

Big Data Management

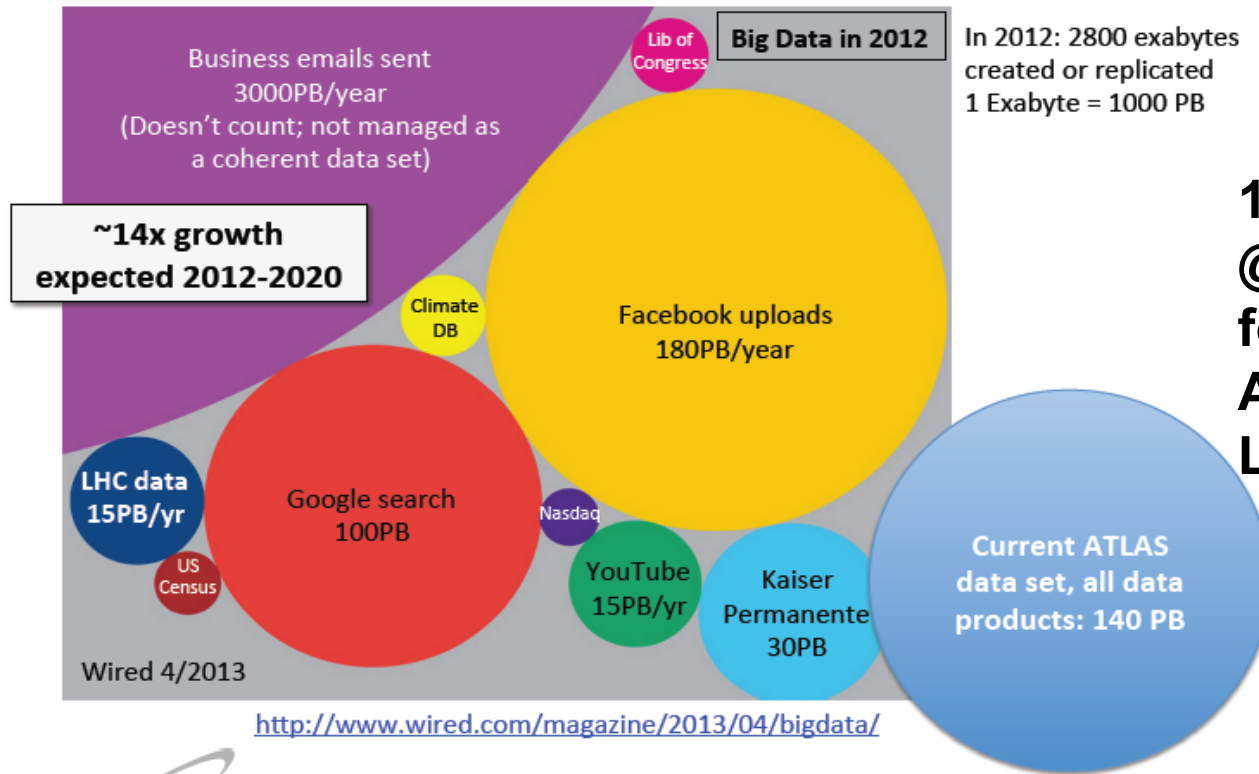
for large Experiments @ DESY

Volker Guelzow
IT-Gruppe
Bruxelles, May 8th, 2014

Big Data in High Energy Physics

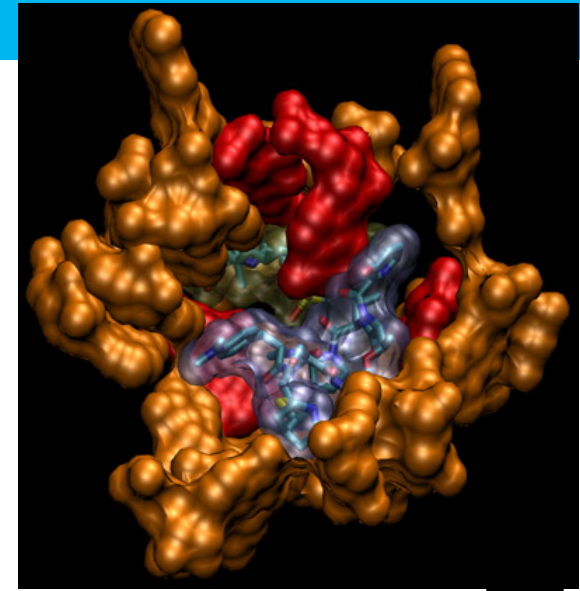
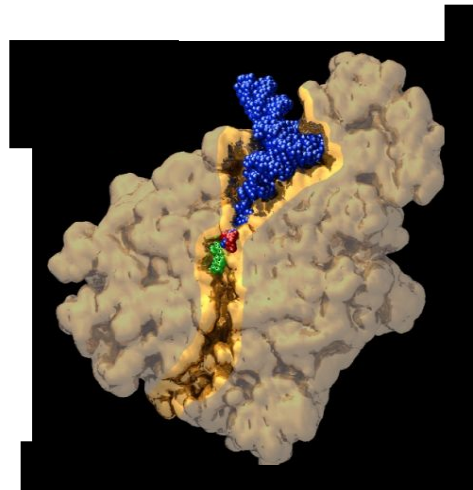
From: Torre Wenaus, CHEP 2013

