# Computing for Research & the Worldwide LHC Computing Grid

Building a global large-scale
ICT infrastructure
for research data processing

Nik[hef

David Groep
DACS & Nikhef
5 November 2024
*KEN3239 rev 1.2*

Peter Higgs and François Englert at the 2013 Nobel prize press conference, Stockholm. Photo: Bengt Nyman, https://www.flickr.com/photos/97469566@N00

# Exploding data? the Large Hadron Collider at CERN

**1964**

**1998 - 2012 … 2028: HL-LHC … 2035+**

P. Higgs, Phys. Rev. Lett. 13, 508:

**16823 characters, 165 kByte PDF**

~50 PiB/year primary data

the LHC obviously looks for a lot more than just the Higgs mechanism. For example Alice looks at the Quark Gluon Plasma, LHCb for CP violation and the matter surplus (and lots more), and ATLAS and CMS look at almost anything. And all look at new BSM physics of course …

Images: ATLAS detector in the cavern at CERN. Source: CERN

# Computing on lots of data – 40 million times/sec

**ATLAS RAW single event**
ROD File
**1.60 MB**

**~60 TByte/s** *(compressed)*

**Trigger system selects 600 Hz ~ 1 GB/s data**

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

*~10k researchers*

*CERN and ~170 institutes*



ATLAS EXPERIMENT

Run Number: 266904, Event Number: 25884805

Date: 2015-06-03 13:41:54 CEST

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)

# Processing at scale for data intensive science



LOFAR

LHC run 2 data
~350 PB 'raw'

Long Term Archive
~ 60 PB

CERN

Library of Congress
5 PB

JS Census
4 PB

Nasdaq ● 3 PB

Google searches
>98 PB

Facebook uploads
>180 PB

*quantities are guesstimates,
as actual data volume
keeps growing as we speak …*

LHC Run 3
>2021
>1000 PB

CERN

HL-LHC
>2028
multiple EB

SKA Phase 1
>~ 2027
~600 PB

# So 'big science' needs some computing …



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

# Scaling computing infra: volume is not the only thing that matters

# Not in one place: the worldwide LHC Computing Grid



~ 1.4 million CPU cores
~ 1500 Petabyte
       disk + archival

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

*largely based on*
*generic e-Infrastructures*
EGI
EuroHPC
PRACE-RI
OpenScienceGrid
ACCESS-CI

# Our journey today …

**let's build some 'scalable' infrastructure for LHC computing, storage, networking, and a global AAI** … *if we make it*
*Using science use cases* from CERN's Large Hadron Collider, the SKA radio telescope, Gravitational Wave detection, structural biochemistry (WeNMR), and more …

**From the bottom up …** of green fields, ships' diesels, and chilly corridors

**Data intensive workflows that drive infrastructure development**

- **why large-scale IT is distributed**: end of faster CPUs, thermal barrier, rise of parallelism

**More than one …**

- **High Performance & High Throughput**: distributed computing, storage and data placement
- **As a service:** herding systems, cloud platforms, containers, and service management

**Networking the systems: linking 'more than one' globally**

- **network design**: elephants vs. mice in shipping large quantities of data … and on cat videos
- *Optical Private Networks* and the *Open Networking Environment LHCone*

**Networking the people**

- **authentication and authorization** technologies
- **multilateral federation:** identity, community management & global trust

**Putting it all together again** *(and maybe an example of a federated anycasted authentication service)*

# Start with …
# a green field approach?

*from field to facility*

# From field to facility





Trekkersveld IV, Zeewolde. From Zeewolde Actueel, https://www.zeewolde-actueel.nl/nieuws/gemeente/254432/bestemmingsplan-trekkersveld-4-ligt-ter-inzage; Microsoft DC Middenmeer, from https://nos.nl/l/2512478,

# Feel the Power



Images: Anton Mors, David Groep, Nikhef

# Converting electricity into … chilled air & heat





Left-side image: frame from a movie by Anton Mors, people replaced by … Adobe Firefly ("without people"?, oh well, this was its best result ☹)

NikhefHousing: a cold aisle

# Where to put large-scale IT: brief look at data centres

- 'tier-1' … 'tier-4' datacenters - increasingly redundant
- all systems are 'lights out', since the DC may be miles away
  - remotely controlled, incl. power-on, remote KVM
- small and large in terms of power and cooling capacity
  - smallish: Nikhef Housing Amsterdam is ~2.5 MW,
  - Meta Zeewolde (now cancelled) would have been 160 MW

- data centre efficiency metric: $PUE = \dfrac{E_{total}}{E_{IT\_equipment}}$



| Current Power | Minimum Power | Peak Power | Average Power | Current / Maximum Power | |
|---|---|---|---|---|---|
| 264 Watt | 264 Watt | 273 Watt | 267 Watt | 264 | 480 Watt |

Reducing cost and impact by improving "Power Unit Efficiency" of the data centre:
- airflow engineering and efficient CRACs
- (free) cooling by changing inflow temperature
- Aquifer Thermal Energy Storage (ATES) to buffer heat (and re-use later for homes)

Typical PUEs vary from 1.03 (in Iceland) to 1.2 for 'good' datacenters in NL

Data centre tiering: Uptime Institute (Tunner, W.P.; Seader, J.H.; Brill, K.G. Tier Classifications Define Site Infrastructure Performance; White Paper)
Remote systems management: IPMI, Redfish and various vendor proprietary solutions – usually dedicated 'out-of-band' network connection, incl. remote KVM

# Every rack should have one

- A bare rack just lacks that nice and warm feeling, so you typically add

- some remotely monitored PDUs
- temperature sensor(s)

- out-of-band management switches
- systems installation net (managed)

- data, storage, and overlay networks
  dual 10/25/100GigE per system
  + optionally a low-latency fabric for HPC,
  like InfiniBand, RoCE, UltraEthernet …

Shown: H234b C06 'SOC' cabinet, Nikhef, front and switches (at back)

# Virtual and cloud services rely on this physical 'stuff'

- HPC systems like the Dutch Snellius, a SuperMUC, LUMI, JUPITER, or Jules Verne,
- data-intensive computing like WLCG, radio astronomy, and so on
- your favourite (or not) typical hyperscalers like AWS, Azure, Google, OVH, Hetzner, …

and all those new AI systems and AI 'factories' that boost Nvidia stock nowadays …



DNI and NL-T1 capacity from 2023 DNI NWO, LOFAR, and WLCG; see https://www.surf.nl/onderzoek-ict/toegang-tot-rekendiensten-aanvragen ; fuse-infra.nl
SURF tape total: ~80 PByte by end 2022; image library at Schiphol Rijk from Sara Ramezani; NikhefHousing: https://www.nikhef.nl/housing/datacenter/floorplan/

# Filling the Data Center

The challenges come
when you have 'more than one'

# Different types of large scale compute resources

- HPC and (computational) cluster computing:
  - modelling for weather/climate, fluid dynamics, but also e.g. QC-simulation

- HTC and data-intensive processing – horizontal scaling:
  - lots of data, as in High Energy Physics (HEP), *omics and protein docking, …
  - conveniently parallel,
    but (intensive) local I/O requirements on memory and scratch storage

- portals and many web applications:
  'horizontal' scaling, often backed by cloud and virtualized resources
  - Cloud-native scaling and containers for 'more of the same, different each time'
  - If it's data at scale: object stores and 'CDN' web-scale caching

HPC: High Performance Computing; HTC: High Throughput Computing; K8S: Kubernetes; CDN: Content Delivery Network

# 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 helps:
  - AMD EPYC Genoa (Zen 4) has 96 cores/die
  - Intel Granite Rapids, Nvidia GraceHopper, …
  - but e.g. Intel Cascade Lake AP was less useful

- CPU design-level performance gains left
  - predictive and 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*



50 Years of Microprocessor Trend Data

Transistors (thousands)
Single-Thread Performance (SpecINT x $10^3$)
Frequency (MHz)
Typical Power (Watts)
Number of Logical Cores

Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2021 by K. Rupp

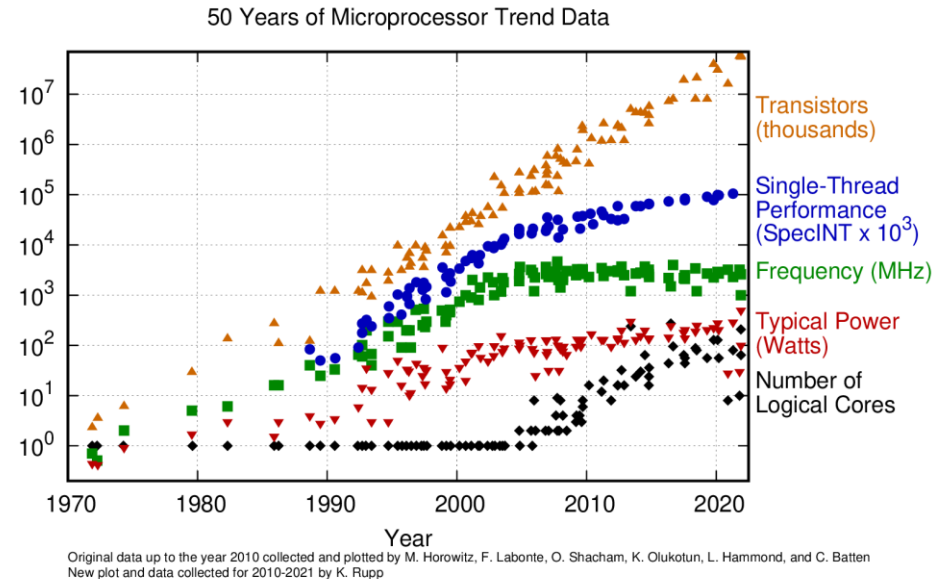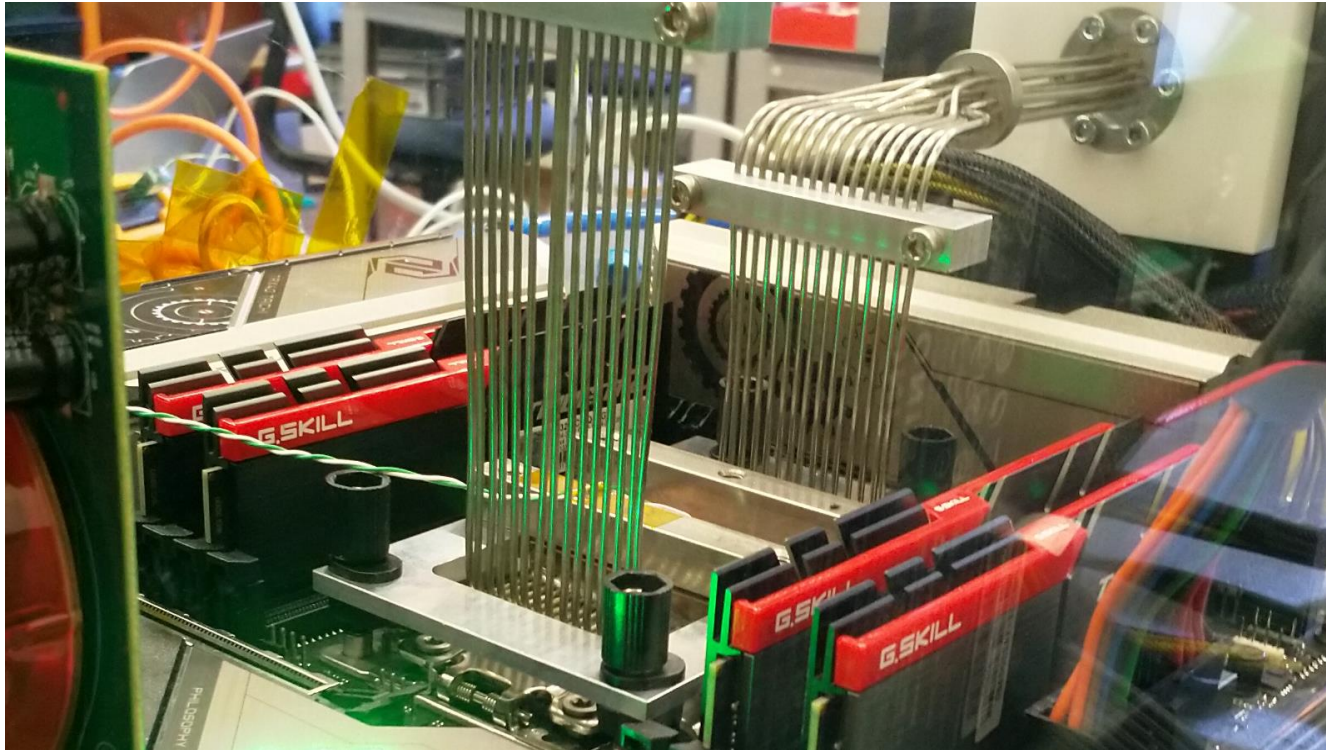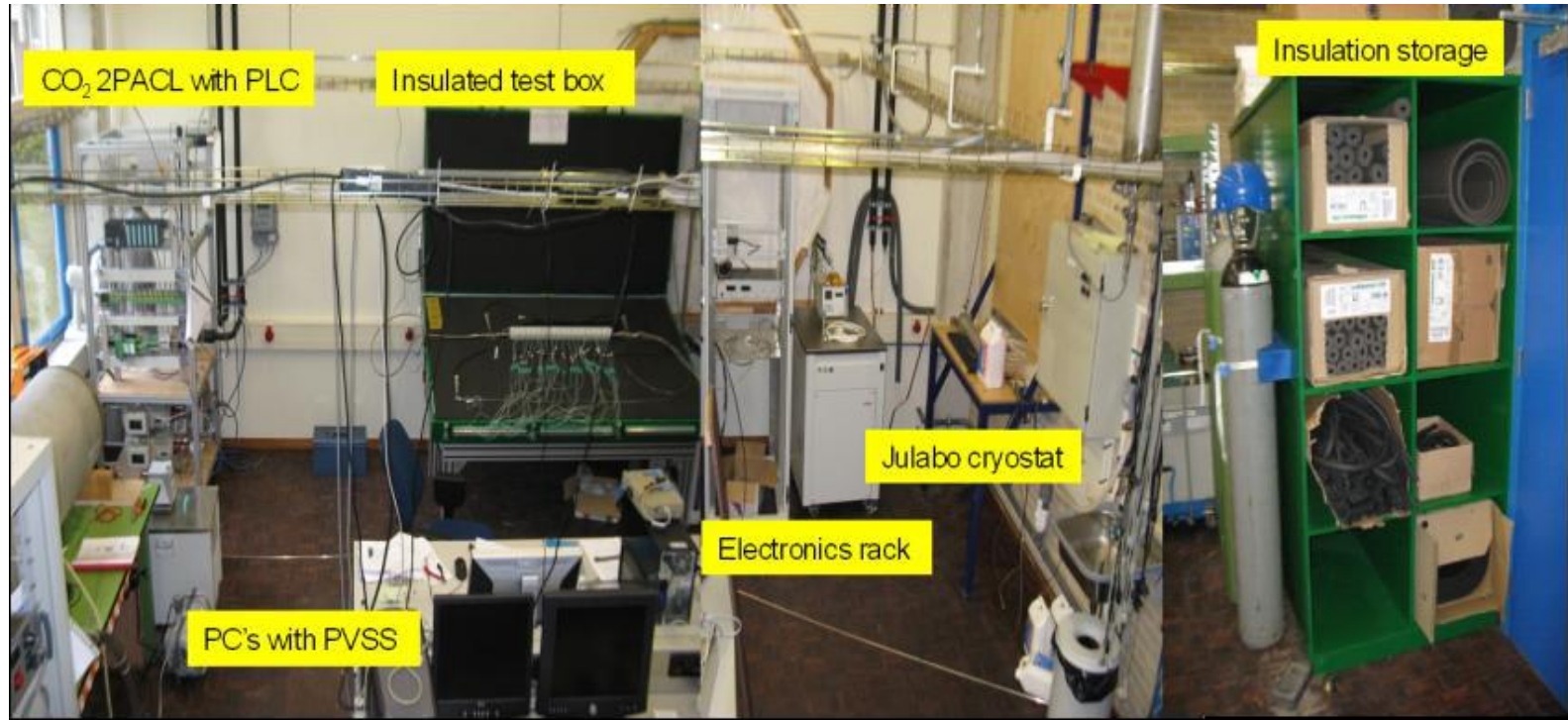Image: K Rupp, https://github.com/karlrupp/microprocessor-trend-data

# 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 (~1.25x nominal speed) sustained over all cores simultaneously, 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

proceed to clusters

# Step one: scale *inside* one system



- 'trivial' step-up is to do multiple sockets in one system
  2-socket, sometimes 4 socket on a motherboard

- to make it appear as a single shared memory system,
  *cache coherency* is required between the CPUs

- useful for tightly coupled parallel applications
  (weather forecasting, fluid dynamics, climate), but
  not needed for 'trivially parallel' high throughput needs

- depending on architecture cache coherency
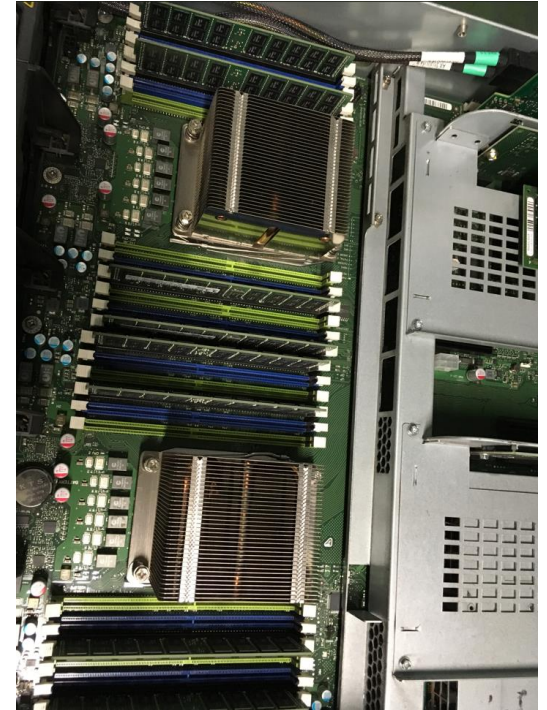  kills single-thread performance (although AMD did lot better here than the Intel *lakes)

Image: dual-socket Fujitsu system at the Xenon experiment site, 2019. source: Tristan Suerink, Nikhef

# CPU design changes may fit application, or not

AMD EPYC effective for applications like WLCG:

- Naples → Rome added shared memory die
- links all cores directly to memory

Rome-Milan improvement?

