# EOS Data Lake

Luca Mascetti, Massimo Lamanna

# "WLCG" DATA LAKE
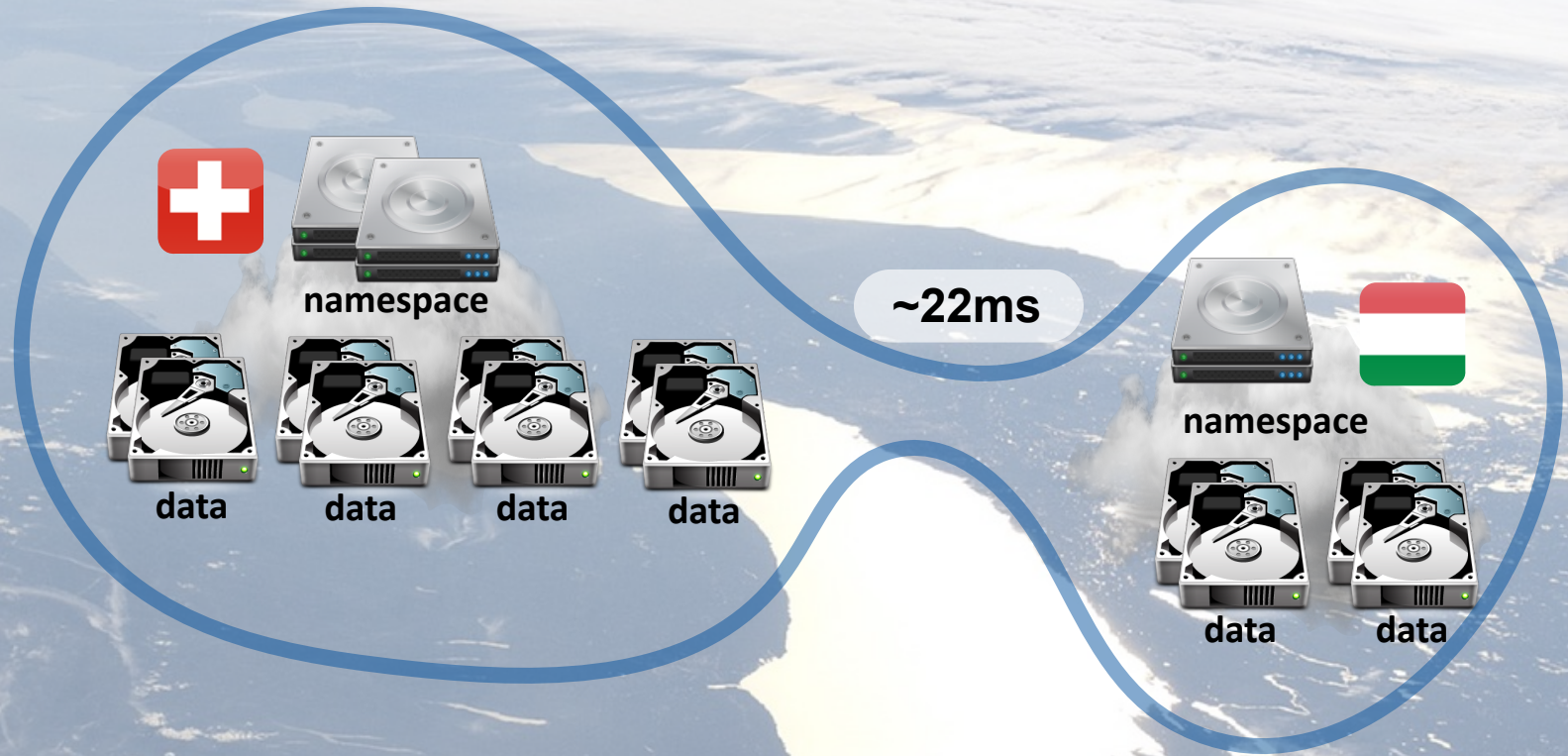
```
n sites
