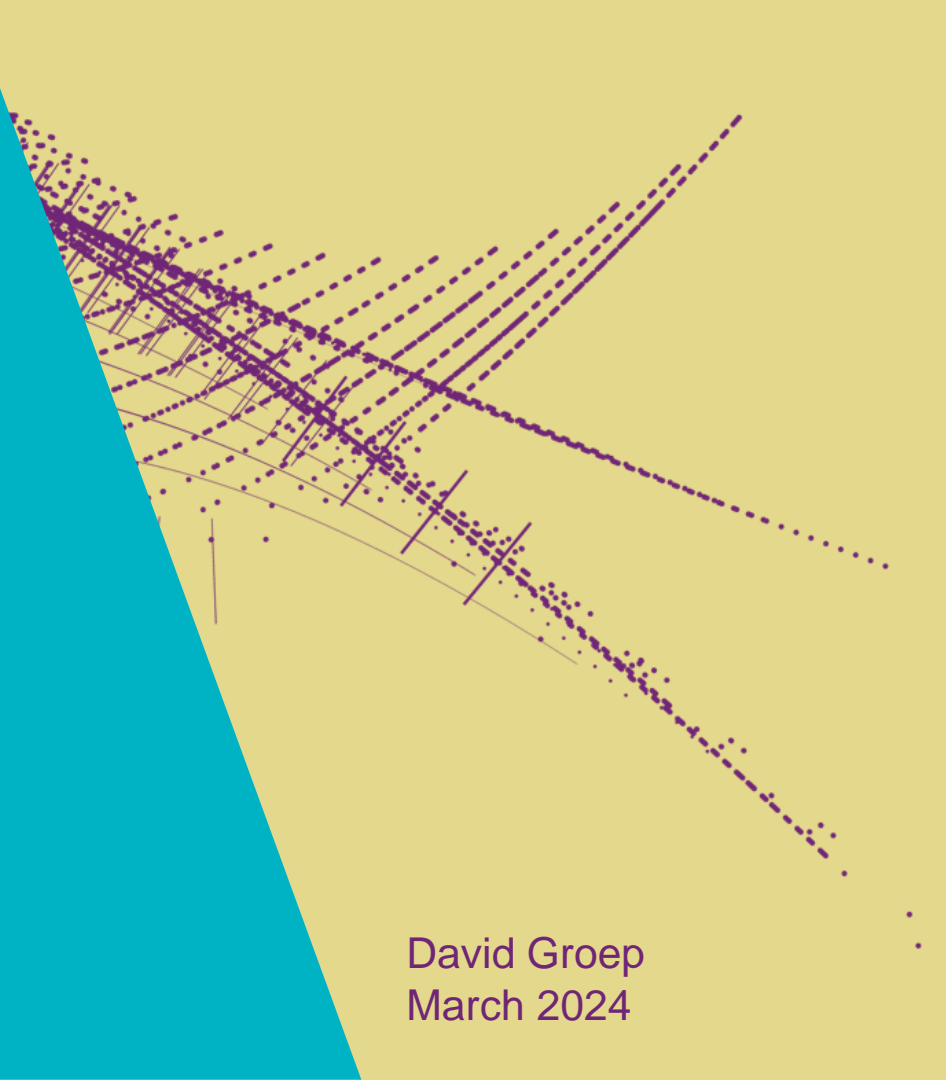




Introduction to Nikhef,
Computing and Networking

Welcome to Nikhef

David Groep
March 2024



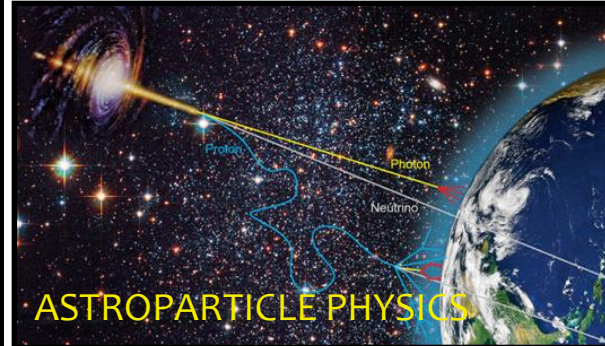
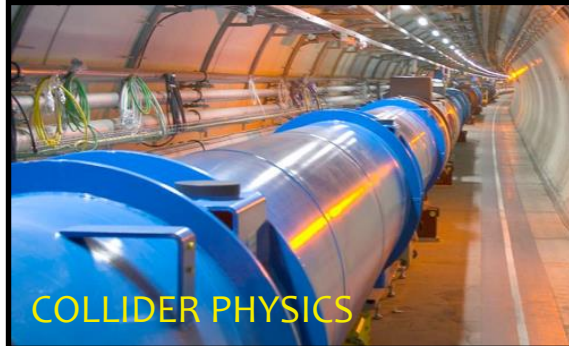
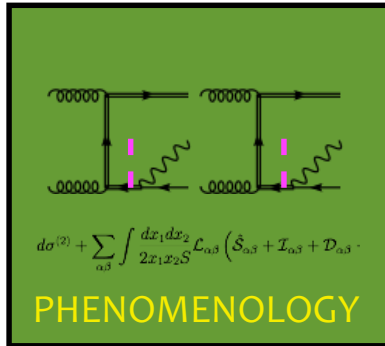
Verleggen van de grenzen van onze kennis

- **Accelerator-based particle physics**

Experiments studying interactions in particle collision processes at particle accelerators, in particular at CERN;

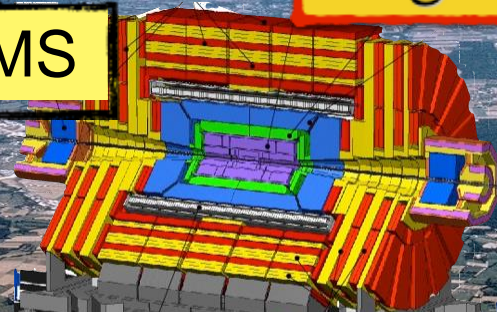
- **Astroparticle physics**

Experiments studying interactions of particles and radiation emanating from the Universe.



Large Hadron Collider

CMS



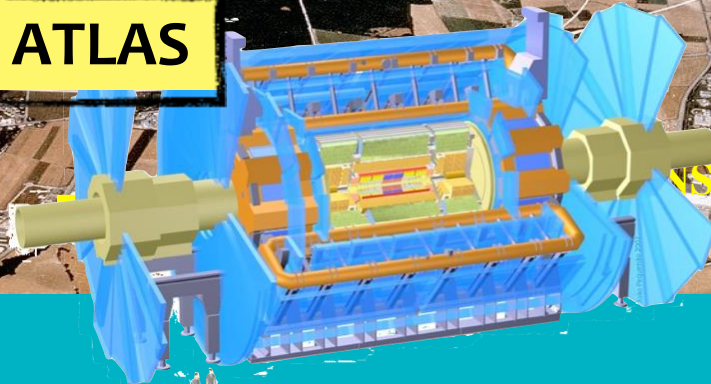
LHCB



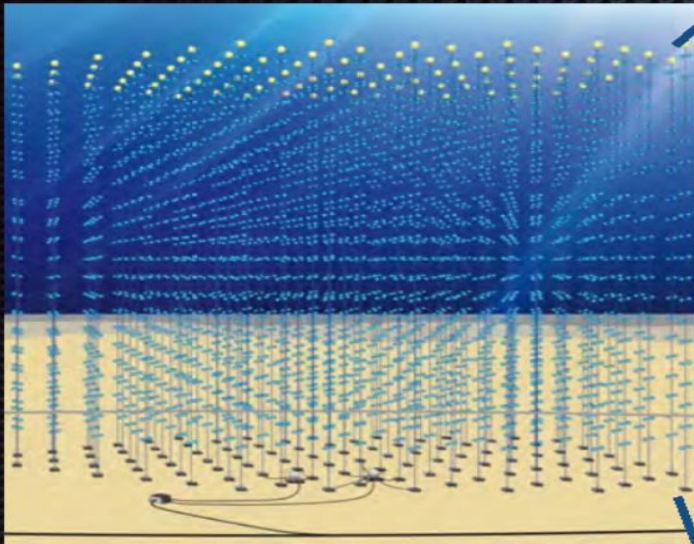
ALICE



ATLAS



Nikhefs neutrino-detector: KM3NeT

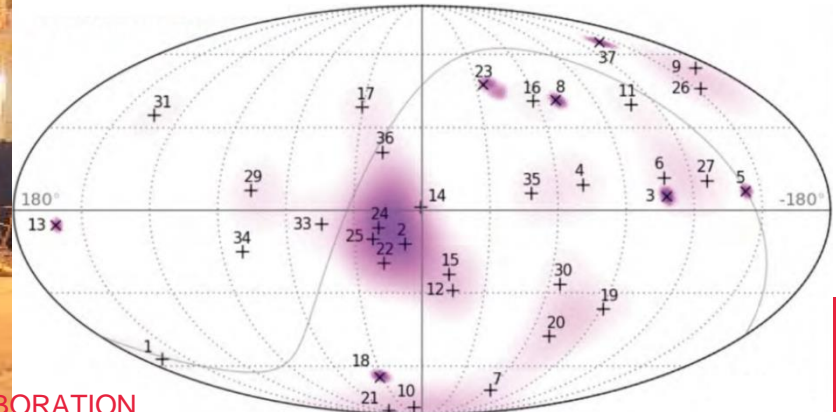


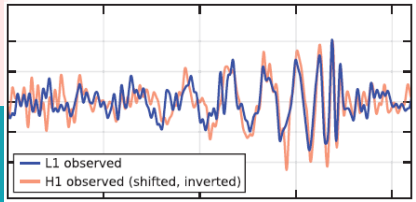
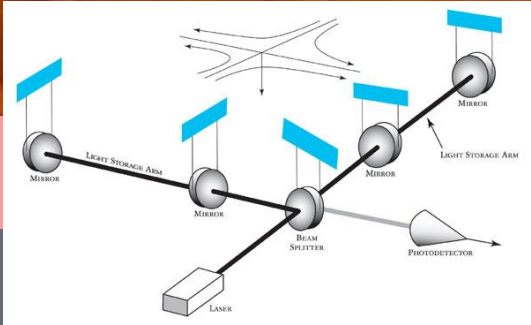


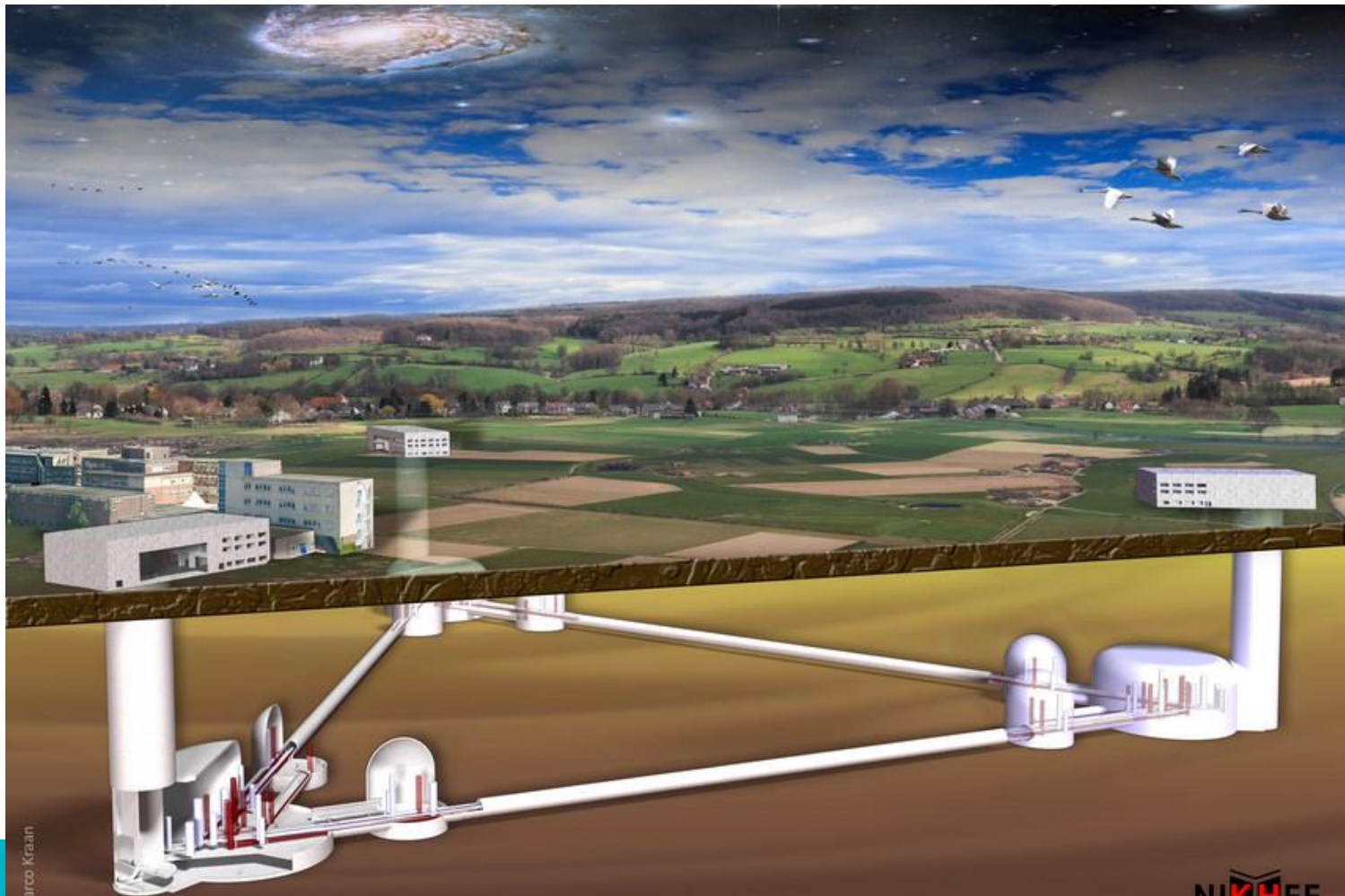
Little white structures prevent the HV bases and cables to touch each other