Data Management Where is LHC in Big Data Terms?



**10 PB HEP Data
@ DESY Tier 2
for
Atlas, CMS,
LHCb**

Big Data in Photon Science



> Examples:

- Simulation of nano structures -> $N^3 \log(N)$
- Protein Crystallography -> 4 MPixels * 10^6 Crystals -> analysis of hundreds TB
- Tomography -> 3D Reconstruction on some 10th TB data/experiment
- Online analysis with some TB/beamline and some Gb/s/beamline
- Offline analysis for some PB/year

Detectors and Data rates

Type	Frame size	Frame rate	Peak rate	Os	Avail.
Medipix	n x 256 x 256 x 2			SL6	Now
PerkinElmer	2048 x 2048 x 2	15 Hz	1.9Gb/s + 8kb/s (log)	Windows-7	Now
Pilatus300k	487 x 619 x 4	~200 Hz	1.56Gb/s	Suse-10	Now
Pilatus1M	981 x 1043 x 4	25 Hz	0.8Gb/s	Suse-10	Now
Pilatus 6M	2463 x 2527 x 4	25 Hz	4.6Gb/s	Ubuntu 10	Now
AGIPD	128 x 512 x 2 x 2 x 352	10 Hz	55 Gb/s	SL6	2015
Eiger	1k x 1k x 2	3 kHz	50 Gb/s	RHEL6	now
Lambda-Si	3 x 1536 x 512 x 2	2 kHz	60 Gb/s	SL6	2013
PCO Edge	2560 x 2160	100 Hz	5.6 Gb/s	Windows-7	now
Percival (1S)	4k x 4k x 2	120 Hz	60 Gb/s	(SL6)	2015
Percival (4S)	8k x 8k	120 Hz	240 Gb/s	(SL6)	late 2015

