

# Local HPC Center usage by Alice : from initial ideas to production

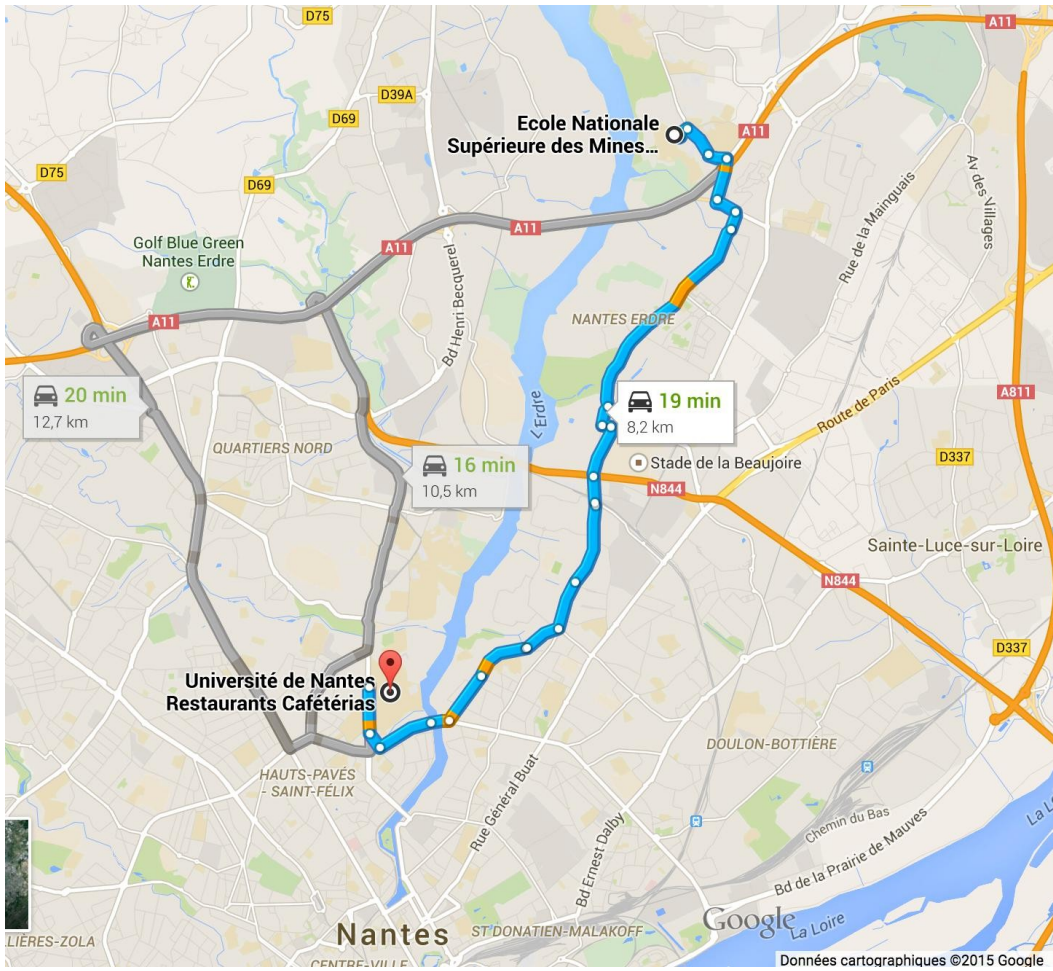
Jean-Michel BARBET, SUBATECH, Nantes

Alice T1/T2 Workshop Derby UK, April 2018

# How it started

- End of 2012, we were made to understand that a fund hunting at the regional level would only be successful if computing resources were to be shared in a common facility
- A joint venture was started with a local HPC Center at the University of Nantes : the CCIPL [1]
- The resources would be hosted by the brand new datacenter at the University

# The CCIPL HPC and the University Datacenter



- Broad range of scientific domains : Chemistry, Physics,...
- 2 clusters, 3440 CPU cores (Sept.2017) CentOS7.3 64bits
- 22 million CPU hours a year
- Omnipath 56Gbits/s interconnect
- 300 TB Fast storage BeeGeeFS
- Staff : 1.8 FTE (1 University, 0.8 CNRS)
- Close to the NREN PoP
- Our share is 200 cores (400 HT)

