



**Consortium des Équipements de Calcul Intensif**  
Funded by F.R.S.-FNRS

[www.ceci-hpc.be](http://www.ceci-hpc.be)



# The new CÉCI common storage

brought to you by Thomas Keutgen, David Colignon, Juan Cabrera, Frédéric Wautelet, Raphael Leplae, Bernard Van Renterghem, Sébastien Skozlowskij and Damien François

2017 CECI scientific day - Louvain-la-Neuve  
[damien.francois@uclouvain.be](mailto:damien.francois@uclouvain.be)  
[www.uclouvain.be/cism](http://www.uclouvain.be/cism)

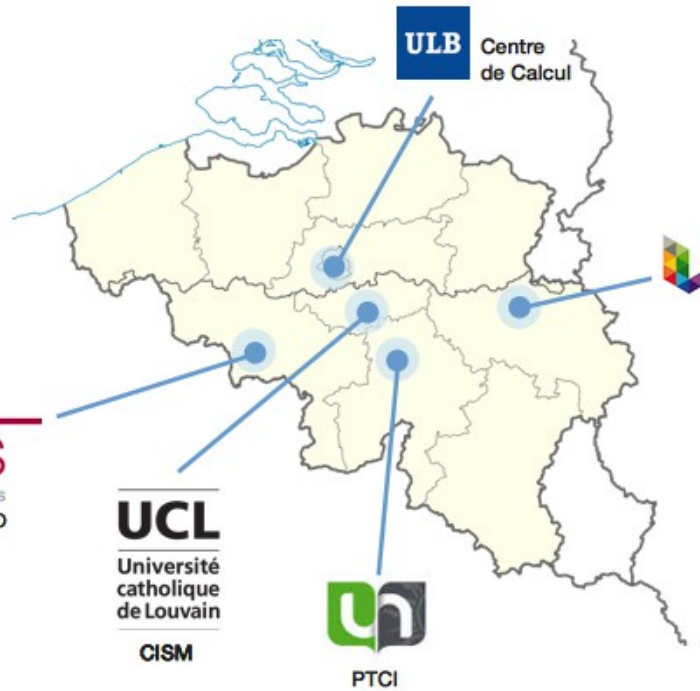


# The new CÉCI common storage Features



6 clusters  
5 sites

Vega



NIC4



Dragon1



Lemaitre2



Hmem



Hercules

new



10Gbps, dedicated, optical network direct links between the sites

# 1. New network

Vega



Dragon1



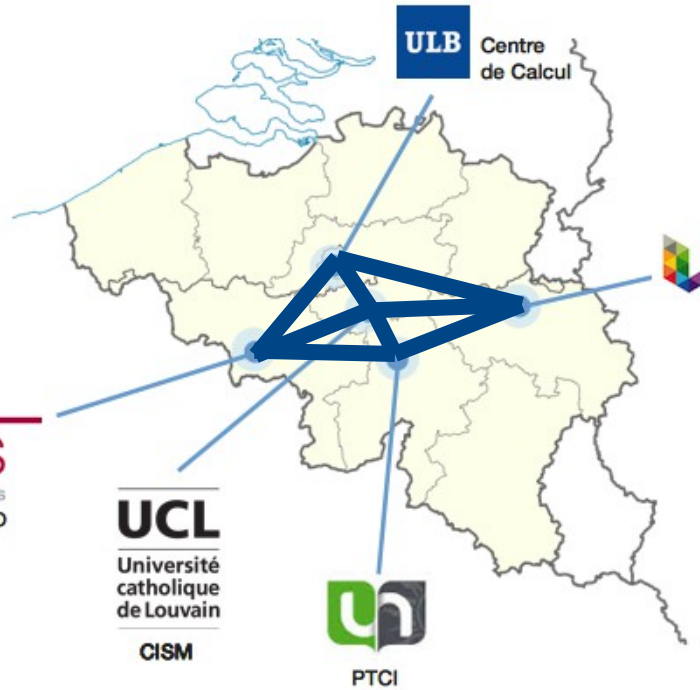
Hercules



Lemaitre2



Hmem



NIC4



New hardware

## 2. New hardware



Two large storage solutions:

- Liège
- Louvain-la-Neuve (DCIII)