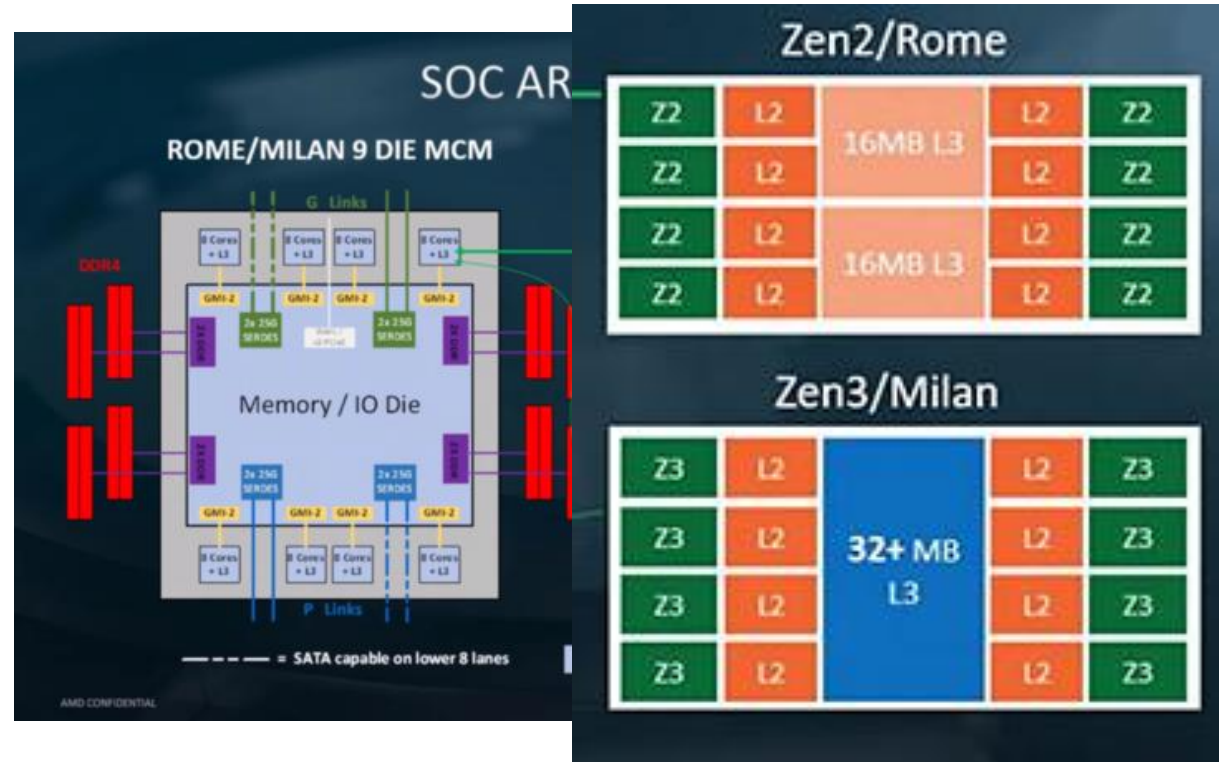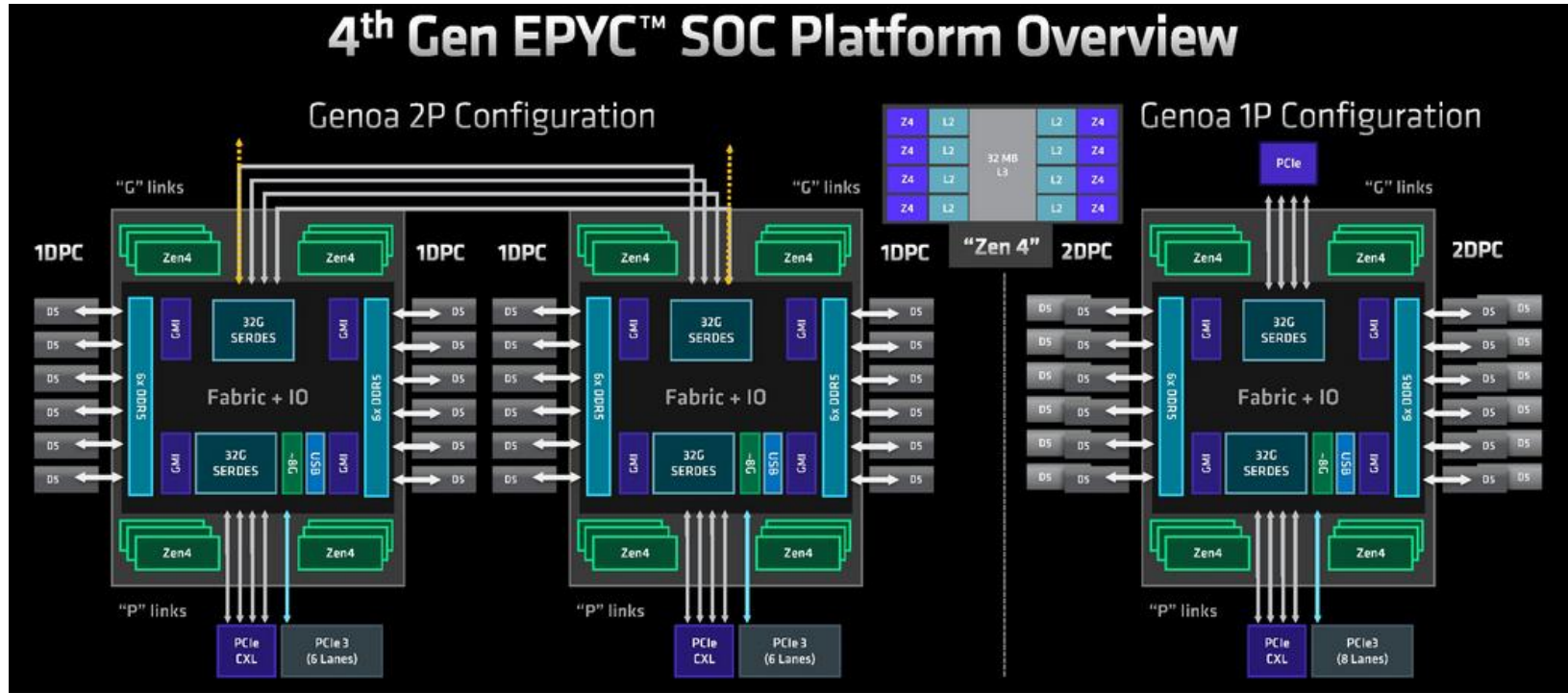- shared L3 cache benefits tightly coupled HPC, but not 'off-die memory' limited HTC
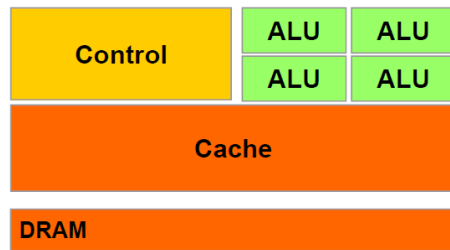


Image source: AMD, retrieved from https://m.hexus.net/tech/news/cpu/135479-amd-shares-details-zen-3-zen-4-architectures/
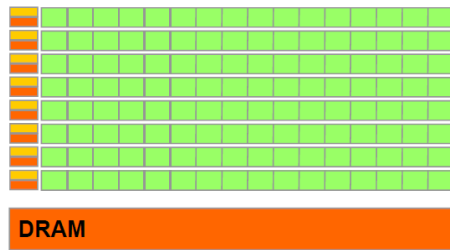
# Scaling up, more examples



AMD EPYC Genoa platform, image from https://www.semianalysis.com/p/amd-genoa-detailed-architecture-makes

# Accelerators – general purpose GPUs



CPU

GPU

*leaving FPGAs out for a moment – but those are particularly useful in guaranteed-latency scenarios!*

- but co-processing comes at a cost of moving data to and from the GPU
- often faster to keep computing and do selection & conditionals later
- computation speed heavily depends on precision (even 4-bit precision is used)
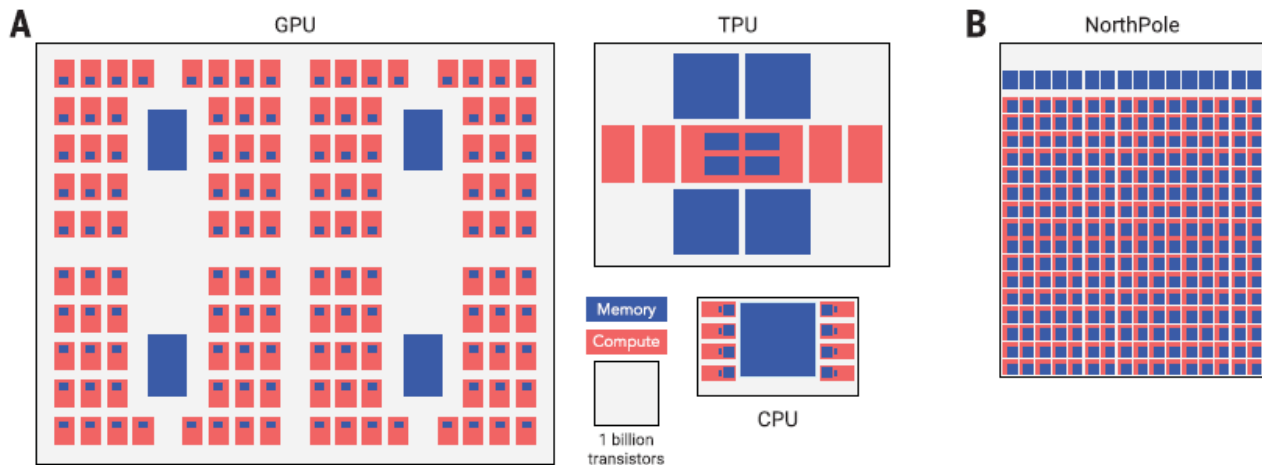- quite power hungry!



Image: 'Massively Parallel Computing with CUDA', Antonino Tumeo Politecnico di Milano, https://www.ogf.org/OGF25/materials/1605/CUDA_Programming.pdf
Floorplan image of die: AMD MI250 GPU, slide source: AMD

# Aiming to remove the data access bottleneck

Separating memory from processing introduces the memory misses that slow down CPU processing as well GPUs due to need for (RDMA) main memory access
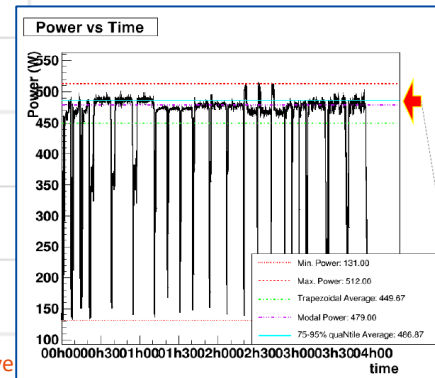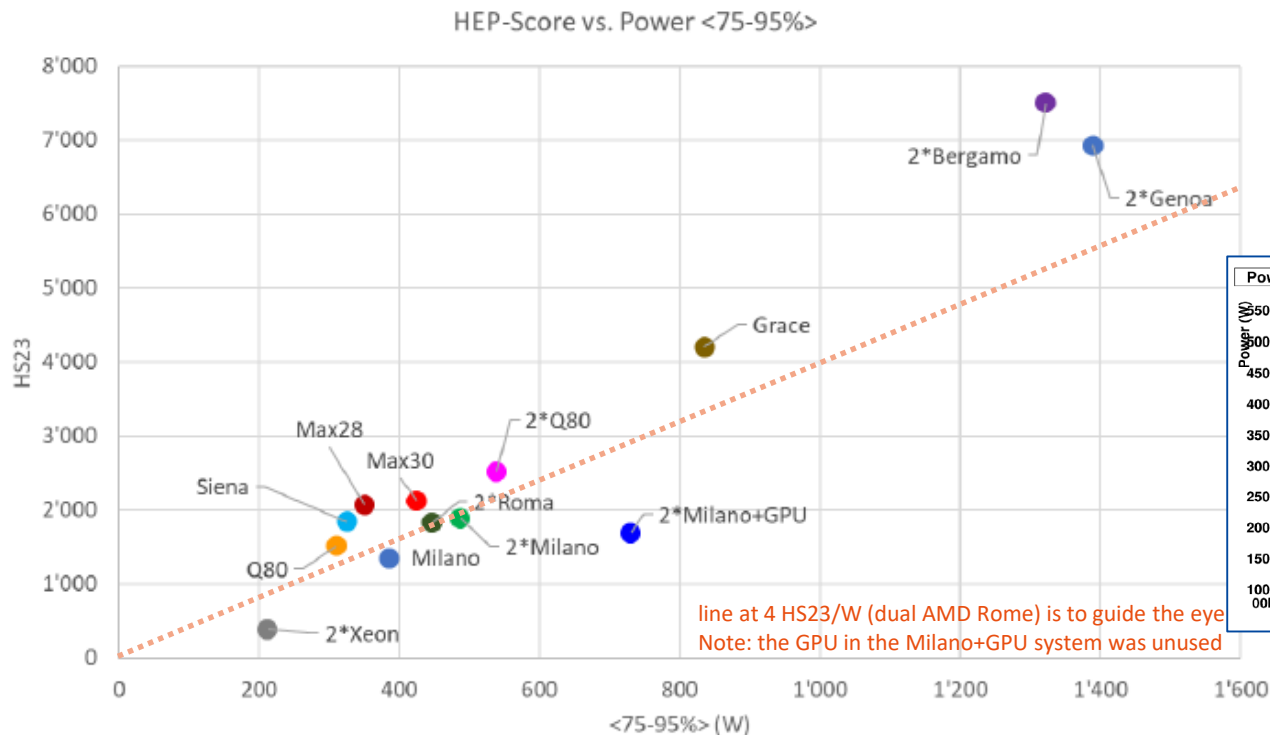
Some very recent designs aim to eliminate this by temporal co-location of program and memory (IBM NorthPole AI, Oct '23) with data-flow driven compute



Physical organization of on-chip memory (blue) and compute (red) are diagrammed for representative processors, scaled to constant transistors per unit area. From Modha's paper Modha et al., *Science* **382**, 329–335 (2023)

