

# **DATA FLOWS, STORAGE AND STORAGE TUNING AND ALL THAT**

RON TROMPERT



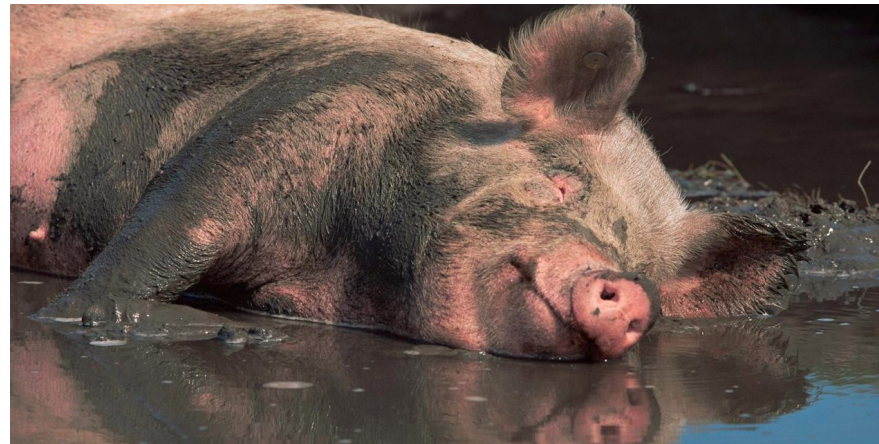
# Contents

- A bit of history



- Things that we do at SURFsara

- Dive in the mud



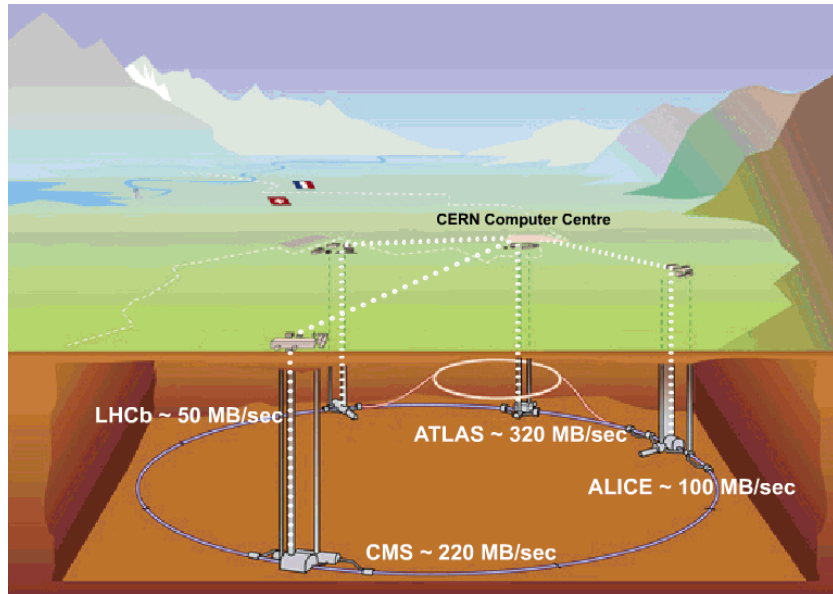


A bit of history

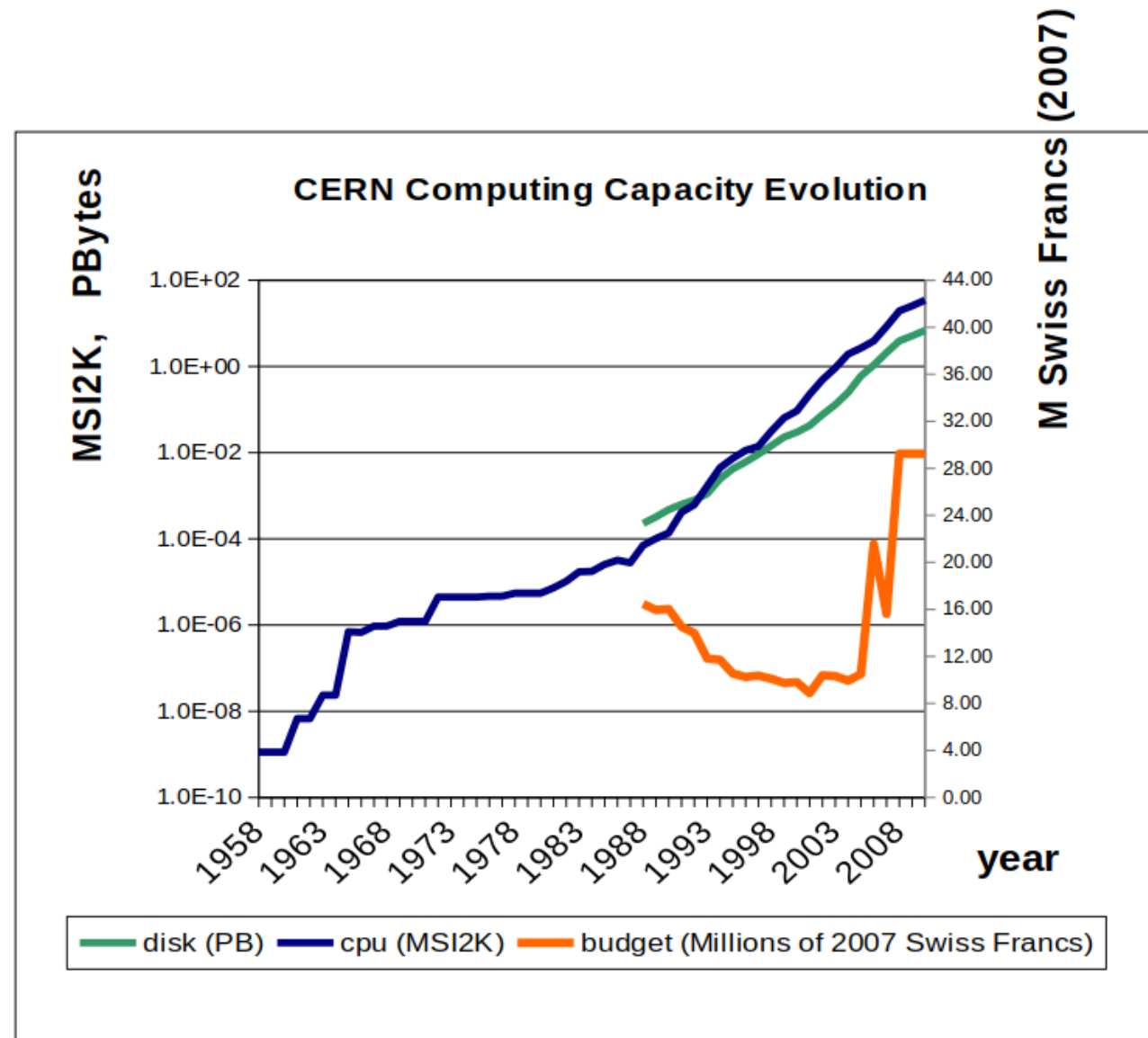




# Once upon a time.....

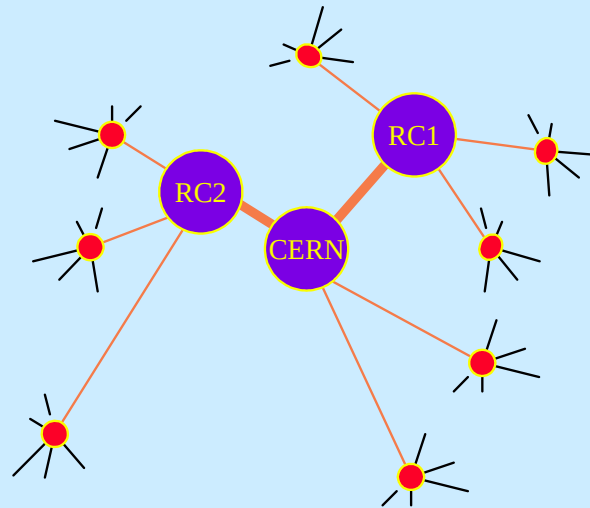


# The LHC challenge

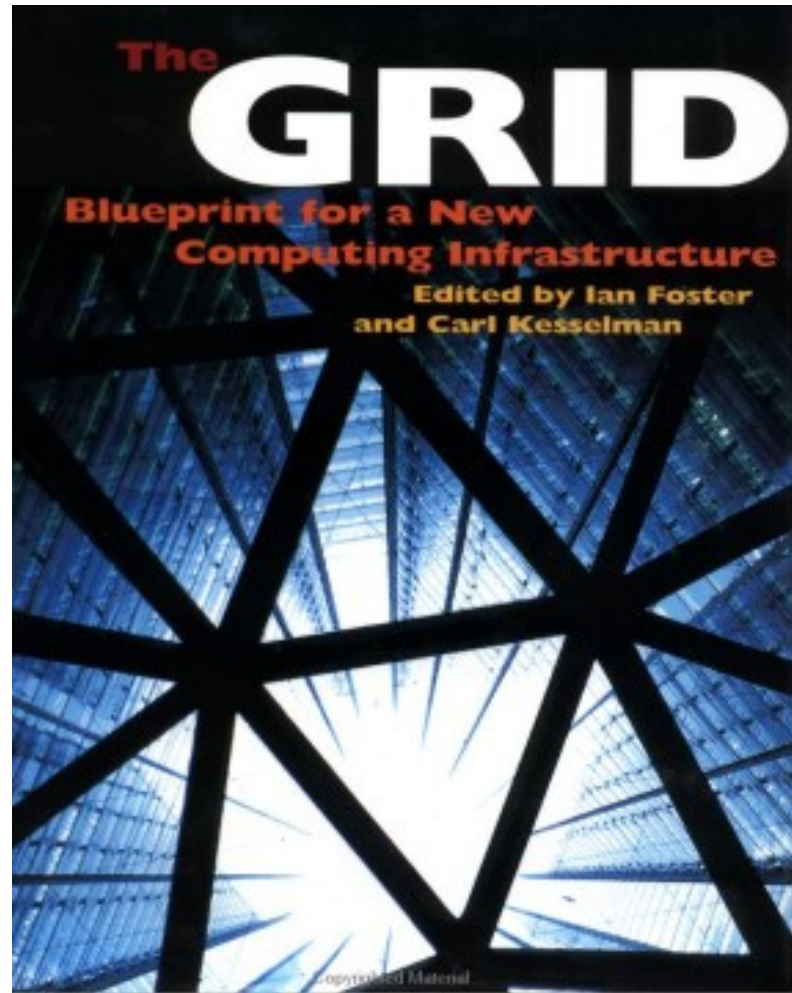


# How to solve this?

- Partially decentralized model
  - replicate the event data at about five regional centres
  - data transfer via network or movable media