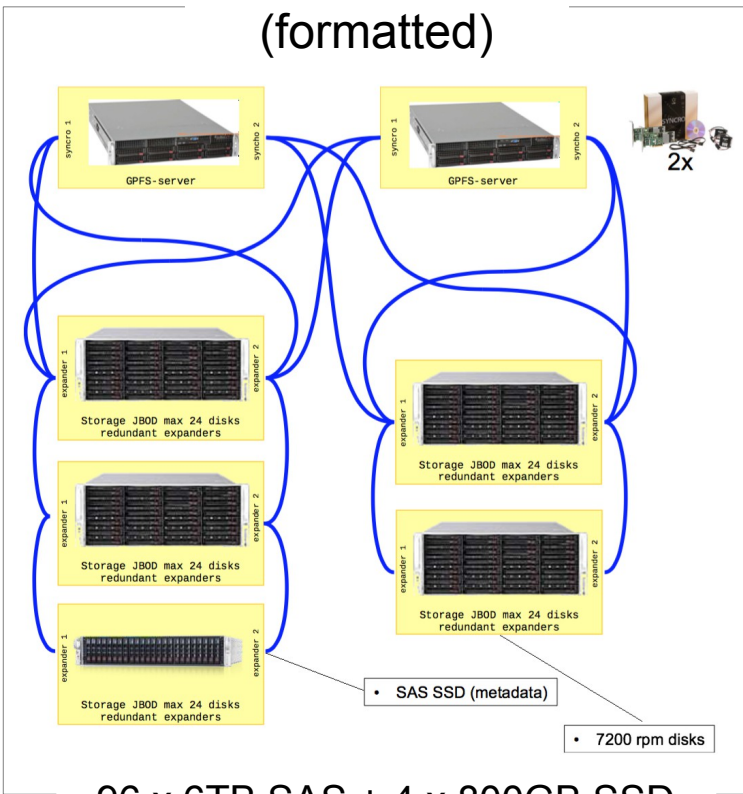


One smaller storage system on each site



New hardware

### 450TB netto (formatted)



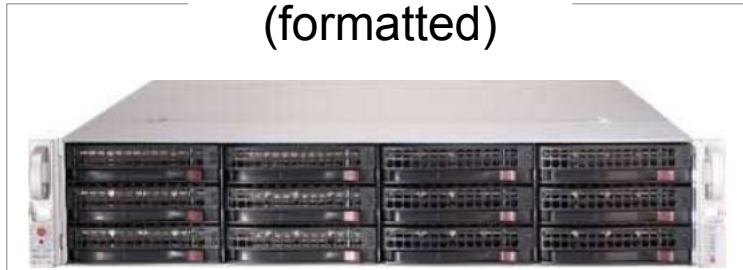
2x

96 x 6TB SAS + 4 x 800GB SSD  
8 x (10+2) RAID6

- Each system is robust to:
- loss of one server
  - loss of one pathway
  - loss of two disks

- Solution is robust to:
- loss of an entire system
  - loss of connectivity

### 55TB netto (formatted)



5x

12 x 6TB SAS + 2 x 480GB SSD  
1 x (10+2) RAID6

- Each system is robust to:
- loss of two disks

- Solution is robust\* to:
- loss of entire disk array
  - loss of network access