Modha *et al.* https://doi.org/10.1126/science.adh1174 or read https://research.ibm.com/blog/northpole-ibm-ai-chip
PCIe card photo from https://www.ibm.com/blogs/solutions/jp-ja/northpole-ibm-ai-chip/

# The energy bottleneck: architecture figure of merit



HEP-Score vs. Power <75-95%>

Data and graphs: Emanuele Simili, Glasgow University, at CHEP2024 (https://indico.cern.ch/event/1338689/contributions/6011562/)
HEPSPEC23 benchmark: https://gitlab.cern.ch/hep-benchmarks/hep-benchmark-suite ('memory-intensive' high throughput processing application benchmark)

# How to get this heat out ... in liquid form, maybe?

- Heat capacity of liquid is much larger than air
- by now (almost) standard for HPC systems

- immersive systems
  look cool, but are 'a bit
  hard' on maintenance



Strongly depends on systems engineering:
when water inlet temperature can be >40
degC, you have almost always free cooling

Image source dual-board system: Lenovo, ThinkSystem SD650
immersive cooling image https://hypertec.com/blog/sustainable-emerging-tech-liquid-immersion-cooling/, PIC T1 centre, Barcelona, ES

# And if large-scale IT does not quite fit ... ahum ...



Image source: https://lambdalabs.com/products/blade

SuperMicro (branded as 'Lambda Blade')
4U chassis, supporting 10 consumer-grade GPUs ...
... with a bump

# but there *is* a serious issue here!



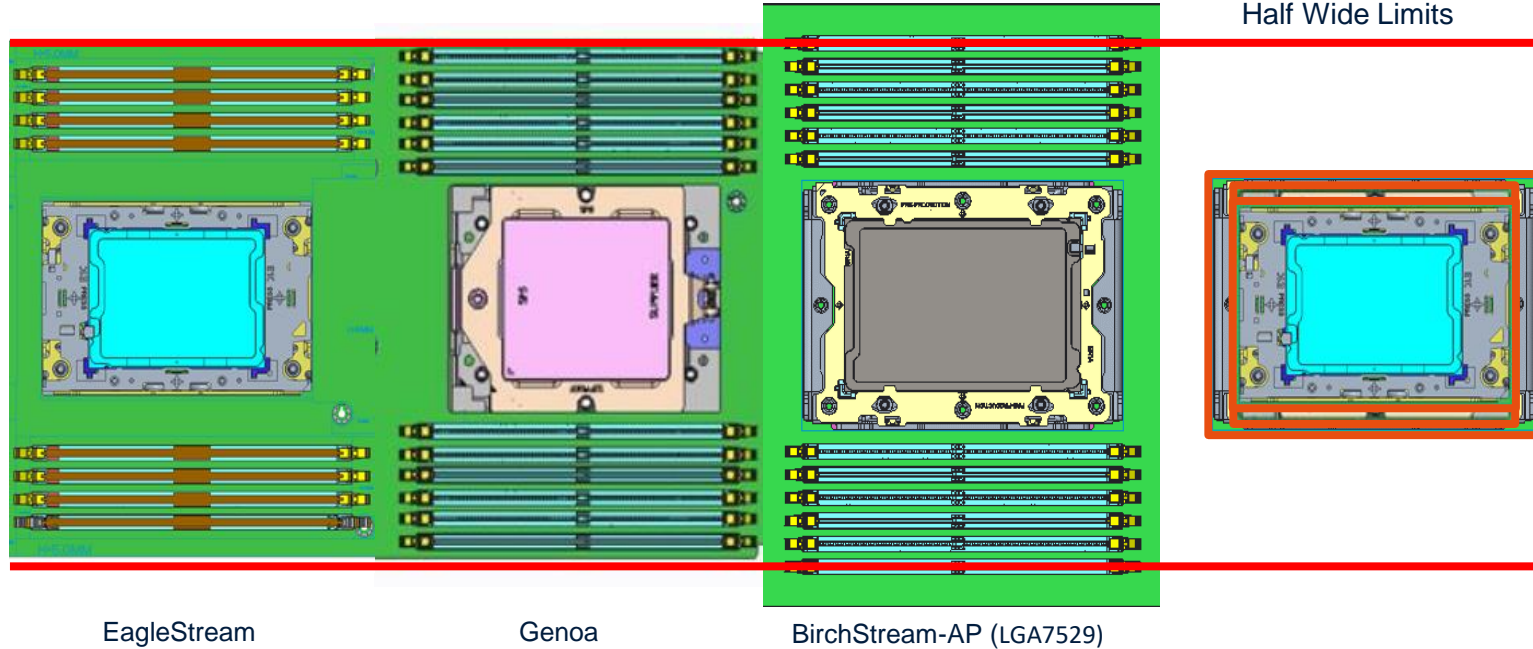Half Wide Limits

EagleStream          Genoa          BirchStream-AP (LGA7529)

Image thanks go to Rick Koopman – Lenovo at the HTCondor Workshop 2024 https://indico.cern.ch/event/1386170/

# Scaling up – beyond one lone motherboard

# 'compute farming': milking computers, in a balanced way

**Data-driven** workloads (like WLCG, SKA, WeNMR) need more than just compute

- **balanced features** for node throughput
  CPU, storage, memory bandwidth
  & latency, NIC & network speed

- **single-socket** multicore systems are fine,
  typically 64-128 cores per system
- **network**: 2x25Gbps (match #cores)
- **memory**: say ~ 8 GiB/core
- **local disk**: 8-16 TB NVME (~100GB/core)
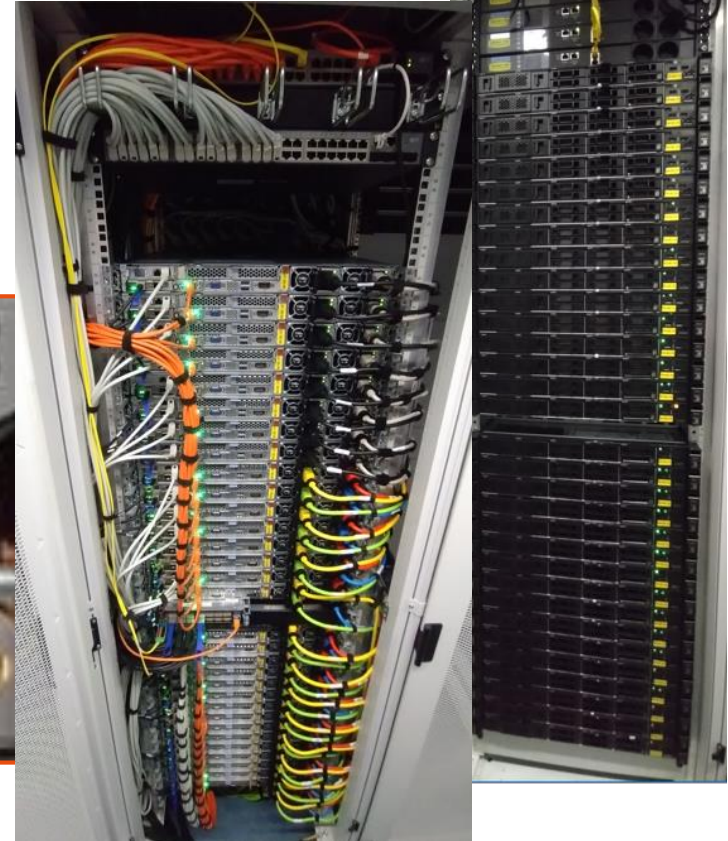- + space (physical + power) to add some **GPU**



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

# But … fancy an interactive console install?



Images: Nikhef Housing H234b NDPF science processing data centre

# Managing multiple system (physical or virtual)

**Fabric (Configuration) Management**
- do you know what is out there?
- update quickly & consistently when vulnerabilities are found?
- versioned repository for rollback?

**note that not all tooling scales in itself**
- **push**: ansible (using ssh logins), or home-brew scripting
- **pull**: each node runs its own actions, e.g. Saltstack, Ansible-agent, Quattor, Chef, …
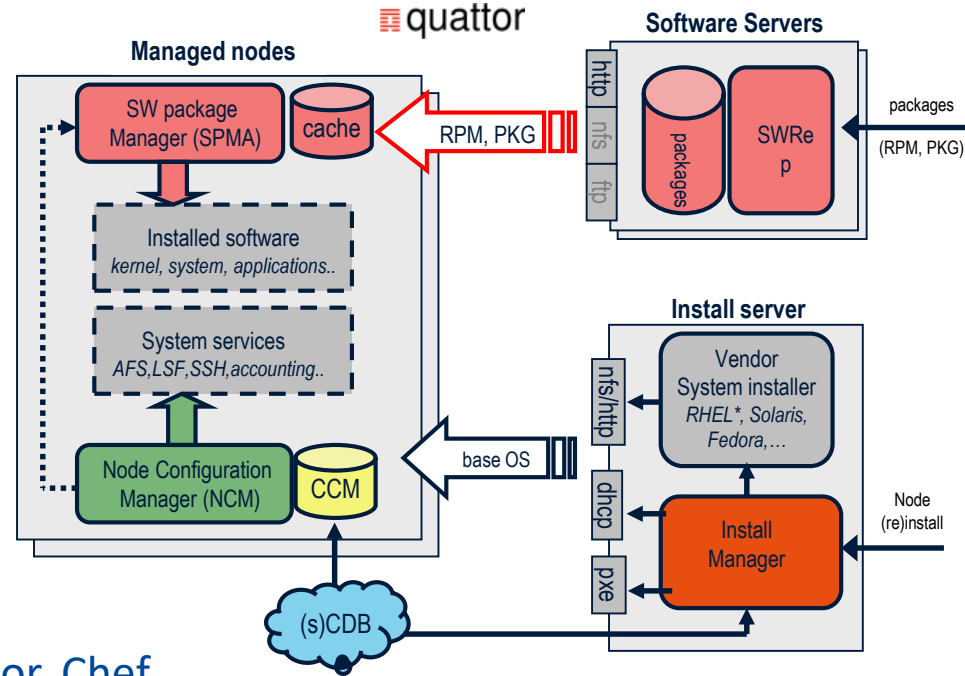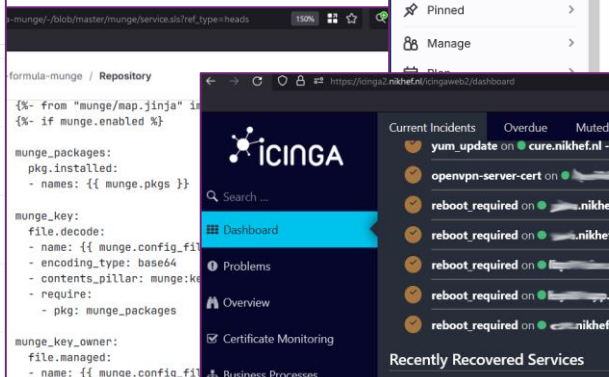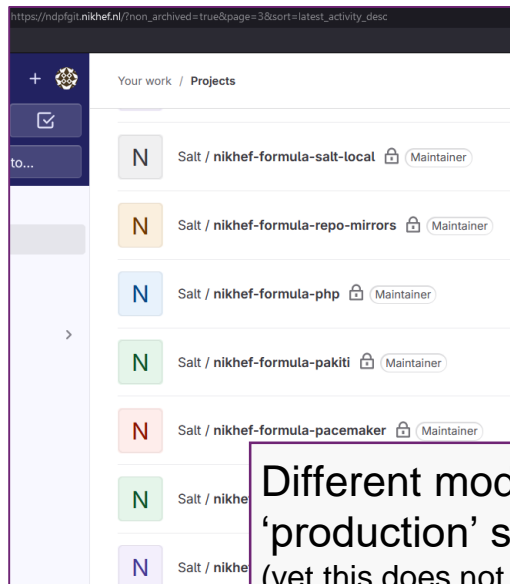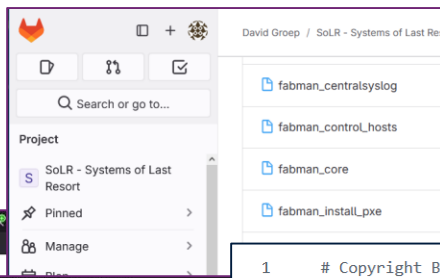


Illustration: German Cancio, CERN, quattor.org, used here as example; see also: ansible.com, saltproject.io, theforeman.org, cfengine.com, puppet.com, …

# Towards 'Software Defined Infrastructure' …



Different modalities are fine as long as all 'production' systems are managed and monitored
(yet this does not apply – for a reason - to the experimental technologies platform and Nationale Speeltuin)

Nikhef NDPF Salt & Reclass (Dennis van Dok, Andrew Pickford, Mary Hester); SoLR Ansible;
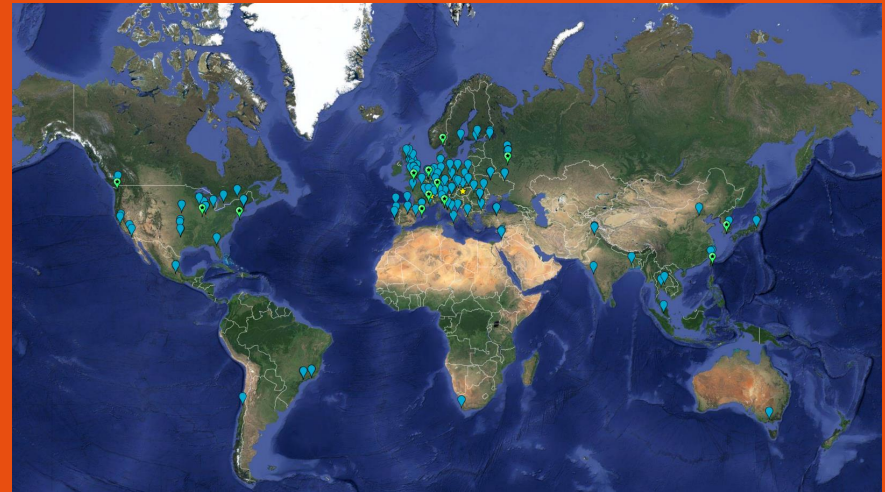Docker Compose for sharemd.nikhef.nl ; example Helm chart from
https://github.com/bitnami/charts/blob/main/bitnami/wordpress/

# More of *more than one* …

The physical layer … and managing software-defined infrastructure



Large-scale IT: worldwide LHC Computing and beyond (2024 ed)

proceed to cloud management ▶

# Cluster computing and 'conveniently parallel' HTC



```
GROUPCFG[auger]      FSTARGET=3     PRIORITY=200    MAXPROC=500     QDEF=augerbig
GROUPCFG[augsgm]     FSTARGET=1     PRIORITY=300    MAXPROC=2       QDEF=augerbig
QOSCFG[augerbig]     FSTARGET=3

# if these are queued, they will generally be of highest priority.
# limit their MAXIJOBs ... we really want two non-ATLAS VOs to be
# of rank higher than ATLAS before we drain the multicore pool.

GROUPCFG[virgo]      FSTARGET=25    PRIORITY=200    MAXPROC=2700 MAXIJOB=10 QDEF
=biggrid
GROUPCFG[ligo]       FSTARGET=23    PRIORITY=200    MAXPROC=2700 MAXIJOB=10 QDEF
=biggrid

# local groups

GROUPCFG[atlas]      FSTARGET=10    PRIORITY=200    MAXPROC=2200    QDEF=niklocal
```

- 'like milking cows' (if you feed them lots of power first)
- parallel access to data comes at a cost of high IOPS

# Batch system platform

Many things in life are *conveniently parallel*
- HEP events & simulation
- ligand matching & drug discovery
- structural biochemistry
- …

**challenge not in parallelism itself**
- we have had HPC systems for ages

**but**
- large numbers of (single-core) jobs
- heterogeneous workloads
  sharing the same set of worker nodes
- computing with concurrent data access

# Scalable submission: HTCondor

Matchmaking based on 'ClassAds'

- both jobs and machines advertise their requirements and capabilities in 'classified advertisements'
- Matchmaking done by the negotiator execution nodes mostly autonomous

helps for scalability and resilience



HTCondor, Miron Livny *et al.*; Compiled from Todd Tannenbaum (2024 HTCondor Workshop) https://indico.cern.ch/event/1386170/contributions/6127903/

# Dutch National e-Infrastructure: High Throughput GINA



Cumulative ncores per VO (SLURM)

**Communities**

**ENMR**: structural biochemistry
**Project MinE**: ALS (health)
**Xenon:** direct DM searches
**TROPOMI**: earth observation
**DUNE**:
- long baseline neutrinos