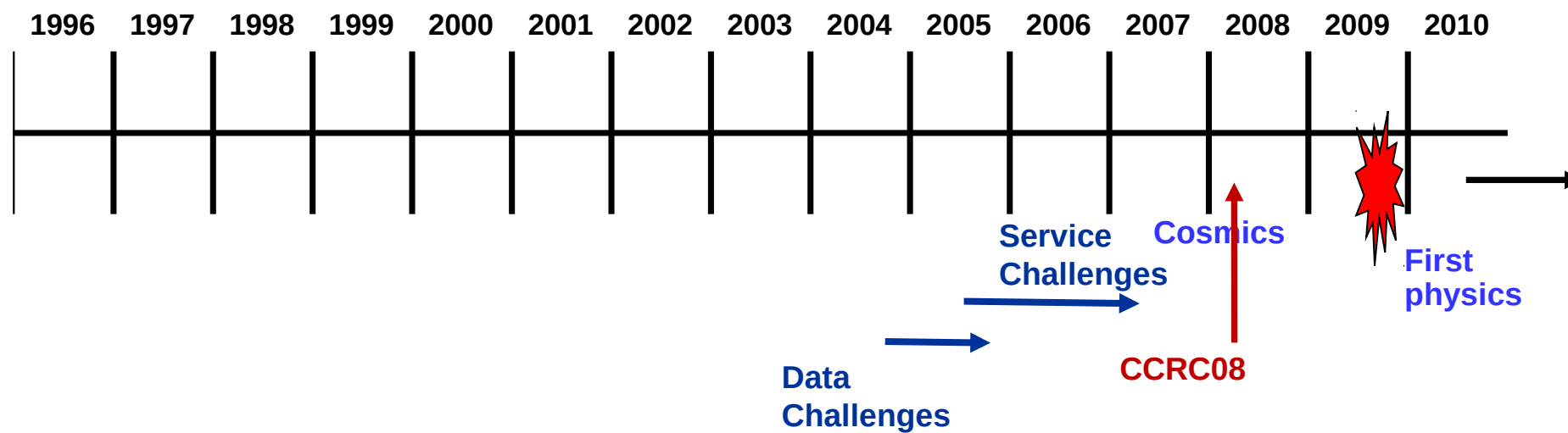
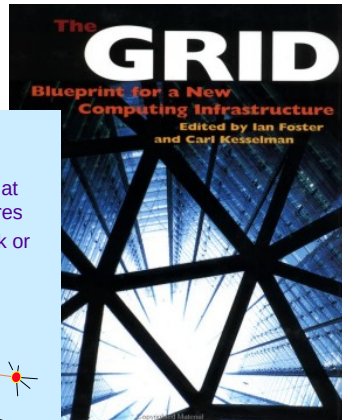
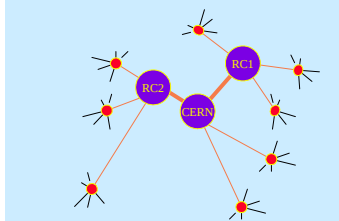
Then this happened



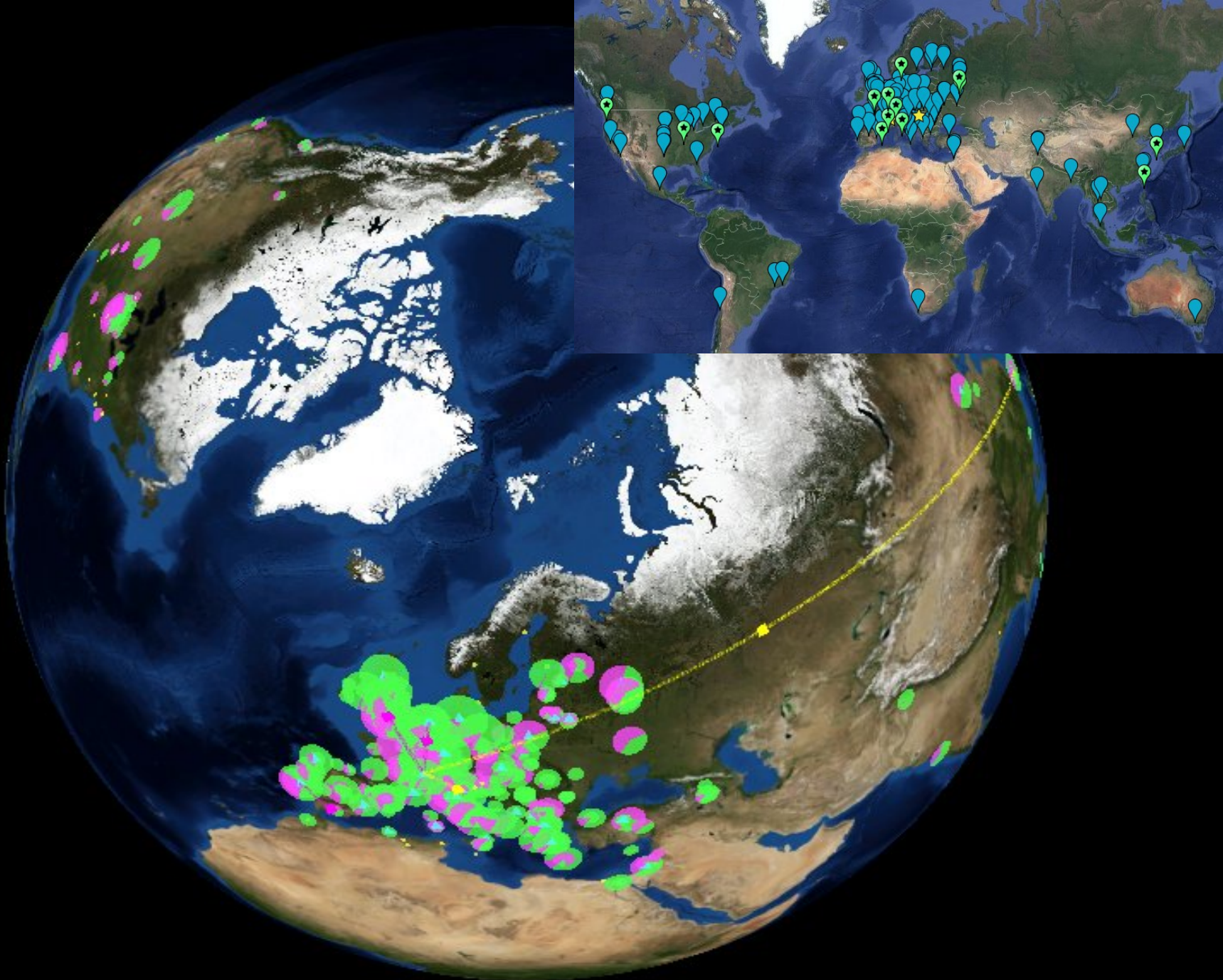


# How did we get there?

- Partially decentralized model
  - replicate the event data at about five regional centres
  - data transfer via network or movable media