k replicas
k<<n
latency > 1ms
```
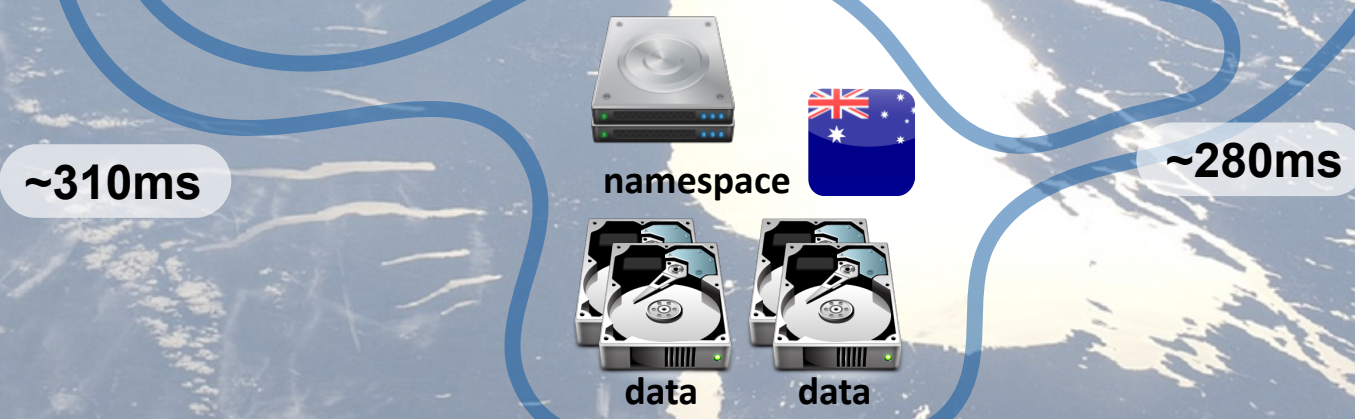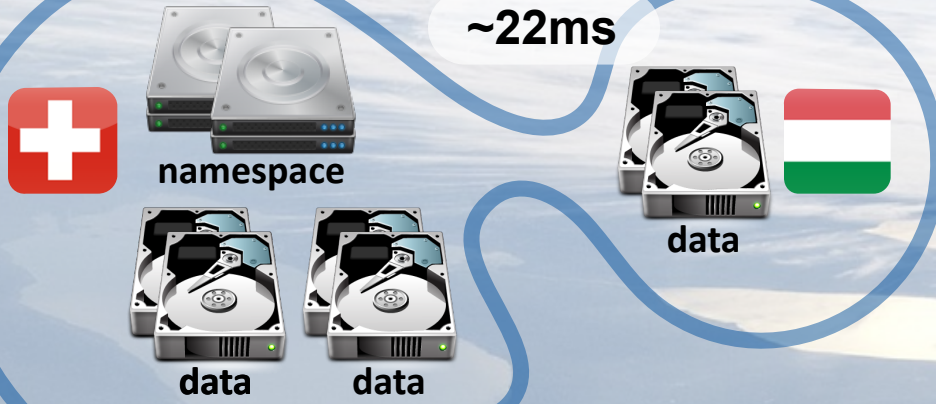
EOS