De Melkweg



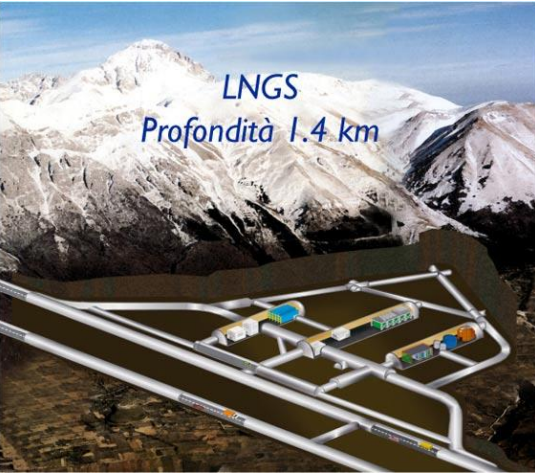
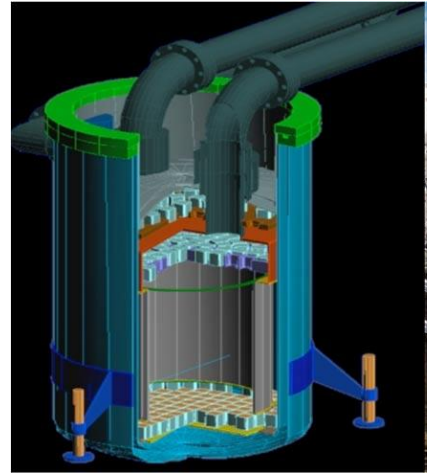
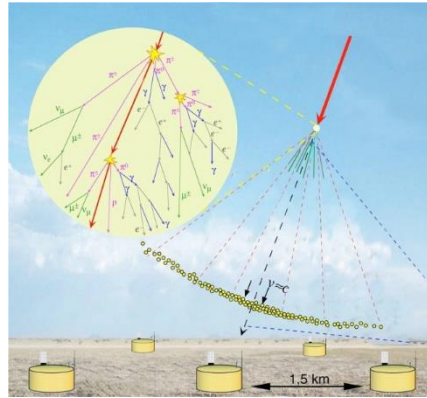
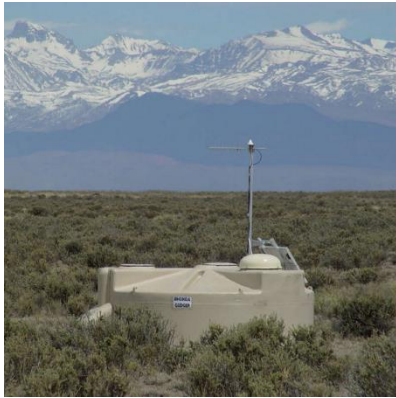




© Marco Kraan

EINSTEIN TELESCOPE – PROJECTED ONTO HEUVELLAND LIMBURG, NL

NIKHEF ikhef

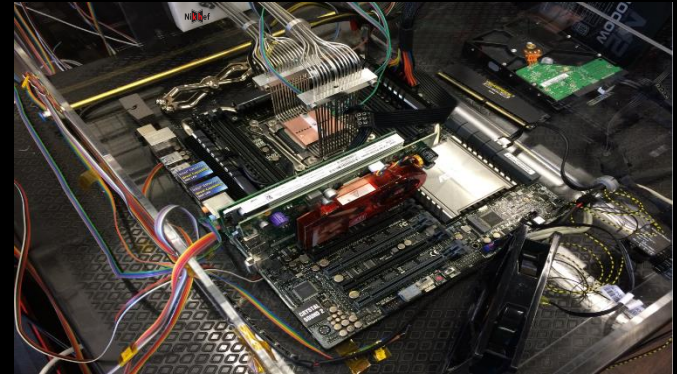
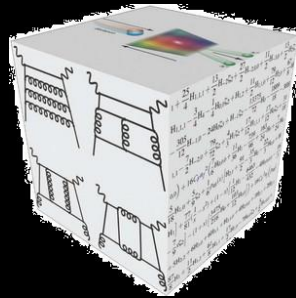


Nikhef collaboration – underpinning programmes



Detector R&D

Theoretical
Physics



Physics Data Processing

PDP 'Physics Data Processing' activities and action lines

Scalable Computing & Algorithms

- algorithms for high-performance software
- data organisation
- accelerator throughput
- rethinking code: at small scale and on large scales

Research Infrastructure & Future Technologies

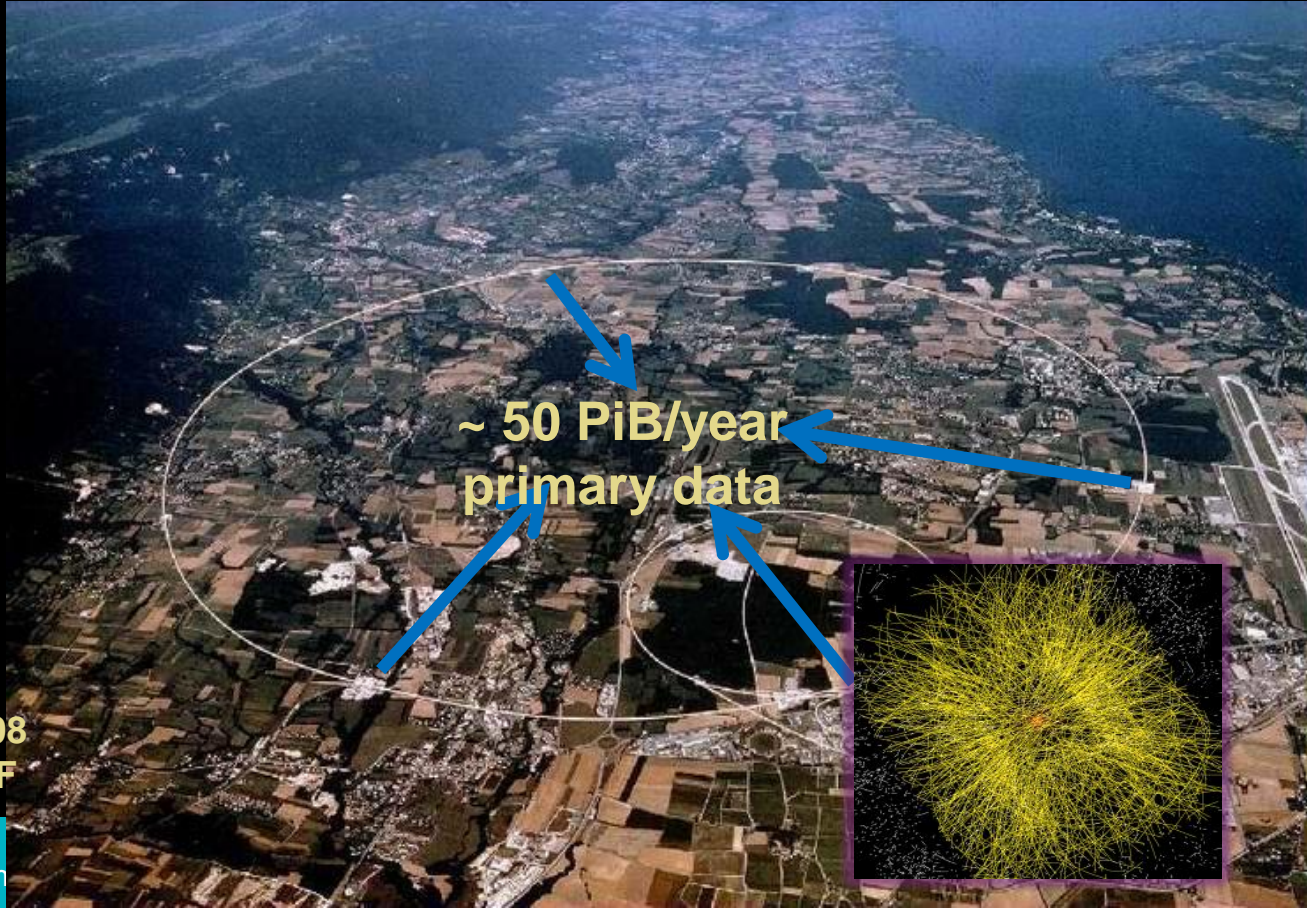
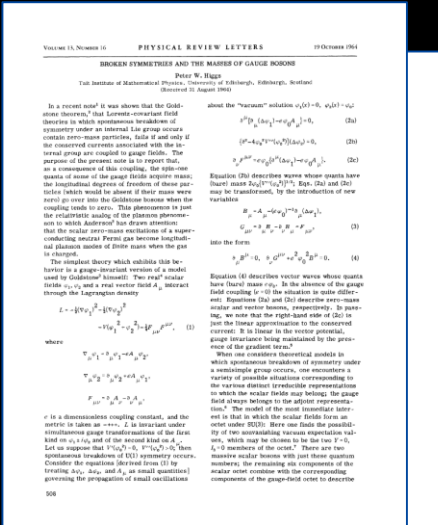
- NDFP processing facility
- National e-Infrastructure *coordinated by SURF*
- experimental next-gen systems engineering
- cross-tier global networks
- stressing & public/private collaborative design

Infrastructure for Secure Collaboration

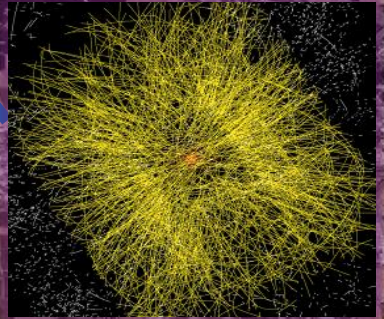
- authentication and authorization protocols
- multi-domain federation
- global trust and identity
- access provisioning
- operational security

Data at the Large Hadron Collider at CERN

1964



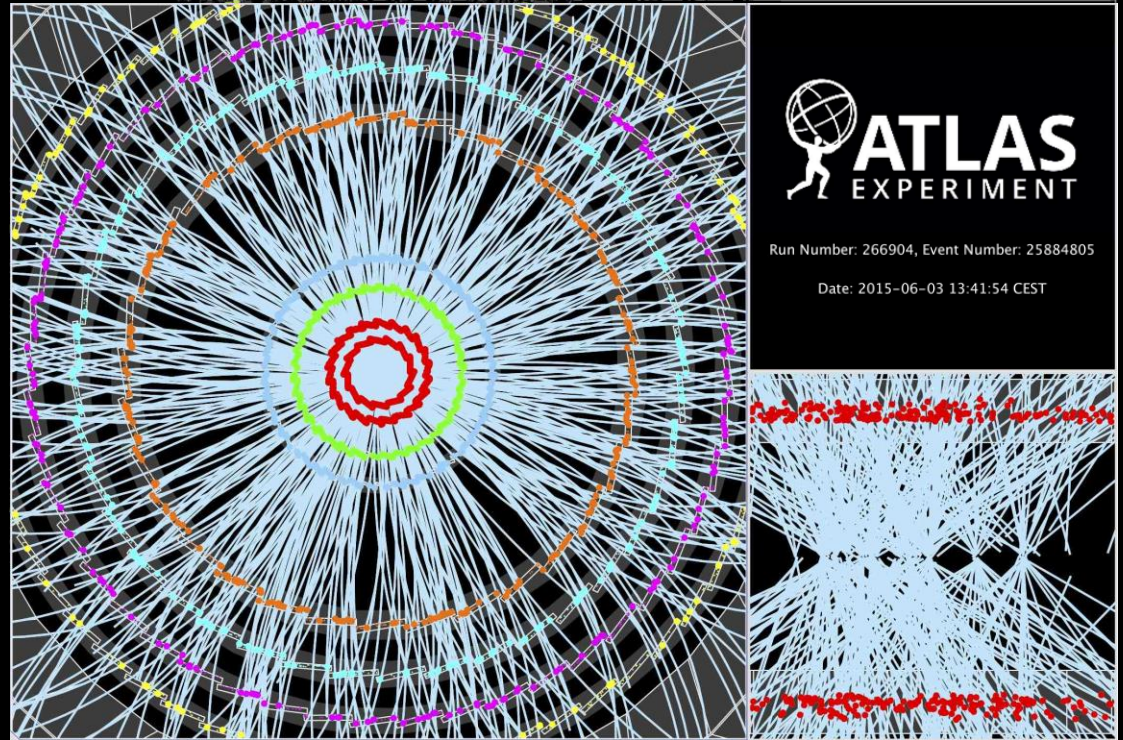
~ 50 PiB/year primary data



P. Higgs, Phys. Rev. Lett. 13, 508
 16823 characters, 165kByte PDF

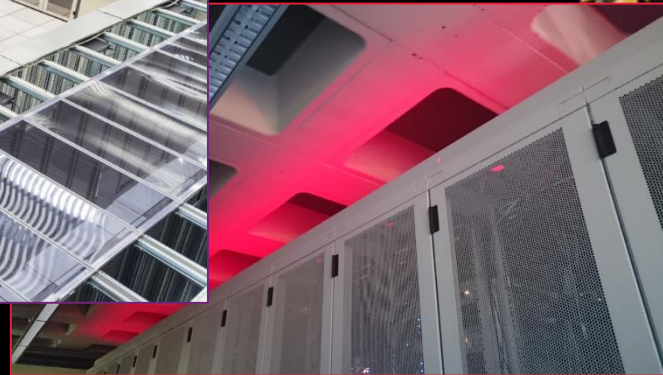
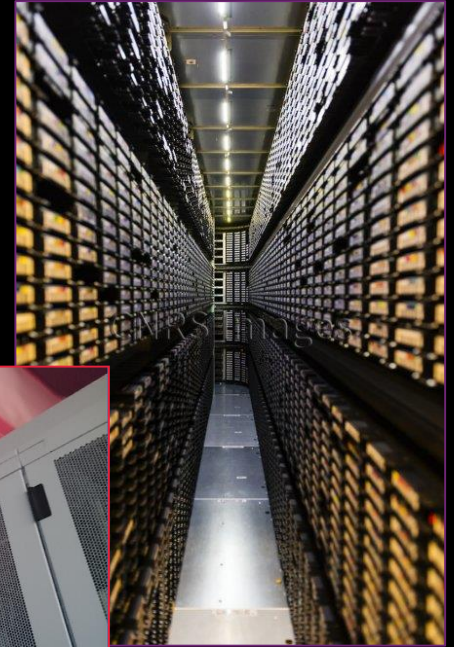
Computing on lots of data – 40 Mevents/sec

~ 10 seconds to compute
a single event at ATLAS
for 'jets' containing ~30
collisions



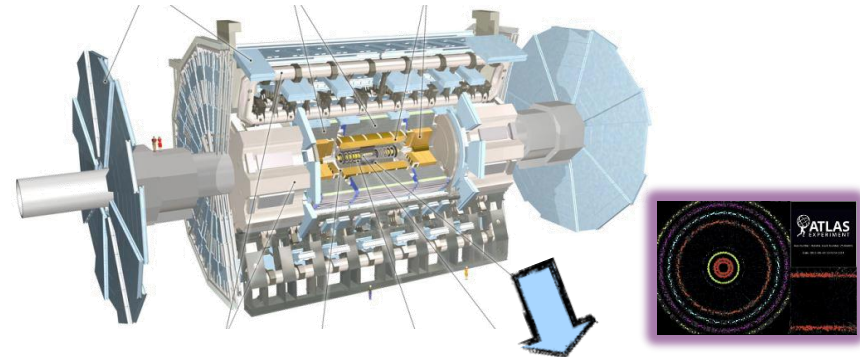
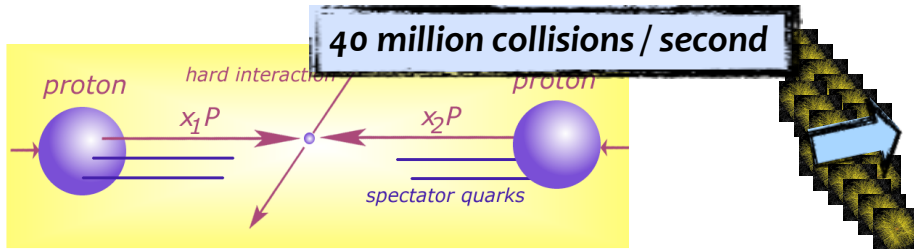
Display of a proton-proton collision event recorded by ATLAS on 3 June 2015, with the first LHC stable beams at a collision energy of 13 TeV;
Event processing time: v19.0.1.1 as per Jovan Mitrevski and 2015 J. Phys.: Conf. Ser. 664 072034 (CHEP2015)

'Big Science' needs some computing ...