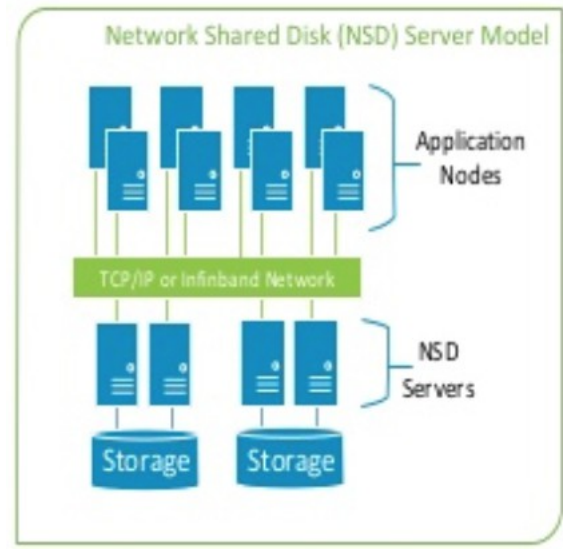


New software

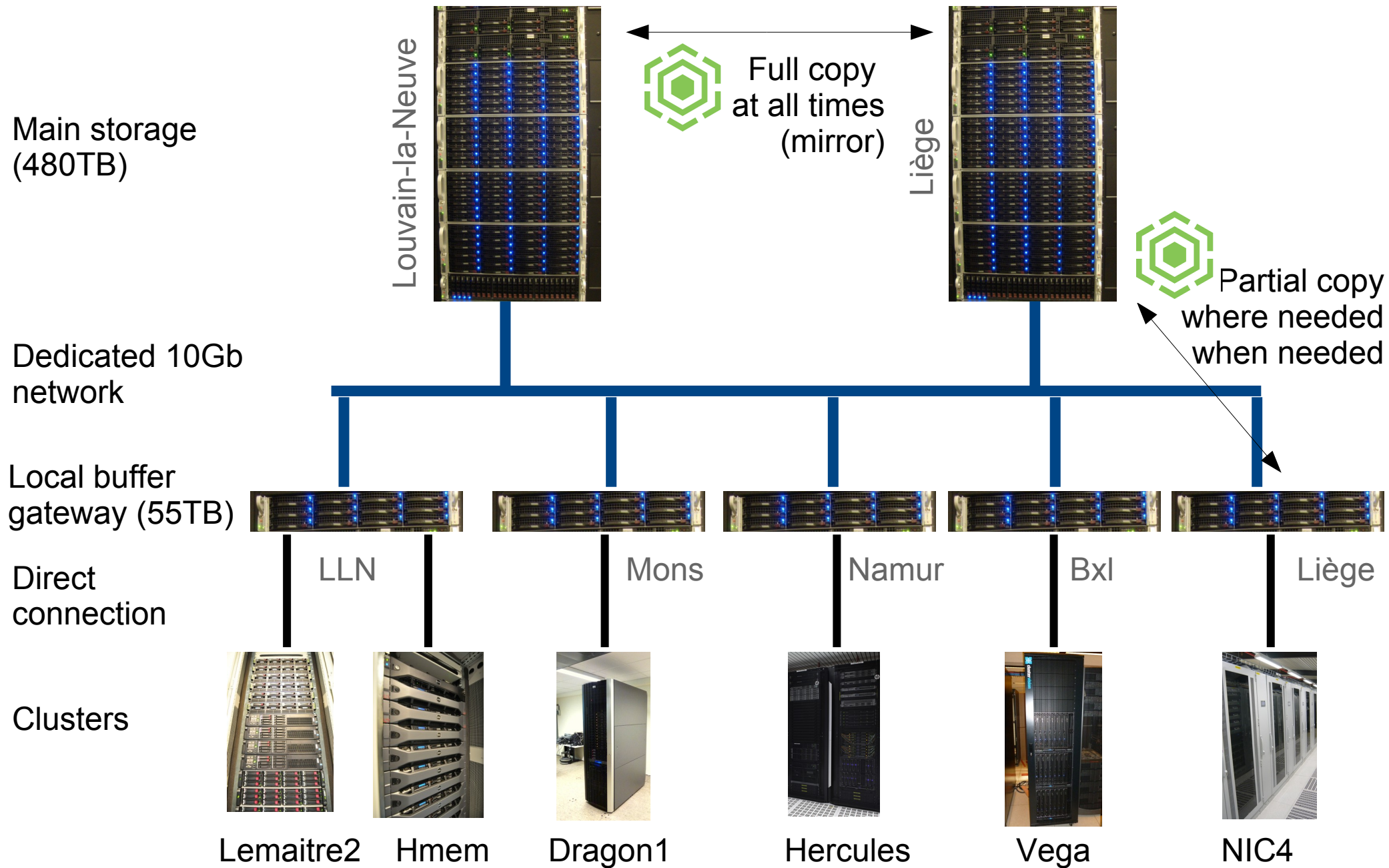
# 3. New software



IBM  
**Spectrum**  
Scale

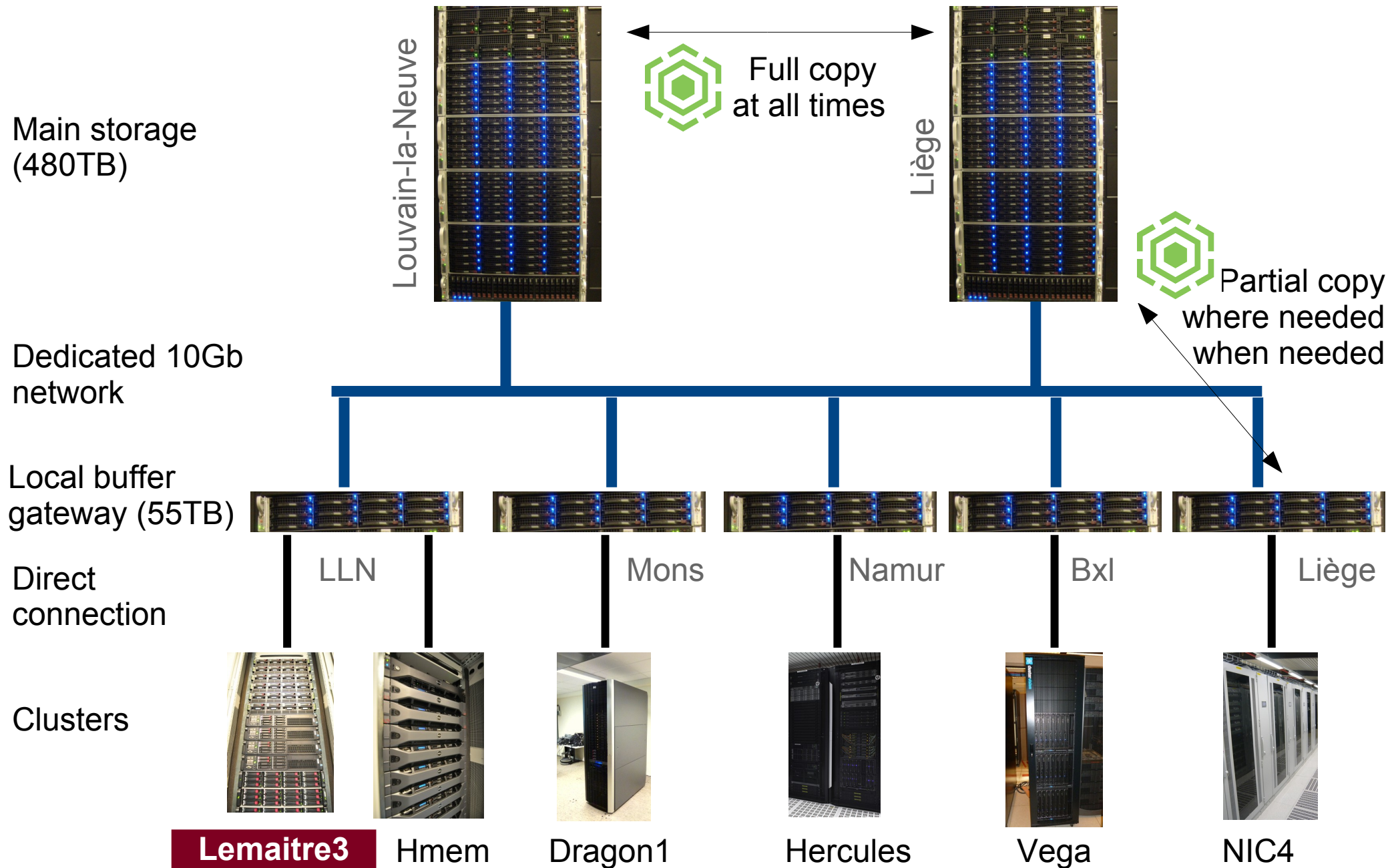


# Setup overview





# Setup overview





The new CÉCI  
common storage  
Short-term benefits

# Four spaces

- /CECI/home
  - Quota 100GB/User
  - Daily snapshots
- /CECI/proj
  - Upon request
  - Quota and duration based on request
- /CECI/trsf
  - Quota per user 100GB soft 10TB hard
  - Automatic purge of files older than 6 months
- /CECI/soft
  - Common software + modules

# Four spaces

- /CECI/home
  - Quota 100GB/User
  - Daily snapshots
- /CECI/proj
  - Upon request
  - Quota and duration based on request
- /CECI/trsf
  - Quota per user 100GB soft 10TB hard
  - Automatic purge of files older than 6 months
- /CECI/soft
  - Common software + modules

# Files written to `$CECIHOME` are visible on all clusters

```
df@ncois:~ $ ssh hmem
df@hmem00:~ $ echo "Hello World" > $CECIHOME/test.txt
df@hmem00:~ $ exit
logout
Connection to hmem.cism.ucl.ac.be closed.
df@ncois:~ $ # I just wrote a file on Hmem in the CÉCI Home
df@ncois:~ $ # I wait a bit
df@ncois:~ $ sleep 30
df@ncois:~ $ # pdsh is a parallel SSH client
df@ncois:~ $ pdsh -g ceci 'cat $CECIHOME/test.txt'
vega: Hello World
dragon1: Hello World
nic4: Hello World
lemaitre2: Hello World
hercules: Hello World
hmem: Hello World
df@ncois:~ $
```

Allow for changes to be propagated

pdsh runs the same command on many hosts at the same time

# Files written to `$CECIHOME` are visible on all clusters

```
HOME
x HOME %1 x HOME %2
Connection to hmem.cism.ucl.ac.be closed.
dfr@ncois:~ $ # I just wrote a file on Hmem in the CÉCI Home
dfr@ncois:~ $ # I wait a bit
dfr@ncois:~ $ sleep 30
dfr@ncois:~ $ # pdsh is a parallel SSH client