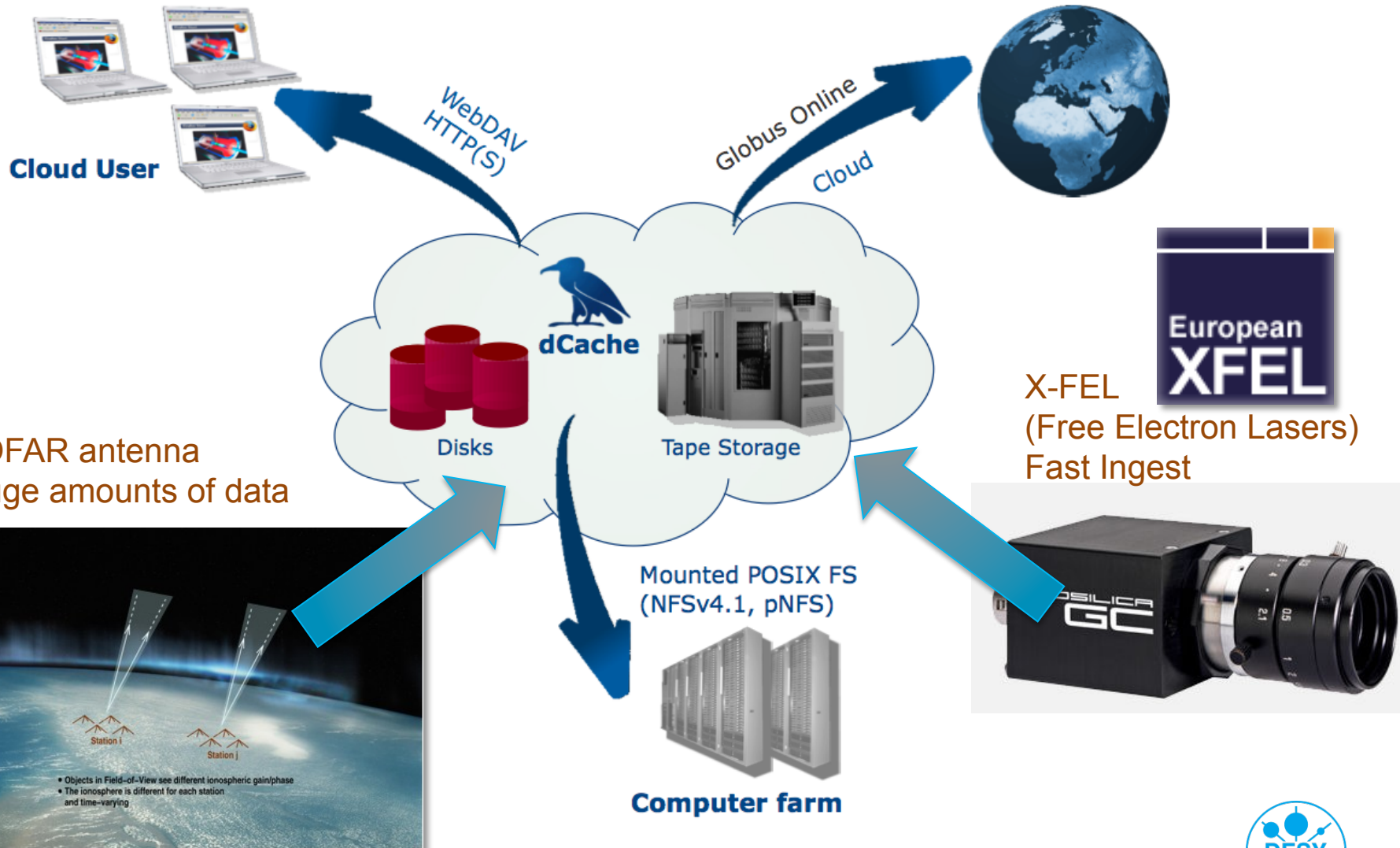


The European XFEL

- Operational 2016
- Expect: 20 PB 2016 → 100 PB (2019) per year
- Big Data Management through dCache
- Tier0 @ DESY
- Distributed Analysis