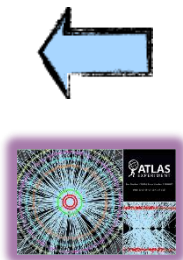
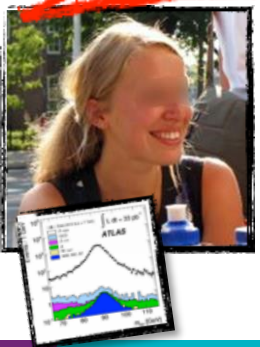


CERN CC B513, image: <https://cds.cern.ch/record/2127440>; tape library: CC-IN2P3 with LHC and LSST data; cabinets: Nikhef H234b

Detector to doctor workflow



Physics analysis by (PhD) students, in papers & analysis notes



Classify particles in collision and their physics properties:

- electrons
- muons
- jets consisting of hadrons
- ...

Trigger system selects 600 Hz ~ 1 GB/s data



diagram adapted from Frank Linde; images: ATLAS collaboration, Nikhef. ... and sorry for the GDPR-blur

Single CPU scaling stopped around 2004

limitation is power, not circuit size

and clock frequency is most 'power-hungry'
still some packages now @ TDP of 400W

multiple cores on the same die helped

AMD EPYC Genoa (Zen 4) has 96 cores on die
Intel Cascade Lake AP looked like a cludge
but now Sapphire Rapids appears better again

CPU design-performance gains left

predictive execution
out-of-order execution
on-die parallelism (multi-core)
pre-fetching and multi-tier caching
execution unit sharing ('SMT')

but at increased risk for security/integrity

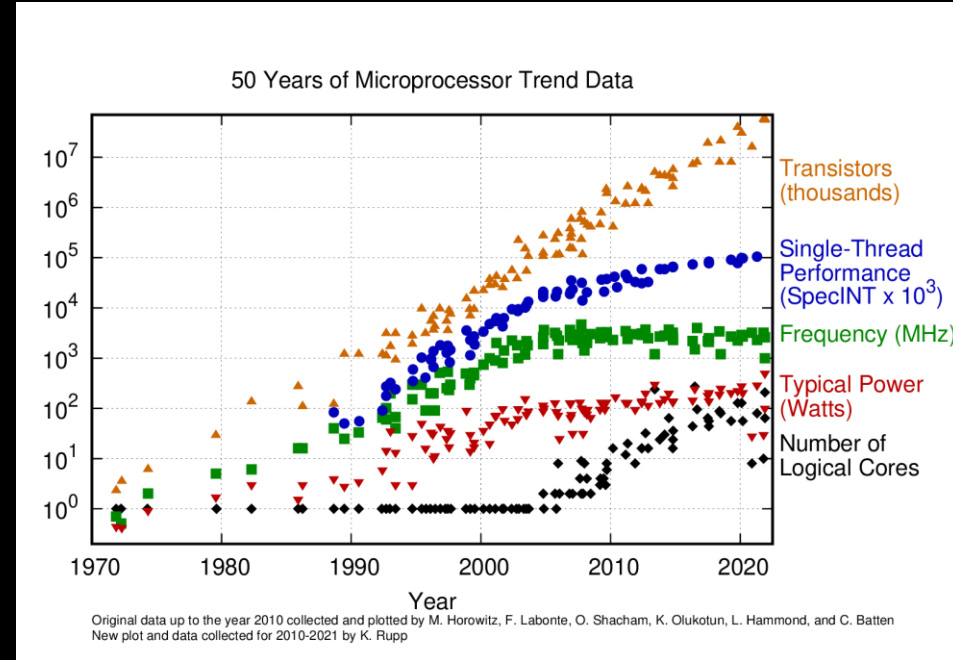
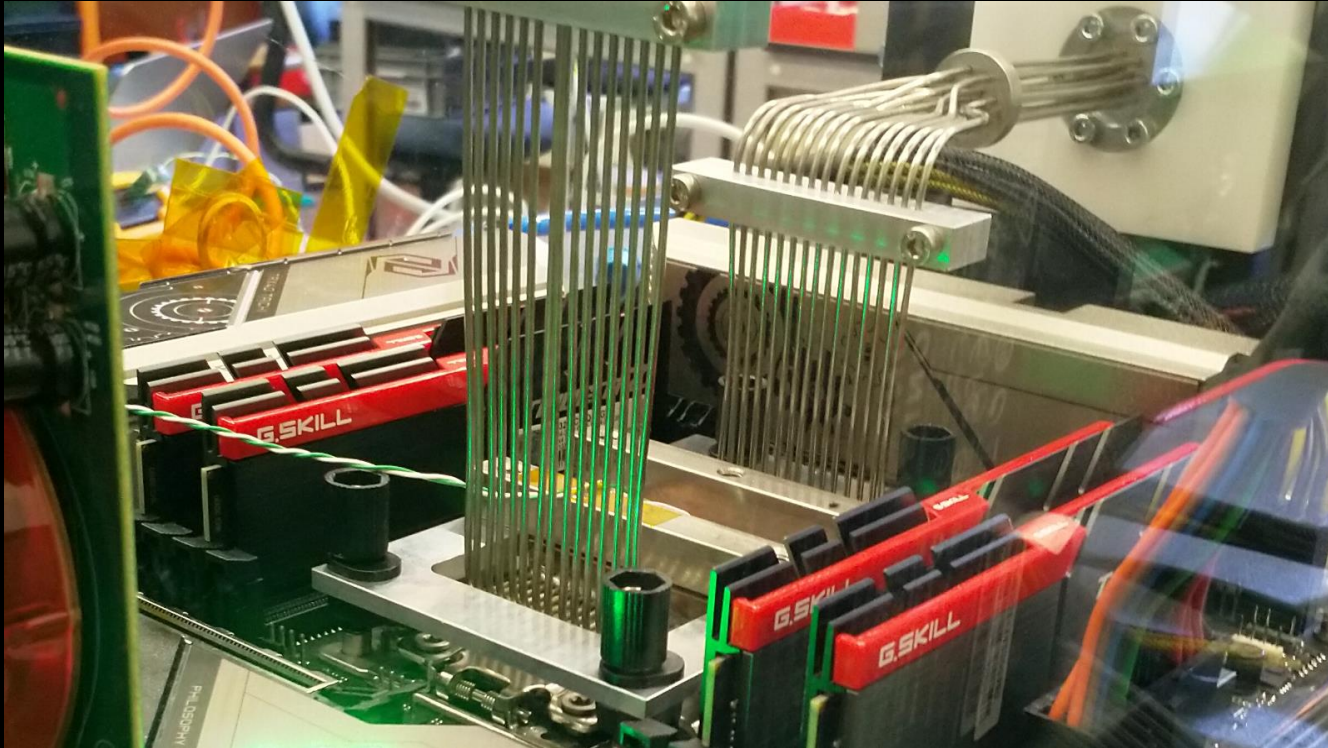


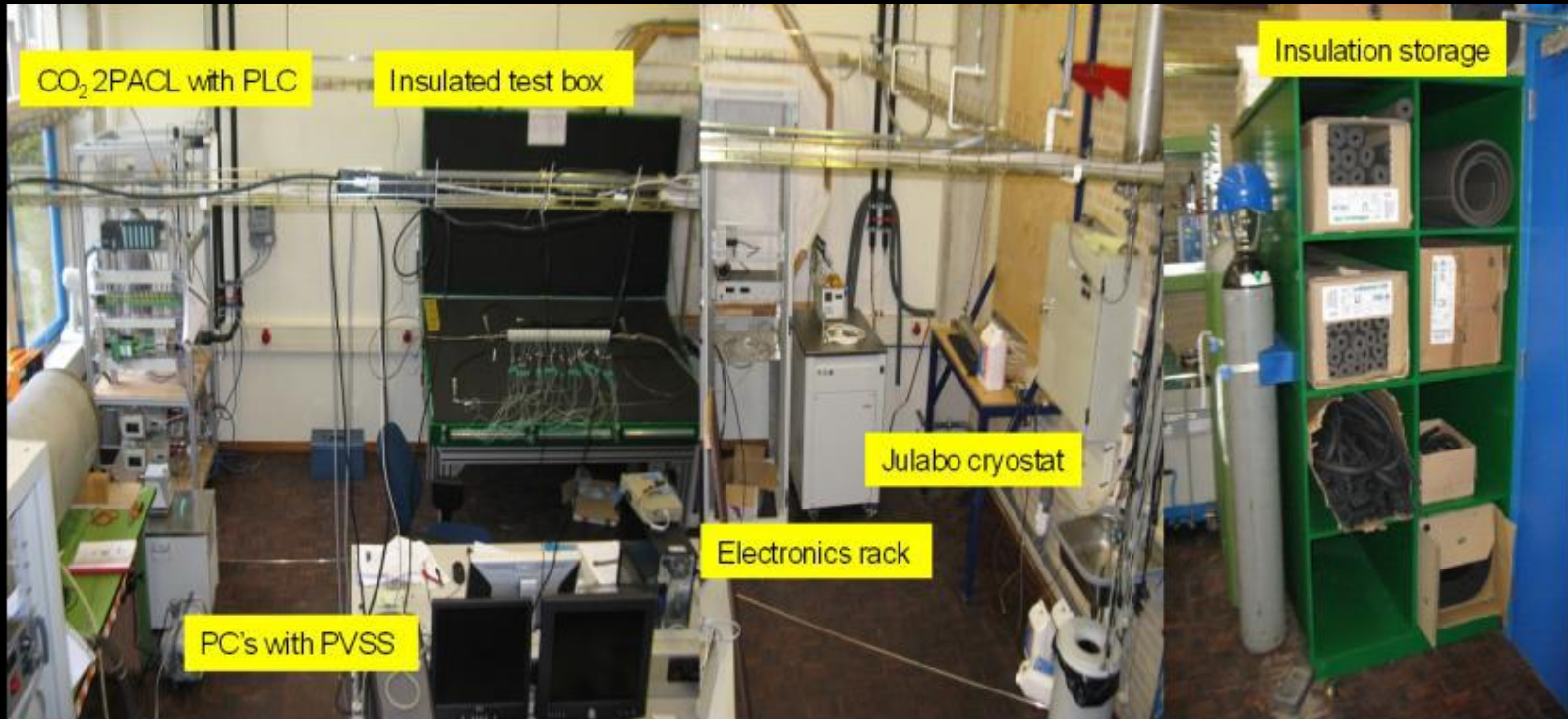
Image: K Rupp, <https://github.com/karlrupp/microprocessor-trend-data>

Fix the thing that didn't scale well, CPU frequency??



LCO₂ cooling of an AMD Ryzen Threadripper 3970X [56.38 °C] at 4600.1MHz processor (~1.5x nominal speed) sustained, using the Nikhef LCO₂ test bench system (<https://hwbot.org/submission/4539341>) - (Krista de Roo en Tristan Suerink)

... since you then need this around it ...



Nikhef 2PA LCO2 cooling setup. Image from Bart Verlaat, Auke-Pieter Colijn *CO2 Cooling Developments for HEP Detectors*
<https://doi.org/10.22323/1.095.0031>

With 20 000+ users, you need something global: WLCG!