dfr@ncois:~ $ pdsh -g ceci 'cat $CECIHOME/test.txt'
vega: Hello World
dragon1: Hello World
nic4: Hello World
lemaitre2: Hello World
hercules: Hello World
hmem: Hello World
dfr@ncois:~ $ ssh hmem 'echo "Goodbye!" >> $CECIHOME/test.txt'
dfr@ncois:~ $ sleep 30
dfr@ncois:~ $ pdsh -g ceci 'cat $CECIHOME/test.txt'
vega: Hello World
vega: Goodbye!
dragon1: Hello World
dragon1: Goodbye!
nic4: Hello World
nic4: Goodbye!
hmem: Hello World
hmem: Goodbye!
lemaitre2: Hello World
lemaitre2: Goodbye!
hercules: Hello World
hercules: Goodbye!
dfr@ncois:~ $
```

Also visible on all  
compute nodes of  
all clusters !



# Example: install your own version of Vim on every cluster

```
dfr@ncois:~ $ ssh hmem
dfr@hmem00:~ $ mkdir -p /dev/shm/dfr/tmp
dfr@hmem00:~ $ # I have created a temp directory in a RAM filesystem
dfr@hmem00:~ $ cd /dev/shm/dfr/tmp
dfr@hmem00:/dev/shm/dfr/tmp $ # I clone the Vim repo
dfr@hmem00:/dev/shm/dfr/tmp $ git clone https://github.com/vim/vim.git
>& /dev/null
dfr@hmem00:/dev/shm/dfr/tmp $ cd vim
dfr@hmem00:/dev/shm/dfr/tmp/vim $ ./configure --prefix=$CECIHOME/.local
dfr@hmem00:/dev/shm/dfr/tmp/vim $ make -j 12 >& /dev/null
dfr@hmem00:/dev/shm/dfr/tmp/vim $ make install >& /dev/null
dfr@hmem00:/dev/shm/dfr/tmp/vim $ ls /CECI/home/ucl/pan/dfr/
bin  ex  rview  rvim  share  view  vim  vimdiff  vimtutor
dfr@hmem00:/dev/shm/dfr/tmp/vim $ echo $PATH
/CECI/home/ucl/pan/dfr/.local/bin:/home/ucl/pan/dfr/.local/bin:/home/uc
l/pan/dfr/.cw/def:/usr/lib64/qt-3.3/bin:/usr/local/bin:/bin:/usr/bin:/u
sr/local/sbin:/usr/sbin:/sbin:/opt/dell/psadmin/bin:/home/ucl/pan/dfr/
CASTEM2015/bin
dfr@hmem00:/dev/shm/dfr/tmp/vim $
```

I work on Hmem

I compile a newer version from sources

I install it in my CECI home directory

I setup the PATH accordingly (on all clusters)