**LIGO/Virgo**:
- Gravitational waves

**Alice, ATLAS, LHCb**
- LHC experiments (in NL)

Graphic: GINA DNI compute service coordinated by SURF

# Estimated Response Time (and predicting it)

- 'Fair share' – distributing load over time in a 'continuous job supply' system



Image: Nikhef NDPF DNI "Grid" cluster. Period: October 6-17, 2022; top-5 communities; GRISview images: Jeff Templon
For work on run time prediction in high-occupancy clusters , see Hui Li *Workload characterization, modeling, and prediction …* https://hdl.handle.net/1887/12574

# For occupancy, intended target audience makes a difference

For organized 'production' computing (planned months in advance in WLCG)
- *predictable* **scheduling** is more important (steady flow of results)
- **maximizing efficiency**: resource cost is the limiting factor in (physics) results
- co-scheduling with data (pre-placement) is required
- community-authorization based access to data sources only

For 'local' users, e.g. students whose progress tomorrow depends on results *today*
- *response time* is more important than efficiency
- fast turn-around/short waiting times by heterogeneous ('competing') user base
- data access must be parallelism-ready, but is 'always' local on-site
- local storage credentials and sharing with desktop and Jupyter environments

*so offering two distinct classes of services is (in this case) intentional*

# NDPF local analysis cluster 'Stoomboot'

period: March 2021 .. October 2022

Running jobs:



Waiting jobs (Week 40, 2022):



Source: NDPF Statistics overview, https://www.nikhef.nl/pdp/doc/stats/ - GRISview images: Jeff Templon for NDPF and STBC

# High *throughput* computing is in the end about data



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

# Can storage support your parallel processing

Basic storage properties

- throughput
- IOPS – I/O Operations per Second
- seek-time

but not many **file systems** support *concurrent parallel access* by many clients

- both data **and** (file system or index) meta-data must be scalably distributed
- typically sacrifice either instant consistency, or (POSIX) semantics,
  (or scalability) in a distributed storage system

Common commercial solutions: GPFS, … but also NetApp, HDS, Dell-EMC, have theirs
Common open source: BeeGFS, gluster, dCache, CephFS, Lustre, …

*… **likely do not use a file system if object storage does the job,** but then you need a catalogue/database*

# Example: client-side managed GlusterFS

- scalable through independence of both clients and servers

- design is stateless: file system meta-data kept in each server's file system

- data itself can be replicated and protected but ... inconsistencies in metadata linger around the corner in case of client failures (e.g. batch system worker nodes)



Image source Gluster community: https://docs.gluster.org/en/main/Quick-Start-Guide/Architecture/

# Example: server-coherent distribution – dCache

- separate client entry points, storage access scheduling, filesystem meta-data (namespaces), and storage
- message layer for eventual consistency
- redirect-based access
  - doors and pools usually on all nodes
  - now also feature of standard NFSv4.1



Images: Tigran Mkrtchyan (DESY, dCache.org), *dCache on steroids - delegated storage solutions*, ISGC 2016, https://dcache.org/manuals/publications.shtml

# dCache: wide area distribution

- can be widely (long latency) distributed
  - Nordic Data Grid Facility: Sweden is quite long (~16ms RTT), and Ljubljana to Umeå is ~30ms RTT (~ 2900km)

- redirect-then-access model limits interactions with any single node across a long-distance links

- at 'cost' of POSIX features like *atime* or concurrent write
  - most distributed applications don't need these anyway
  - but indeed it's not a good backing store for databases ☺



The NDGF dCache instance spans datacentres across Scandinavia and Slovenia, but is administered and used as a single instance.

# Structure of application data placement impacts storage (hardware) systems design

pre-staging all data locally allows for **latency hiding**, posix-style access with lseek(2), and a fast, local, '$TMPDIR'
*e.g. why there are Data Transfer Nodes (DTNs) in the 'Science DMZ' concept*



**but**, nowadays, pre-staging started coming at a cost, when using **SSDs** as local 'scratch' area … because of their hardware characteristic 'endurance'

Photo HGST nVME from: Dmitry Nosachev on Wikimedia Commons CC-BY-SA; Image Science DMZ and Data Transfer Nodes: ESnet fasterdata.es.net

# Especially with *WORN* storage: Write Once Read Never

Frequency distribution of **read-back vs. write** volume, observed on local scratch for NDPF execution nodes for *outside ('grid') access (blue) vs local access (orange)*

**Access pattern is rather different. But why?**

- external users pre-stage, because it is built into data management frameworks (like DIRAC, Athena),
- 'local' users stream output data (dCache with NFSv4) and use $TMPDIR mainly for merging partial results

Different types of workload (here analysis vs processing) determine the choice of systems hardware



Data: NDPF execution nodes, based on SSD SMART data, integrated over total device lifetime; plot shows number of local analysis nodes scaled to DNI-WLCG count; collected using smartctl on 2020-10-28 – in total 97 'DNI' and 34 'STBC' SSDs were used in the analysis

# As a service!

'Cloud' Services and
Service Management



There is NO CLOUD, just other people's computers

Maastricht University | Department of Advanced Computing Sciences

# Scaling things '… as a service'

**The managed servers usually are not physical**

- although there is lots of 'fixed' virtualization of systems, network and (block) storage

When scale, or environment, must be flexible, you get *software defined infrastructure*

- IaaS: Infrastructure as a Service
- PaaS: Platform as a Service
  (containers, but also a batch system …)
- SaaS: Software as a Service
  (like the science application portal like WeNMR)

powerful tools, but also easy to get wrong (i.e. having plain-text secrets in the version control system to automate redeployment). And abstractions are *leaky*!



Image from CERN's OpenShift, A Lossent *et al* 2017 *J. Phys.: Conf. Ser.* **898** 082037 https://doi.org/10.1088/1742-6596/898/8/082037

# Moving the management boundary

## Infrastructure-as-a-Service

| |
|---|
| Application |
| Data |
| Runtime environment |
| Middleware |
| Operating system |
| Virtualisation layer |
| Physical server |
| Storage devices |
| Network |

Guest

Hyper visor
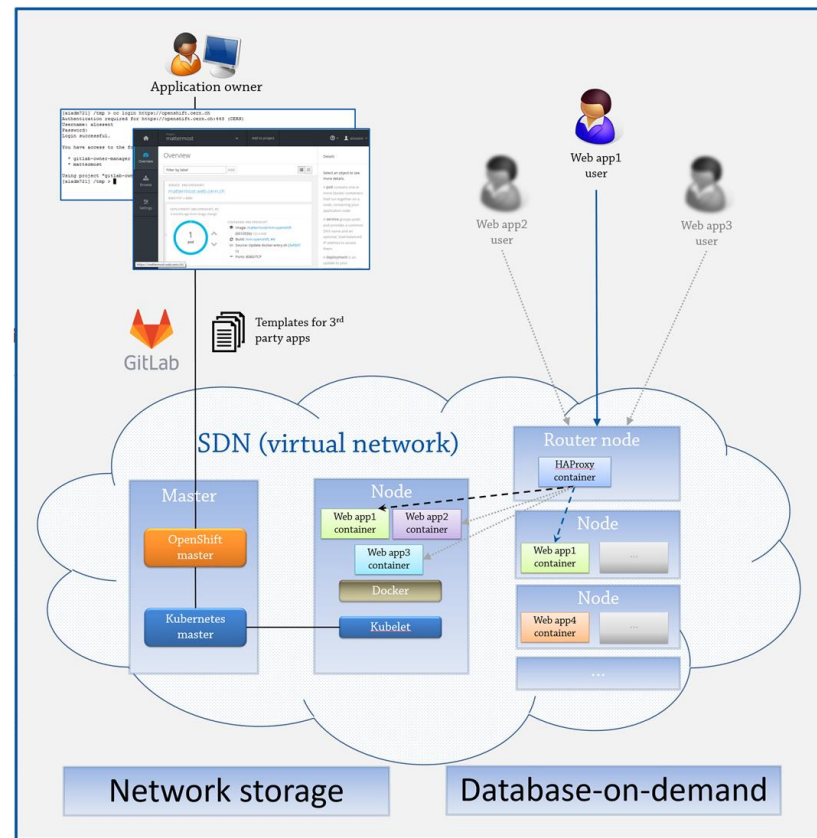
Host

## Platform-as-a-Service

| |
|---|
| Application |
| Data |
| Runtime environment |
| Middleware |
| Operating system |
| Virtualisation layer |
| Physical server |
| Storage devices |
| Network |

## Software-as-a-Service

| |
|---|
| Application |
| Data |
| Runtime environment |
| Middleware |
| Operating system |
| Virtualisation layer |
| Physical server |
| Storage devices |
| Network |

Astronomy catalogue: https://vizier.cds.unistra.fr/

IaaS: openstack.com, Oracle OCI; PaaS: dsri.maastrichtuniversity.nl, apptainer.org, cvmfs.readthedocs.io,  kubernetes.io, slurm.schedmd.com; SaaS: Jupyter.org

# 'Cloudification' eases systems management …



OpenShift (OKD) system at CERN (accessible for CERN users only) – at Maastricht use the DSRI infrastructure: https://dsri.maastrichtuniversity.nl/

# Common interfaces to the different clouds?

'protocol hourglass'



hourglass image: Alessio Merlo in The Condor on the Grid: state of art and open issues,

# Standard interfaces for compute and data?

hourglass model 'kind-of' worked for IP and web with http as common standard

- a very simple stateless interface

protocols for higher-level services never quite reached this level of global interop

- requirements too complex and stateful
- use cases were usually scoped

slowly changing now but only for similarly simple things, like on-line object storage

Is distributed computing too bespoke …?



Interoperable cloud? Compare OGF's OCCI WG GFD.221 (https://www.ogf.org/documents/GFD.221.pdf) with e.g. Amazon S3 API or the OwnCloud CS3 interfaces

# DIRAC: spanning heterogeneous resource models

Add a scheduling layer!

'any (IT) problem can be solved by adding an extra level of indirection'

*DIRAC is just one example*



Image: DIRAC project, A. Tsaregorodtsev *et al.* CPPM Marseille, from https://dirac.readthedocs.io/ ; CVMFS (CERN VM File System) is a common software distribution platform using distributed signed data objects in a cached hierarchy using CDN techniques, see https://cernvm.cern.ch/fs/

# An overlay network of containers

*Nobody wants a cloud per-se … what folk want is a solution …*



'alien containers' HPC integration - container computing, using curated application images

Image sources: NDPF JupyterHub service "Callysto";  SLATE: Service Layer At The Edge – Rob Gartner (UChicago), Shawn KcMee (UMich) *et al.* – slateci.io

# Containerised workloads: between 'PaaS' and 'SaaS'



Images: Oksana Shadura et al (UNebraska Lincoln), Brian Bockelman (Morgridge Institute) at CHEP2023 https://indico.jlab.org/event/459/contributions/11610/

# Services: serving the users, not the IT department

What actually drives your IT architecture
and service management system?

- strategic requirements of your organisation?
- what are 'appropriate' service levels
  and impact in case things go differently?
- balancing stability, innovation, and engagement

Potentially separate 'service classes'

- 'enterprise' services
- research computing services (the 'primary business')?
- experimental services (innovation, future strategy)?

Alternative approach to 'IT landscape' architecture https://go.nikhef.nl/principles-of-digitalisation

# Service portfolios – catalogues are nice, up to a point



Catalogues from Nikhef, European Open Science Cloud EU Node (free VMs for 'all' researchers, subject to https://open-science-cloud.ec.europa.eu/system/files?file=2024-10/EOSC-EU-Node-User-Access-Policy-v1.0.pdf)

# Example: FitSM – Federated IT Service Management

Structuring service management
- ISO 20k
- https://www.fitsm.eu/
- ITIL (now at ITIL v3)

and a whole bunch of others, like COBIT, AgileSM, …

Slide with PR list from https://www.fitsm.eu/

FitSM: ITSM process framework

1. Service portfolio management (SPM)
2. Service level management (SLM)
3. Service reporting management (SRM)
4. Service availability & continuity management (SACM)
5. Capacity management (CAPM)
6. Information security management (ISM)
7. Customer relationship management (CRM)
8. Supplier relationship management (SUPPM)
9. Incident & service request management (ISRM)
10. Problem management (PM)
11. Configuration management (CONFM)
12. Change management (CHM)
13. Release & deployment management (RDM)
14. Continual service improvement management (CSI)

**Core management processes for any IT service**

14

# Putting 'more than one' thing together

Connecting the data:
The Internet Is Not Enough!

Large-scale IT: worldwide LHC Computing and beyond (2024 ed)

# 'Elephant streams in a packet-switched internet'

*'You may have plenty of shovels,*
*but where to leave the sand?'*

- wheelbarrow works fine in your garden
- want to send it to different places?
  Use waggons on a train,
  or ships with containers
- always from A-to-B?
  A conveyer belt will do much better!

…  although you still need
   a hole to dump it in …



Image conveyor belt tunnel near Bluntisham, Cambridgeshire by Hugh Venables, CC-BY-SA-4.0 from https://www.geograph.org.uk/photo/4344525

# A quick look at internet routing …

network paths
from various places
in Western Europe

towards an IP address at CERN



Data: RIPE NCC Atlas project, TraceMON IPmap, atlas.ripe.net, measurement 9249079

# Many paths to Rome … i.e. to your server

- From a home connected to Freedom Internet to *spiegel.nikhef.nl*

```
[root@kwark ~]# traceroute -6 -A -T gierput.nikhef.nl
traceroute to gierput.nikhef.nl (2a07:8500:120:e010::46), 30 hops max, 80 byte packets
 1  2a10-3781-17b6.connected.by.freedominter.net (2a10:3781:17b6:1:de39:6fff:fe6b:4558) [AS206238]  0.810 ms  1.052 ms  1.330 ms
 2  2a10:3780::234 (2a10:3780::234) [AS206238]  7.460 ms  7.655 ms  7.705 ms
 3  2a10:3780:1::21 (2a10:3780:1::21) [AS206238]  8.868 ms  9.054 ms  9.103 ms
 4  et-0-0-1-1002.core1.fi001.nl.freedomnet.nl (2a10:3780:1::2d) [AS206238]  10.017 ms  9.934 ms  10.263 ms
 5  as1104.frys-ix.net (2001:7f8:10f::450:66) [*]  10.898 ms  11.744 ms  11.797 ms
 6  gierput.nikhef.nl (2a07:8500:120:e010::46) [AS1104]  11.502 ms  7.800 ms  7.357 ms
```
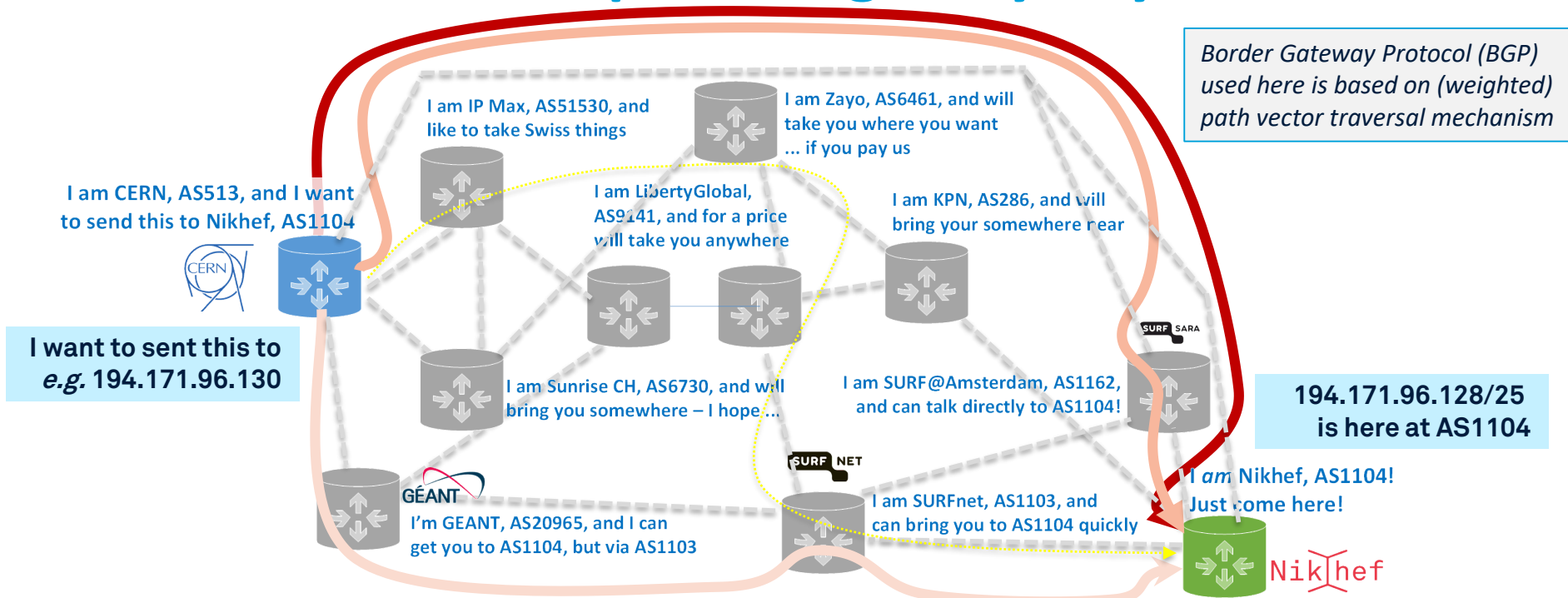
- but from Interparts in Lisse, NH:

```
[root@muis ~]# traceroute -6 -A -I gierput.nikhef.nl
traceroute to gierput.nikhef.nl (2a07:8500:120:e010::46), 30 hops max, 80 byte packets
 1  2a03:e0c0:1002:6601::2 (2a03:e0c0:1002:6601::2) [AS41960]  1.380 ms  1.371 ms  1.369 ms
 2  2a02:690:0:1::b (2a02:690:0:1::b) [AS41960]  1.305 ms  1.312 ms  1.312 ms
 3  et-6-1-0-0.asd002a-jnx-01.surf.net (2001:7f8:1::a500:1103:2) [AS1200]  1.957 ms  2.000 ms  2.052 ms
 4  ae47.asd001b-jnx-01.surf.net (2001:610:e00:2::49c) [AS1103]  2.443 ms  2.505 ms  2.507 ms
 5  irb-4.asd002a-jnx-06.surf.net (2001:610:f00:1120::121) [AS1103]  2.041 ms  2.138 ms  2.138 ms
 6  nikhef-router.customer.surf.net (2001:610:f01:9124::126) [AS1103]  8.977 ms  7.957 ms  7.951 ms
 7  gierput.nikhef.nl (2a07:8500:120:e010::46) [AS1104]  7.922 ms  8.093 ms  8.081 ms
```