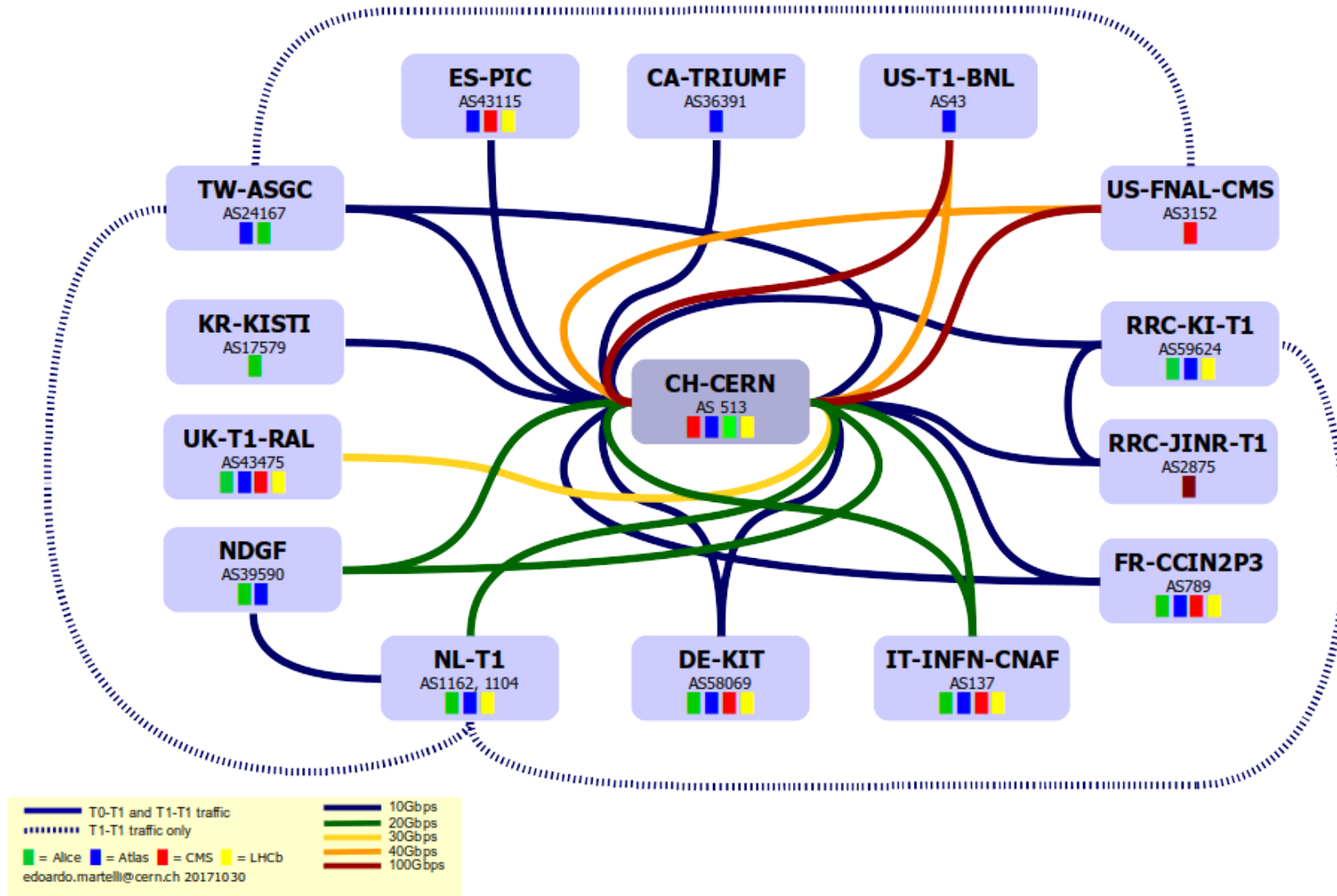


**WLCG**  
Worldwide LHC Computing Grid

- 42 countries
- 170 datacenters
- 50-70PB of data per year
- 800000 jobs per day
- 500PB disk
- 750PB tape

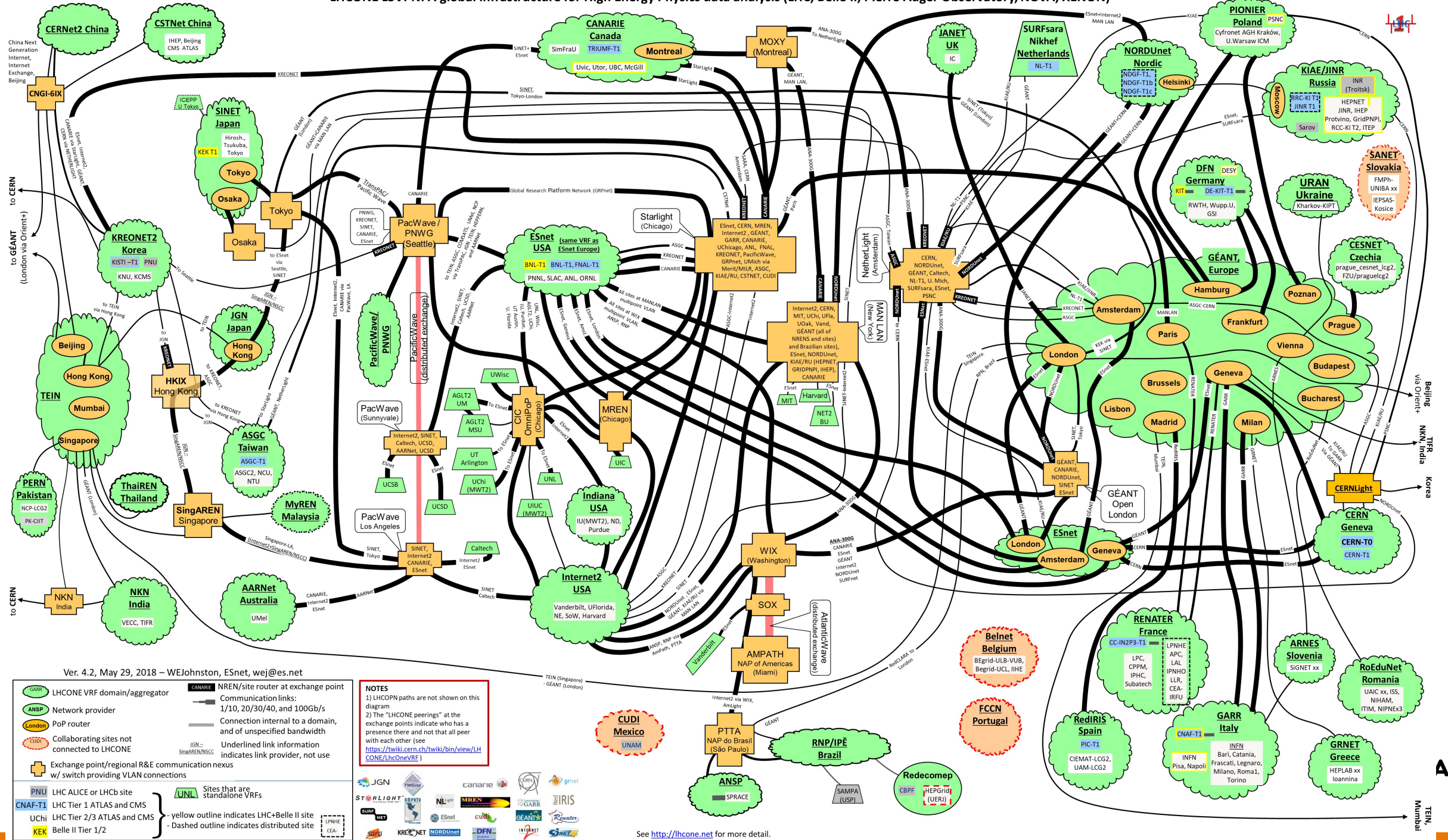


# LHCOPN





# LHCONE L3VPN: A global infrastructure for High Energy Physics data analysis (LHC, Belle II, Pierre Auger Observatory, NOVA, XENON)