# A catch though...

By default, compilers will tune the code for the CPU of the machine they run on.

		run on					
		Hmem	Lemaitre2	Dragon1	Hercules	Vega	NIC4
compile on	Hmem	ok	sub-opt	sub-opt	sub-opt	sub-opt	sub-opt
	Lemaitre2	crash	ok	sub-opt	sub-opt	sub-opt	sub-opt
	Dragon1	crash	crash	ok	ok	sub-opt	ok
	Hercules	crash	crash	ok	ok	sub-opt	ok
	Vega	crash	crash	crash	crash	ok	crash
	NIC4	crash	crash	ok	ok	sub-opt	ok

# Mitigation:

## **GCC:**

- march=core2 to build binaries running everywhere
- mtune=[westmere|sandybridge|bdver1]  
to tune for the cluster you use the most

## **Intel** (multiple code paths):

- xSSE2 to build binaries running everywhere
- axSSE4.2,AVX,CORE-AVX2  
to add additional paths optimized for each cluster

# Mitigation:

## GCC: “Function multi-versioning”

Toggle line numbers

```
1 __attribute__ ((target ("default")))
2 int foo ()
3 {
4     // The default version of foo.
5     return 0;
6 }
7
8 __attribute__ ((target ("sse4.2")))
9 int foo ()
10 {
11     // foo version for SSE4.2
12     return 1;
13 }
14
15 __attribute__ ((target ("arch=atom")))
16 int foo ()
17 {
18     // foo version for the Intel ATOM processor
19     return 2;
20 }
21
22 __attribute__ ((target ("arch=amdfam10")))
23 int foo ()
24 {
25     // foo version for the AMD Family 0x10 processors.
26     return 3;
27 }
28 int main ()
29 {
30     int (*p)() = &foo;
31     assert ((*p) () == foo ());
32     return 0;
33 }
```

In the above example, 4 versions of function foo are created. The first version of foo with the target attribute "default" is the default version. This version gets executed when no other target specific version qualifies for execution on a particular platform. A new version of foo is created by using the same function signature but with a different target string. Function foo is called or a pointer to it is taken just like a regular function. With the new support, GCC takes care of doing the dispatching to call the right version at runtime.

# Mitigation:

## Intel compiler: “Manual processor dispatch”

### Example

```
#include <stdio.h>
// need to create specific function versions for the following processors:
__declspec(cpu_dispatch(core_2nd_gen_avx, core_i7_sse4_2, core_2_duo_ssse3, generic ))
void dispatch_func() {}; // stub that will call the appropriate specific function
version

__declspec(cpu_specific(core_2nd_gen_avx))
void dispatch_func() {
    printf("\nCode for 2nd generation Intel Core processors with support for AVX goes
here\n");
}

__declspec(cpu_specific(core_i7_sse4_2))
void dispatch_func() {
    printf("\nCode for Intel Core processors with support for SSE4.2 goes here\n");
}

__declspec(cpu_specific(core_2_duo_ssse3))
void dispatch_func() {
    printf("\nCode for Intel Core 2 Duo processors with support for SSSE3 goes here\n");
}

__declspec(cpu_specific(generic))
void dispatch_func() {
    printf("\nCode for non-Intel processors and generic Intel processors goes here\n");
}

int main() {
    dispatch_func();
    printf("Return from dispatch_func\n");
    return 0;
}
```