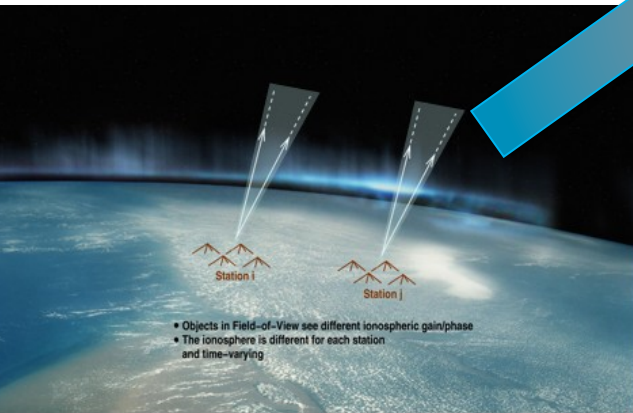


dCache Big Data Cloud



LOFAR antenna
Huge amounts of data

X-FEL
(Free Electron Lasers)
Fast Ingest

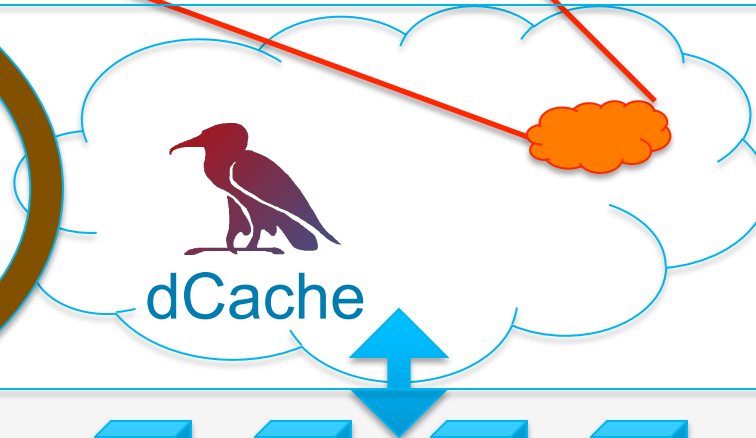