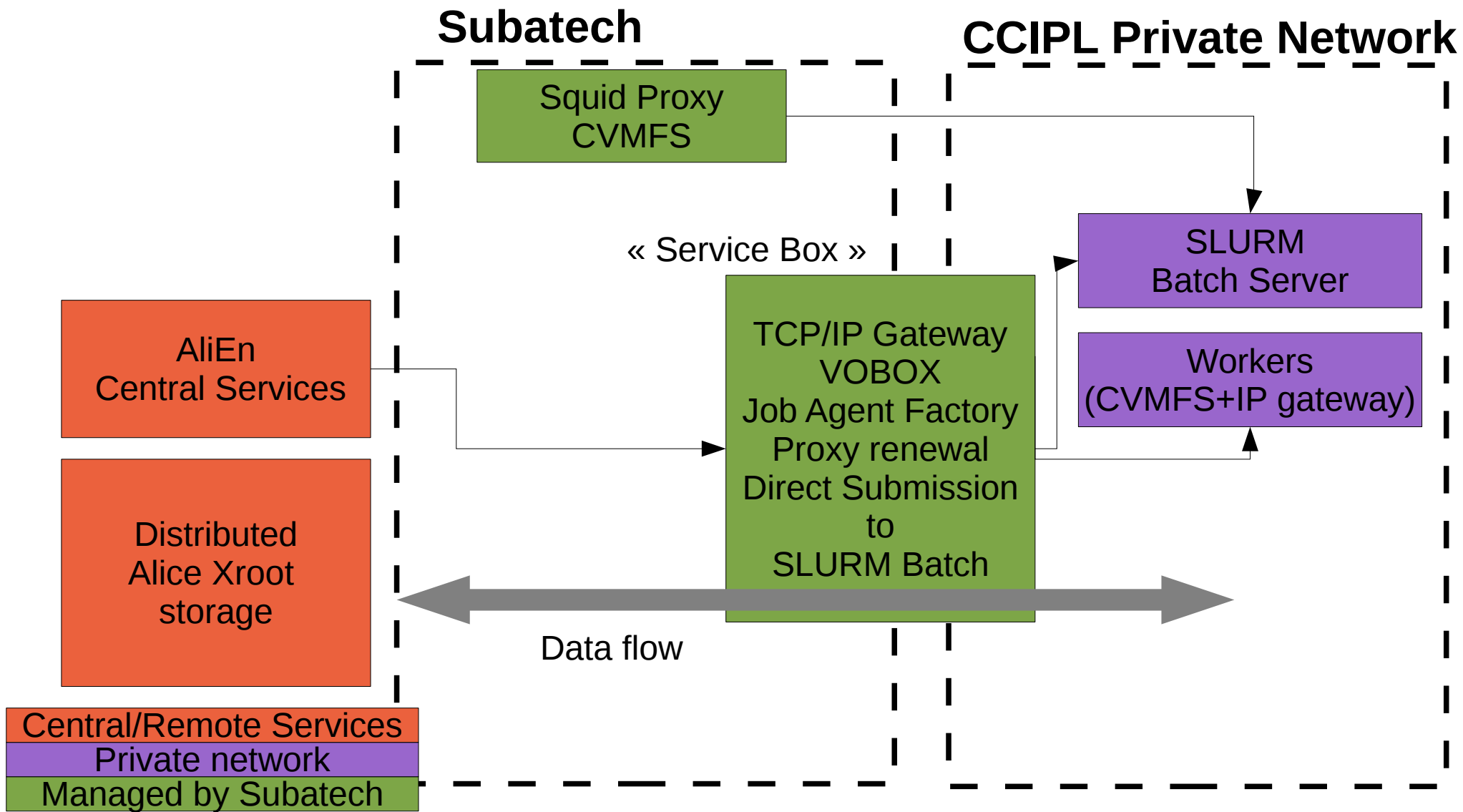
# Our requirements

- As a partner, we had the opportunity to state our requirements and we asked for :
  - A « service box » with 2x10Gbits/s network adapt.
  - External network access via the service box (gateway)
  - CVMFS operational on the worker nodes
  - Local disk on the WN (CVMFS cache + ...)
  - Enough RAM/core to run Alice jobs
- We managed to have the storage installed at Subatech (+1PB under EOS in Nov 2016)