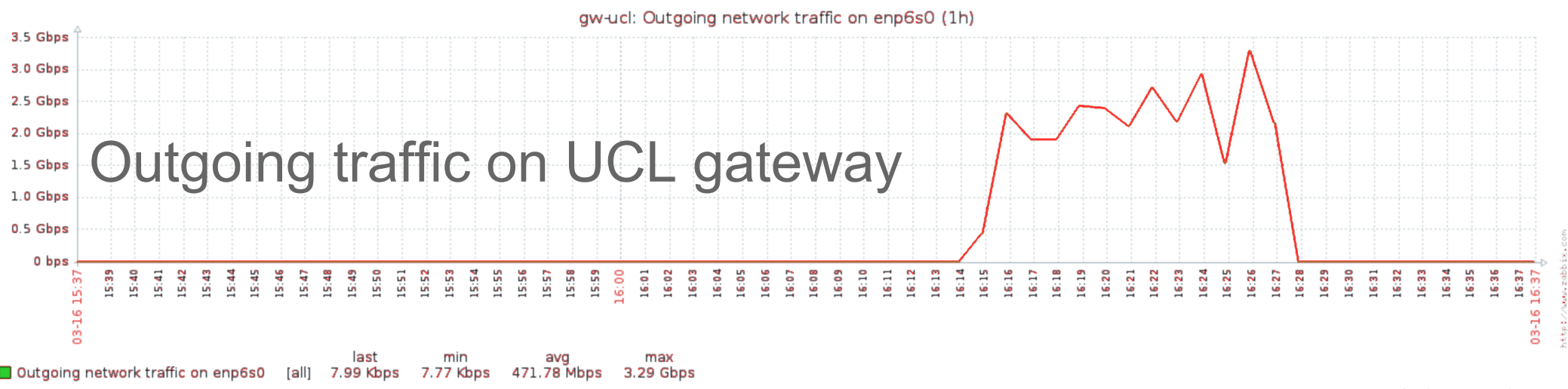
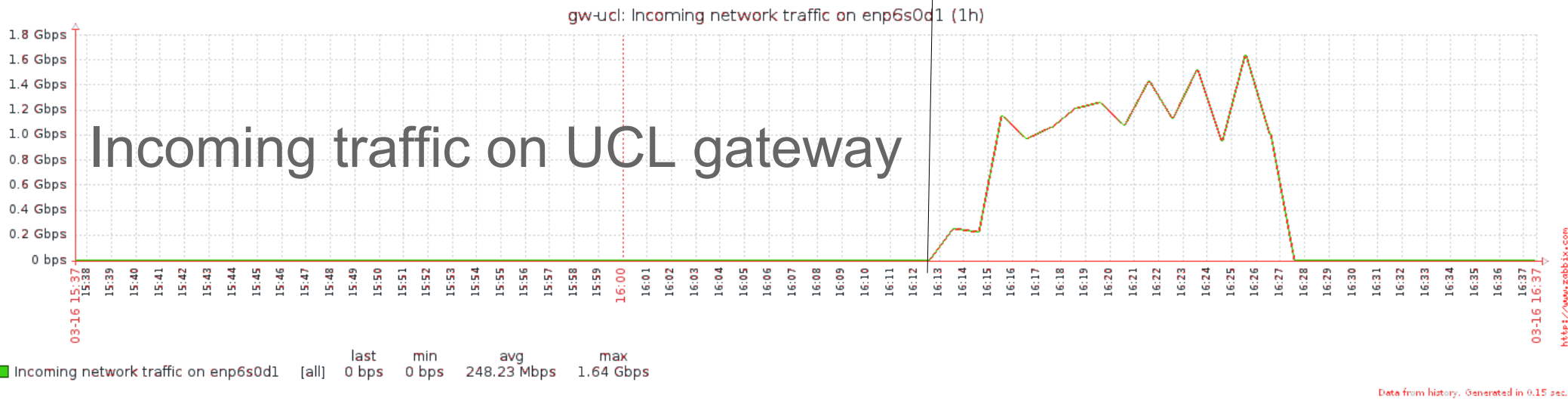
# Four spaces

- /CECI/home
  - Quota 100GB/User
  - Daily snapshots
- /CECI/proj
  - Upon request
  - Quota and duration based on request
- /CECI/trsf
  - Quota per user 100GB soft 10TB hard
  - Automatic purge of files older than 6 months
- /CECI/soft
  - Common software + modules