~ 1.4 million CPU cores
~ 1500 Petabyte
disk + archival

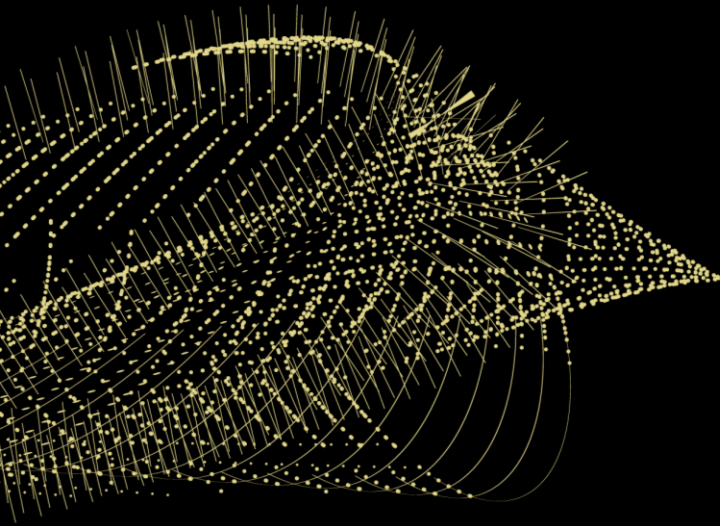
170+ institutes
40+ countries
13 'Tier-1 sites'
NL-T1:
SURF & Nikhef

e-Infrastructures
EGI
PRACE-RI
EuroHPC
OpenScienceGrid
XSEDE (ACCESS)

Earth background: Google Earth; Data and compute animation: STFC RAL for WLCG and EGI.eu; Data: <https://home.cern/science/computing/grid>
For the LHC Computing Grid: wlcg.web.cern.ch, for EGI: www.egi.eu; ACCESS (XSEDE): <https://access-ci.org/>, for the NL-T1 and FuSE: fuse-infra.nl, <https://www.surf.nl/en/research-it>

Global distribution of computing and data placement

WLCG and EGI Advanced Computing for Research

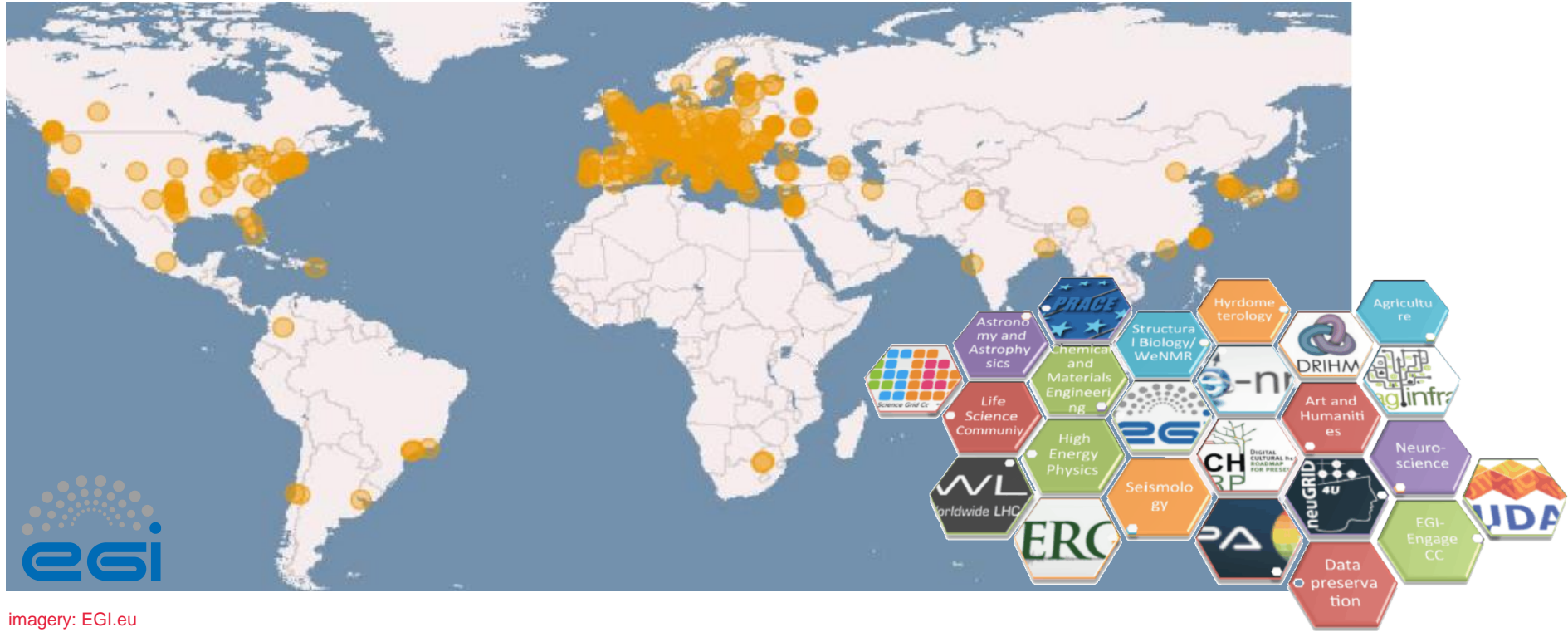


- > 1 user group
- > 1 organization
- > 1 system
- > 1 site



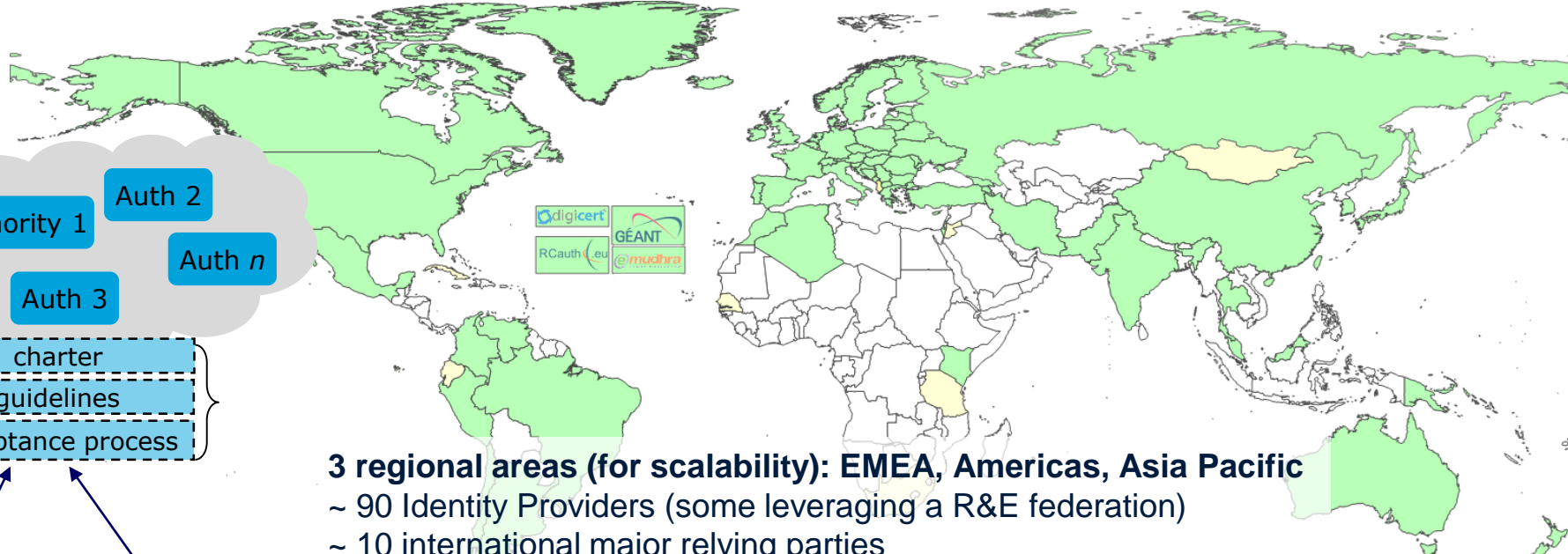
More of more than one ...
@Nikhef

E-INFRASTRUCTURES: EGI, EUDAT, GEANT, PRACE, ...



imagery: EGI.eu

Policy-bridged global federations for research computing



3 regional areas (for scalability): EMEA, Americas, Asia Pacific

~ 90 Identity Providers (some leveraging a R&E federation)

~ 10 international major relying parties

~ 60 countries / economic areas / international treaty orgs

> 1000 relying service provider collaborations

relying party 1

relying party n

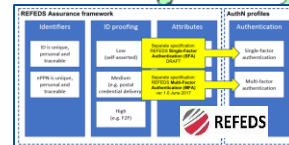


Image: Interoperable Global Trust Federation IGTF, <https://igtf.net/>; REFEDS Assurance Framework RAF: <http://refeds.org/assurance>, <https://refeds.org/profile/mfa>

OpenID Connect Federation

OIDC endpoints + trust policy data for registration can be federated in a meta-data feed

- makes OIDC 'federatable' (plain oidc is single OP)
- as for PKIX, can be technical or policy bridge
- delegated metadata makes 'OIDC-fed' scale in webscale scenarios

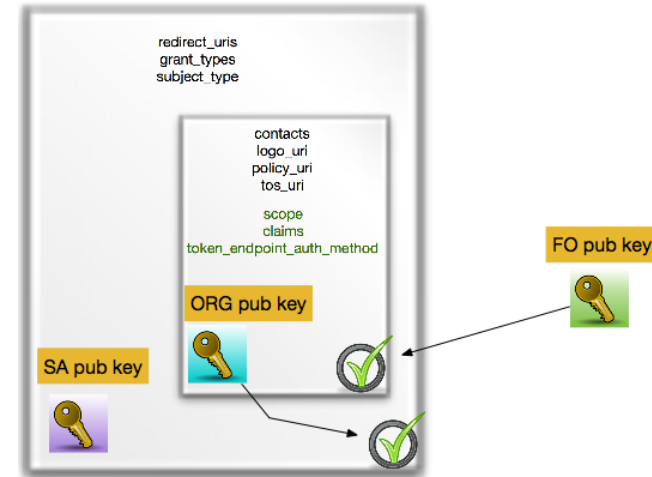
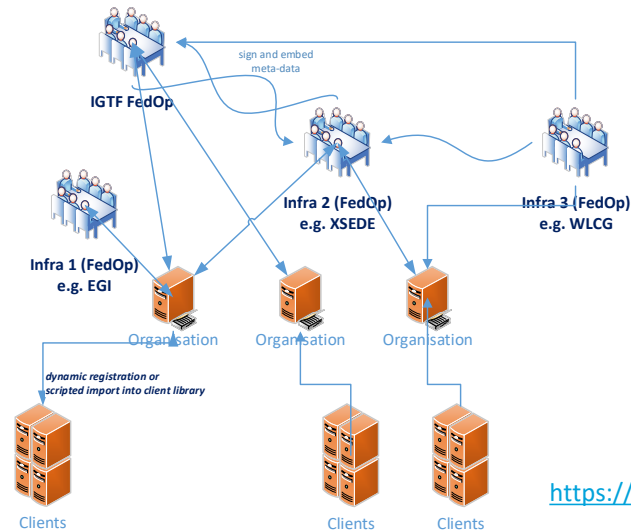
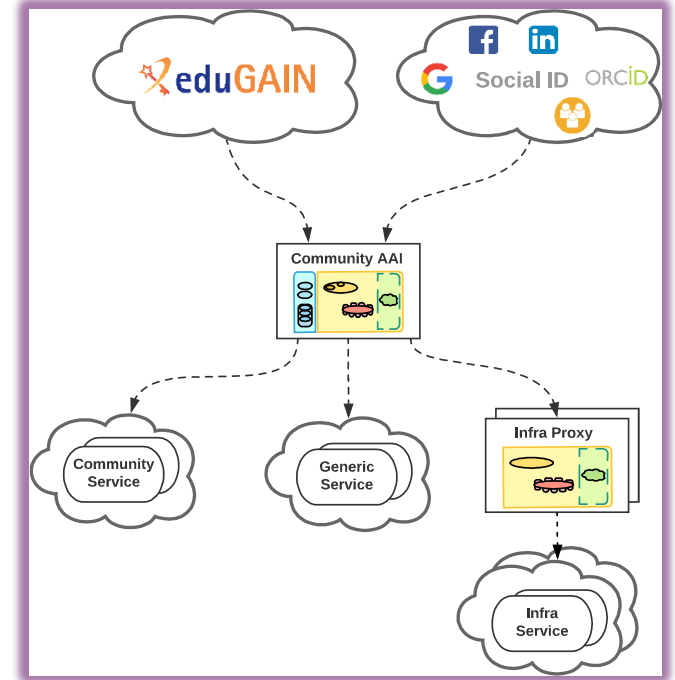
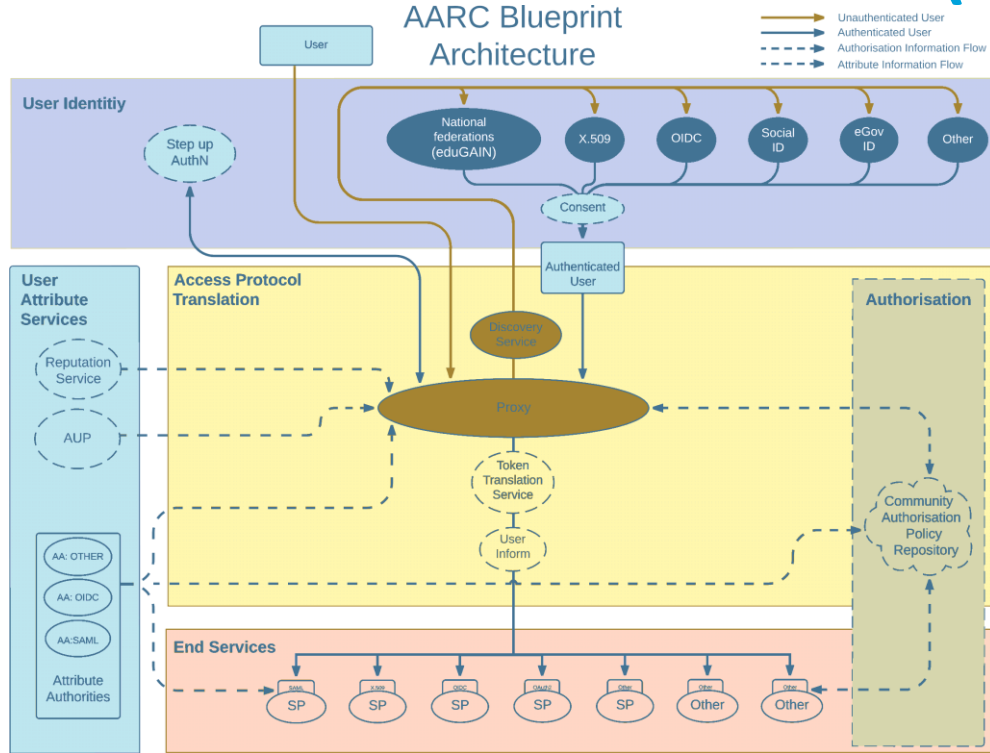


Image: Roland Hedberg, University of Umeå
OpenID Connect Federation:

https://openid.net/specs/openid-connect-federation-1_0.html

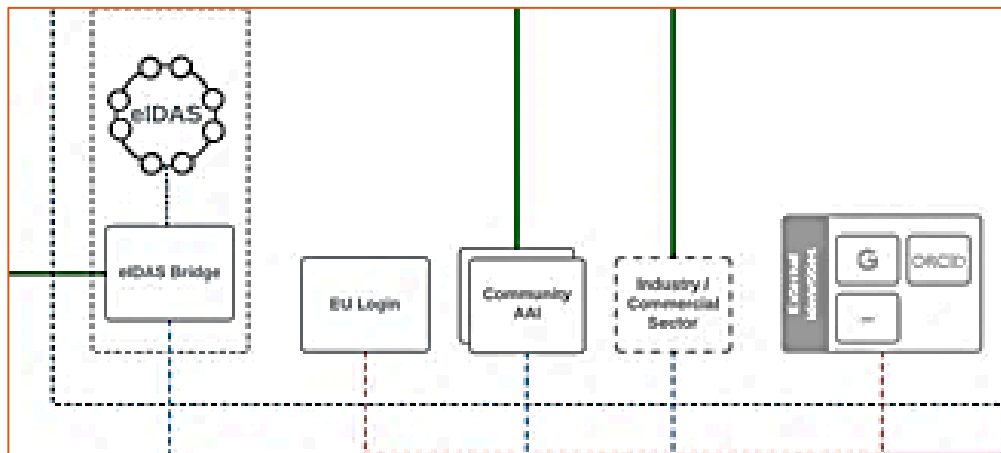
Most trust flows from the (research) community