# Timeline

- Dec 2013 : Creation of COS (steering comitee) of the CCIPL : application for funding
- Mar 2014 : Presentation[2] at Tsukuba Alice T1/T2 workshop
- Dec 2014 : Funding secured (CPER-2014-2020) : 300K€ for Subatech (storage+CPU)
- Jan 2015 : Work on specifications for the procurement bid started
- Feb 2015 : Presentation [3] at Torino Alice T1/T2 workshop
- Jun 2015 : First demonstrator in Subatech (2 WNs in a private network with Torque)
- Dec 2015 : Specifications ready for the procurement bid
- May 2016 : 2nd demonstrator in Subatech (change batch system from torque to HT-condor)
- Jun 2016 : Bid opening, best bidder selected (Slurm, CentOS7, Trinity,...)
- Aug 2016 : 3rd demonstrator in Subatech (change batch system from HT-condor to SLURM)
- Nov 2016 : Acceptance tests, validation, benches
- Dec 2016 : Network link between Subatech and CCIPL
- Mar 2017 : Installation of the « vobox » also being a network gateway
- Mar 2017 : The new CCIPL HPC cluster is officially opened to users
- May 2017 : CVMFS available on Wns, greenlight to start with Alice production,  
reconfigure the AliEn site Subatech\_CCIPL to use the vobox/gateway.
- 30 May 2017 : First jobs running and ...done !**
- 1st April 2018 : WLCG Accounting starts !

# Current setup



# ALICE setup

- Activities limited to MC production (user « aliproduct »)

**cerequirements**

```
other.user=="aliproduct" || other.user=="barbet"
```

- SLURM commands in AliEn LDAP :

<b>statusarg</b>	-a
<b>statuscmd</b>	queue
<b>submitarg</b>	--partition=Subatech --nodes=1 --ntasks-per-core
<b>submitcmd</b>	sbatch
<b>TTL</b>	172800
<b>type</b>	SLURM <small>required</small>
<b>killcmd</b>	scancel

# Local monitoring of the Alice-box

**Current Network Status**  
 Last Updated: Fri Mar 30 14:01:24 CEST 2018  
 Updated every 90 seconds  
 Nagios® Core™ 3.2.3 - [www.nagios.org](http://www.nagios.org)  
 Logged in as *jean-michel*

[View History For This Host](#)  
[View Notifications For This Host](#)  
[View Service Status Detail For All Hosts](#)