AS41960: Interparts; AS1200: AMS-IX route reflector; AS1103: SURFnet; AS1104: Nikhef; AS206238: Freedom Internet – on the FrysIX there is direct L2 peering

# Where do internet packets go anyway?



Border Gateway Protocol (BGP) used here is based on (weighted) path vector traversal mechanism

I am IP Max, AS51530, and like to take Swiss things

I am Zayo, AS6461, and will take you where you want ... if you pay us

I am CERN, AS513, and I want to send this to Nikhef, AS1104

I am LibertyGlobal, AS9141, and for a price will take you anywhere

I am KPN, AS286, and will bring your somewhere near

I want to sent this to *e.g.* 194.171.96.130

I am Sunrise CH, AS6730, and will bring you somewhere – I hope ...

I am SURF@Amsterdam, AS1162, and can talk directly to AS1104!

194.171.96.128/25 is here at AS1104

I'm GEANT, AS20965, and I can get you to AS1104, but via AS1103

I am SURFnet, AS1103, and can bring you to AS1104 quickly

I *am* Nikhef, AS1104! Just come here!

grey-dash lines for illustration only: may not correspond to actual peerings or transit agreements; red lines: the three existing LHCOPN and R&E fall-back routes; yellow: public internet fall-back (least preferred option)

# Announcing routes: the Border Gateway Protocol

```
davidg@deelqfx-re0> show route receive-protocol bgp 192.16.166.21 table LHCOPN

LHCOPN.inet.0: 316 destinations, 344 routes (316 active, 0 holddown, 0 hidden)
  Prefix                    Nexthop              MED      Lclpref     AS path
* 109.105.124.0/22          192.16.166.21        10                   513 39590 I
* 117.103.96.0/20           192.16.166.21        10                   513 24167 I
* 128.142.0.0/16            192.16.166.21        10                   513 I
* 130.199.48.0/23           192.16.166.21        10                   513 43 ?
* 130.199.185.0/24          192.16.166.21        10                   513 43 ?
* 130.246.176.0/22          192.16.166.21        10                   513 43475 I
```
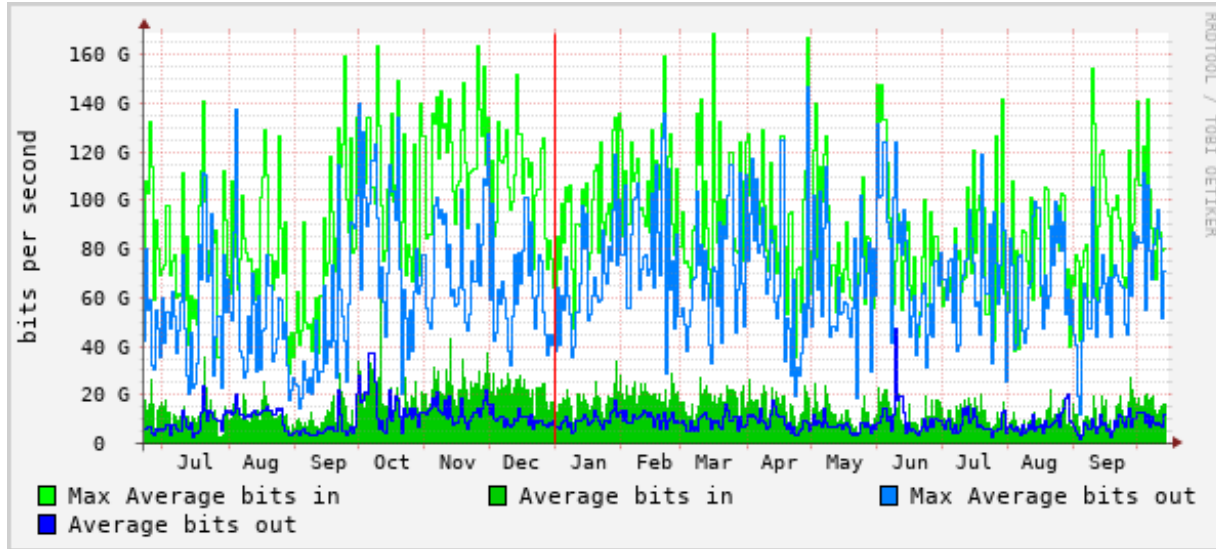
```
davidg@deelqfx-re0> show route advertising-protocol bgp 192.16.166.21 table LHCOPN

LHCOPN.inet.0: 316 destinations, 344 routes (316 active, 0 holddown, 0 hidden)
  Prefix                    Nexthop              MED      Lclpref     AS path
* 192.16.186.160/30         Self                                      I
* 194.171.96.128/25         Self                                      I
* 194.171.98.112/29         Self                                      I
```

IPv4 routes advertised from AS513/CERN (for all sites on LHCOPN) to AS1104/Nikhef (top), and the routes announced by AS1104/Nikhef to CERN, on 5 Nov 2022

# Typical data traffic to and from the processing cluster

# Network is more than just what it says on the tin

More network bandwidth does
not mean your *data* gets there faster

- memory requirements (since TCP needs a capability to re-transmit)

- tcp 'slow start'
- congestion control algorithms

TCP throughput calculator

**Theoretical network limit**
rough estimation: rate < (MSS/RTT)*(C/sqrt(Loss)) [ C=1 ] (based on the Mathis et.al. formula)
network limit (MSS 9000 byte, RTT: 150.0 ms, Loss: $2.304*10^{-11}$ ($2*10^{-09}$%)) : **100000.00 Mbit/sec.**

**Bandwidth-delay Product and buffer size**
BDP (100000 Mbit/sec, 150.0 ms) = **1875.00 MByte**
required tcp buffer to reach 100000 Mbps with RTT of 150.0 ms >= **1831054.7 KByte**
maximum throughput with a TCP window of 1831054 KByte and RTT of 150.0 ms <= **100000.00 Mbit/sec.**

Useful sources: https://www.switch.ch/network/tools/tcp_throughput/, https://fasterdata.es.net/
tcp slow-start graphic from Abed et al, *Improvement of TCP Congestion Window over LTE- Advanced Networks* IJoARiC&CE  2012

# The cat video that destroyed it all …

- TCP protocol sensitive to packet loss
  - 3 lost packets is enough to trigger this

- different congestion avoidance algorithms exists (~20 by now)

- loss severely impacts links w/large 'bandwidth-delay-product' (BDP)
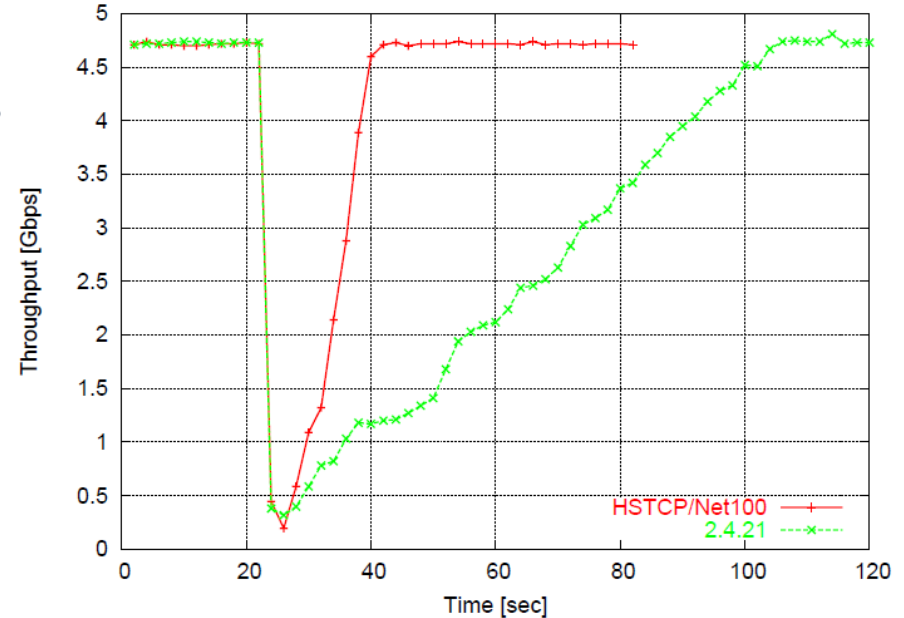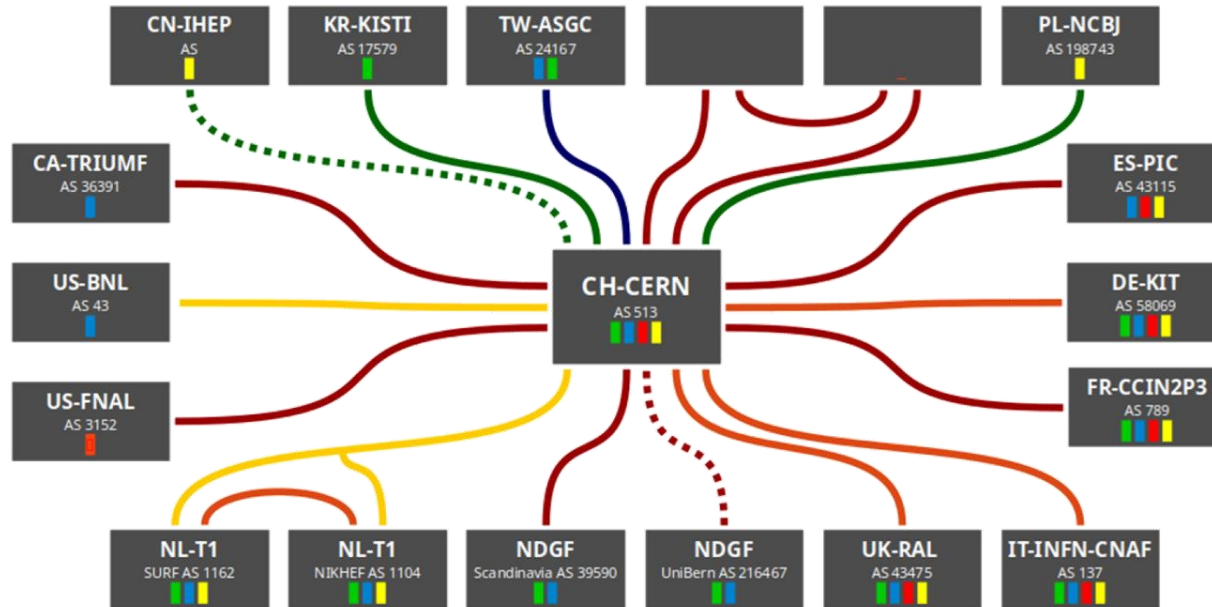
- NL: ~3 ms, US East: 150ms



Figure 10: HSTCP versus stock TCP recovery time

source: Catalin Meirosu et al. *Native 10 Gigabit Ethernet experiments over long distances* in FGCS, doi:10.1016/j.future.2004.10.003 – aka. ATL-D-TN-0001

# LHCOPN – distributing raw data



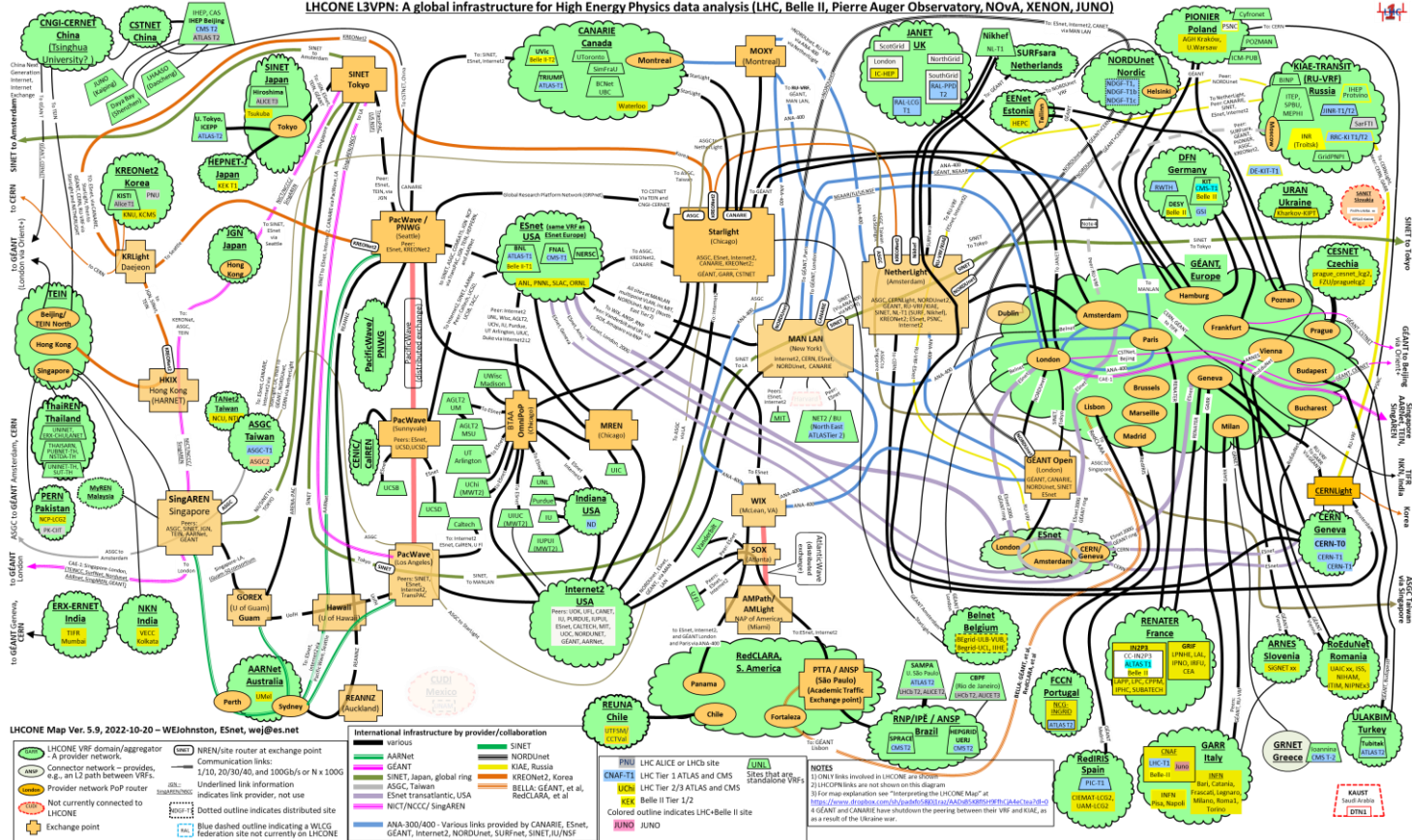Image source: Edoardo Martelli, CERN, https://lhcopn.web.cern.ch/

