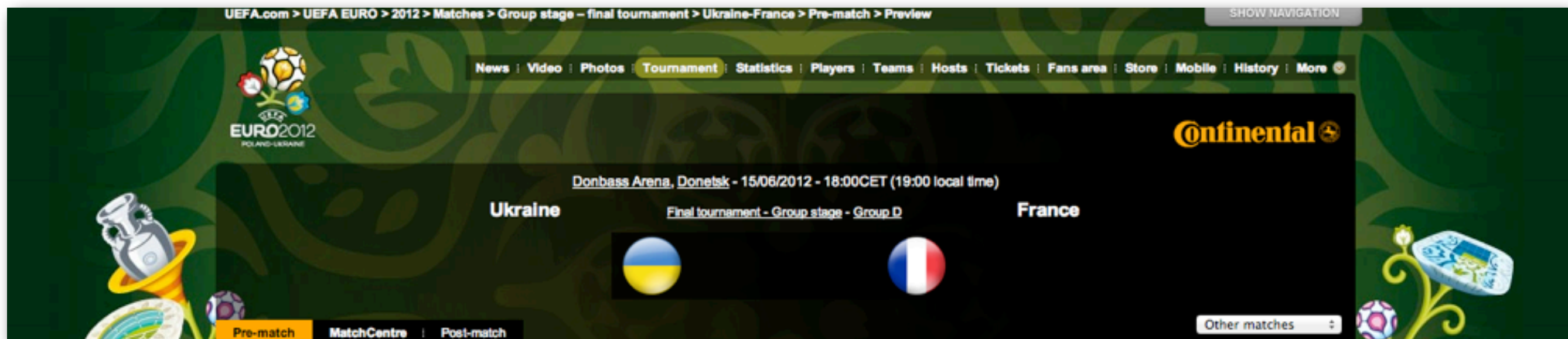
FR	3570	06-15 10:01	0	18	0	2326	18403	0	27	11805	342	7256	23271	3650	36	14%	
FR sites	Pilots	Latest	pending	defined	waiting	assigned	activated	sent	starting	running	holding	transferring	finished	failed	cancelled	cores	%fail
ALL			0	18	0	2326	18403	0	27	11805	342	7256	23271	3650	36	14%	
BEIJING	84	06-15 10:01	0	0	0	0	5	0	0	58	3	206	583	161	0	22%	
BEIJING_MCORE (brokeroff)			0	0	0	0	0	0	0	0	0	0	0	0	0	8	
CPPM	127	06-15 10:01	0	2	0	147	911	0	0	456	7	461	1114	81	2	7%	
GRIF-IRFU	188	06-15 10:01	0	1	0	150	679	0	0	475	2	69	3164	9	2	0%	
GRIF-LAL	152	06-15 10:01	0	1	0	0	2154	0	0	386	23	698	433	40	4	8%	
GRIF-LPNHE	95	06-15 10:01	0	0	0	9	506	0	0	360	7	404	872	64	1	7%	
IN2P3-CC	746	06-15 10:01	0	1	0	0	2902	0	0	2015	37	0	4011	1642	1	20%	
IN2P3-CC-T2	223	06-15 10:01	0	0	0	0	573	0	0	315	3	0	899	16	2	2%	
IN2P3-CC_MCORE (brokeroff)	1	06-15 10:01	0	0	0	0	0	0	0	0	0	0	0	0	0	16	
IN2P3-CC_SGE_VL	751	06-15 10:01	0	2	0	0	3437	0	0	2021	121	794	2089	1325	1	36%	
IN2P3-LPSC	99	06-15 10:01	0	0	0	147	332	0	0	394	23	219	785	57	3	7%	
LAPP	114	06-15 10:01	0	0	0	3	838	0	0	410	6	416	385	20	6	5%	
LPC	126	06-15 10:01	0	2	0	32	206	0	26	334	8	373	795	90	2	100%	
ROMANIA02	39	06-15 10:01	0	1	0	63	172	0	0	136	5	166	486	11			
ROMANIA07	114	06-15 10:01	0	1	0	1007	554	0	0	965	14	875	1608	44			
ROMANIA14	34	06-15 10:01	0	2	0	42	242	0	0	206	9	134	558	2			
ROMANIA16	66	06-15 10:01	0	2	0	168	388	0	1	387	14	334	1275	6			
TOKYO	152	06-15 10:01	0	1	0	1	2100	0	0	929	9	398	1629	81			
TOKYO-extra	459	06-15 10:01	0	2	0	557	2404	0	0	1958	51	1709	2585	1	2	0%	
IT	2040	06-15 10:01	0	17	0	159	13860	0	25	6007	385	2660	26217	1307	36	5%	