### Host Status Totals

Up	Down	Unreachable	Pending
1	0	0	0
All Problems		All Types	
0		1	

### Service Status Totals

Ok	Warning	Unknown	Critical	Pending
12	1	0	0	0
All Problems		All Types		
1		13		

### Service Status Details For Host 'nanlcg12'

Host ↑↓	Service ↑↓	Status ↑↓	Last Check ↑↓	Duration ↑↓	Attempt ↑↓	Status Information
<a href="#">nanlcg12</a>	<a href="#">ALIEN-CE-PROXY</a>	WARNING	03-30-2018 13:58:06	0d 5h 18m 18s	1/1	SERVICE STATUS: WARNING Proxy Lifetime 110657 secondes
	<a href="#">CVMFS-ALICE</a>	OK	03-30-2018 13:57:16	52d 12h 54m 22s	1/1	SERVICE STATUS: OK CacheUsed:41% FileDescUsed:0%
	<a href="#">NTP</a>	OK	03-30-2018 13:55:35	219d 5h 51m 13s	1/4	NTP OK: Offset -0.0006240362418 secs
	<a href="#">Network-load</a>	OK	03-30-2018 13:57:13	52d 9h 35m 39s	1/3	OK: ens2f0 10000Mb/s link up RX/TX=33.59/1382.26 kb/s
	<a href="#">SSH</a>	OK	03-30-2018 13:58:06	52d 3h 48m 33s	1/3	SSH OK - OpenSSH_6.6.1 (protocol 2.0)
	<a href="#">Unix CPU</a>	OK	03-30-2018 13:58:15	52d 5h 51m 19s	1/3	OK: CPU is 0% used
	<a href="#">Unix DISK</a>	OK	03-30-2018 13:59:44	52d 3h 36m 53s	1/3	DISK OK - free space: / 30629 MB (59% inode=99%):
	<a href="#">Unix LOAD</a>	OK	03-30-2018 13:58:24	302d 21h 6m 39s	1/3	OK - Charge moyenne: 0.16, 0.19, 0.19
	<a href="#">Unix PROCS</a>	OK	03-30-2018 13:58:14	52d 10h 48m 30s	1/3	PROCS OK: 458 processus
	<a href="#">Unix RAM</a>	OK	03-30-2018 13:57:23	52d 4h 14m 10s	1/3	OK: 0% Used Memory - Total: 128661 MB, used: 0 MB, free: 128661 MB
	<a href="#">Unix TRIPWIRE</a>	OK	03-30-2018 13:59:33	1d 10h 21m 53s	1/3	Tripwire OK - violation found: 0
	<a href="#">Unix ZOMBIE</a>	OK	03-30-2018 13:57:14	52d 12h 51m 20s	1/3	PROCS OK: 0 processus avec ETAT = Z
	<a href="#">clustertier2running</a>	OK	03-30-2018 13:59:24	1d 17h 52m 0s	1/3	Running jobs on CCIPL OK - Working: 400 : Queued: 20