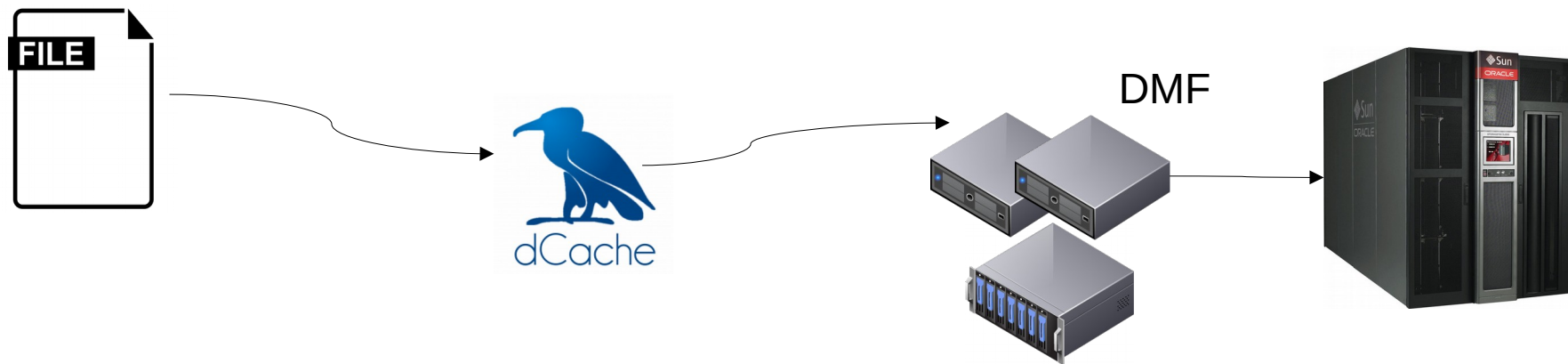




# dCache

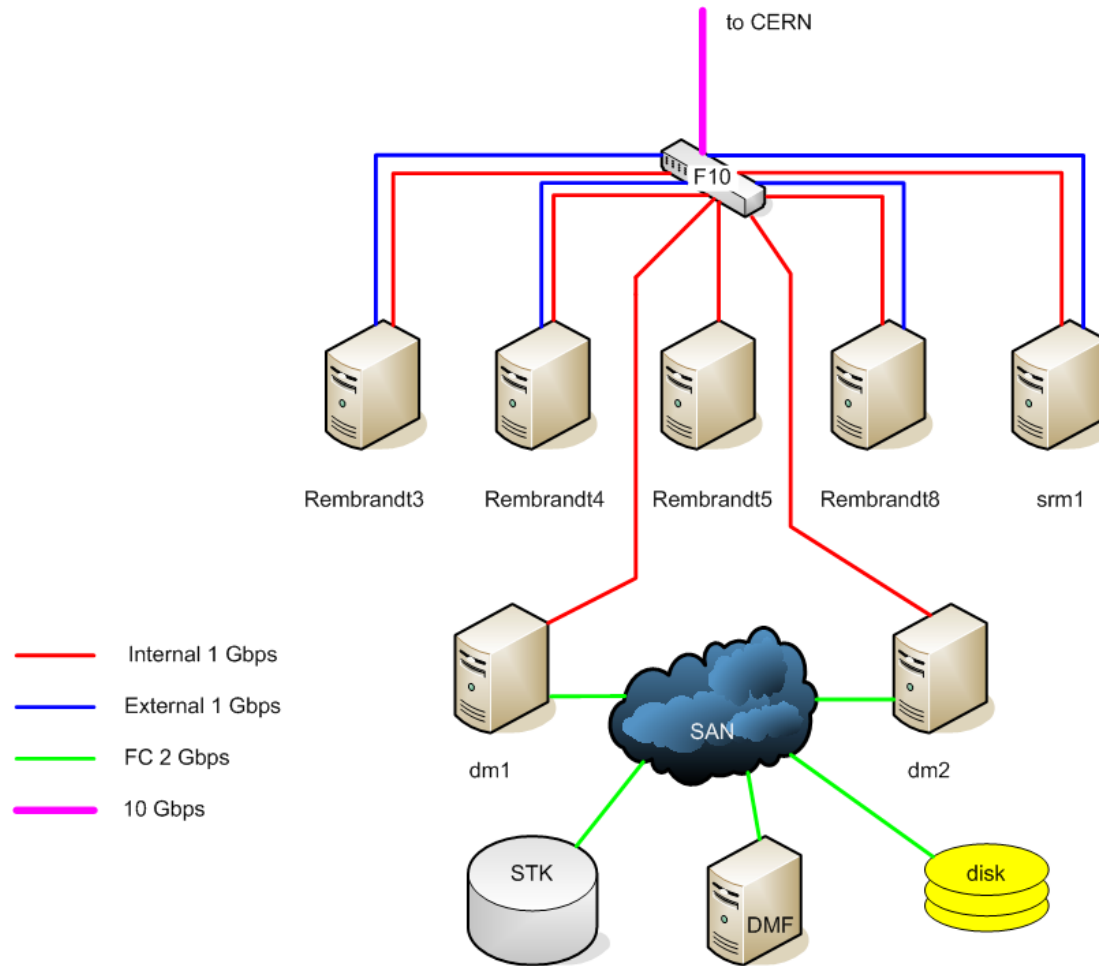


- Manages a large number of disk servers (NAS)
- Put, Get, Delete
- Data is immutable
- Seamless integration with tape systems

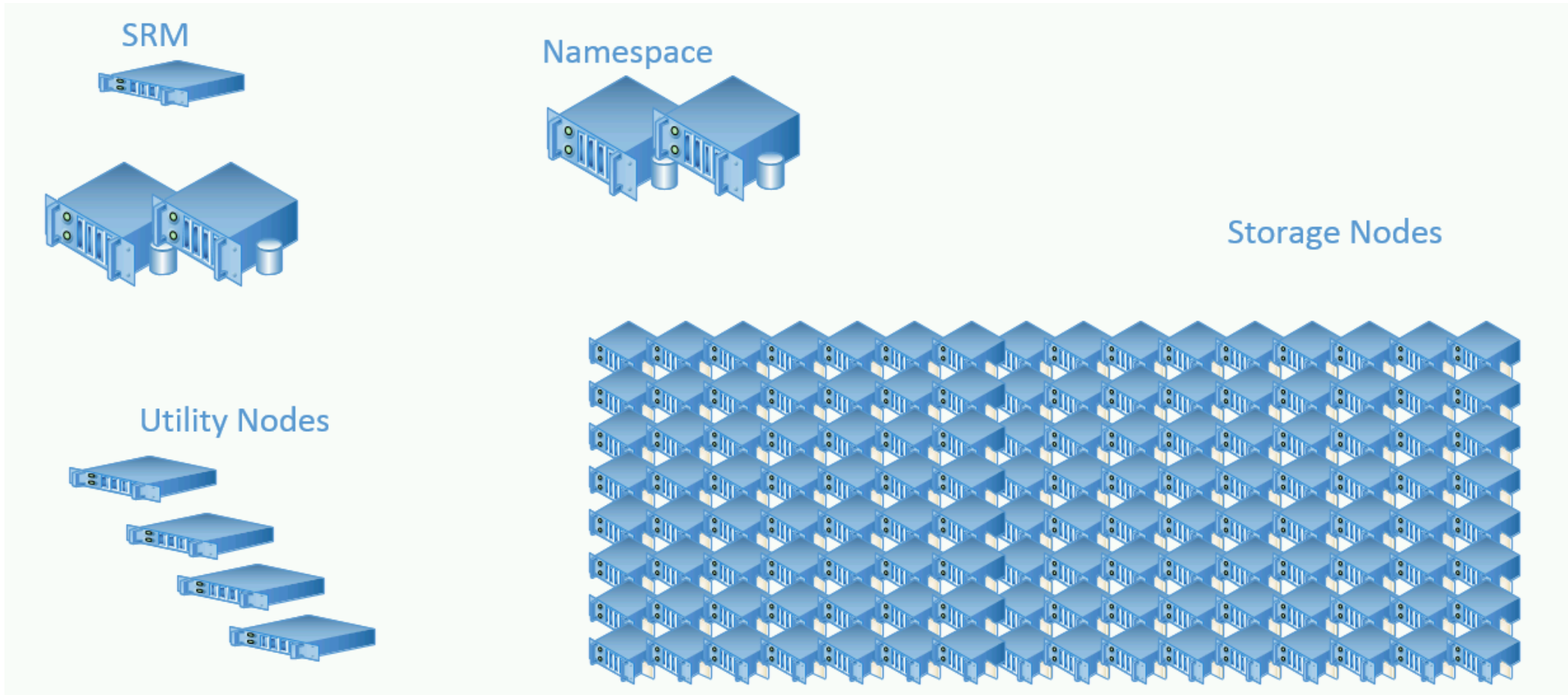




# 2005



# 2018



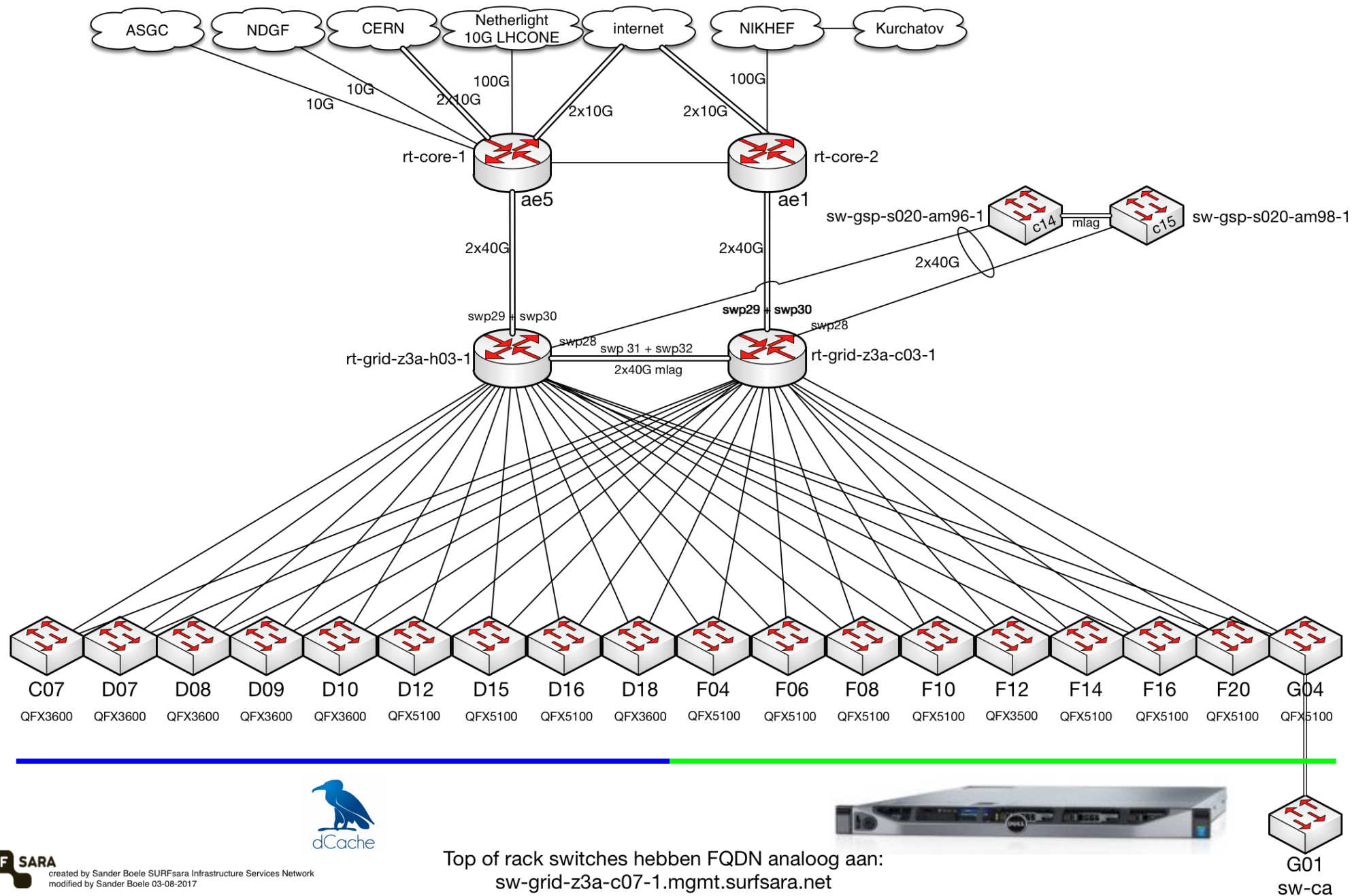


# dCache in numbers



- 14PB storage capacity
- Manages 49PB on disk and tape
- Last year: 117 PB read, 22PB written, 7PB stored to tape and 8PB restored from tape
- 128 disk pool nodes with 80-400TB each