AARC Blueprint Architecture (2019) AARC-G045 <https://aarc-community.org/guidelines/aarc-g045/>; stacked proxies: EOSC AAI Architecture EOSC Authentication and Authorization Infrastructure (AAI), ISBN 978-92-76-28113-9, <http://doi.org/10.2777/8702>

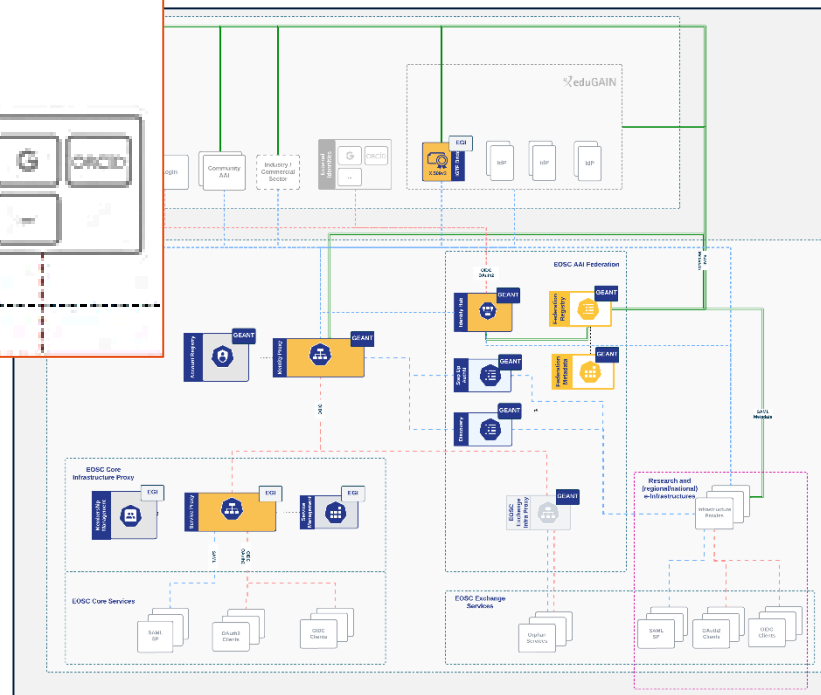
EOSC AAI Federation

Identity assurance brings the true value: authenticators are aplenty, and 'MFA' far less interesting than vetted identities. But HEI home IdPs seem reluctant to provide it ...



user identity comes 'with the user' from outside, mediated by the research community, ORCID, or from the home member state involved

Image: EOSC AAI for the EOSC Core and Exchange Federation for the EOSC European Node by Christos Kanellopoulos, Nicolas Liampotis, David Groep (June 2023)



Bridges: token translation example 'SAML' to PKIX



Community Science Portal

GSIFTP demo

Info Browse Proxy info User info Logged in as david@nikhef.nl

gsiftp://prometheus.desy.de: /

dr-x-----	1	david	david	512 Feb 7 06:00	lost+found
dr-x-----	1	david	david	512 Feb 7 06:01	VOS
dr-x-----	1	david	david	512 Feb 7 06:01	Users
dr-x-----	1	david	david	512 Feb 7 06:02	UTF-8
dr-x-----	1	david	david	512 Feb 7 06:03	Music
dr-x-----	1	david	david	512 Feb 7 06:04	Video
dr-x-----	1	david	david	512 Feb 7 11:21	upload

Delete selected entry Browse... No file selected. Upload file Create directory

Remote name:

dCache EGI AARC

RCaAuth.eu The white-labeled Research and Collaboration Authentication CA Service for Europe

RCaAuth.eu Online CA consent page

The Master Portal below is requesting access to your personal information and to act on your behalf.

If you approve, please accept, otherwise, cancel.

Details on which attributes are released, why, to whom, and how they are processed can be found in the RCaAuth.eu OASIS CA privacy policy. For further information on the CA see the OASIS CA glossary.

Approve

No, continue No, cancel

Master Portal information:

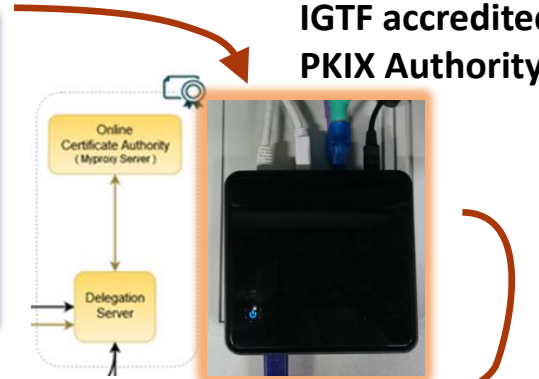
Name: EGI Master Portal
Description: EGI Master Portal
URL: https://caauth.eu/portal/online-ca.php

Information that will be sent to the Master Portal:

url: https://caauth.eu/portal/online-ca.php
uri: https://caauth.eu/portal/online-ca.php
uri_authn_req: https://caauth.eu/portal/online-ca.php
uri_authn_resp: https://caauth.eu/portal/online-ca.php
uri_authz_req: https://caauth.eu/portal/online-ca.php
uri_authz_resp: https://caauth.eu/portal/online-ca.php
uri_authn_req: https://caauth.eu/portal/online-ca.php
uri_authn_resp: https://caauth.eu/portal/online-ca.php
uri_authz_req: https://caauth.eu/portal/online-ca.php
uri_authz_resp: https://caauth.eu/portal/online-ca.php
uri_authn_req: https://caauth.eu/portal/online-ca.php
uri_authn_resp: https://caauth.eu/portal/online-ca.php
uri_authz_req: https://caauth.eu/portal/online-ca.php
uri_authz_resp: https://caauth.eu/portal/online-ca.php

Infrastructure Master Portal Credential Store

IGTF accredited PKIX Authority



REFEDS R&S Sirtfi Trust

RCaAuth.eu The white-labeled Research and Collaboration Authentication CA Service for Europe

English | Nederlands | Español | Français | Deutsch

You have previously chosen to authenticate at Nikhef

Log in at Nikhef

Research and e-Infrastructures | Common | UK | Netherlands | Sweden | Switzerland | Other countries | Miscellaneous

EGI AAI Checkin
ELIXIR research Infrastructure AAI

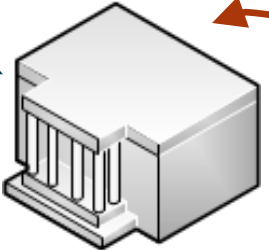
The RCaAuth.eu WebUI is provided by RCaAuth.eu. For support, please contact the help desk of your own home organisations. Service built on OpenProxy and MyProxy.

Policy Filtering WAYF to eduGAIN

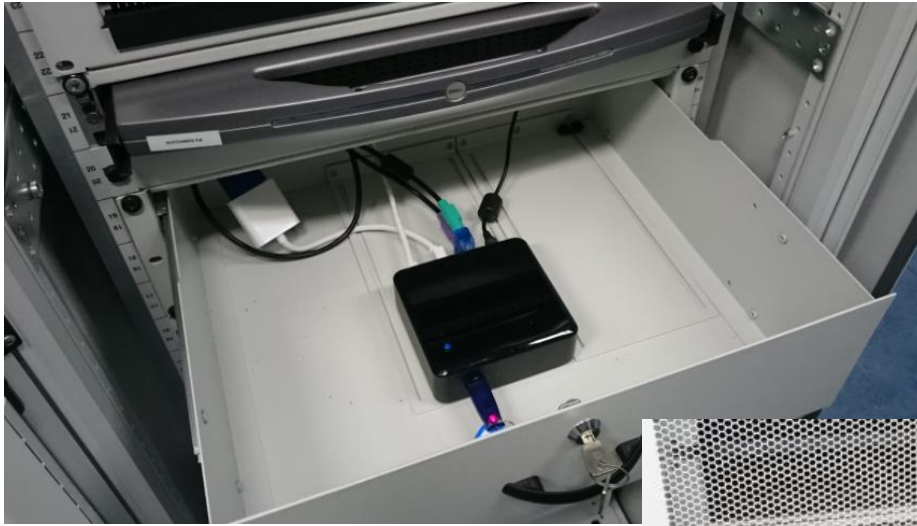


User Home Org or Infrastructure IdP

see also <https://rcdemo.nikhef.nl/>



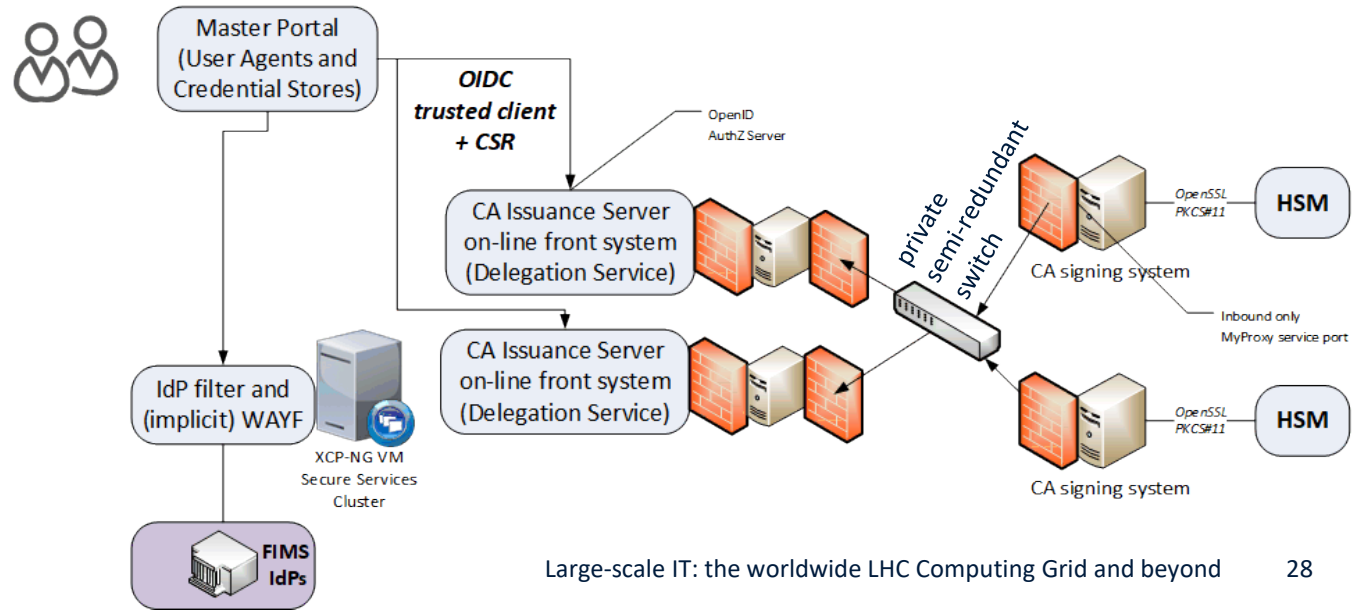
With a single, yet fully compliant, 'Heath Robinson' CA



A single-site locally-highly-available RCauth at Nikhef Amsterdam

- Most 'fault-prone' components are
 - Intel NUC (single power supply)
 - HSM (can lock itself down, and the USB connection is prone to oxidation)
 - DS front-end servers (physical hardware, albeit with redundant disks and powersupplies)

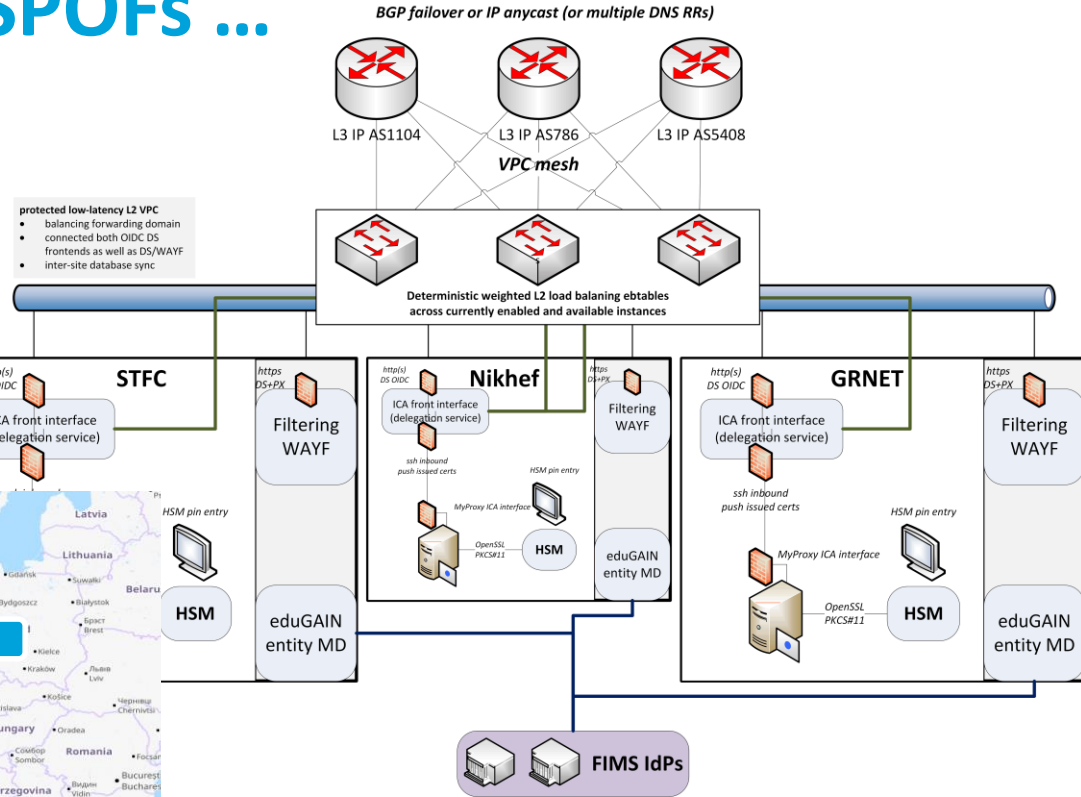
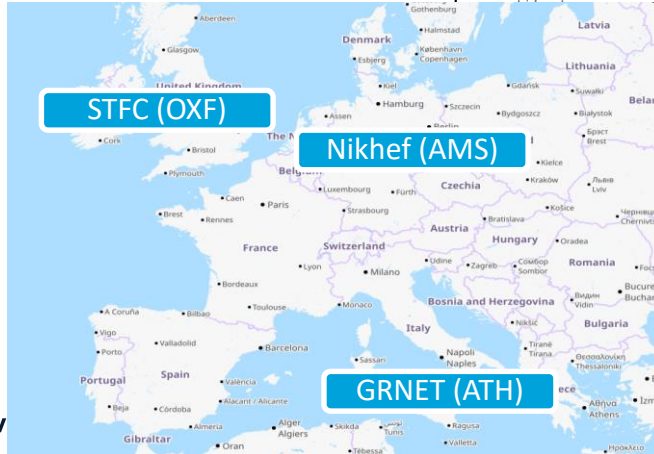
Eliminated
SPOFs first
using 'local HA'



Since we do not like SPOFs ...