namespace

~22ms

namespace

data    data    data    data

data    data

~22ms

namespace

data

data

data

~300ms

data

~310ms

namespace

~280ms

data

data

data

5

# Storage Pools Overview

1GB File

client

```
/eos
    /asia
        /taiwan
    /australia
        /melbourne
    /europe
        /geneva
        /budapest
    /dualcopy
        /gva-bud
        /mel-gva
        /mel-bud
    /triplecopy
        /mel-gva-bud
        /mel-gva-tpe
```

550 MB/s (noauth)
470 MB/s (auth)

3 nodes GVA
2 nodes BUD
6 nodes MEL
2 nodes TPE

45 MB/s (auth)

30-35 MB/s (auth)

Storage pools were created with filesystems from all four sites.
Files were replicated according to the different configured policy
(e.g. 3 replicas: MEL-GVA-TPE).

CERN

# Make data access easy
# Make Analysis simple
# Facilitate Science



- Scale-out filesystem underneath the ownCloud app, using the eosd fuse interface for file IO
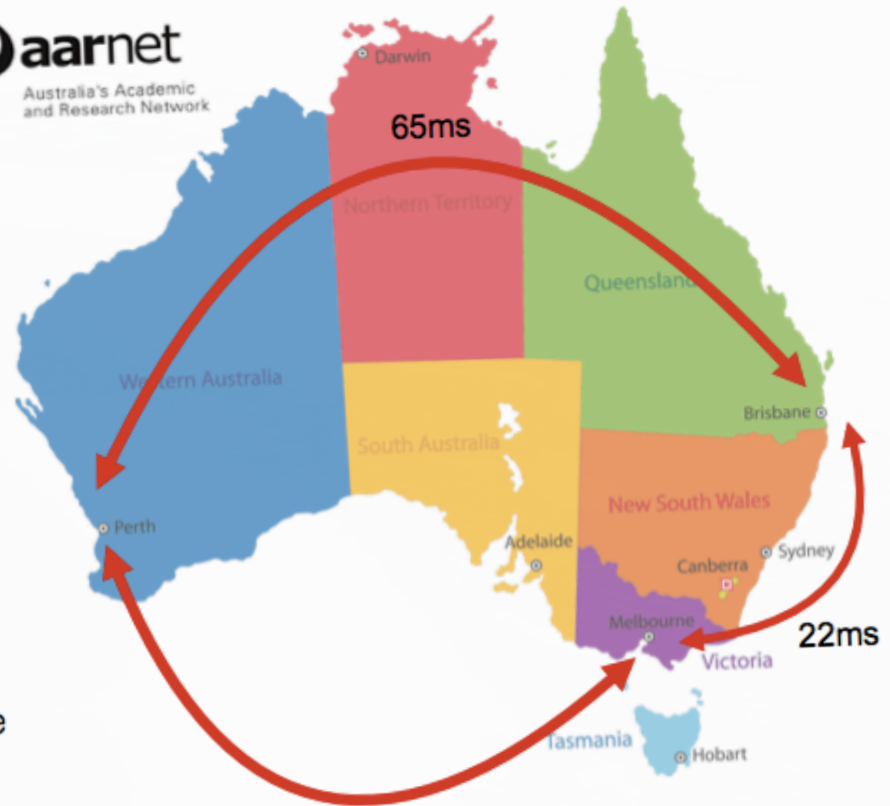
- Geo-distributed setup: Brisbane, Melbourne, Perth
  - ~1PB (scale to ~20PB next year)

- Australian National University, in Acton Canberra: mirror archives of both genome sequences and open or freely available software distributed among three sites