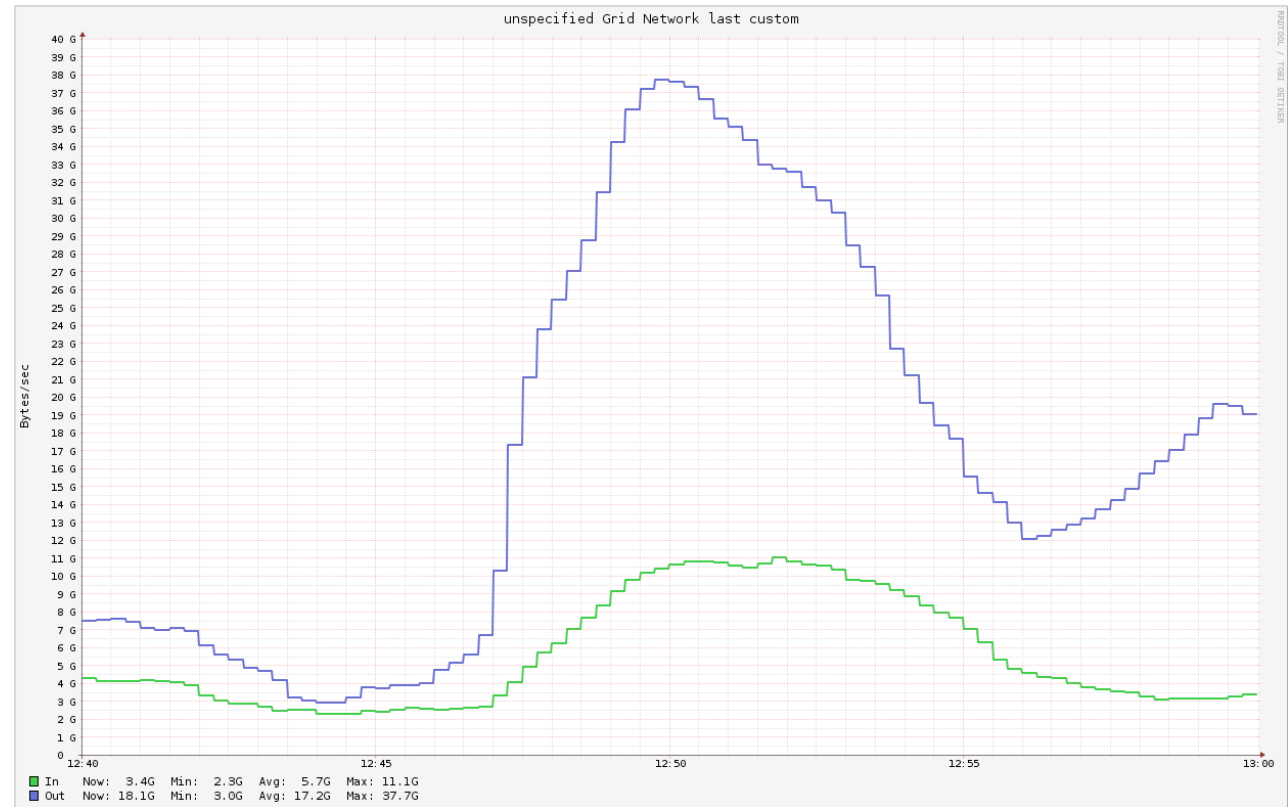
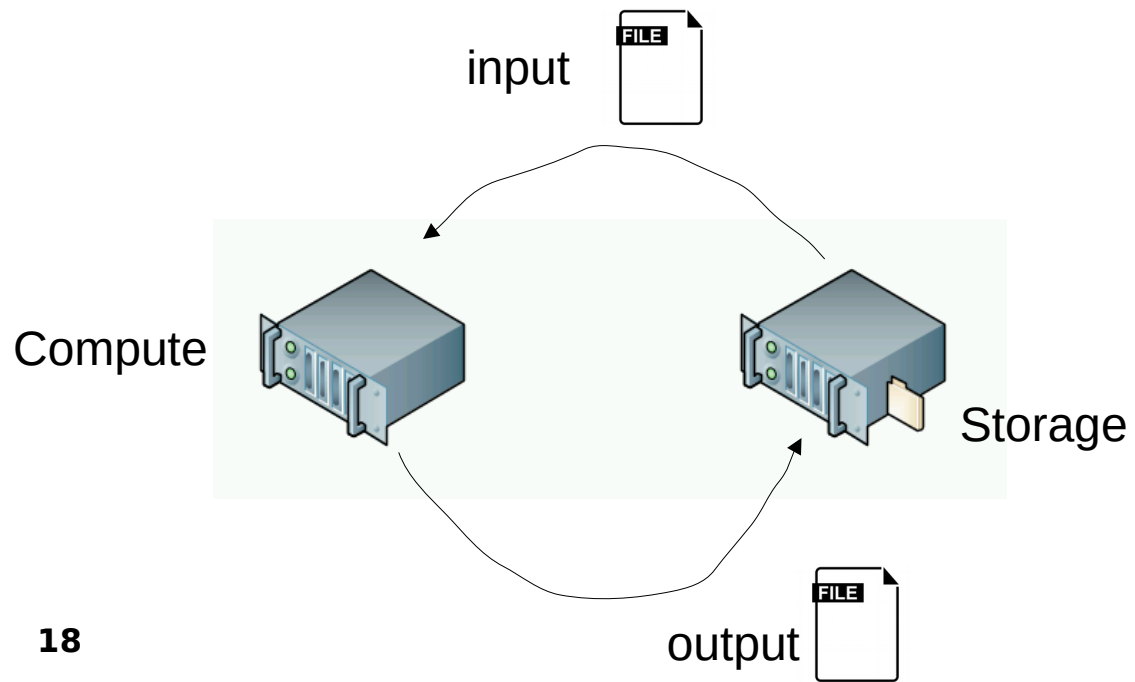
# GRID network