# LHCOPN – traffic levels for data transfer (DC24)



From Lassnig, M., & Wissing, C. et al. (2024). WLCG/DOMA Data Challenge 2024: Final Report. Zenodo. https://doi.org/10.5281/zenodo.11444180

# LHCone



LHCone ("LHC Open Network Environment") – visualization by Bill Johnston, ESnet version: October 2022 – updated with new AS1104 links
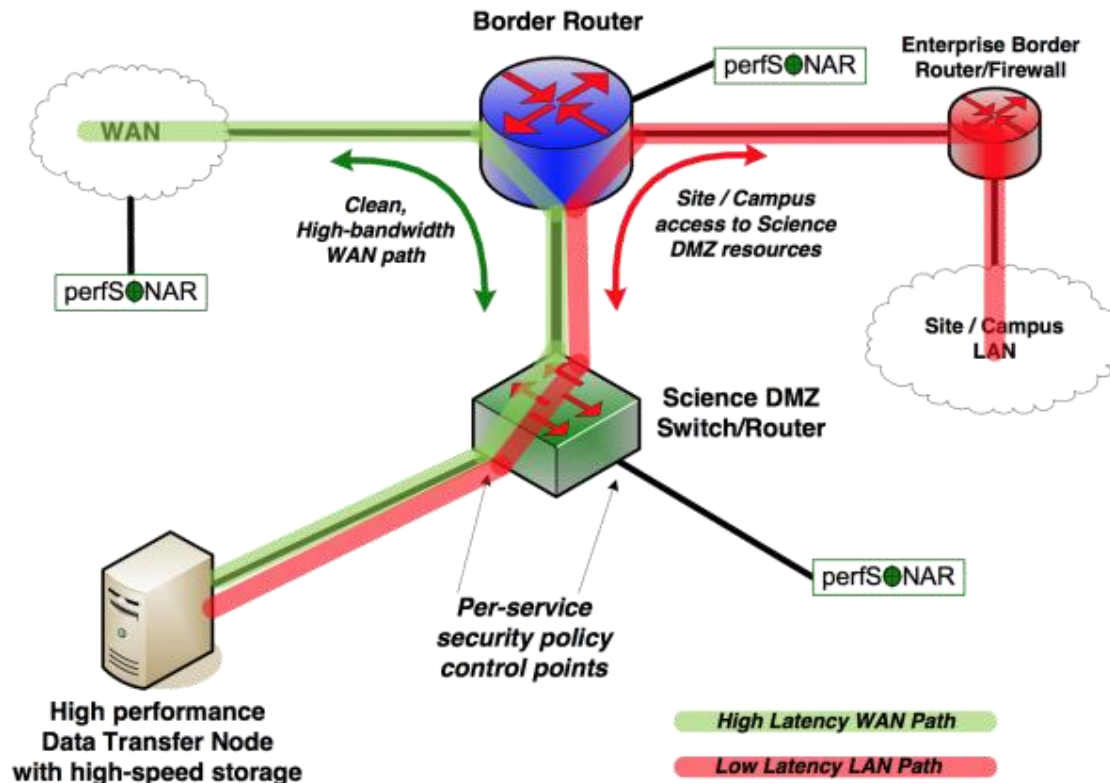
# 'ScienceDMZ'

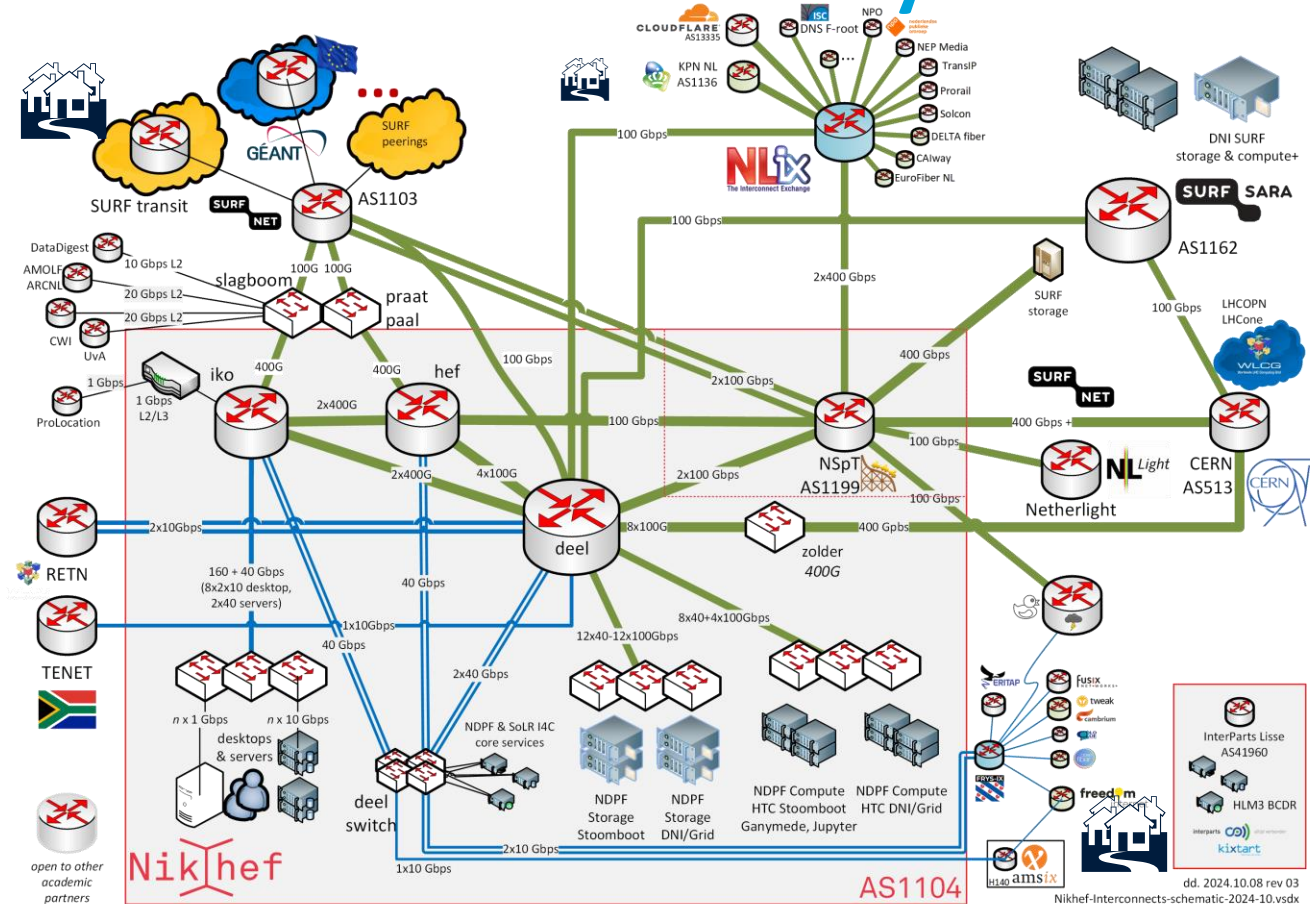**Predicable performance and data access for research**

'**where research services, data, and researchers meet**'

- latency hiding through caching
- **security zoning/segmentation** protects specific data sets
- **outside any enterprise perimeter**

Image and 'ScienceDMZ' concept promulgated by ESnet (see fasterdata.es.net)

# Just one random autonomous system: AS1104



AS1104
state as of Oct 2024
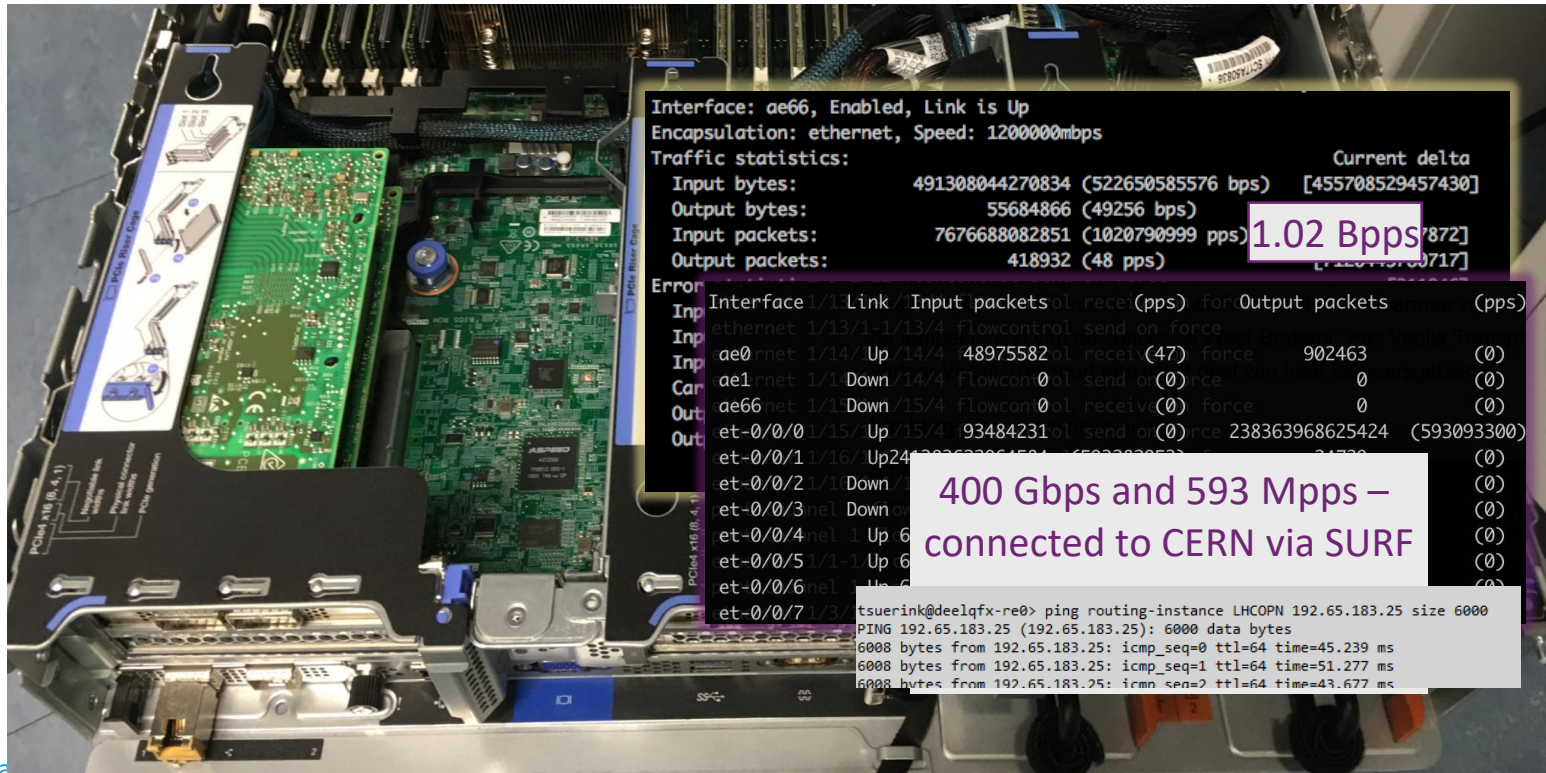
# Exercising the network – sensor data and events



Image: Ballenballknoknk, Tristan Suerink

# Scaling data access: 'system-aware design' at application layer

Reading data 'scattered' in a file - simply using POSIX-like IO - when done over the network severely exposes latency

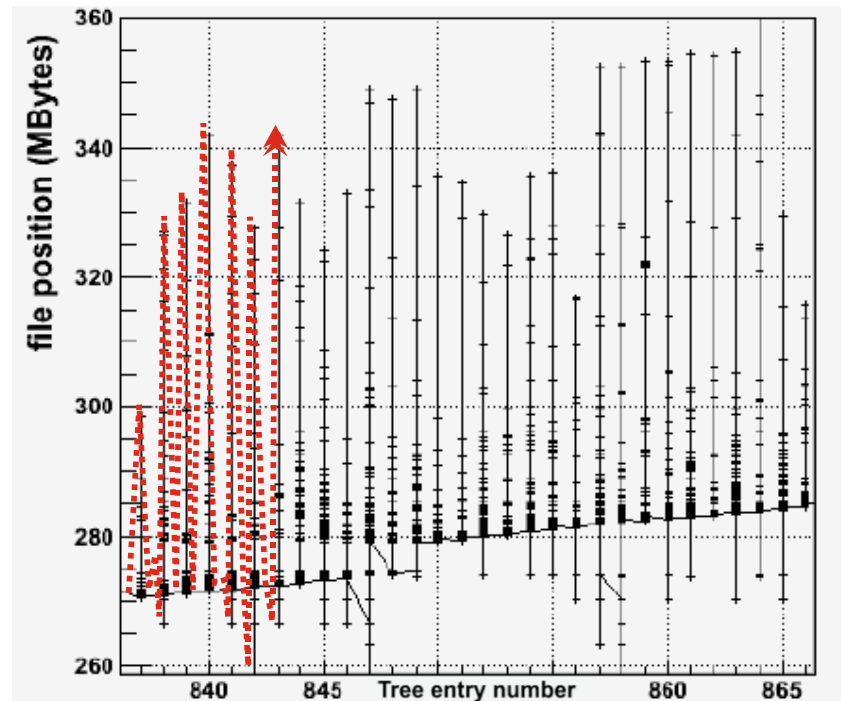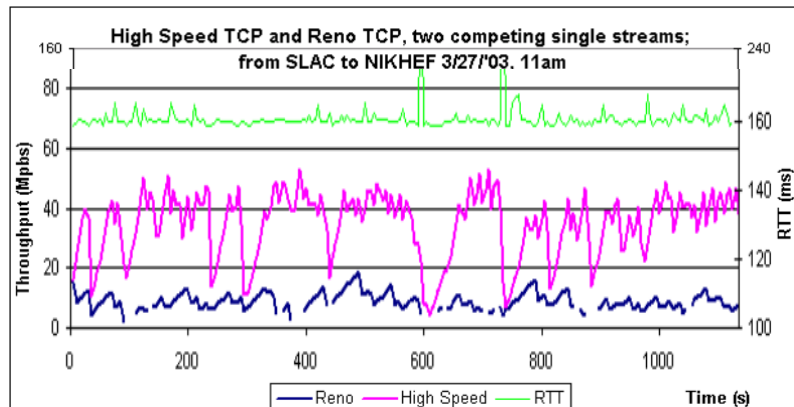*and TCP slow-start makes that even worse*



Image of TCP slow-start and packet loss impact (in Mpps): Antony Antony et al., Nikhef, for DataTAG, 2003(!)
Right: base graphic: Philippe Canal "Root I/O: the fast and the furious", CHEP2010 Access pattern reflects Root versions < 5.28, before Ttree caching and 'baskets'