Distributed High Availability setup

- across the 3 sites
- design for minimal effort
- readily-available techniques
 - L3 VPN (OpenVPN) or L2 VPC
 - Linux HAProxy

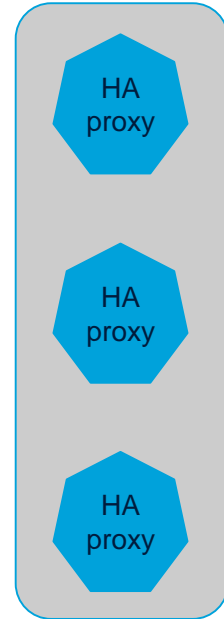


work supported by the EOSC Hub and EOSC Future Horizon Europe projects

A transparent multi-site setup is needed for the user

User

- connects to HA proxy at **{wayf,pilot-ica-g1}.rcauth.eu**
- HA proxy sends users to “**closest**” working service
- primarily **forward to its own DS** when available



If a HA loses its backend DS, can still route to another DS over VPC/VPN backend

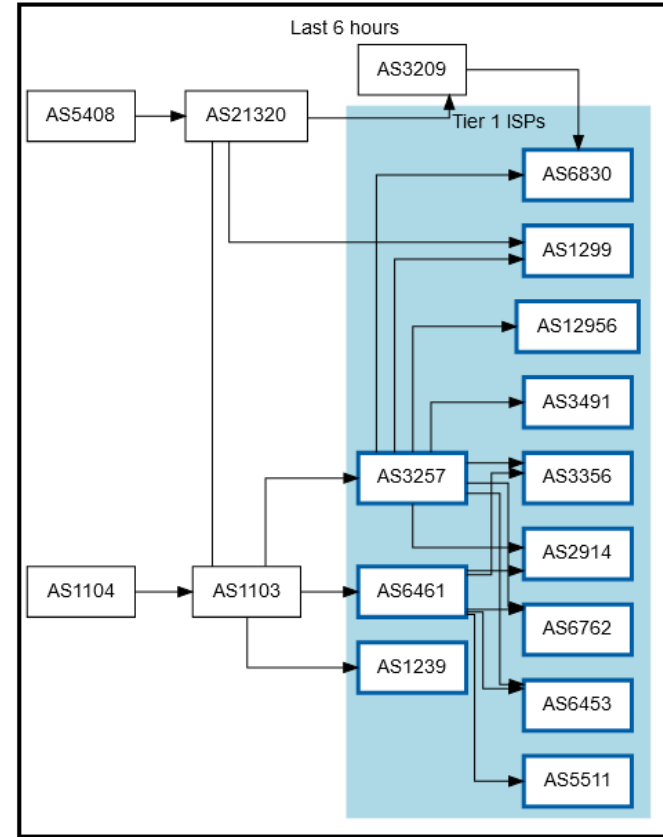
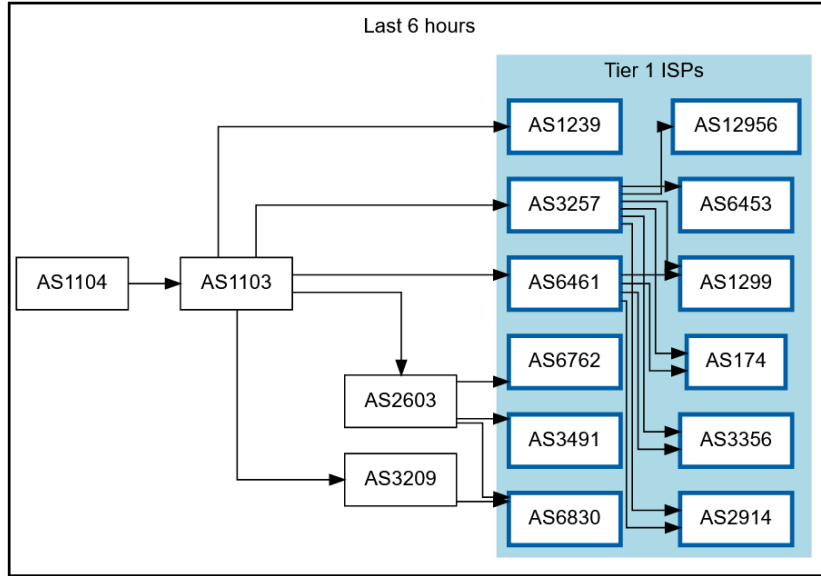
Straightforward proven solution is IP anycast

wherever the user is, the service is at

- **2a07:8504:01a0::1**
- or for legacy IP users at 145.116.216.1

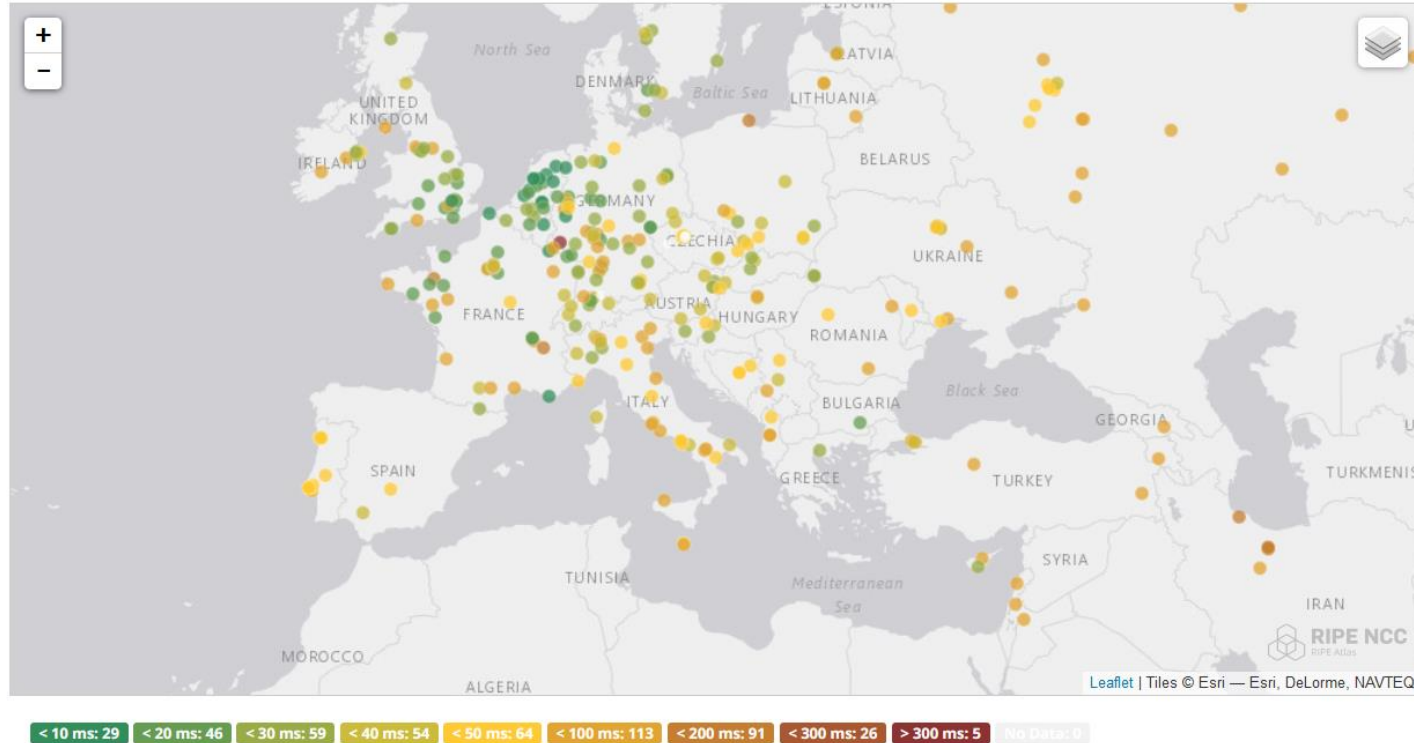
selected imagery: Mischa Sallé, Jens Jensen, Nicolas Liampotis

Getting 2a07:8504:1a0::/48 out there



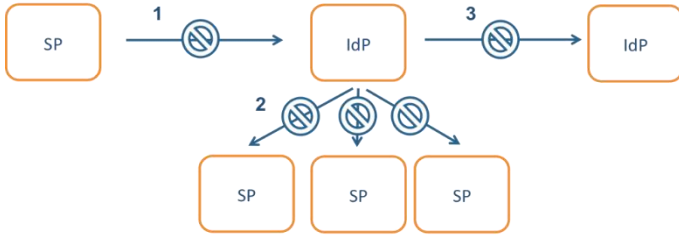
route maps: bgp.tools for 2a07:8504:1a0::/48 – IPv4 for 145.116.216.0/24 is similar – imagery from November 2022

And you get reasonable load balancing in Europe for free



map: RIPE NCC RIPE Atlas - 500 probes, distributed across Europe (<https://atlas.ripe.net/measurements/50949024/>)

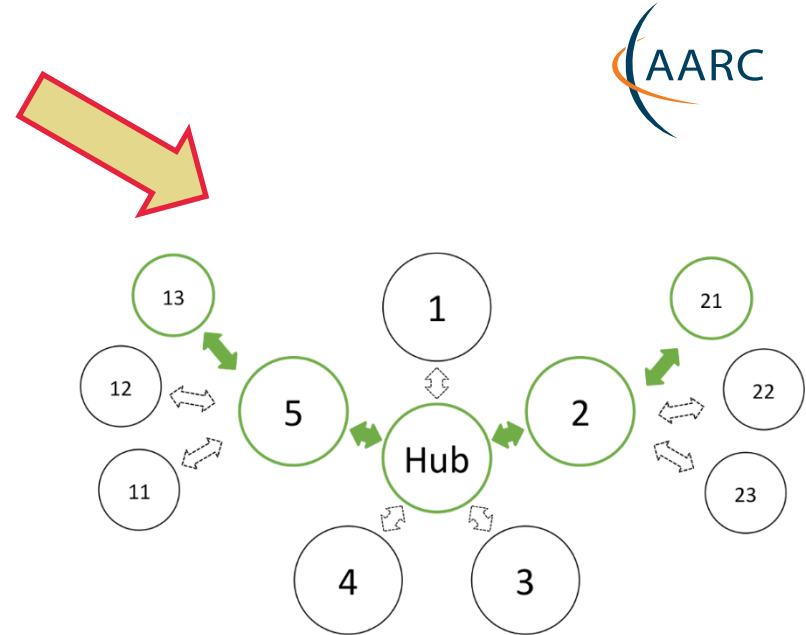
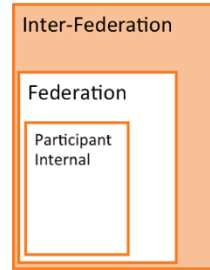
MANY PARTIES, SHARED SECURITY CHALLENGES



Incident Response Communication, communication blocks

Challenges

- IdP appears outside the service's security mandate
- Lack of contact or lack of trust in the IdP which to the SP is an unknown party
- IdP **fails to inform other affected** SPs, for fear of leaking data, of reputation, or just lack of interest and knowledge
- No established channels of communication, esp. not to federations themselves!



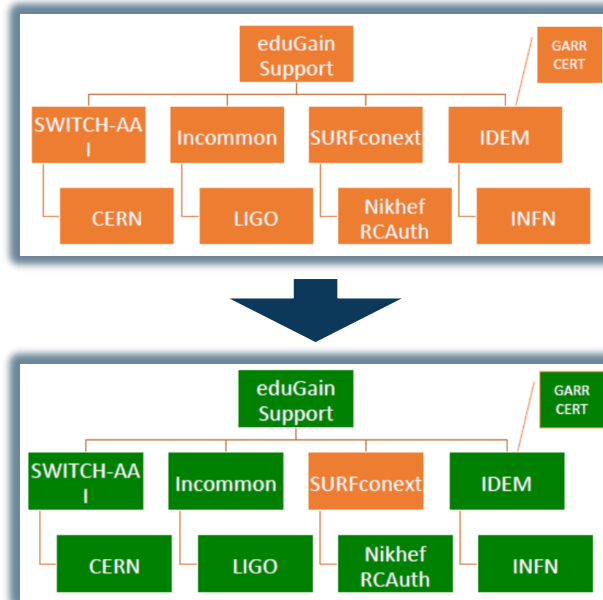
Inter-Federation Incident Response Communication



EXERCISES – COMMUNICATIONS AND ACTIONS



parties involved in response challenge



More than one: Nikhef Science Data Centre ('234b')



Physical farms: selecting the ‘worker nodes’

For HTC applications
– like WLCG, IGWN,
but also SKA, WeNMR – typically
balanced features for node throughput
(CPU, storage, memory bandwidth, network)

single-socket multicore systems are fine,
typical: 64-128 cores per system
network: 2x25Gbps
(+ ‘out of band’ management like IPMI)
memory: 8 GiB/core
local disk: 4TB NVME PCIe Gen4 x4
+ space (physical + power) to add **GPU**