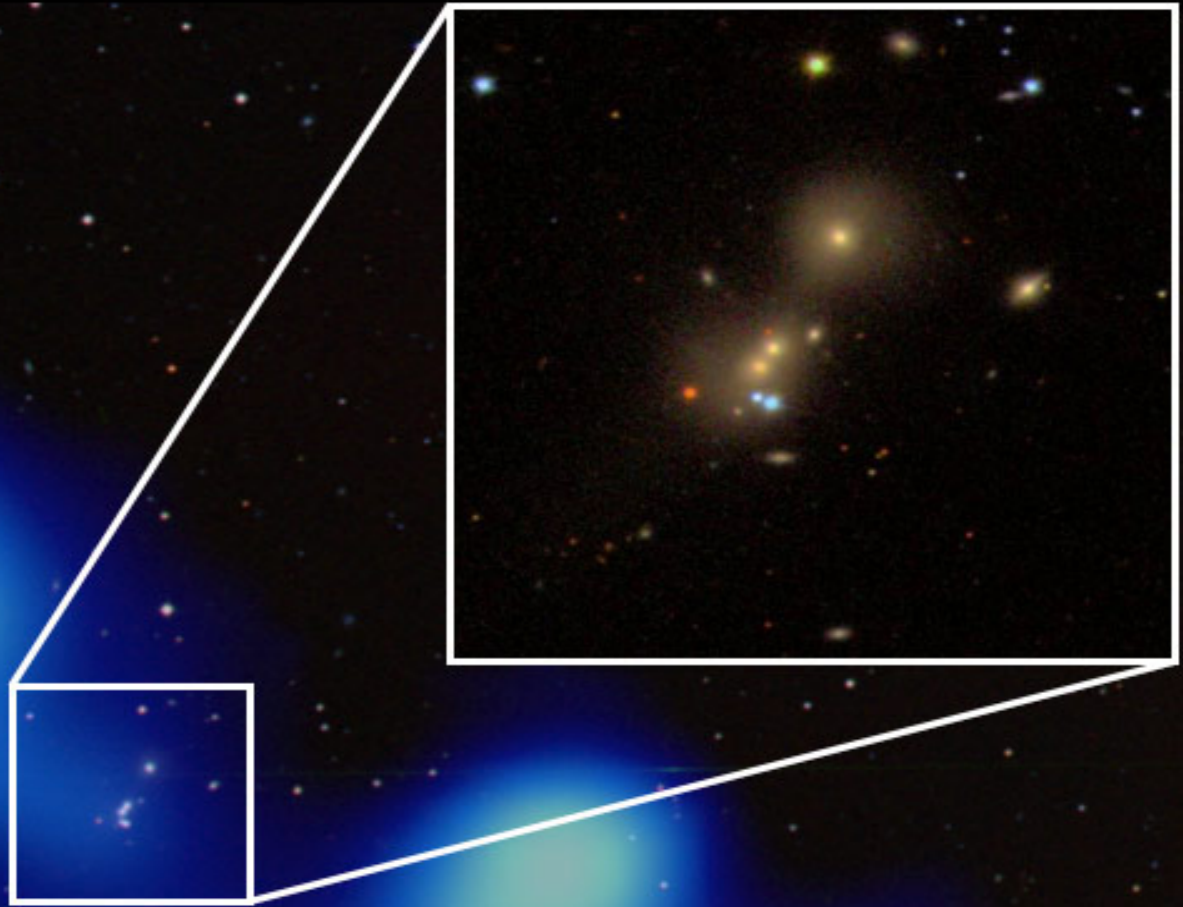




# Data processing

- Thousands of running jobs processing data

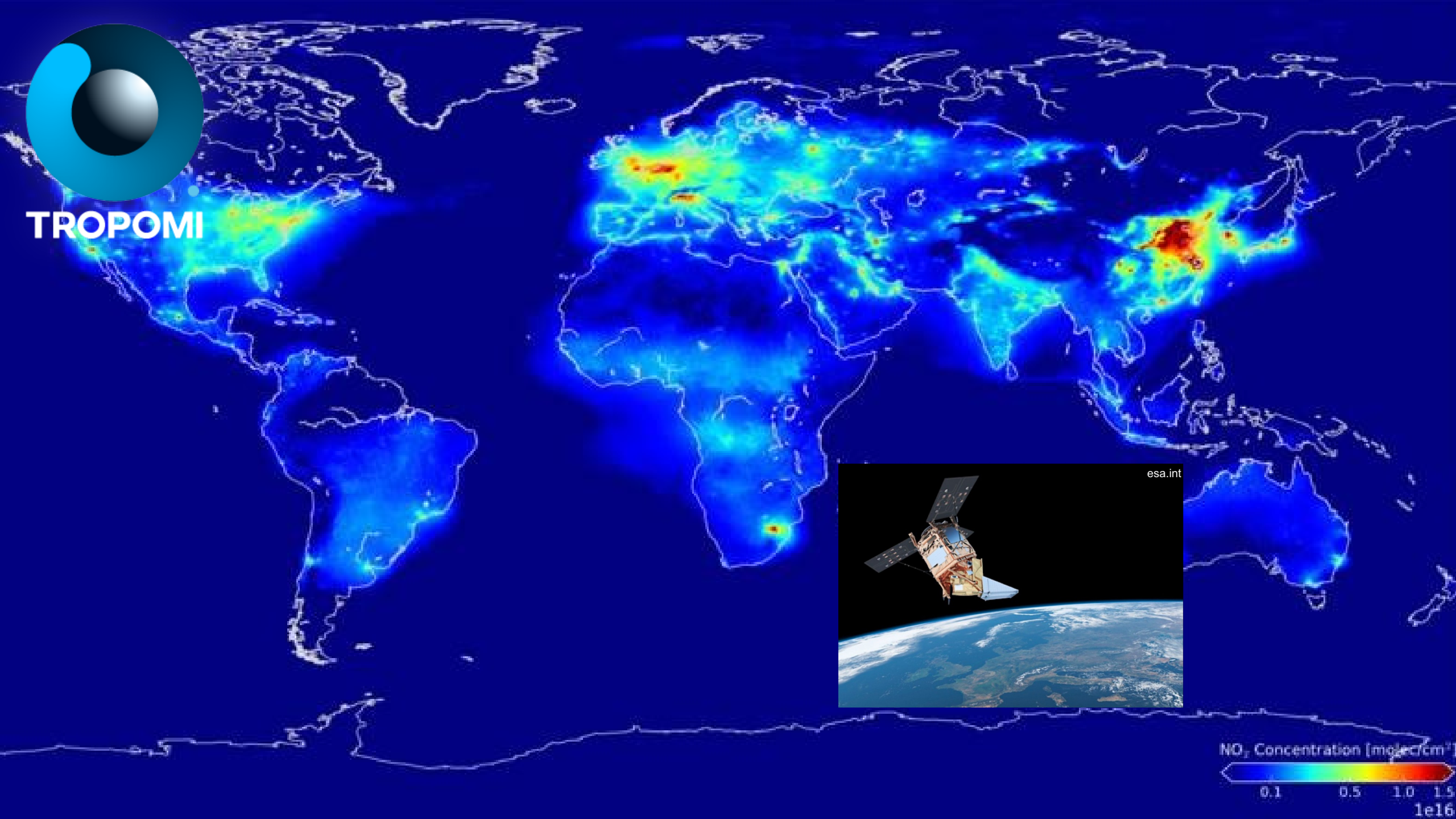






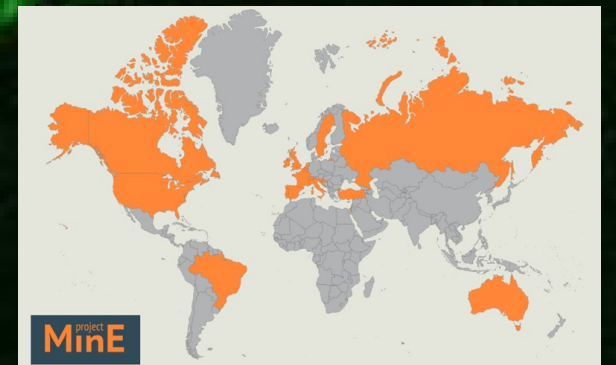
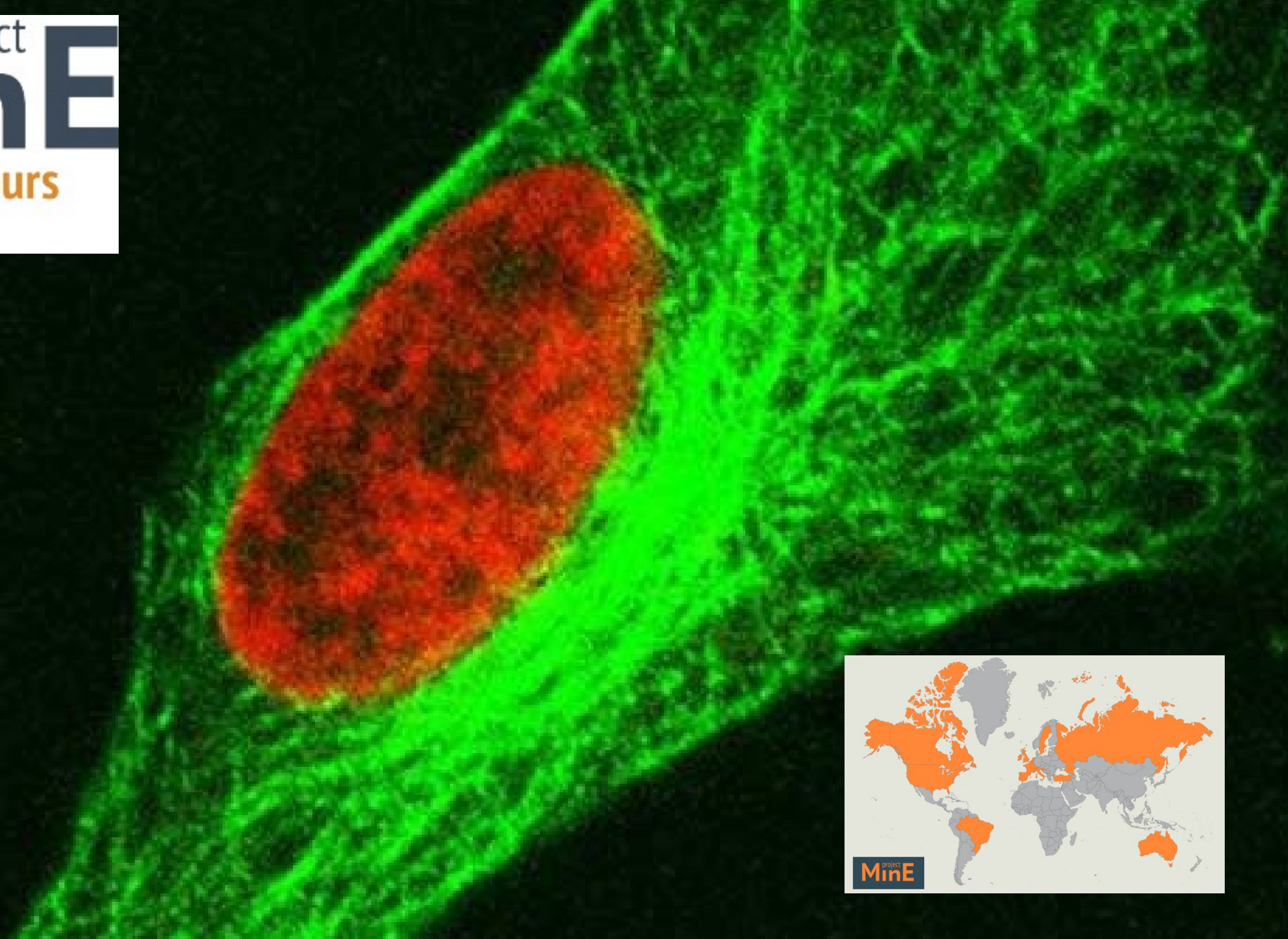