# And sometimes traffic is triggered by researchers scaling up 'accidentally' from a laptop to a cluster without too much thought
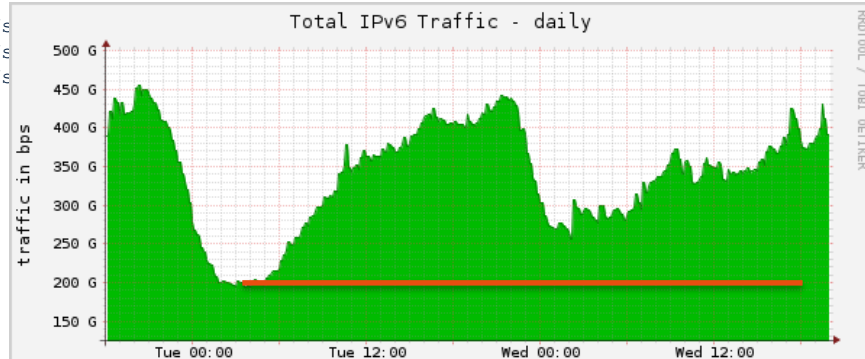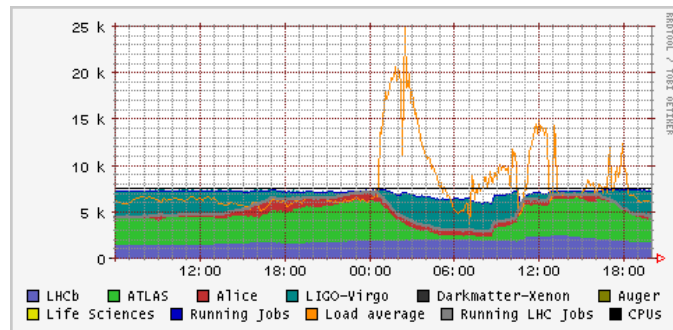
A researcher doing mass creation of containers, rebuilding their python 'virtual env' for each job, running on >> 4000 cores

```
[root@wn-pep-002 ~]# top
top - 09:40:47 up 71 days, 12:17,  2 users,  load average: 110.38, 101.43, 106.3
Tasks: 700 total,   7 running, 666 sleeping,   0 stopped,  27 zombie
%Cpu(s): 17.0 us,  2.0 sy,  0.0 ni, 81.0 id,  0.0 wa,  0.0 hi,  0.0 si,  0.0 st
KiB Mem : 39462902+total, 23514457+free, 10406320 used, 14907812+buff/cache
KiB Swap: 67108860 total, 66841340 free,   267520 used. 37964784+avail Mem

  PID USER      PR  NI    VIRT    RES    SHR S  %CPU %MEM     TIME+ COMMAND
82661 ligo000   20   0 5618756 396356    924 R 360.0  0.1   5:14.43 mksquashfs
72615 ligo000   20   0 5626336 248516    816 R  90.0  0.1   5:44.11 mksquashfs
83257 ligo000   20   0 5611608 219300    852 S  90.0  0.1   1:17.66 mksquashfs
...
```





Pulling the python packages at line rate and downloading public python repositories ultimately *will* trigger Cloudflare and flood SURFnet

June 28th, 2023, data from Nikhef NDPF stats & cricket (top),
SURFnet asd001b-jnx-01 to asd001b-jnx-04 (left),
AMS-IX SFlow https://stats.ams-ix.net/sflow/index.html (bottom)

# For example for HL-LHC, or SKA, more is needed > 2028 …

- 'Typical' network is now mixed 400G-100G
- Push experiments to 800Gbps in metro area, and a local (AMS) loop has been demonstrated
- next: 800 → 1600G AMS-GVA ☺



**BTG**

Home | BTG | BTG Services | INTUG | Innovatielab | Activiteiten | Lobby & Opinie | Publicaties

**Minister Adriaansens opent testomgeving voor volgende generatie netwerktechnologieën**

in Amsterdam is door minister Micky Adriaansens van Economische Zaken en Klimaat ...erotonde is een testomgeving waar SURF en Nikhef gaan experimenteren met nieuwe ...ng beschikt over een internetsnelheid van 800 Gbit/s, wat meer dan 1000 keer sneller ...n gemiddeld huishouden in Nederland. De innovatierotonde stelt Nederlandse ...e doen naar de volgende generatie netwerktechnologieën.

...en onderzoek naar bandbreedte op het internet groeit. Onderzoekers willen steeds meer ...over de landsgrenzen heen met elkaar delen. De bandbreedte van het netwerk speelt ...ote hoeveelheden data snel te kunnen verwerken, is de verwachting dat 800Gbit/s ... De innovatierotonde maakt het mogelijk om te experimenteren met nieuwe

798.49 Gb/s

Web screenshot: btg.org,
Images Nokia 7750-SR1x in Nikhef AMS H234b: Tristan Suerink

# Research data traffic looks like ... a DDoS to others ☺



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

# Access, Trust & Identity

More than one user, *from*
more than one organizational domain, *in*
more than one country

Large-scale IT: worldwide LHC Computing and beyond (2024 ed)

# WLCG: when we met a global trust scaling issue



- 170 sites
- ~50 countries & regions
- ~20000 users

so … just *how* many interactions ??



people photo: a small part of the CMS collaboration in 2017, Credit: CMS-PHO-PUBLIC-2017-004-3; site map: WLCG sites from Maarten Litmaath (CERN) 2021

# Scaling issues – credentials at each site does not work



state of Grid and the LHC computing in 2000

# Authentication – proving *who* are you

Authenticating to a *single service* is relatively simple

- per-service identity (username) and secrets (e.g. password or TOTP token)
- server-side: list of valid users and (hashed and hopefully salted) secrets

```
[root@kwark ~]# cat /etc/passwd
root:x:0:0:root:/root:/bin/bash
bin:x:1:1:bin:/bin:/sbin/nologin
daemon:x:2:2:daemon:/sbin:/sbin/nologin
adm:x:3:4:adm:/var/adm:/sbin/nologin
lp:x:4:7:lp:/var/spool/lpd:/sbin/nologin
sync:x:5:0:sync:/sbin:/bin/sync
shutdown:x:6:0:shutdown:/sbin:/sbin/shutdown
halt:x:7:0:halt:/sbin:/sbin/halt
```

```
root:$6$s8ciAG5gLuv2bPQS$6EcskgtKvQ.rHbif
davidg:$6$nDYcIez2Uaufbtlg$R1hS/Qjn0qYQZk
marianne:$6$p3CeevG6jfNDqZjl$HKHqUTnt2fEqQfkA/m5J3oAOA0zSvgLCKOSQhPS
```

# Authorization – what you are allowed to do

soon needs specifying **access rights** to resources, based on an access **policy**

- might be implicit or ad-hoc

- be in formal policy language
  like XACML (*example: Argus PDP)*

- or be service-specific
  *example: Linux sssd config*

```
resource "http://cern.ch/authz/ce1" {
    action "http://cern.ch/authz/actions/ce-submit" {
        rule permit {
            vo="atlas"
            pilot-job="true"
        }
        rule deny {
            pilot-job="true"
        }
    }
}
```

*simplified Argus policy language – can map directly to XACML*

```
ldap_access_order = filter,authorized_service
ldap_access_filter = (|(memberOf=cn=gridSrvAdministrators,ou=DirectoryGroups,dc=farmnet,
dc=nikhef,dc=nl)(memberOf=cn=gridMWSecurityGroup,ou=DirectoryGroups,dc=farmnet,dc=nikhef
,dc=nl)(memberOf=cn=nDPFPrivilegedUsers,ou=DirectoryGroups,dc=farmnet,dc=nikhef,dc=nl))
```

Policy example: Argus system, https://argus-documentation.readthedocs.io/en/stable/misc/examples.html; service-specific: sssd.conf ldap auth_provider

# Authorization and access control

Access control is ultimately enforced by the service provider
*(unless data-level encryption is used, where the data owner retains some control)*



policy overlap diagram by Olle Mulmo, KTH for EGEE-I JRA3, policy pie: OpenGrod Forum OGSA working group and Globus Alliance

# Authorization policy subjects

AuthZ policies need subject attributes ('claims')

- **bound to an verifiable identity** statement
  - e.g. visa are strongly linked to a specific entity, and asserted by a trusted party (by the service)
- be a **bearer token**
  - scoped to a relying party, a service, or an action
- **self-asserted**
  - quite useless unless backed by verifiable evidence, like in self-sovereign identity schemes

Transport mechanisms (see also RFC2903)

- pushed alongside the service access,
- pulled from the source as needed, or
- pushed by the attribute source as an agent

USA visa image source: https://2009-2017.state.gov/m/ds/rls/rpt/79785.htm ; RATP bearer token, issued for the Paris public transport system

# Access control in a single domain

- Dedicated to each service
  where you need access

- Usually strongly linked to authorization:
  at times even
  different accounts for different roles

- In a multi-organizational system becomes

$$\mathcal{O}(n_{sites} * n_{services}) * \mathcal{O}(n_{users})$$



Image: AARC NA2 training module "Authentication and Authorisation 101" - https://aarc-community.org/training/aai-101/

# Authentication and Authorization Infrastructure



Image: AARC NA2 training module "Authentication and Authorisation 101" - https://aarc-community.org/training/aai-101/

# Federation

portability of identity information across otherwise autonomous administrative domains



Shibboleth IdP image and SAML2 auth flow by SWITCH (CH) – see also https://refeds.org/ on federation structure and (assurance and security) guidelines

# One simple federation you know: eduroam

*service-specific* trust
between organisations
globally

hierarchical RADIUS servers based
an 802.1x secure exchange
over TLS or EAP-TTLS
tunneling your credentials
back to your home institution

RADIUS server then instructs WiFi access point/controller



eduroam: Klaas Wieringa et al., image from https://eduroam.org/how/, GEANT ; RADIUS: RC2865 https://www.rfc-editor.org/rfc/rfc2865; see also freeradius.org

# Multipurpose federation with SAML: SURFconext & eduGAIN



Images: SURFconext IdP dashboard by SURF, showing some services tagged with REFEDS R&S; eduGAIN map: GEANT, https://technical.edugain.org/status

# Your favourite federated service?



https://surfspot.nl/

# SAML federation

| Attributes | Values |
|---|---|
| E-mail | davidg@nikhef.nl |
| Affiliation | • employee<br>• member<br>• faculty |
| Targeted ID | https://sso.nikhef.nl/sso/saml2/idp/metadata.php!https://attribute-viewer.aai.switch.ch/shibboleth!b9f858169ea28dc68b6753baa10 84d8c039e36a7 |
| Common Name | David Groep |
| Display Name | David Groep |
| Principal Name | davidg@nikhef.nl |
| Home organization (international) | nikhef.nl |
| Home organization type (international) | urn:mace:terena.org:schac:homeOrganizationType:int:other |



SAML2.0 auth flow

Try at https://attribute-viewer.nikhef.nl/ and select "Login via a global authentication SAML source"
Firefox: use F12, and SAML tracer https://addons.mozilla.org/nl/firefox/addon/saml-tracer/ (by Tim van Dijen of SimpleSAMLphp fame)

SAML WebSSO flow image: SWITCH, CH

# Under the hood, sends a (signed) XML document

```
<saml:Subject>
    <saml:NameID Format="urn:oasis:names:tc:SAML:2.0:nameid-format:persistent">xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx</saml:NameID>
    <saml:SubjectConfirmation Method="urn:oasis:names:tc:SAML:2.0:cm:bearer">
      <saml:SubjectConfirmationData NotOnOrAfter="2022-10-21T18:16:40Z"
        Recipient="https://attribute-viewer.aai.switch.ch/Shibboleth.sso/SAML2/POST"
        InResponseTo="_64c10a60c382bdaeb328653d9d25951c" /></saml:SubjectConfirmation>
</saml:Subject>
<saml:Conditions NotBefore="2022-10-21T18:11:39Z"
                 NotOnOrAfter="2022-10-21T18:16:40Z">
    <saml:AudienceRestriction>
      <saml:Audience>https://attribute-viewer.aai.switch.ch/shibboleth</saml:Audience>
    </saml:AudienceRestriction>
</saml:Conditions>
<saml:AuthnStatement AuthnInstant="2022-10-21T17:33:29
                     SessionNotOnOrAfter="2022-10-22T0
                     SessionIndex="_90f745f18f712b6a56
  <saml:AuthnContext>
      <saml:AuthnContextClassRef>urn:oasis:names:tc:SAM
      <saml:AuthenticatingAuthority>https://sso.nikhef.
  </saml:AuthnContext>
</saml:AuthnStatement>
```

```
<saml:AttributeStatement>
    <saml:Attribute Name="urn:mace:dir:attribute-def:cn"
                    NameFormat="urn:oasis:names:tc:SAML:2.0:attrname-format:uri">
      <saml:AttributeValue xsi:type="xs:string">David Groep</saml:AttributeValue>
    </saml:Attribute>
    <saml:Attribute Name="urn:oid:2.5.4.3"
                    NameFormat="urn:oasis:names:tc:SAML:2.0:attrname-format:uri">
      <saml:AttributeValue xsi:type="xs:string">David Groep</saml:AttributeValue>
    </saml:Attribute>
    <saml:Attribute Name="urn:mace:dir:attribute-def:eduPersonAffiliation"
                    NameFormat="urn:oasis:names:tc:SAML:2.0:attrname-format:uri">
      <saml:AttributeValue xsi:type="xs:string">employee</saml:AttributeValue>
      <saml:AttributeValue xsi:type="xs:string">member</saml:AttributeValue>
      <saml:AttributeValue xsi:type="xs:string">faculty</saml:AttributeValue>
    </saml:Attribute>
    <saml:Attribute Name="urn:oid:1.3.6.1.4.1.5923.1.1.1.1"
    ...
```

# Different tech, similar concept: X.509 RFC5280 client certificates

```
Version: 3 (0x2)
Serial Number:
    34:f3:e3:5f:c0:53:0b:a6:ef:2b:4a:79:01:b5:50:3b
Signature Algorithm: sha384WithRSAEncryption
Issuer: C = NL, O = GEANT Vereniging, CN = GEANT eScience Personal CA 4
Validity
    Not Before: Apr  2 00:00:00 2022 GMT
    Not After : May  2 23:59:59 2023 GMT
Subject: DC = org, DC = terena, DC = tcs, C = NL, O = Nikhef, CN = David Groep davidg@nikhef.nl
Subject Public Key Info:
    Public Key Algorithm: rsaEncryption
        RSA Public-Key: (4096 bit)
        Modulus:
            00:f0:0d:c0:ff:ee:f0:0d:f0:0d:c0:ff:ee:f0:0d:
            ...
            ff:50:6d
        Exponent: 65537 (0x10001)
X509v3 extensions:
    X509v3 Key Usage: critical
        Digital Signature, Key Encipherment
    X509v3 Basic Constraints: critical
        CA:FALSE
    X509v3 Extended Key Usage:
        E-mail Protection, TLS Web Client Authentication
    X509v3 Certificate Policies:
        Policy: 1.2.840.113612.5.2.2.5
```

You should be able to get an 'IGTF-DOGWOOD' assurance certificate from RCauth.eu.
Go to https://rcdemo.nikhef.nl/ and select the 'Basic demo' and use 'run non-VOMS' to get and view your short-lived certificate

are back-channel interactions

| run non-VOMS demo |

# Certificates chains & constraint proxy identity delegation

- PKIX certificates are ASN.1 structures in a distinguished binary encoding (DER format)
- contains the tuple (issuer, subject, serial) + validity period + key material + **extensions**
- within it is the message digest (hash), signed with private key of the issuer
- Verifiable using the issuer's public key



RFC3820 'proxy' certificates extend this concept to (policy-constraint) identity delegation

To get an RFC3820 proxy certificate using your own federated identity, use RCauth.eu – see https://rcdemo.nikhef.nl/ and use the "Basic Demo" option

# OpenID Connect and OAuth2

- Quite .well-known
  (used by lots modern 'non-enterprise' SSO)

- shows in its initial design: one source of
  identity (Openid Provider, 'OP'), and many
  services (Relaying Parties, 'RP')



| Show OpenID Connect Client | |
|---|---|
| Name | hekel.nikhef.nl |
| Description | Hekel using mod_auth_openidc |
| Client id. | _f6bfe81892e680e4ecfc3b41ecf1a15d141c0d106b |
| Client secret | _ |
| Auth. source | saml2 |
| Redirect URI | https://hekel.nikhef.nl/rp/redirect_uri |
| Scopes | openid
profile
email
assurance |

Shown is the 'implicit flow', other flows possible. Image source: AARC NA2 training on AAI 101
See https://openid.net/ for protocols and standardization work

# Federation: different technologies, same idea



**SAML - Security Assertion Markup Language and WebSSO ('SAML2Int')**
- XML-formatted 'attribute statements' over web transport (usually POST)
- SAML-Metadata: list of entities with description of bindings with entityAttributes

**PKI - Public Key Infrastructures**
- trusted third party (a *certification authority* a.k.a. *CA*)
  signs X.509 formatted certificates with name, issuer, serial number, and extensions
- CAs can sign end-entities as well as other CAs (hierarchically or by cross-signing)
- *bridge CAs* render a technical implementation of a shared policy (assurance)
- *policy-bridges* don't sign anything, but curate *distribution*
  (like browsers and operating systems based on CA/BF requirements, IGTF for research infras)

**OpenID Federation – Federating OpenID Connect parties**
- federate end-points for OIDC Providers and Relying Parties (or OAuth2), with similar models

*note federation based on 'ultimate trust' domains (e.g. cross-realm Kerberos) also exists …*

# Federation: technological …
# … or policy bridge

trust remains with the relying party
can be *bridged* by either cross-signing (left)
or by policy agreements (right)



*Policy-bridge trust federation: EGI.eu infrastructure leveraging the IGTF federation*



Left-hand image: 4 Bridges Forum, source: Scott Rea (then: Dartmouth University)
Images: cabforum.org, WebTrust logo: from DigiCert.com; image MS root store, https://learn.microsoft.com/en-us/security/trusted-root/program-requirements

# 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

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

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

# Federation: technology, interoperability, policy



Image from SWITCH (CH) and edugain.org

# Managing complexities of federation & identity



WebFTS prototype
'FIM4R' in wLCG
Romain Wartel et al.