*"This system is presently running 0.3.187, and has been so trouble free that I keep forgetting it's there." David Jericho -AARNetSolutions Architect*
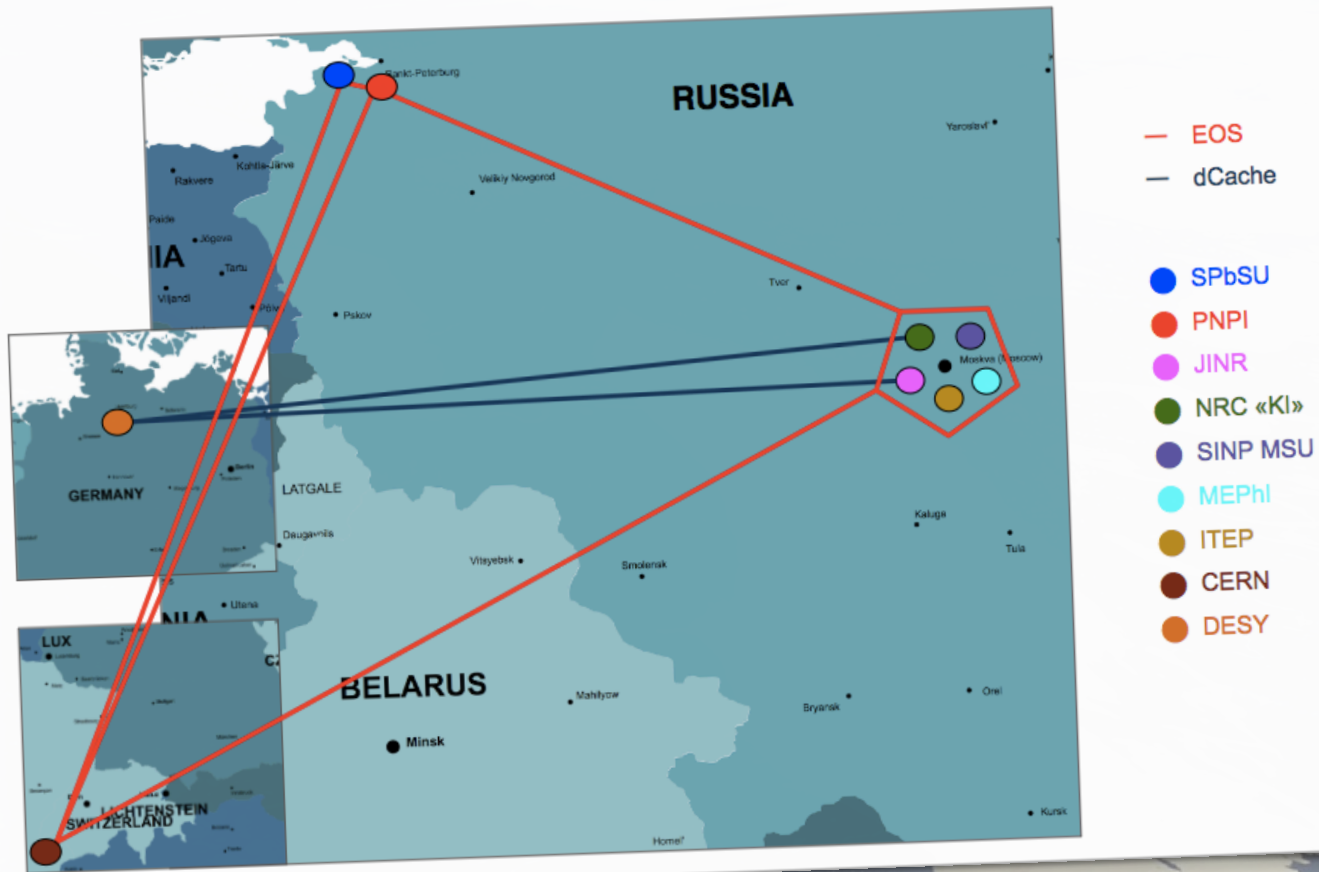
**aarnet** Australia's Academic and Research Network