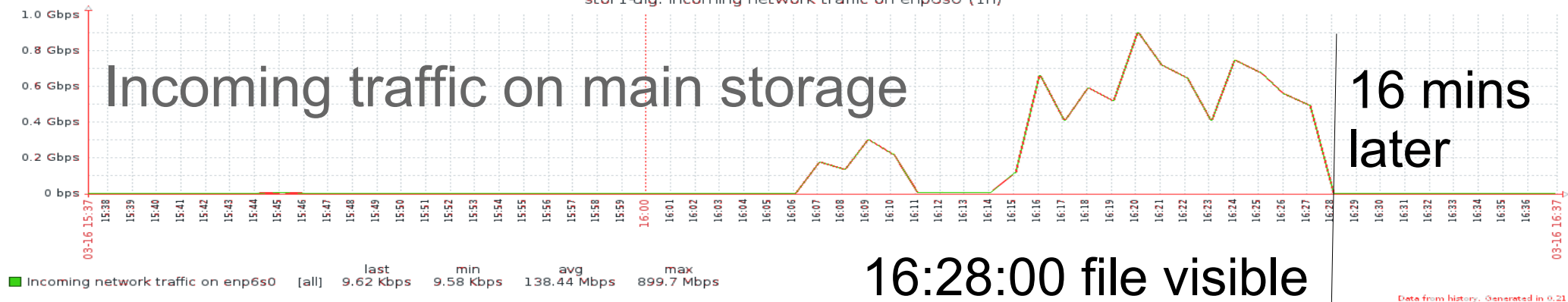
# Copy of 100GB file

from the scratch of Lemaitre2 to scratch of NIC4

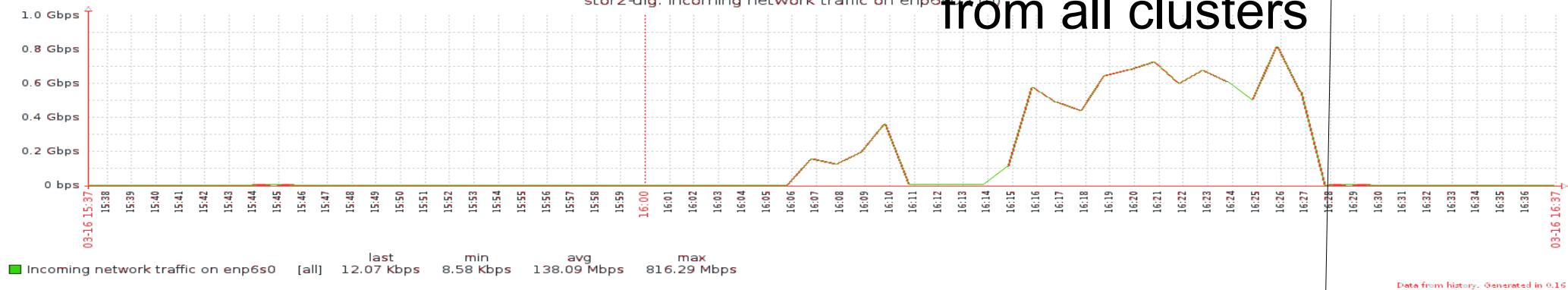
16:12:31 command issued



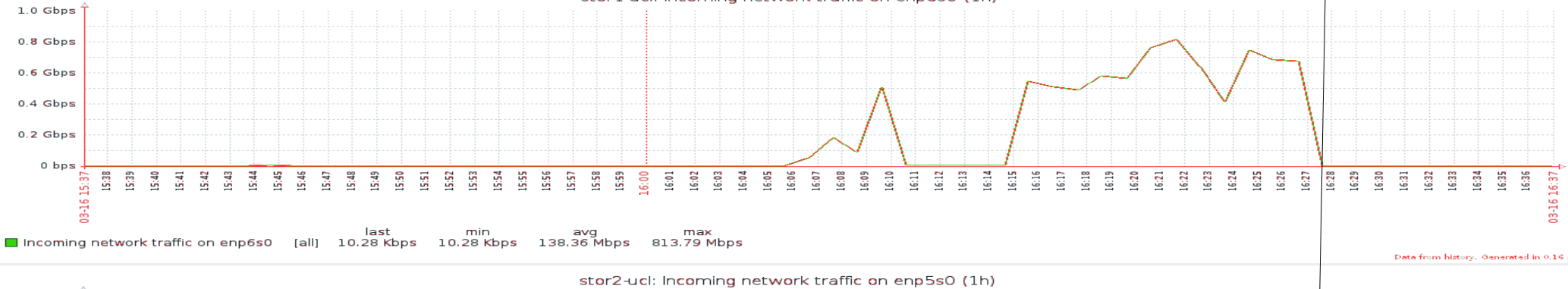
stor1-ulg: Incoming network traffic on enp6s0 (1h)



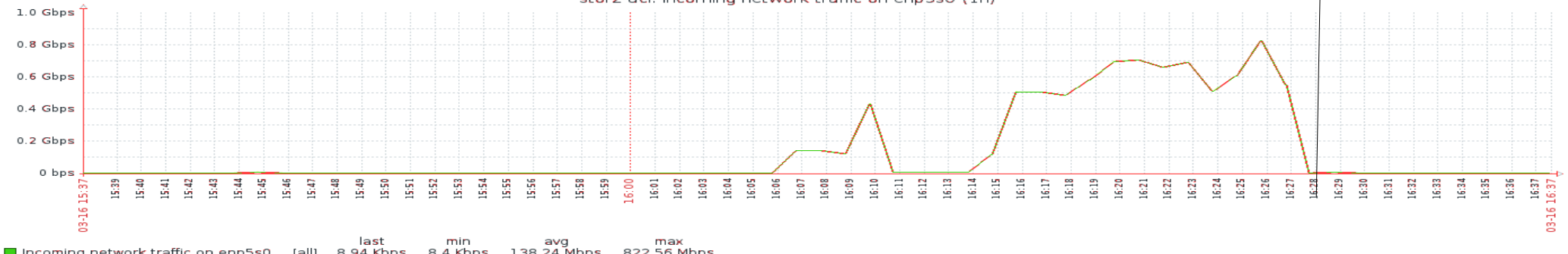
stor2-ulg: Incoming network traffic on enp6s0 (1h)



stor1-uc1: Incoming network traffic on enp6s0 (1h)

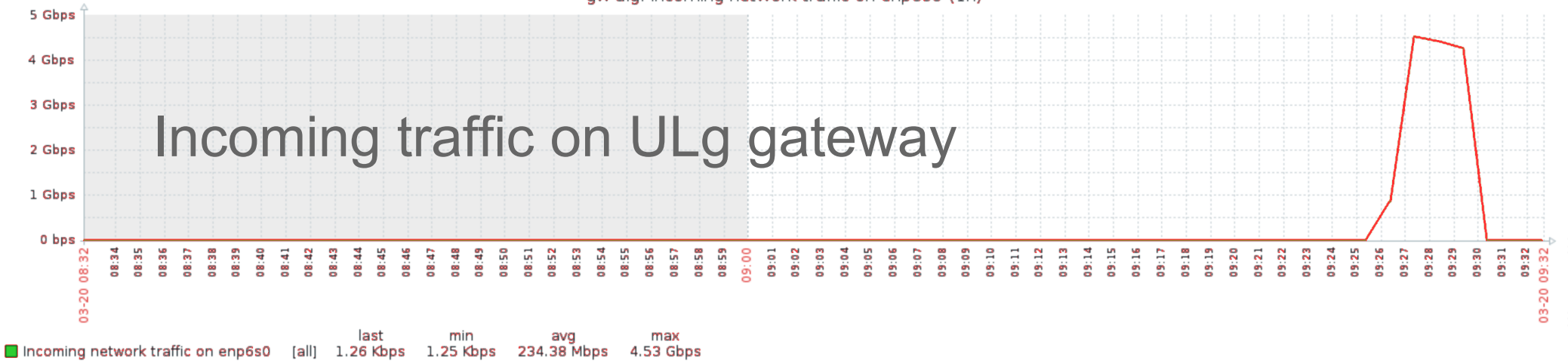


stor2-uc1: Incoming network traffic on enp5s0 (1h)