dCache – OwnCloud Data Management

WEB 2.0



Unlimited hierarchical
Storage Space
NFS 4.1
CDMI

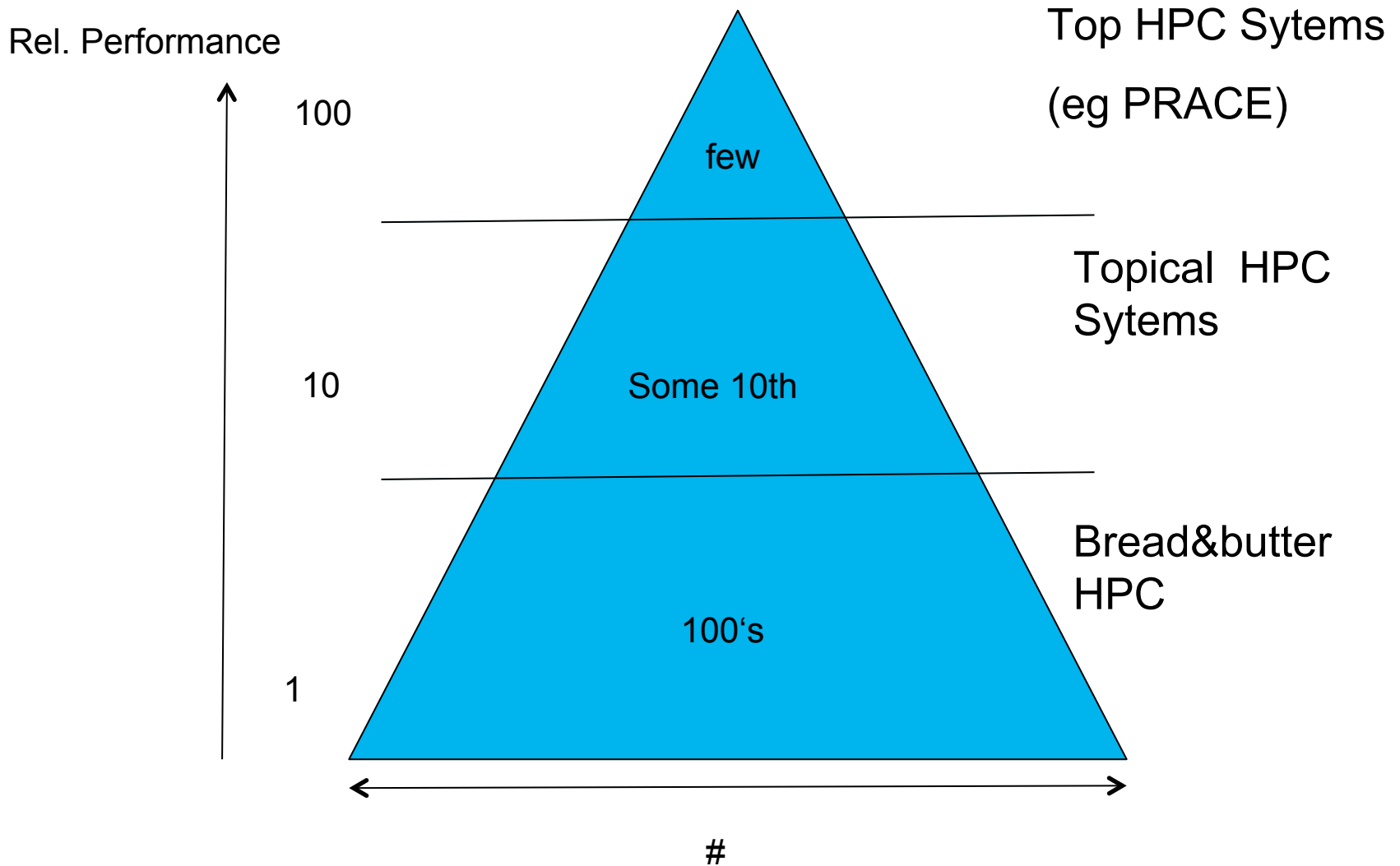


SSDs

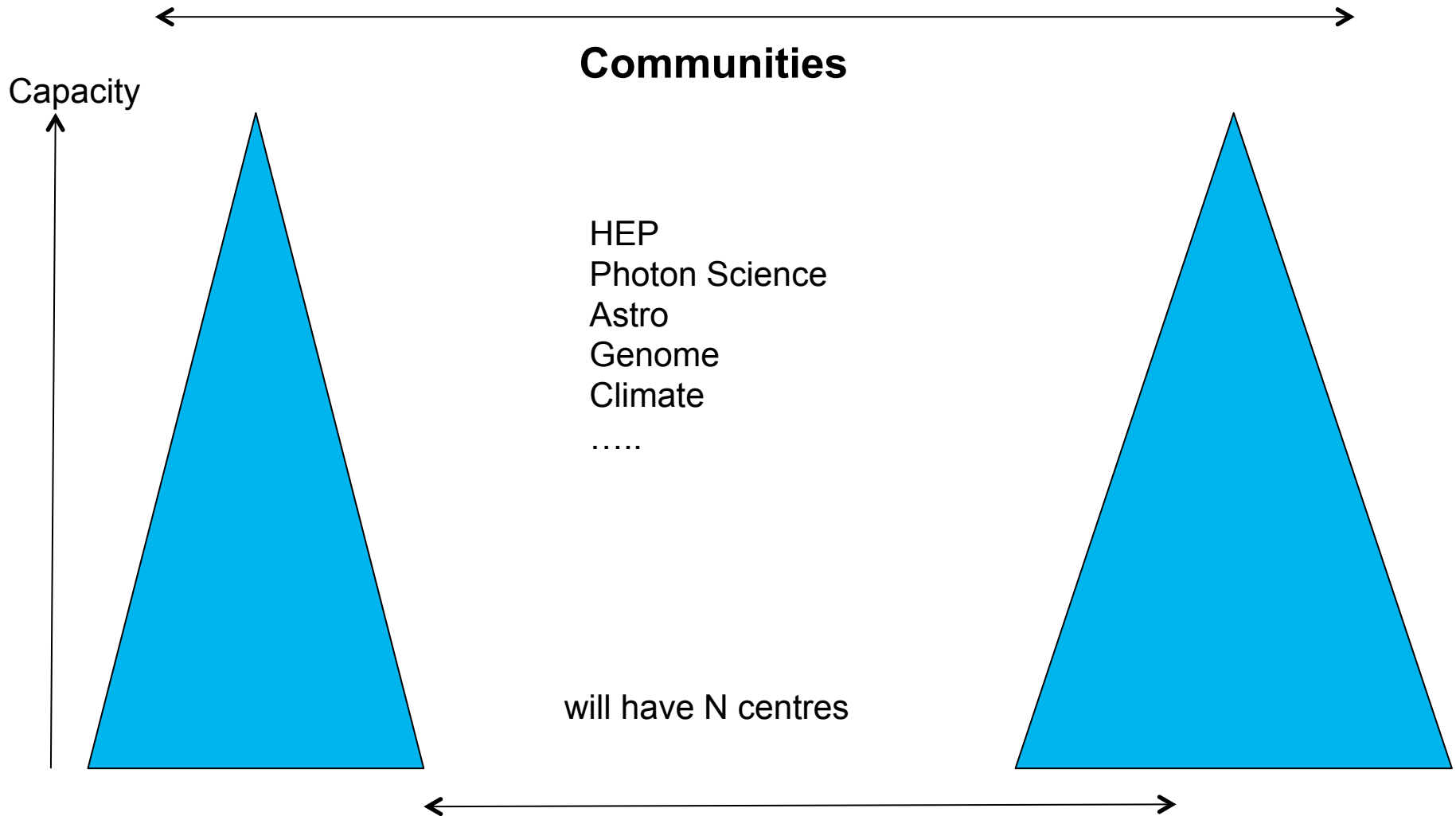
Spinning Disks

Tape, Blue Ray ...