**aspera** **RSYNC**

Australian National University

65ms
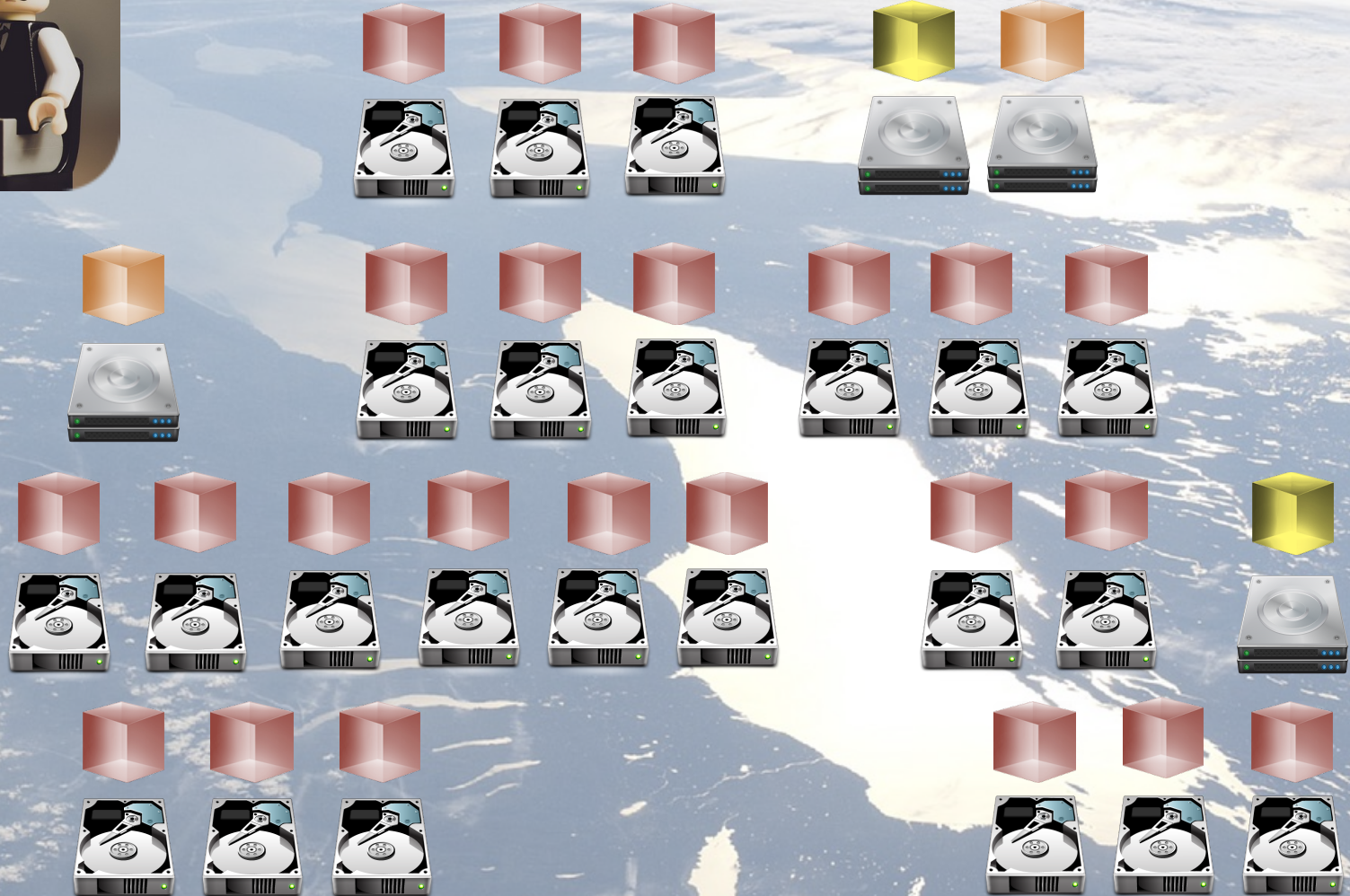
22ms

Background image: Shutterstock

# Storage Management
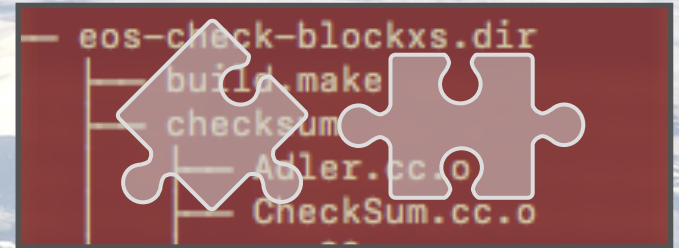
Are we able to reduce the complexity
and the expertise globally required?

# EOS

**Decoupling storage software and hardware maintenance**

```
├── eos-check-blockxs.dir
│   ├── build.make
│   ├── checksum
│   │   ├── Adler.cc.o
│   │   ├── CheckSum.cc.o
```

API*

* Already in use by CERN repair team and sysadmins team

# Network Evolution

in-lake
connectivity

Current Grid Model

Resources Provided        Expertise Required

Resource and Manpower Consolidation

Data Lake Model

Resources Provided        Expertise Required

# Thanks for the attention!



www.cern.ch