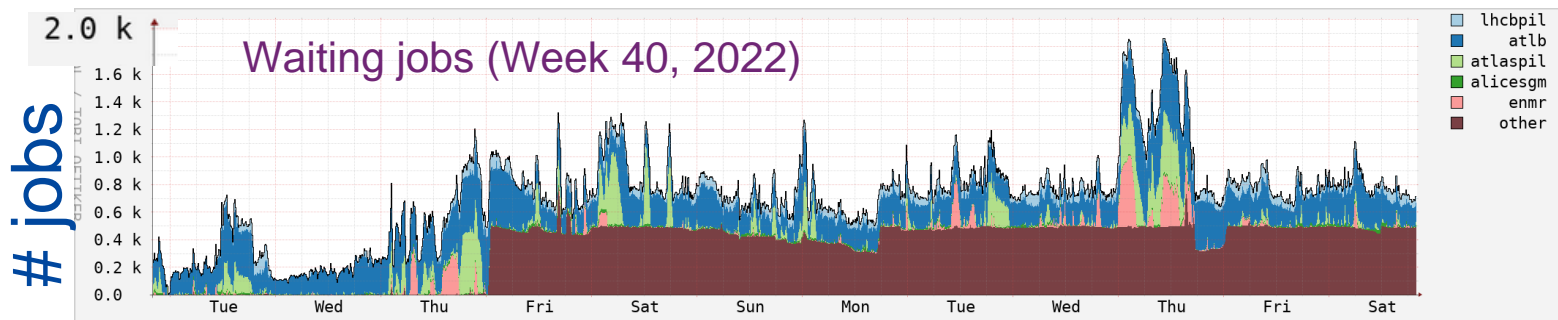
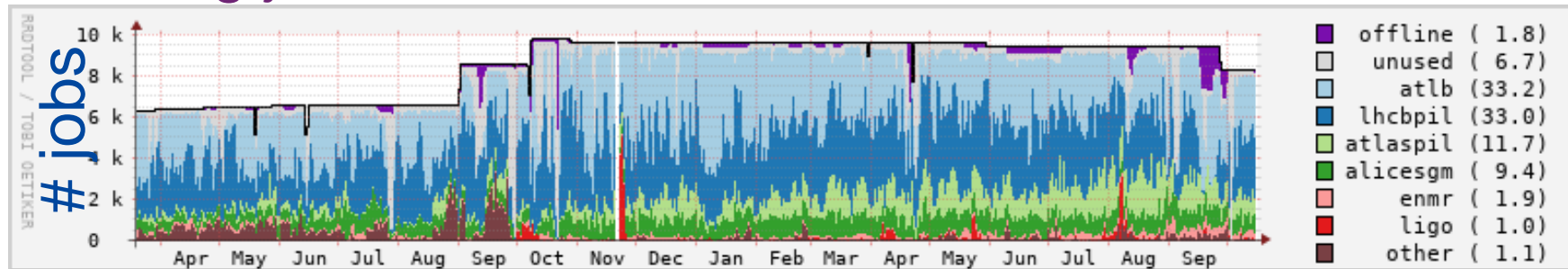


Image: Cluster ‘Lotenfeest’ at the Nikhef NDPF, acquired March 2020. Lenovo SR655 with AMD EPYC 7702P 64-Core single-socket

NDPF 'WLCG and Dutch National Infra' cluster

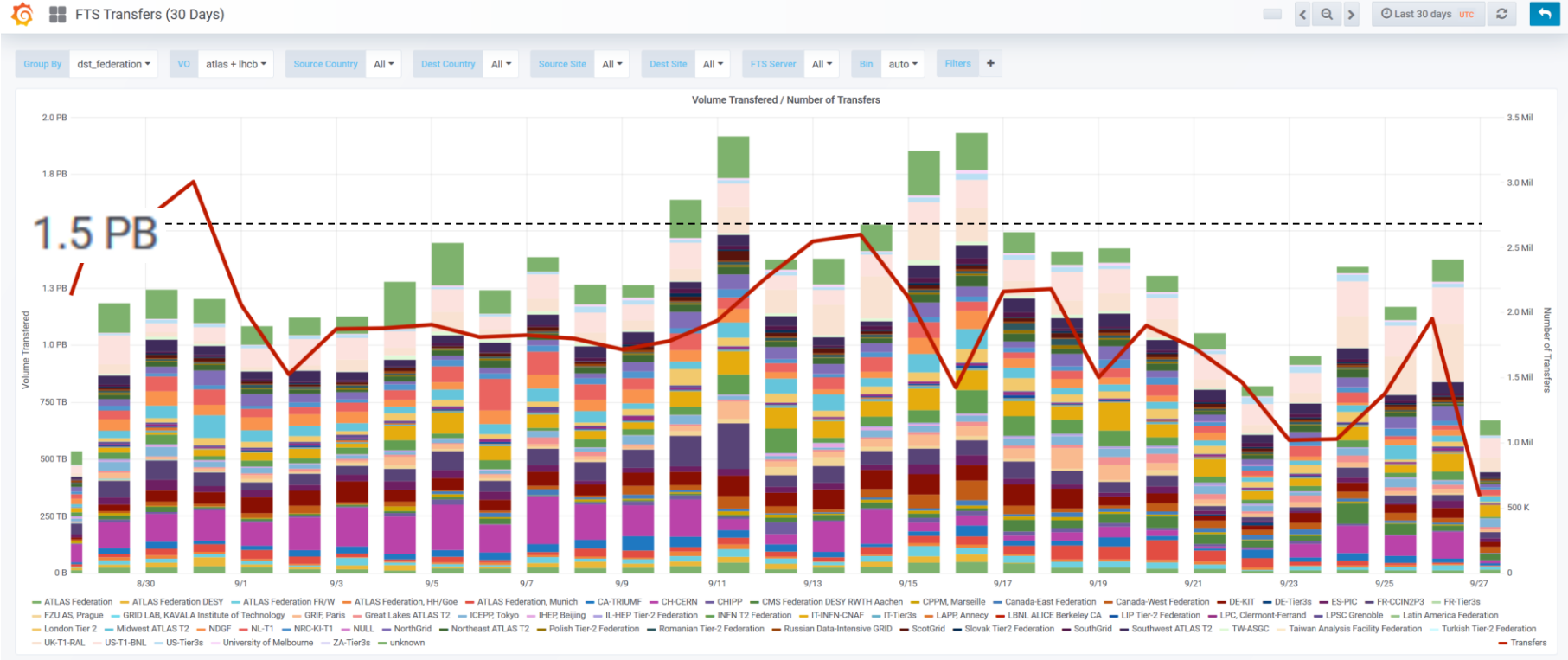
Running jobs:

period: March 2021 .. October 2022



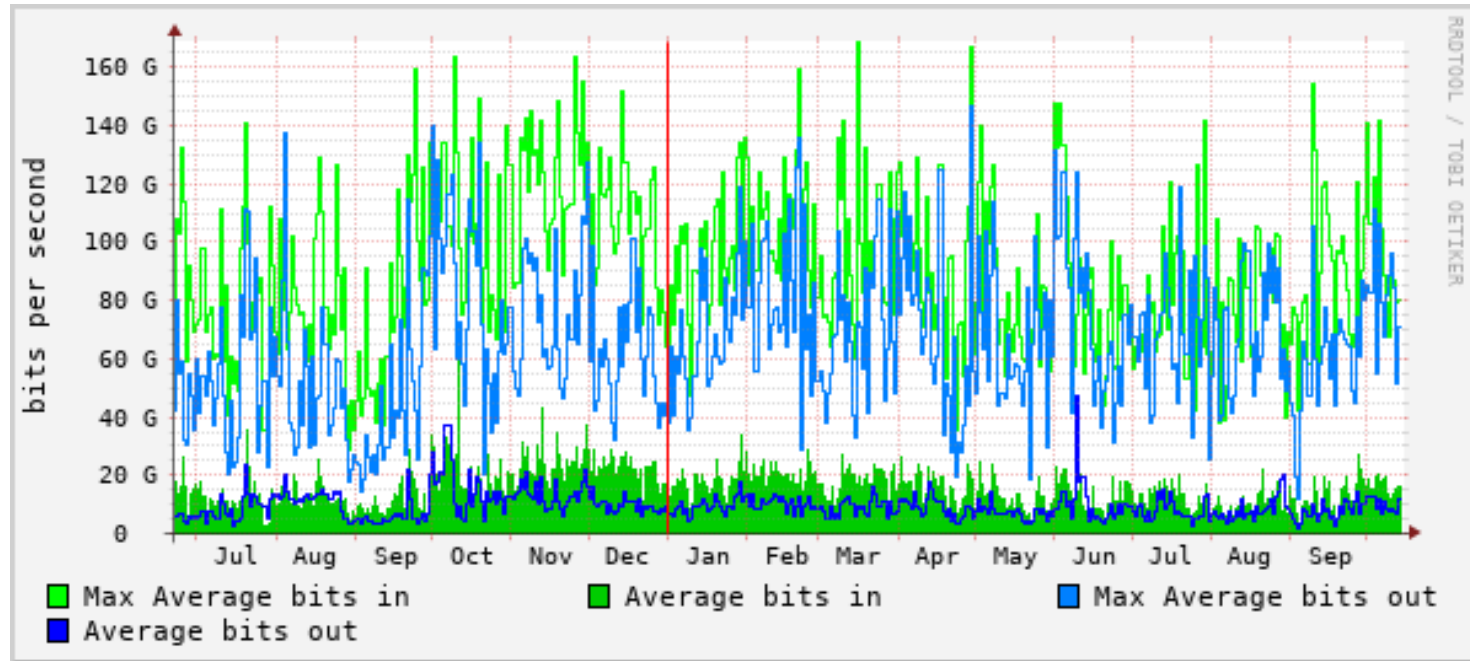
drainage event on Sept 27 are nodes being moved to the LIGO-VIRGO specific cluster; Source: NDPF Statistics overview, <https://www.nikhef.nl/pdp/doc/stats/>
'other' waiting jobs are almost all for the Auger experiment - GRISview images: Jeff Templon for NDPF and STBC

Global high throughput computing needs moving of data

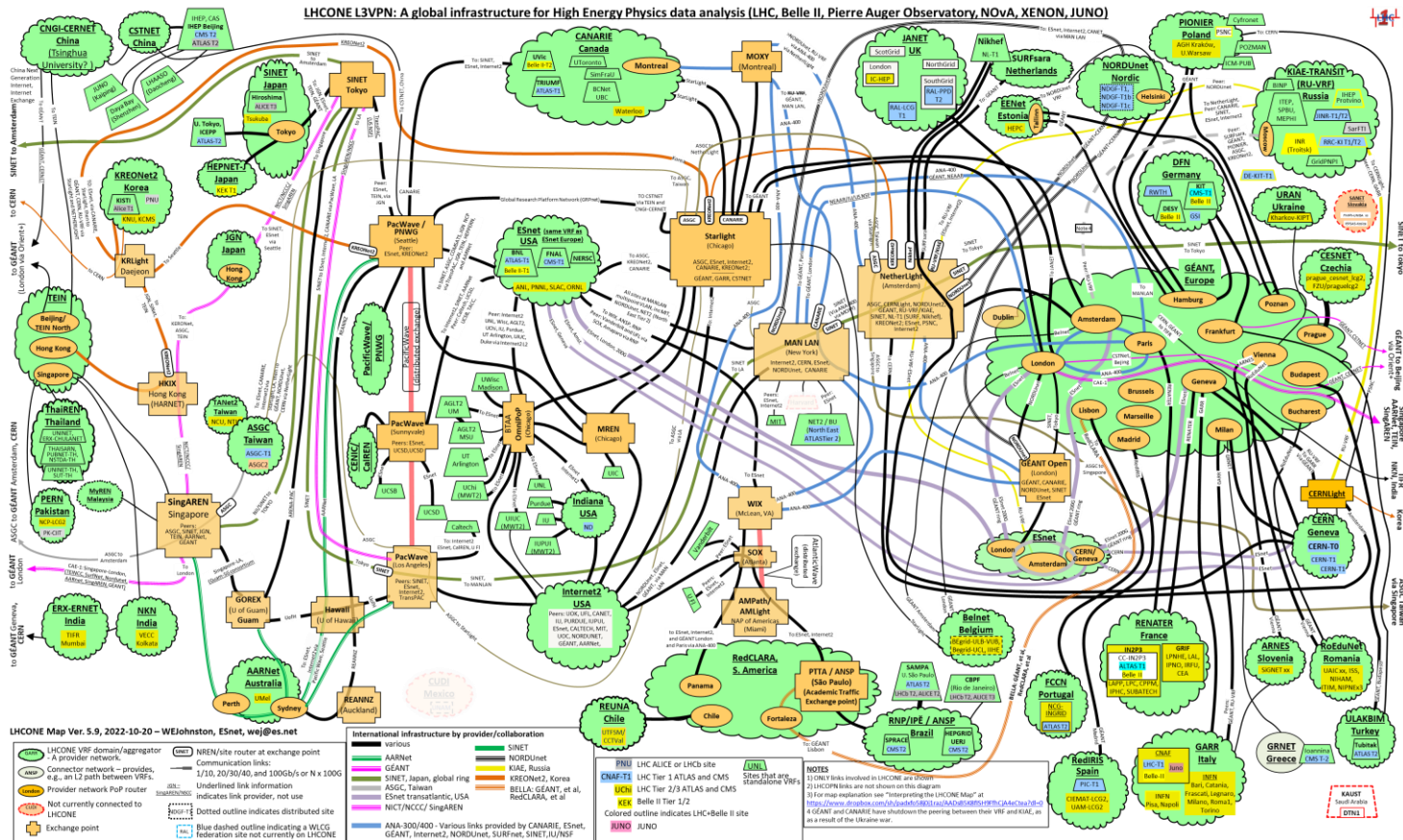


source: <https://monit-grafana.cern.ch/d/000000420/fts-transfers-30-day> ; data: November 2020 ; CERN FTS instance WLGC: daily transfer volume ATLAS+LHCb