TROPOMI



esa.int







Dive in the mud

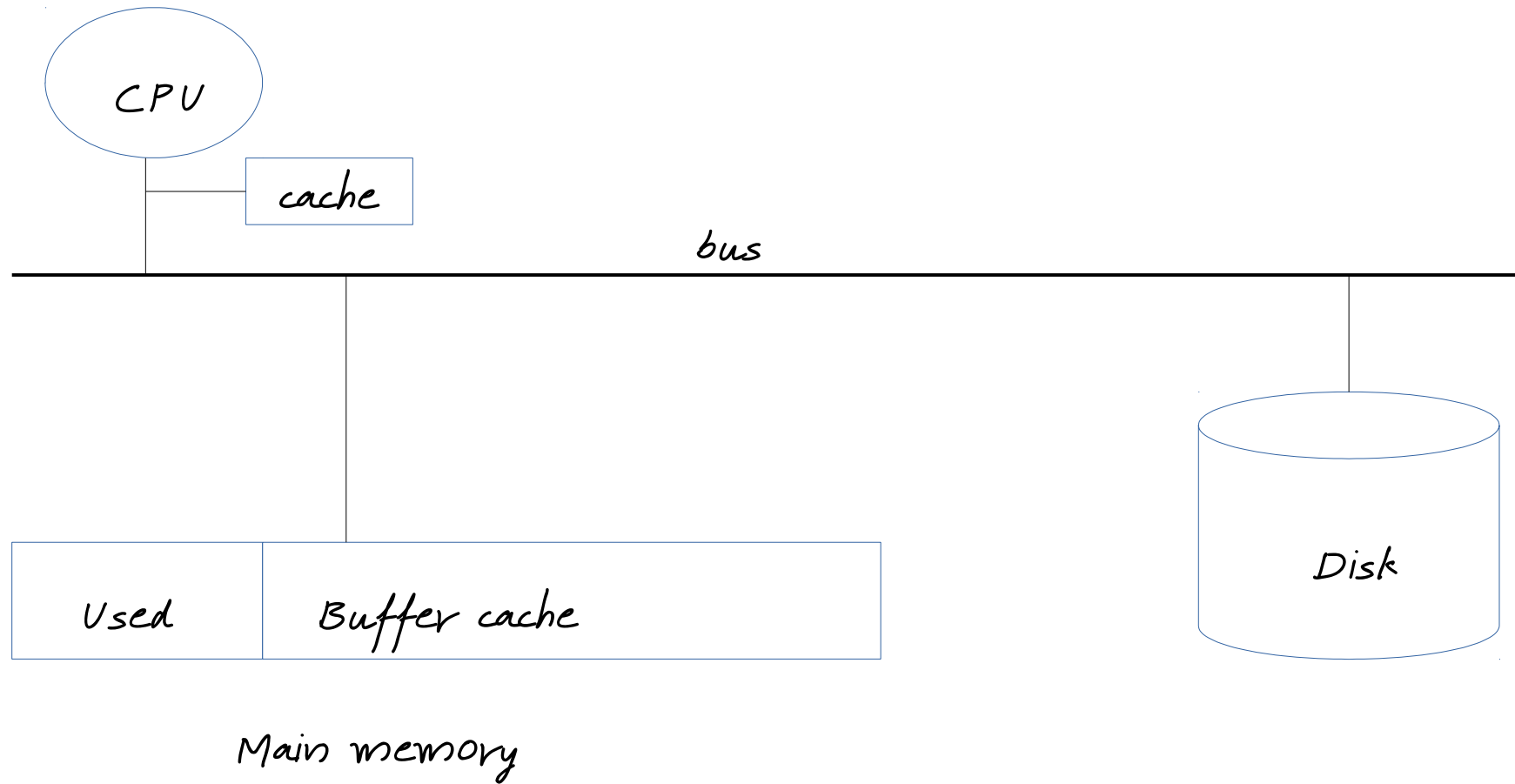


# Testing

- Tools like Bonnie++ or iozone are fine but the best test is doing what your users are going to do.
- What applications are going to run to transfer data?
- How big are files typically
- What block sizes are used to read/write the data?
- Do people do random I/O or not?
- There is one other thing,...

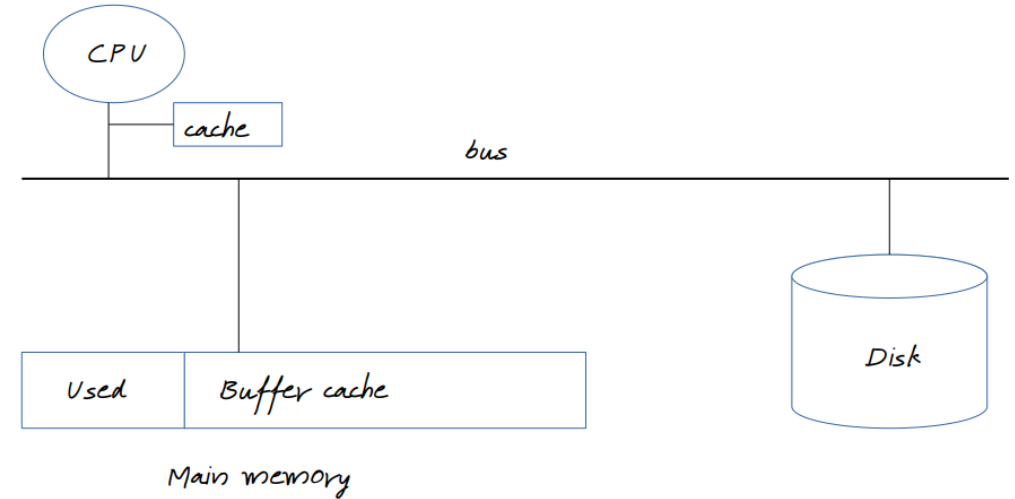


# The Buffer Cache



# Testing

- Reading from disk is a lot slower than from memory
- Caching in memory is great
- Written files are often read back soon after
- Save source code in your editor and start compiling it
- Save a pdf you have got by mail and opening it
- Files are often read more than once in a short time span
- Forward an email
- Reduce the effects of the file system buffer cache as much as possible while testing. You want to test the disk subsystem, network card and CPU, not memory





# Testing: dd

- Test only the disk system

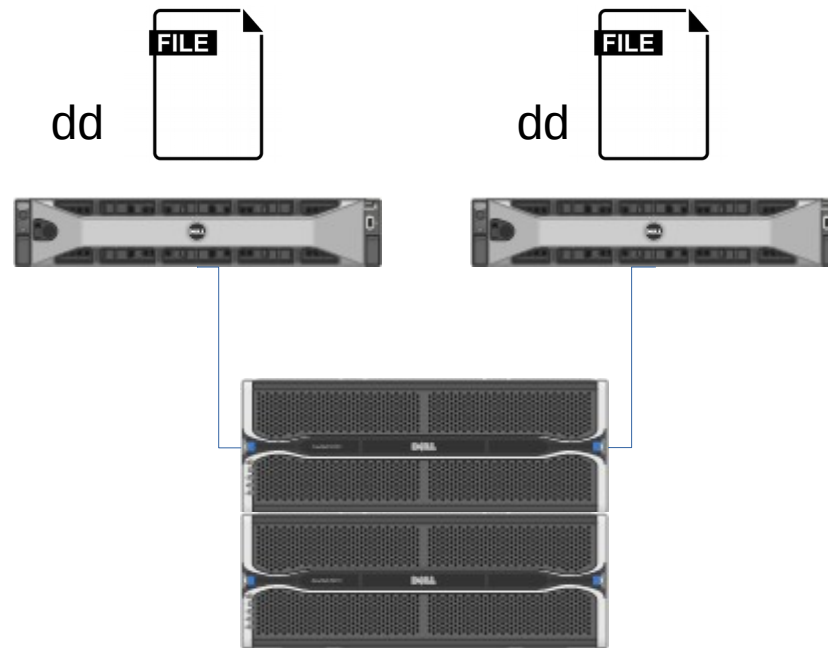
- Receipe:

- Create a number of testfiles. The size of the testfiles should be what you will expect in production.
- In total, the amount of space taken by the testfiles should be considerably larger than the amount of memory in the machine. We use 25x more than there is memory.
- Pick at random files to read for testing. Chances are slim that you pick a file that is still in memory from a previous test.
- Spawn a bunch of concurrent processes that are continously reading and writing files using dd. Reading and writing at the same time.
- Run dd's with the same block size as your application would use.
- Let the tests run until the average throughput stabilises.



# Testing: dd

- When multiple nodes are attached to the same disk system, then run the tests on all of them