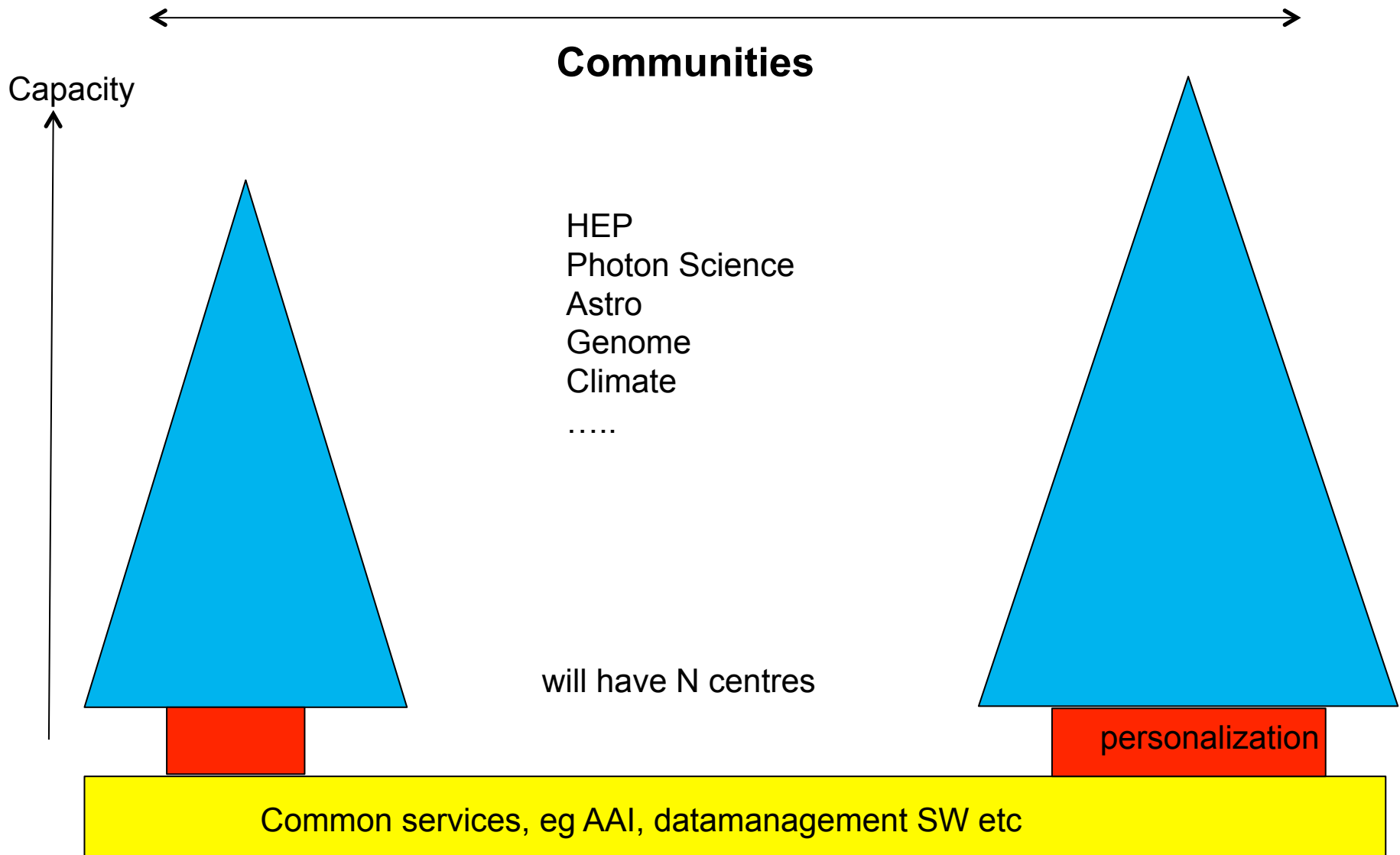
The HPC Pyramide, a european model for „big data“?



The federated data infrastructure model for science



The federated data infrastructure model for science



User interface layer

- > AAI: access to federated resources
(adapting existing solutions from other projects
like EduGain based solutions)
- > Security, policies, life cycle management
(like „how long? Ownership? Who has access? What kind of
data?..)
- > Portals for scientific and industrial users
(like access to resources, virtual accounting, industrial usage,..)
- > Open Access



- > Cloud solutions

(open stack, Helix Nebula, Dirac-like systems)

- > Provisioning of an advanced computing facility (GPU, parallel, green programming) in a federation
- > Distributed competence through simulation Lab's
- > Cooperation with PRACE and potentially others



Data management

- > Advanced (distributed) baseline storage solutions
(f.i. high performance for ingest und access,...)
- > Advanced data management software solutions
- > Data cloud solutions
- > Metadata handling
- > Distributed competence through data Lab's
- > Data preservation (extension of DPHEP)



Summary

- > Setting up federated structures for scientific users
- > Making Distributed competence accessible
- > Setting up advanced Network strategies and developments (like LHCone) with NREN's, Dante, Terena
- > Making advanced ICT technology available to science
- > Cooperation through RDA