Typical data traffic to and from the processing cluster



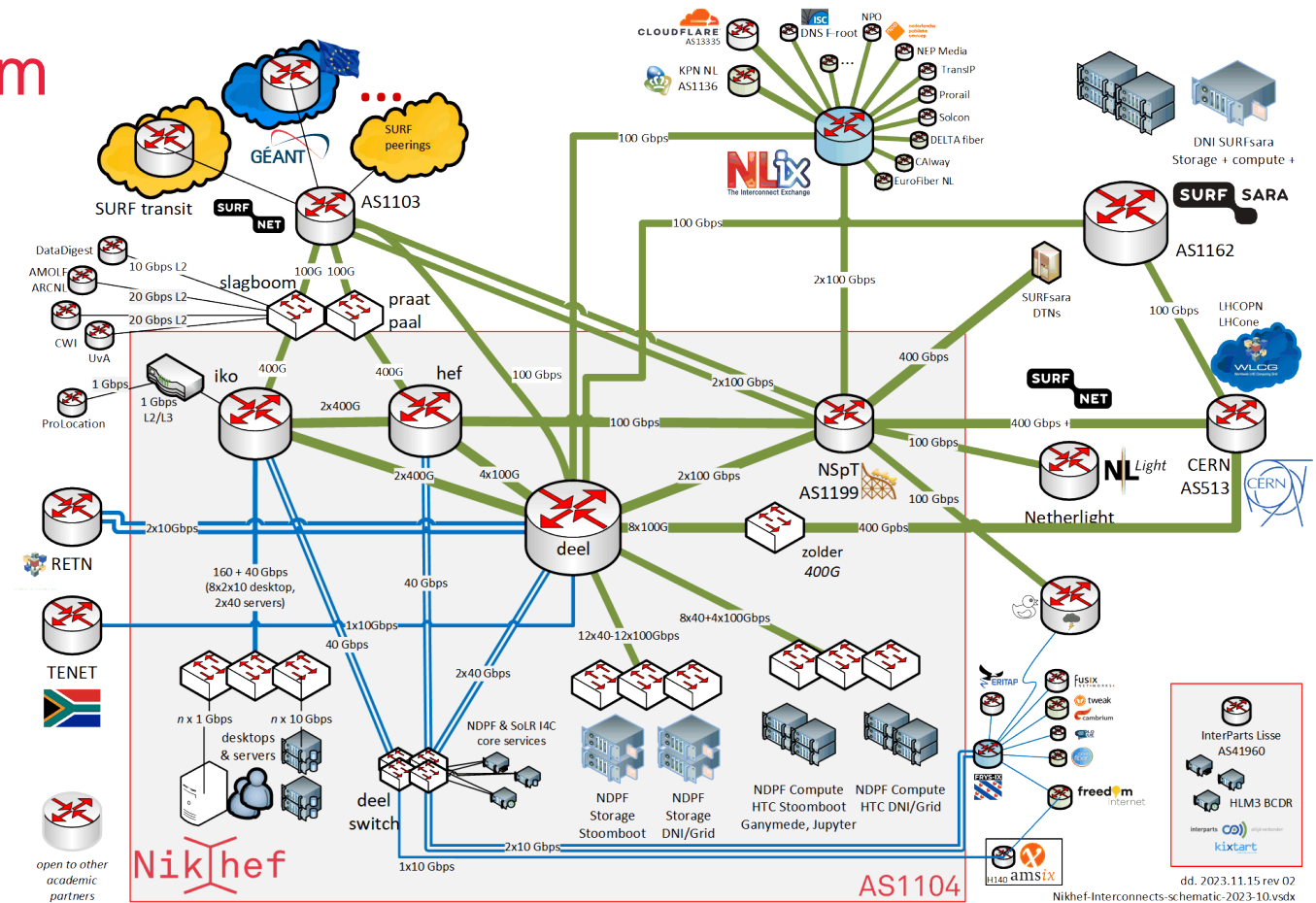
Source: Nikhef cricket graphs period June 2021 – October 2022 – aggregated (research) traffic to external peers from deelqfx – <https://cricket.nikhef.nl/>



LHCone (“LHC Open Network Environment”) – visualization by Bill Johnston, ESnet version: October 2022 – updated with new AS1104 links



Just one random (smallish) autonomous system



AS1104

AS1104

- InterParts Lisse AS41960
- HLM3 BCDR
- interparts
- kixtart

dd. 2023.11.15 rev 02
Nikhef-Interconnects-schematic-2023-10.vsdX



Exercising the network – sensor data and events

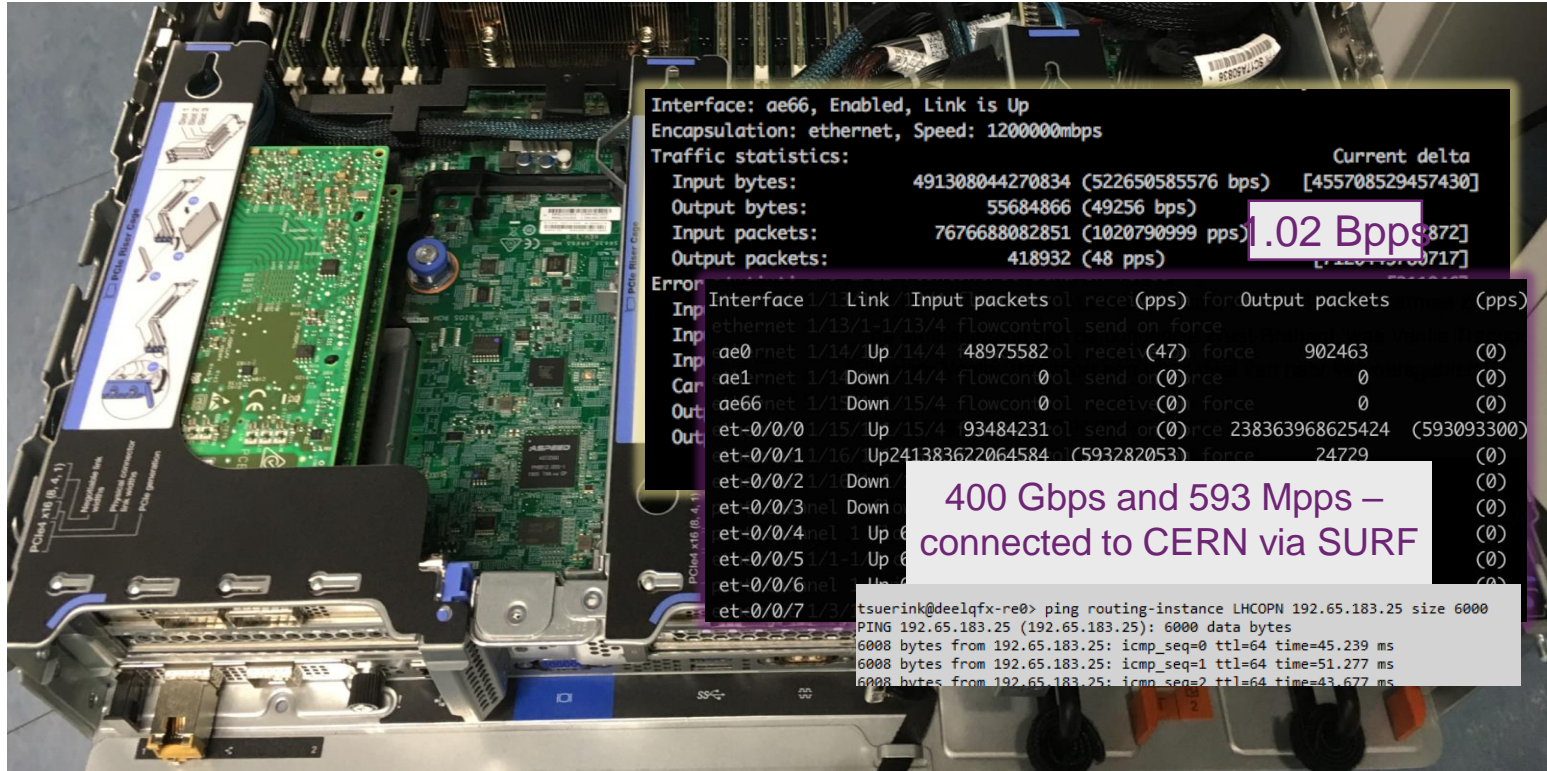
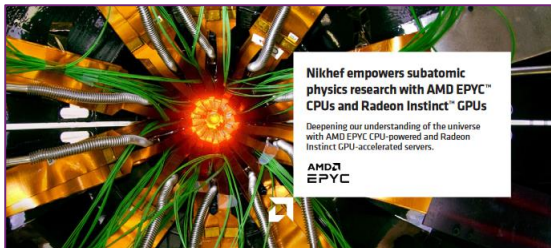


Image: ballenbak.nikhef.nl, Tristan Suerink

Innovation on infrastructure



Nikhef empowers subatomic physics research with AMD EPYC™ CPUs and Radeon Instinct™ GPUs

Deepening our understanding of the universe with AMD EPYC CPU powered and Radeon Instinct GPU-accelerated servers.



Many of the latest scientific discoveries are as much about the computing power used to analyze experimental data as they are about the theories behind them. At the forefront of advancing the processing capabilities for subatomic physics research is Nikhef, the Dutch National Institute concentrating on this area. Nikhef has provided computing that has helped with the discovery of gravitational waves in 2016, the Higgs boson, and the fundamental physics in between, including confirmation that many of the heavy elements in the universe are produced in neutron star mergers.

"The institute performs blue-sky research to learn more about the nature of the universe and the building blocks of matter," explains Toed Aaij, Scientific Staff Member at Nikhef. "The fundamental goal of this institute is to find the big universal box of building blocks everything is made from," adds Tristan Smeik, IT Architect at Nikhef. The more computing power that the institute can draw at this specific, the more that can be discovered. This led the team to AMD EPYC™ processors and Radeon Instinct™ GPUs, which delivered the performance Nikhef's workloads required and the solution price that aligned with their budget.

Data-hungry science
Nikhef is involved in many different experiments, but all of them require a considerable level of computing power. "About 100 scientific staff work at Nikhef," explains Aaij. "These staff usually work on one (or sometimes more than one) of the experiments Nikhef is involved in.

Three of these experiments are at CERN, the ATLAS, LHC, and ALICE experiments. There are several astroparticle physics experiments. One is the Pierre Auger experiment, covering several thousand square kilometers of Pampa in Argentina. The area is equipped with detectors to search for air showers caused by extremely high energy particles that arrive from the universe. Then there is the neutrino physics experiment OPERA, and dark matter research with the XENON experiment. Finally, there is a large gravitational waves physics group that is a member of the LIGO-Virgo experiment collaboration."

"If there's one thing all these experiments have in common, it's the increasing amounts of data that the experiments produce. "The scientists always want more data," says Smeik. "I think there are few experimental physics papers that do not end with 'we need more data.' And in this field of physics, to get more data you build a more sensitive experiment." In the case of the Large Hadron Collider (LHC) at CERN, the box of data produced will be particularly huge.

"In about four years the LHC will increase the number of collisions detected by about a factor of 10," says Aaij. "This means that the experiments will start producing a similarly increasing amount of data. If we look at the growth of storage space and compute capacity over time, then we do not expect to open get close to a factor 10 in increase of performance for a flat budget. We need to deal with that, because we need to process the data. Otherwise, we can't do science with it." This is where AMD EPYC processors and GPU acceleration have offered the best solutions to satiate the hunger for growing data processing ability.

FUNGIBLE

NIKHEF, SURF AND FUNGIBLE SET NEW BENCHMARK FOR THE WORLD'S FASTEST STORAGE PERFORMANCE

Companies Double Current Performance Record, Setting the New Bar at 6.55 Million Read IOPS



CUSTOMER
Nikhef

INDUSTRY
Subatomic Physics

CHALLENGES
Increasing data throughput with higher I/O and memory bandwidth

SOLUTION
Diverse AMD EPYC™ processors and Radeon Instinct™ GPUs, and AMD Radeon Instinct™ M50 GPUs

RESULTS
Faster processing and the ability to harness GPU-accelerated machine learning to cope with rapidly expanding experimental data volume

AMD TECHNOLOGY AT A GLANCE
AMD EPYC™ 7002G processors with 32 cores
AMD EPYC™ 7002P processors with 64 cores
AMD Radeon Instinct M50 GPUs

TECHNOLOGY PARTNER
Lenovo

AMD + NIKHEF CASE STUDY



Image: Minister of Economic Affairs M. Adriaansens launched the Innovation Hub with Nikhef, SURF, Nokia and NL-ix, January 2023. Composite image from <https://www.surf.nl/nieuws/minister-adriaansens-lanceert-testomgeving-voor-supersnelle-netwerktechnologie>



Our science data flows are somebody else's DDoS attack



Het begon in 2018. Een bijzondere samenwerking tussen overheden om ervoor te zorgen dat de dienstverlening overal wordt.

Het 'red team' is verantwoordelijk voor de aanvallen, het 'blue team' voor de verdediging. Een van de partijen die aan de avond meedoet is [Nikhef](#). Tristan, IT architect bij Nikhef, geeft aan "dat zij dit

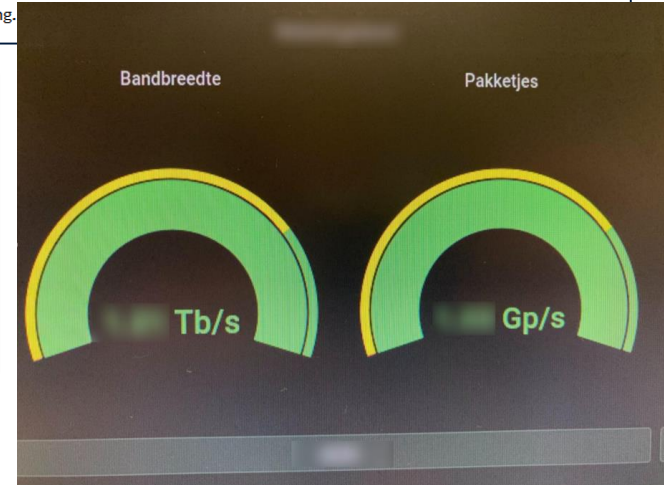


Image sources: belastingdienst.nl, rws.nl, nu.nl, werkentegennederland.nl

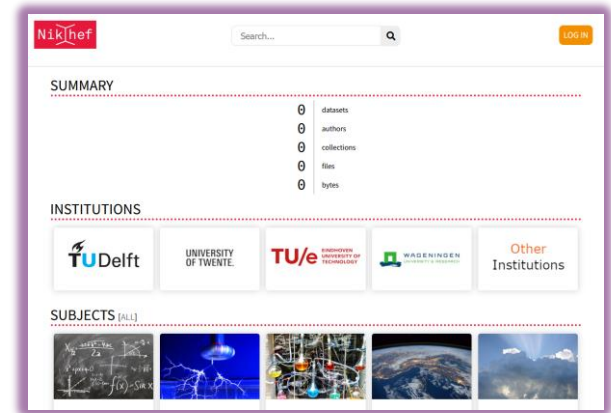
But what about some other digital competences?

We have the network (physical and in terms of expertise) for processing and analysis, but other elements are missing

- accessible analysis preservation for our on-prem experiments and R&D projects
- longer-term 'local' re-use in programmatic research
- complexity of managing ongoing large analyses
- data and software of non-collaboration-based outputs

Now preparing for **analysis preservation** and RDM through managed 'snapshots' as part of the analysis pipelines

- *linking Stoomboot dCache to a new institutional Research Data Management*



Power usage and efficient data centres



Nikhef scientific data centre (the 'glass box') designed for 400kW total use + cooling in 47 racks



De snelste CPU/GPU is voor ons niet altijd de beste (*sorry gamers & miners!*), want 5 jaar energie en beheer zijn even kostbaar als de server zelf

WKO: Warmte Koude Opslag

21% van het vermogen is nodig om te koelen, maar: we mogen 3500GJoule/jaar (~112 kWjaar, ~982 000 kWh) aan studenten tegenover leveren om ze warm te houden !



Let's go on tour!



David Groep

davidg@nikhef.nl

<https://www.nikhef.nl/~davidg/presentations/>

 <https://orcid.org/0000-0003-1026-6606>



Maastricht University

Nikhef