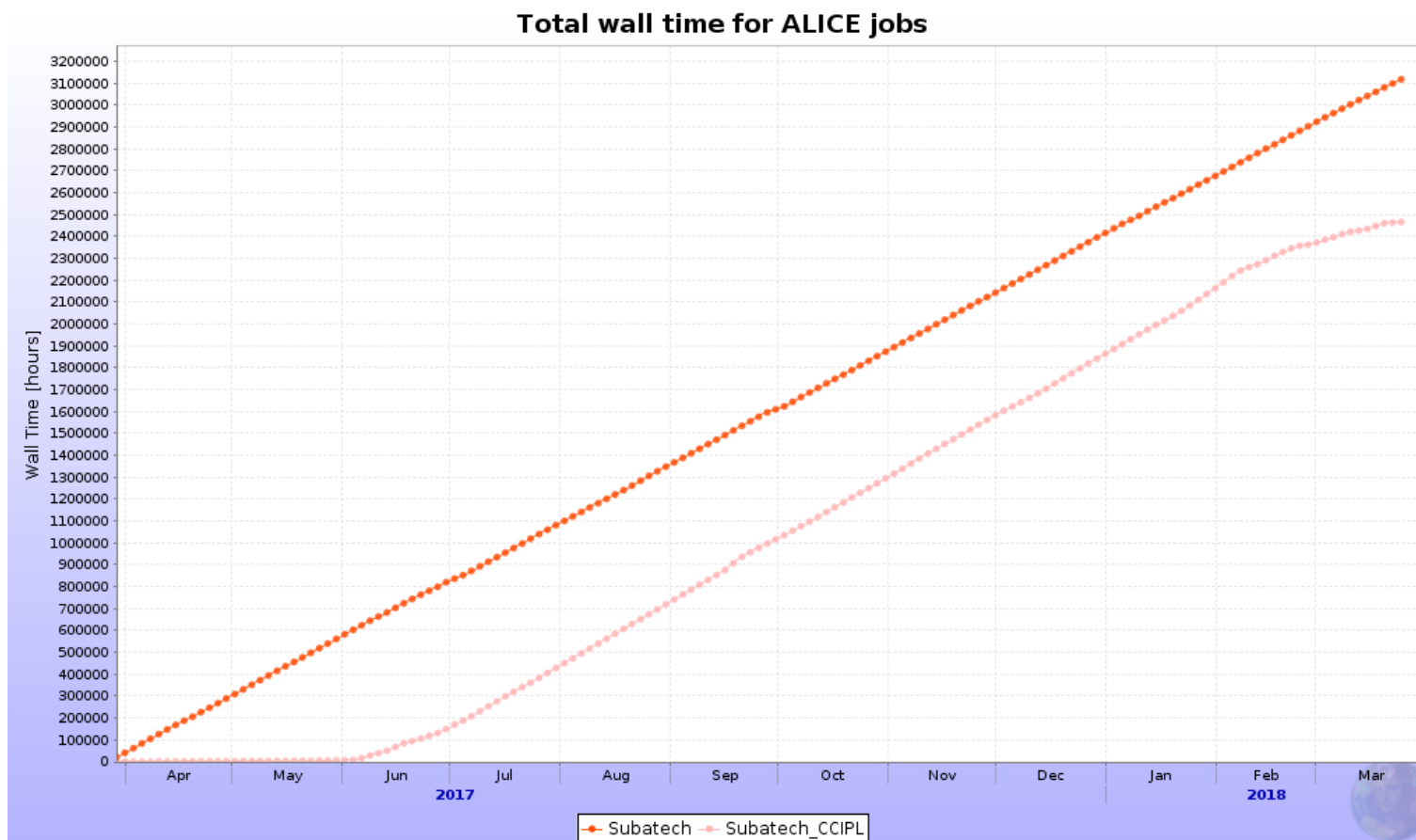
# Local monitoring with Nagios



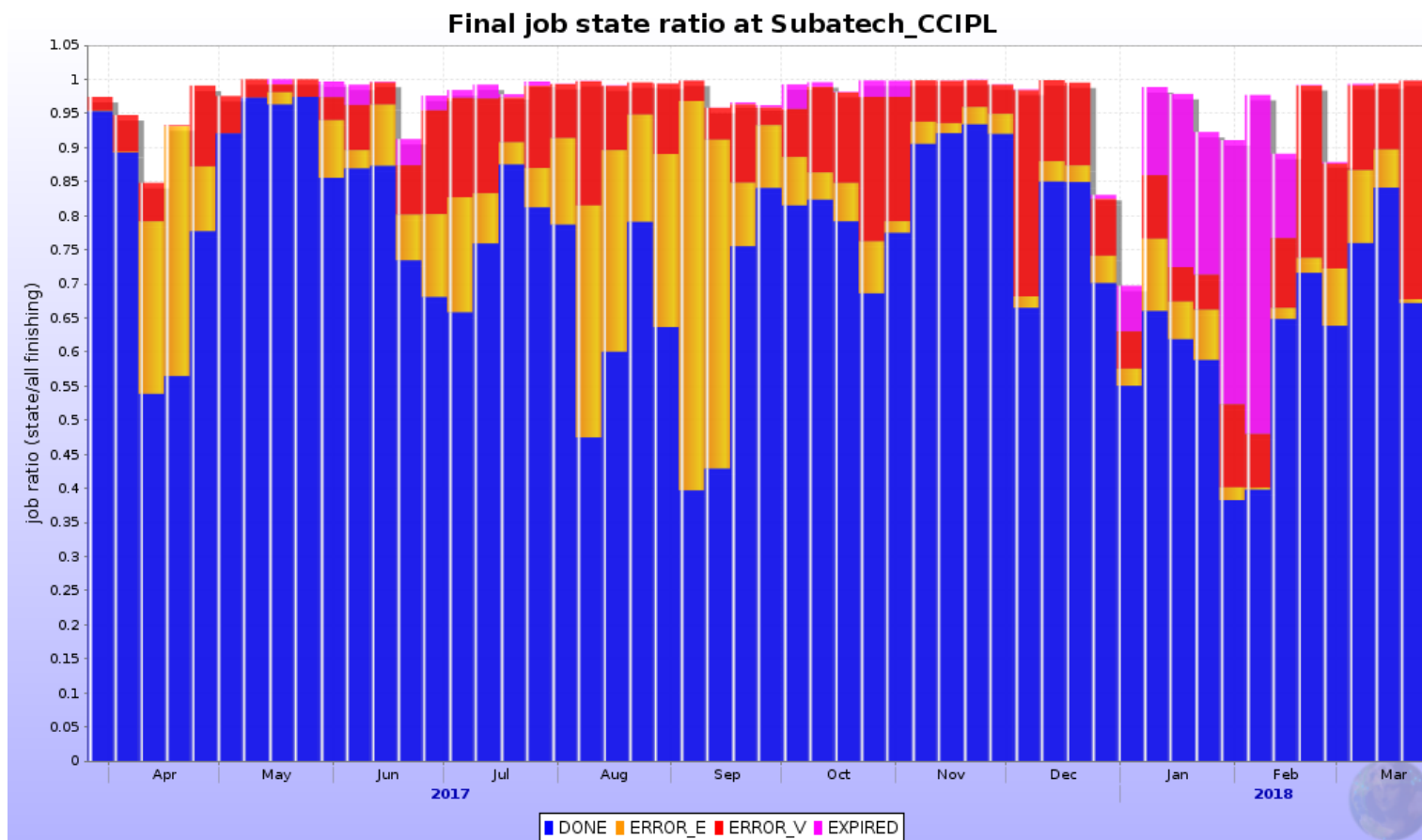
# Solved issues

- Overload of the CCPIIL network filesystem (due to the high number of small files)
  - BeeGeeFS for /home
  - NFS for /scratch
  - => Alien WORK\_DIR = /tmp/ALICE
- Forwarding of network traffic through the Alice box is restricted (firewalld). This requires identification of network flows from the Wns
  - Dcache port numbers not easy to guess

# Some graphs from Monalisa



# Some graphs from Monalisa



# WLCG Accounting 1/3

- We were forced to provision in an HPC-Center
- We intend to pledge these resources to have this contribution recognised
- For this we need to publish accounting data in the WLCG accounting database
- No on the shelf solution, but part of recent work on publishing from an ARC-CE could be reused

# WLCG Accounting 2/3

- F.Schaer from CEA-IRFU France wrote a script to send ARC-CE UsageRecords to the national MonBox using APEL loader (using MQ).
- If we can create similar UsageRecords from SLURM, the same setting can be reused
- Shell script to generate previous day accounting data from SLURM sacct
- Python script to create xml UsageRecords
- APEL loader sends data to IRFU MonBox
- Data is sent to APEL repository by the MonBox

# WLCG Accounting 3/3

- Started to publish accounting data this way on April 1st (corresponds to the start of a pledge period)
- Problem : Collecting correct CPUtime !
  - Is this a SLURM issue ?
  - Or is it related to how the batch shell is built by AliEn (SLURM.pm) ?
- How to check in APEL ?
  - By DN ? (Resource center view, restricted)