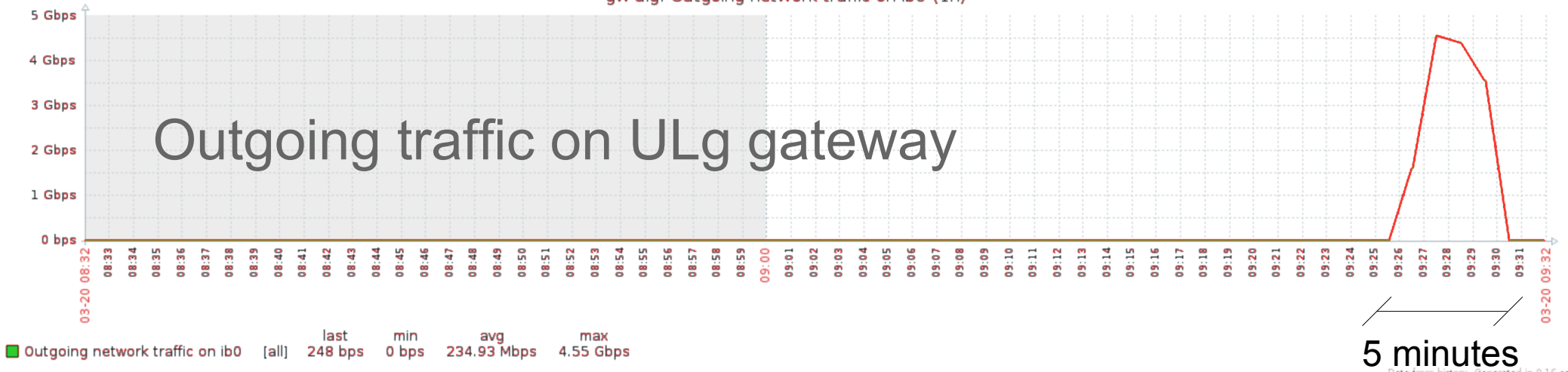




gw-ulg: Incoming network traffic on enp6s0 (1h)



gw-ulg: Outgoing network traffic on ib0 (1h)



Total 21 minutes  
vs.

```

x dfr@nic4 361
dfr@nic4:~ $ scp dfr@lemaitre2:$GLOBALSCRATCH/bigfile.dat $
GLOBALSCRATCH
bigfile.dat 0% 79MB 6.3MB/s 4:25:25 ETA

```

UNamur (Université de Namur) su

Using the common filesystem — CÉCI

Docs

CÉCI

Search docs

QUICK START - FIRST STEPS

- Creating an account
- Connecting to the clusters
- Copying files
- Editing files
- Slurm Quick Start Tutorial

ACCESSING THE CLUSTERS

- Troubleshooting and frequent mistakes

MANAGING FILES

- Transferring files to and from the clusters

Using the common filesystem

- Home
- Trfs
- Proj
- Soft

CÉCI v: latest

Docs » Using the common filesystem

## Using the common filesystem

All CÉCI clusters are connected to a central storage system that is visible to all compute nodes of all clusters. This system runs on a fast, dedicated, network. It will become the home of the users in the near future, but in this first phase, it is set up as an additional home besides the default, cluster-specific, home.

This storage system is installed at two CÉCI locations and data are replicated synchronously on both locations to ensure data safety and a certain level of high availability. Moreover, on each site, a local cache is setup to mask the latencies of the network and make sure the user experience is as smooth as possible. Those caches are replicated asynchronously with the central storage, meaning that files that are written there will appear after some delay on the other clusters. It also means that if you modify the same file from two different clusters, the result is undefined.

**Warning**

Do not write to the same file from two different clusters at the same time. This would corrupt the file.

The storage is split into four distinct directories:

# Four spaces

- /CECI/home
  - Quota 100GB/User
  - Daily snapshots
- /CECI/proj
  - Upon request
  - Quota and duration based on request
- /CECI/trsf
  - Quota per user 100GB soft 10TB hard
  - Automatic purge of files older than 6 months
- /CECI/soft
  - Common software + modules

# Gotcha's

- Using the `du` command on a cluster will give you **your usage on the gateway/cache connected to that cluster.** and not your usage on the main storage.
- Use the `quota -usa` command on a cluster to get your usage on the main storage.
- Some clusters mount the storage via *automount* so the `quota` command **might not display the CÉCI storage.** In that case, run `ls $CECIHOME` to trigger the mounting and run `quota` again.
- Quota are enforced at the gateway/cache level (100GB) and at on the main storage (200GB).
  - More than 100GB on a cluster => you cannot write from that cluster anymore
  - More than 200GB on the main storage => data are not transferred from the clusters anymore

# What's next...

- **Fine-tune** configuration

Your feedback: your local CÉCI system administrator

- Setup **procedure to request** group space
- Build **common software** installation

All further information through  
the CÉCI users mailing list



**Consortium des Équipements de Calcul Intensif**  
Funded by F.R.S.-FNRS

[www.ceci-hpc.be](http://www.ceci-hpc.be)



**Try the new CÉCI  
common storage!**