On test

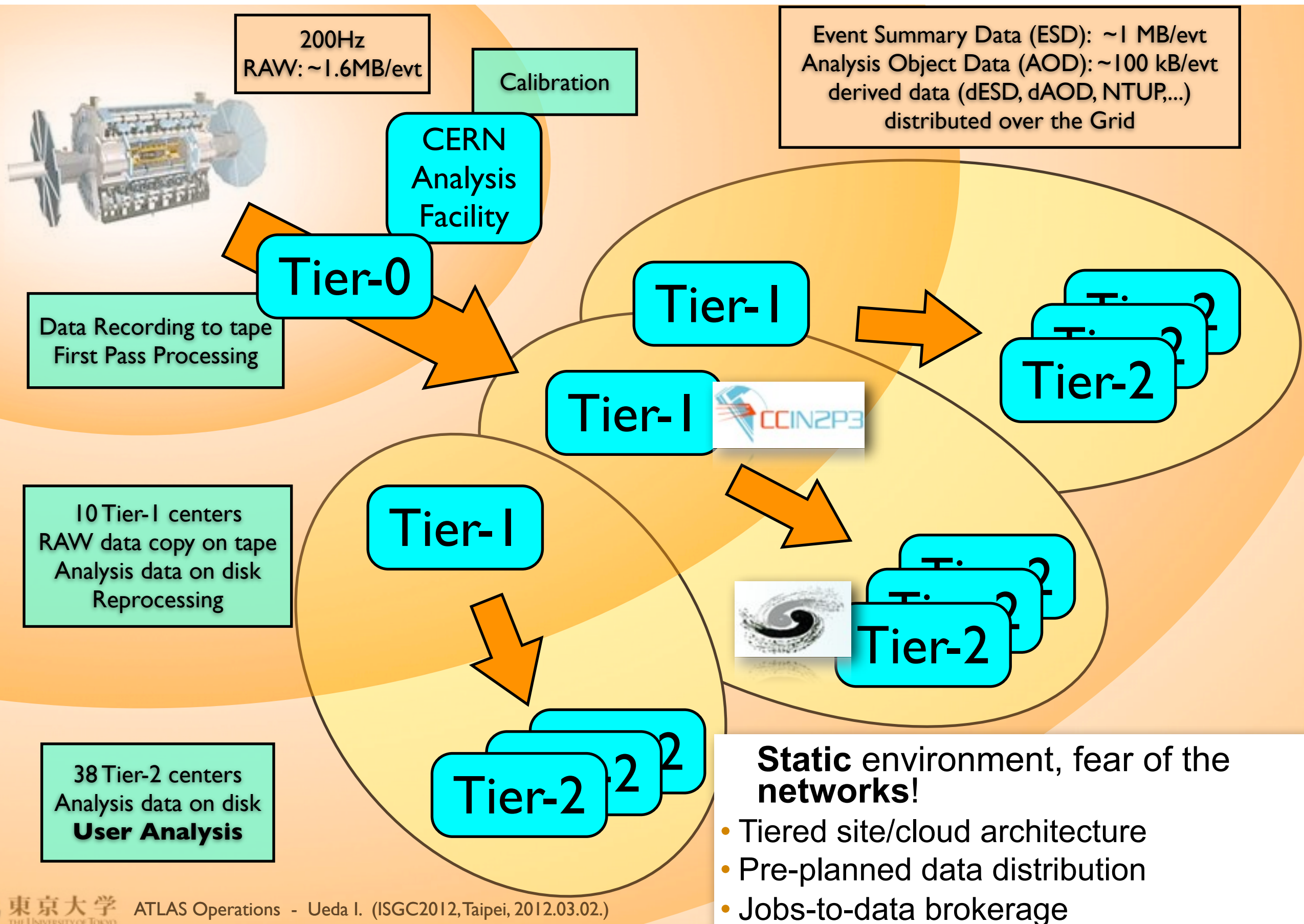
widely used in 2013?

On test
widely used in 2013?

And now...



ATLAS Computing Model: T0 Data Flow



ATLAS Computing Model: MC Data Flow