# Possible enhancements

- Opportunistic usage of free WNs outside our partition
  - Would require configuring also other worker nodes with CVMFS + the alicebox as a network gateway
  - Jobs could be preempted any time
  - Is it worth the effort and potential issues ?



# Security : risks

- Now we have a server physically located in another site
- Sharing worker nodes with local users
- We are not responsible for applying security fixes

# Security : measures

- Establish minimum Incident response :
  - Contacts, emergency procedures
- Trust and close relations with HPC admins
- Basic risk analysis :
  - Understand the risk scenarii and have security measures to make them acceptable.

# Références

- [1] CCIPL : Centre de Calcul Intensif des Pays de la Loire  
<http://www.ccipl.univ-nantes.fr/>
- [2] Opportunistic usage of local computing resources : ideas to discuss  
Alice T1/T2 Workshop, Tsukuba, Mar 2014, JM Barbet  
<https://indico.cern.ch/event/274974/contribution/95/material/slides/0.pdf>
- [3] Collaboration project with HPC cluster in Nantes  
Alice T1/T2 Workshop, Turin, Fev 2015, L. Aphecetche  
<https://indico.cern.ch/event/354209/contribution/34/material/slides/0.pdf>
- [4] Mutualisation de ressources CPU dans un centre de calcul régional  
LCG-France Meeting, IPNO, Orsay, June 2015, JM Barbet

# Thank you

## Credits :

Yann Dupont administrator of the CCIPL and the IT services at Nantes University  
Khalil Chawoshi (head) and the Subatech local IT team  
Jérôme Bernier and the CCIN2P3 Network team  
Frédéric Schaer CEA-IRFU for parts of the accounting solution  
Laurent Aphecette Subatech  
Miguel, Costin, the Alice team and people on the alice-lcg-task-force