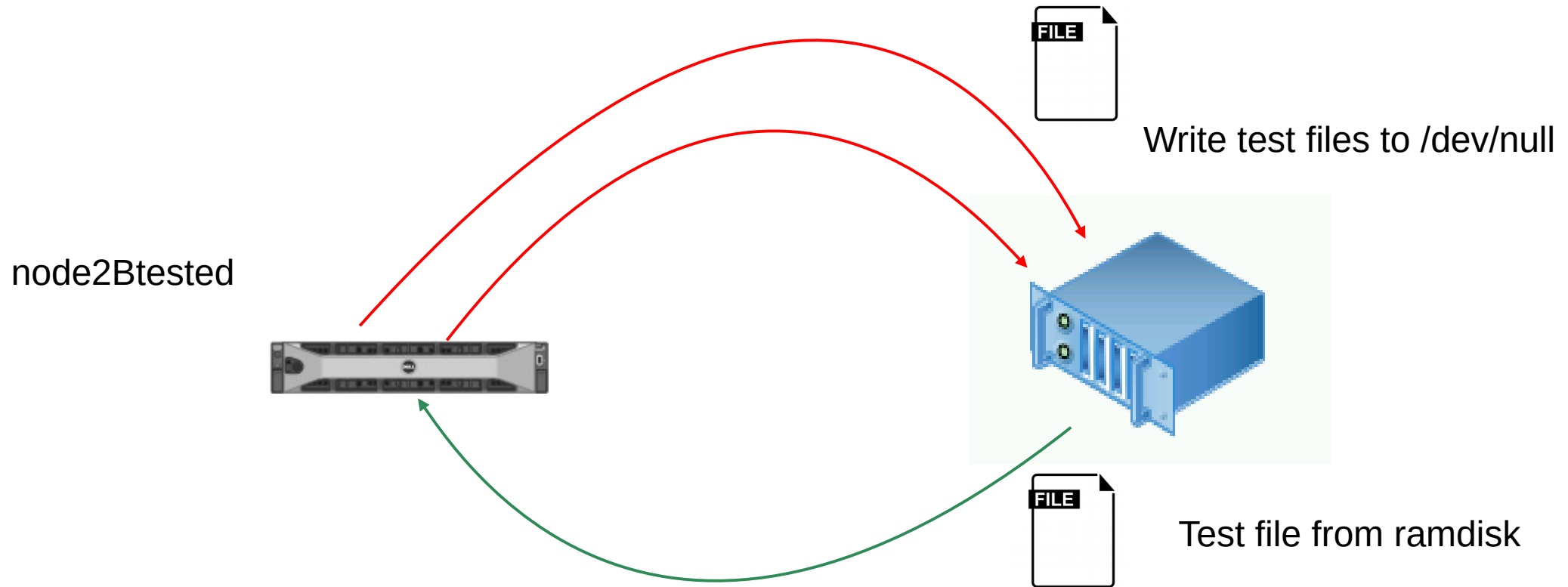




# Testing: gridftp

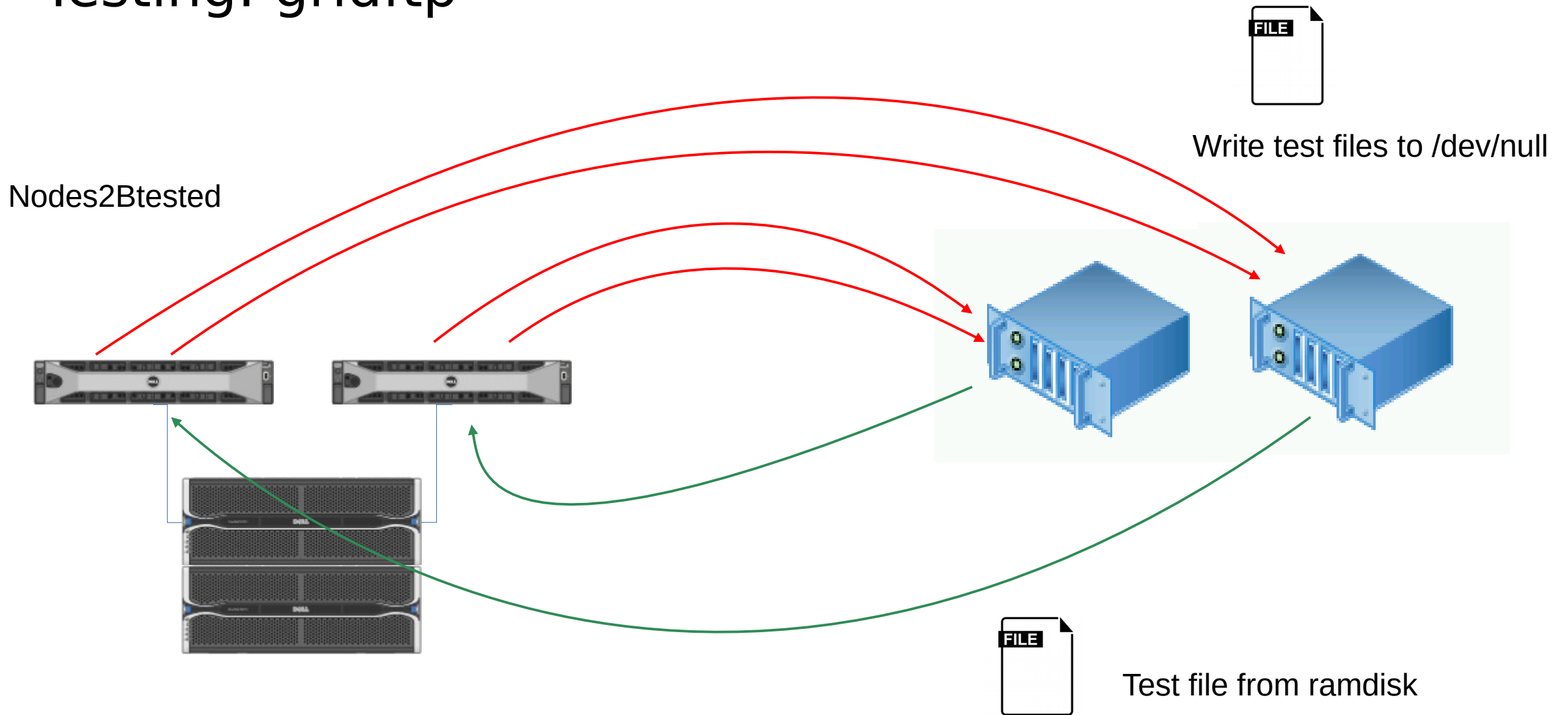
- Test disk and network
- Receipe:
- Same as dd tests.
- Setup gridftp server on the node that is tested and on another node on the same LAN having an equally performant NIC. The another node should not become a bottleneck. Avoid testing the disk system of the other node.
- Spawn multiple concurrent processes that transfer files to and receive files from another machine.
- We are talking about gridftp here but other protocols are fine as well.

# Testing: gridftp





# Testing: gridftp



# Testing

- [https://github.com/SURFsara-e-infra/dd\\_tests.git](https://github.com/SURFsara-e-infra/dd_tests.git)
- [https://github.com/SURFsara-e-infra/gridftp\\_tests.git](https://github.com/SURFsara-e-infra/gridftp_tests.git)



# Tuning

- Kernel IO scheduler
- `cat /sys/block/sd???.queue/scheduler`
- To our experience, when you use a hardware RAID controller, then this controller is probably a lot smarter than you are in controlling what HDDs do. We use noop.
- [https://access.redhat.com/documentation/en-us/red\\_hat\\_enterprise\\_linux/7/html/performance\\_tuning\\_guide/chap-red\\_hat\\_enterprise\\_linux-performance\\_tuning\\_guide-storage\\_and\\_file\\_systems#sect-Red\\_Hat\\_Enterprise\\_Linux-Performance\\_Tuning\\_Guide-Considerations-Generic\\_tuning\\_considerations\\_for\\_file\\_systems](https://access.redhat.com/documentation/en-us/red_hat_enterprise_linux/7/html/performance_tuning_guide/chap-red_hat_enterprise_linux-performance_tuning_guide-storage_and_file_systems#sect-Red_Hat_Enterprise_Linux-Performance_Tuning_Guide-Considerations-Generic_tuning_considerations_for_file_systems)

# Tuning

- Set MTU to 9000 (normally MTU=1500)
- Better performance on high bandwidth networks
- Faster recovery from packet loss
- MTU=9000 all the way, but what if it isn't?
- Fragmentation leading to overhead and performance degradation
- DF (don't fragment) bit set, then no fragmentation, packet is dropped and ICMP packet send back so sender can adjust MTU
- Sometimes ICMP is blocked → MTU black hole
- `net.ipv4.tcp_mtu_probing = 1`
- Disables MTU probing and turns it on when a black hole is detected
- `net.ipv4.tcp_congestion_control = htcp`
- <http://fasterdata.es.net/host-tuning/linux/>



# Tuning

- Send and receive buffer sizes
- `net.ipv4.tcp_rmem, net.ipv4.tcp_wmem`
- Modern kernels are getting better in auto tuning but it may be worthwhile to look at these autotuning settings
- Send and receive buffer sizes
- Example: `net.ipv4.tcp_rmem='4096 87380 8388608'`
  - 4096= minimum buffer for each tcp connection
  - 87380=default buffer size
  - 8388608=maximum receive buffer for tcp connection
- [https://wwwx.cs.unc.edu/~sparkst/howto/network\\_tuning.php](https://wwwx.cs.unc.edu/~sparkst/howto/network_tuning.php)
- <https://paulgrevink.wordpress.com/2017/09/08/about-long-fat-networks-and-tcp-tuning/>

# Tuning

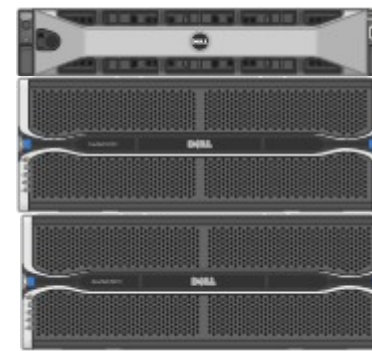
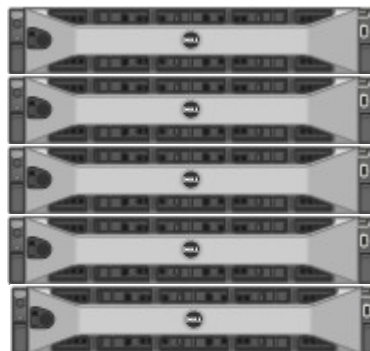
- Linux vm tuning
- vm.dirty\_background\_ratio, vm\_dirty\_ratio, vm.dirty\_expire\_centisecs, vm.dirty\_writeback\_centisecs
- Fast disk system with battery backedup RAID controller:
  - vm.dirty\_background\_ratio = 5
  - vm.dirty\_ratio=10
- [https://lonesysadmin.net/2013/12/22/better-linux-disk-caching-performance-vm-dirty\\_ratio/](https://lonesysadmin.net/2013/12/22/better-linux-disk-caching-performance-vm-dirty_ratio/)
- Max connections
- net.core.somaxconn (default 128)



# Concluding remarks

- Hardware considerations

- Set of smaller nodes
- Nodes with separate disk system



- HDDs, SSDs, NVMe

- For what applications?
- €€€?

- Scaling up

- Data migration scenario when hardware is decommissioned (outages? And for how long)

- I/O requirements

- We use 12MiB/s/TiB

**NOT SURE IF PRESENTATION WAS SO GOOD  
NO ONE HAD ANY QUESTIONS**



**OR NO ONE WAS PAYING ATTENTION**