ELIXIR reference
architecture 2016
Mikael Linden et al.

communities had either invented
their own 'proxy' model to abstract complexity

or they were composed of many services
each of which had to manage federation complexity

Community images: Romain Wartel, CERN; Mikael Linden, CSC; Lukas Hammerle, SWITCH

# Multiple sources of authority: the community

- authorization assertion providers (attribute authorities) use the identifier(s) from authentication in their membership services

- *source of authority* for attributes is distributed

for example:
- community membership from an experiment
- affiliation status from home organisation

*may be jointly needed to access sensitive data*
*that is subject to medical-ethical clearance*

# 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

# Composite AAIs: proxies beyond just the research infrastructures

Proxy model harmonizes IdPs from many sources

- **eduID**-style identifiers
  - 'life-long learning' identifiers
  - independent student identifier (the ESI) for mobility & Erasmus-without-papers
  - eduGAIN-alignment, but also a 'provider of last resort'

- **eIDAS** and government eID (e.g. DigID)
  - identity assurance step-up

- **ORCID** provides identifier portability through linking
  - provides name linking and persistent attribution
  - since it persists, also very useful to allow access *independent of home organisation* throughout a carrer

Composite AAI image source: Christos Kanellopoulos (GEANT), Marcus Hardt (KIT)

# When many proxies from different groups come together

When collaborations cross different domains (or an industry sector with lots of mergers and spin-offs)

- proxies with each group
- inter-federate SP/IdP interfaces
- each federation can add own policy and entity filtering

Example
European Open Science Cloud (EOSC)
AAI based on federations and proxies



Christos Kanellopoulos (GEANT) for the EOSC AAI Federation in "The EOSC Core", https://eoscfuture.eu/wp-content/uploads/2022/04/EOSC-Core.pdf

# European Open Science Cloud (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)

# Same blocks underlie e.g. the Fenix and Puhuri HPC ecosystem



Fenix image via Christos Kanellopoulos, diagram via Anders Sjöström (NeIC, Puhuri) at the TNC23 workshop

# Also the basic blocks for your student identity& Erasmus+



## MyAID Architecture

- Provides an Authentication Proxy for the core Erasmus+ services (Online Learning Agreement, Dashboard, PhD Hub and the Erasmus+ App).

- Supports authentication via eduGAIN, eIDAS and Google

Christis Kanellopoulos (GÉANT) for the Erasmus+/Erasmus Without Papers programme

# What value does our university ID bring in a life-long learning environment? Time to think less institution-centric?

## EBSI Wave 2 (15 MS, 20 HEIs, 2 EUA)

### Study

**01** A student gets a diploma with a list of course units validated from Erasmus (Transcript of Records Credential) (ES/BE/IT)

**02** A student applies for a PhD with a Bachelor / Master degree from a foreign country (Bachelor/Master Diploma Credential) (RO/GR/FR)

**03** A student gets access to local discounts using student credential (European Student IDentity) (BE/ES)

**04** A refugee presents an EQPR to a European Italian University to apply for a Master (EQPR - CoE Refugee Passport) (IT/DE)

### Work

**05** A graduated citizen applies for a job with a Degree from a foreign country (License to Practice Credential) (GR/CY)

### Grow

**06** A PhD student applies for specific courses in a foreign country (Cross-border Micro-credentials) (FI/LT)

GEANT Association

Stichting Internet Domeinregistratie Nederland

SURF BV

Vezcozo BV

Dienst Uitvoering Onderwijs – Dutch Education Ministry

Images from Lluís Ariño, for the DC4EU project. See e.g. https://www.dc4eu.eu/consortium/#netherlands

# Putting it back together again

Common patterns in scalability

# A global infrastructure of EGI, OSG and WLCG, …



'an infrastructure with components matched to application need'
- systems architecture: compute (HTC clusters), networking, storage, and application structure
- in a balanced and {energy,cost}-efficient setup

BerkeleyDB Information System for EGI, from top-level BDII at ldap://bdii03.nikhef.nl:2170/o=grid; Earth visualization: https://dashb-earth.cern.ch/, Google Earth

**Maastricht University**  | DACS

# European Open Science Cloud (EOSC) ecosystem example



and many more systems and 'data spaces' besides EOSC: *e.g.* Copernicus EO data, GAIA-X, sectoral spaces, …

EOSC: https://eoscfuture.eu/wp-content/uploads/2022/04/EOSC-Core.pdf; data spaces image: https://digital-strategy.ec.europa.eu/en/library/building-data-economy-brochure

# Looking for a common pattern?

- It's all about *balanced* systems
  - systems are like congested highways: no use solving just *one* bottleneck
  - and the bottlenecks may be inside the system as well as in interconnects

- Horizontal scaling, and be as stateless as possible
  - although persistent storage obviously has to retain some state ☺
  - edge scales horizontally, and scaling from 2+ is much easier than from 1→ 2

- You can move problems around, but it's hard to actually *solve* them
  - e.g. lack of a single common interface implies one needs adaptors and plugins

- Scaling *collaboration and trust* federation is as complex as scaling systems
  - composing services across administrative domains is ubiquitous
  - but beyond a certain size, $\mathcal{O}(100)$, you will find need for some policy and review

**… since some things are fun, but not quite *that* scalable …**

Liquid $CO_2$ cooling test bench,
24.33% overclocked
using CineBench R20
best sustained, i.e. without LN2…
In a Nikhef-AMD collaboration

| | SCORE | USER | | FREQUENCY | HARDWARE | COOLING | HW | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| 1. | 23323 pts | | Splave | 5400.2 MHz | AMD Ryzen Threadripper 3970X | LN2 | 0pts | | 0 | |
| 2. | 23081 pts | | Alex@ro | 5375 MHz | AMD Ryzen Threadripper 3970X | LN2 | 0pts | | 1 | |
| 3. | 22064 pts | | Hiwa | 5050.6 MHz | AMD Ryzen Threadripper 3970X | LN2 | 0pts | | 0 | |
| 4. | 21601 pts | | keeph8n | 5000.4 MHz | AMD Ryzen Threadripper 3970X | LN2 | 0pts | | 0 | |
| 5. | 20022 pts | | Nikhef | 4600.1 MHz | AMD Ryzen Threadripper 3970X | SS | 0pts | | 0 | |

T Suerink, K de Roo: https://hwbot.org/submission/4539341_nikhef_cinebench___r20_with_benchmate_ryzen_threadripper_3970x_20022_pts

The 9 kinds of physics seminar

- The "Typical"
- The "Ideal"
- The "Unprepared Theorist"
- The "Unprepared Experimentalist"
- The "Well-meaning Undergrad"
- The "Guest From Another Department"
- The "Nobel Prize Winner"
- "Poetry"
- The "Politician"

http://manyworldstheory.com/2013/10/03/the-9-kinds-of-physics-seminar/

# More Q&A time!

David Groep, davidg@nikhef.nl
*https://www.nikhef.nl/~davidg/presentations/*
*https://orcid.org/0000-0003-1026-6606*

David Groep, davidg@nikhef.nl
*https://orcid.org/0000-0003-1026-6606*

Nik|hef

Maastricht University | Department of Advanced Computing Sciences

# Distributed collaborative services
## *a more technical example with RCauth.eu*

Credential translation in the AARC BPA
    … building RCauth.eu
Leveraging federation and collaboration
    for ubiquitous research credentials

# *Bridges and Token Translation Services*
# TCS - for users that manage to grasp the idea



**TCS is a SAML Service Provider** (today by Sectigo)
to eduGAIN: where eligible authenticated users obtain
client certificates for access to many research services
**A globally recognized identity for all employees & students** (they are automatically eligible!).

GEANT Trusted Certificate Service - https://ca.dutchgrid.nl/tcs/,
https://cert-manager.com/customer/surfnet/idp/clientgeant, https://www.geant.org/Services/Trust_identity_and_security/Pages/TCS.aspx

# Seamless in-line token translation services from 'SAML' to PKIX



user facing | hidden back-end

**Community Science Portal**

**IGTF accredited PKIX Authority**

**Infrastructure Master Portal Credential Store**

**User Home Org**
*or Infrastructure IdP*

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

REFEDS R&S
Sirtfi Trust

**Policy Filtering WAYF to eduGAIN**

Built on CILogon and MyProxy, see www.cilogon.org   CILogon   125

# Unique certificated from FIM via eduPerson and REFEDS R&S

Sources of naming and uniqueness, that work *today*

- **eduPersonPrincipalName** – scoped point-in-time unique identifier, which could be, but usually is not, privacy preserving: "davidg@nikhef.nl", "P70081609@maastrichtuniversity.nl"
- **eduPersonTargetedID** – scoped transient non-reassigned identifier, like urn:geant:nikhef.nl:nikidm:idp:sso!*27c8d63ed42c84af2875e2984*
- **subject-id** - a scoped persistent non-reassigned identifier, which should be privacy-preserving: 44f7751265a6e8b228f9@nikhef.nl

Plus the (domain-name based) schacHomeOrganisation and a '**representation of the real name**'

**/DC=eu/DC=rcauth/DC=rcauth-clients/O=*orgdisplayname*/CN=*commonName +uniqeness***

uniqueness will added to commonName via hashing of *ePPN, ePTID, subject-id*, so that an enquiry via the issuer allows unique identification of the vetted entity"

# The 'back side' of a typical RCauth portal data flow



**Parsed ID Token:**

```
stdClass Object
(
    [typ] => JWT
    [kid] => E01796EA0367564935B0981731B9B116
    [alg] => RS256
)
stdClass Object
(
    [sub] => P70081609@unimaas.nl
    [idp] => http://login.maastrichtuniversity.nl/adfs/services/trust
    [eduPersonTargetedID] => http://login.maastrichtuniversity.nl/adfs
    [idp_display_name] => Maastricht University
    [cert_subject_dn] => CN=Groep\, David (DACS) KWwWAnhI4psmiGTw 1,O=
    [name] => Groep, David (DACS)
    [eduPersonPrincipalName] => P70081609@unimaas.nl
    [given_name] => David
    [family_name] => Groep
    [email] => david.groep@maastrichtuniversity.nl
    [iss] => https://aai.egi.eu/mp-oa2-server
```

**Proxy information:**

```
subject    : /DC=eu/DC=rcauth/DC=rcauth-clients/O=maastrichtuniversity.nl/CN=Groep, David (DACS) KWwWAnhI4psmiGTw 1/CN=208760481/CN=466908503
issuer     : /DC=eu/DC=rcauth/DC=rcauth-clients/O=maastrichtuniversity.nl/CN=Groep, David (DACS) KWwWAnhI4psmiGTw 1/CN=208760481
identity   : /DC=eu/DC=rcauth/DC=rcauth-clients/O=maastrichtuniversity.nl/CN=Groep, David (DACS) KWwWAnhI4psmiGTw 1/CN=208760481
type       : RFC compliant proxy
strength   : 2048 bits
path       : /tmp/x509up_uiI4TkF
```

**Maastricht University**

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

Maastricht University

# 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'**

Maastricht University

# 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

Maastricht Univ

# 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



**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

HA proxy

HA proxy

VPC or VPN interconnect

HA proxy

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

selected imagery: Mischa Sallé, Jens Jensen, Nicolas Liampotis

**Maastricht University** | DACS

# Anycast: when the same place exists many times



**So we used**

- 3 (for now: 2) sites
- one VM at each site exposing 2a07:8504:01a0::1
- smallest v6 subnet (/48)
- bird + a service probe
- each site's own ASN
- some IRR DB editing
- IPv4 is similar, with a /24

*and some monitoring*

routing image: SIDNlabs - https://www.sidnlabs.nl/en/news-and-blogs/the-bgp-tuner-intuitive-management-applied-to-dns-anycast-infrastructure

# 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



| < 10 ms: 29 | < 20 ms: 46 | < 30 ms: 59 | < 40 ms: 54 | < 50 ms: 64 | < 100 ms: 113 | < 200 ms: 91 | < 300 ms: 26 | > 300 ms: 5 | No Data: 0 |

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

# Shortest path, also when mixing with the default-free zone

```
[root@kwark ~]# traceroute -IA 145.116.216.1
traceroute to 145.116.216.1 (145.116.216.1), 30 hops max, 60 byte packets
  1  cmbr.connected.by.freedominter.net
       (185.93.175.234) [AS206238]
  2  connected.by.freedom.nl
       (185.93.175.240) [AS206238]
  3  et-0-0-0-1002.core1.fi001.nl.freedomnet.nl
       (185.93.175.208) [AS206238]
  4  as1104.frys-ix.net (185.1.203.66) [*]
  5  parkwachter.nikhef.nl
       (192.16.186.141) [AS1104]
  6  gw-anyc-01.rcauth.eu
       (145.116.216.1) [AS786/AS5408/AS1104]
```

*rcauth.eu HA proxy*

Route from home to RCauth.eu, from my home Freedom Internet ISP

*me, at home*

# RSA Crypto

Just in case … you cannot factor '55'

# Establishing trust at a distance

Remote trust needs cryptography in some way

**Client authentication**
- pre-shared secrets, may be salted hashed on service side
- required: secure one-way hash function
- need a **protected channel** between identifiable end-points

**Mutual authentication**
- eithers need a lot of shared keys, a trusted third party (TTP), or mesh validation (WoT)
- with the TTP and multiple services comes the need for crypto
- across administrative domains, *key distribution* is the larger challenge

The cryptography used can be either *symmetric* or *asymmetric*, 'public key'

# Asymmetric crypto: RSA interlude needed?



$(d,n)$

$(d,e,p,q)$

$(e,n)$

$n = pq$

Alice

$(e,n)$

$c$

$D_{d,n}(c) \rightarrow m$

$c = E_{e,n}(m)$

$m$

$E_{e,n}(m) = m^e \bmod(n)$
$D_{d,n}(c) = c^d \bmod(n)$
$m = D(E(m)) = E(D(m))$    (*reversibility*)

if and only if   $de = 1 \bmod(\phi(p,q))$
where       $\phi(p,q) = (p\text{-}1)(q\text{-}1)$
and         $(p\text{-}1)$ prime relative to $e$

Bob

Rivest, Shamir and Adleman, Communications of the ACM 21 (2), 120-126

# 6-bit RSA (note: this might be broken quickly …)

- Take a (small) value $e$ = **3**
- Generate a set of primes ($p,q$), each with a length of $k$/2 bits, with ($p$-1) prime relative to $e$.
  ($p,q$) = **(11,5)**
- $\phi(p,q)$ = (11-1)(5-1) = **40**; $n=pq=$**55**
- find $d$, in this case **27** [3*27 = 81 = 1 mod(40)]

- Public Key: **(3,55)**
- Private Key: **(27,55)**

$E_{e,n}(m) = m^e \bmod(n)$
$D_{d,n}(c) = c^d \bmod(n)$
$m = D(E(m)) = E(D(m))$
if a.o. if $\quad de = 1 \bmod(\phi(p,q))$
where $\quad \phi(p,q) = (p\text{-}1)(q\text{-}1)$

# Message exchange

Encryption:

- Bob thinks of a plaintext $m(<n)$ = **18**
- Encrypt with Alice's public key **(3,55)**
- $c$=E$_{3;55}$(18)=$18^3$ mod(55) = 5832 mod(55) = **2**
- send message **"2"**

Decryption:

- Alice gets **"2"**
- she knows private key **(27,55)**
- E$_{27;55}$(2) = $2^{27}$ mod(55) = **18** !

$E_{e,n}(m) = m^e \bmod(n)$
$D_{d,n}(c) = c^d \bmod(n)$
$m = D(E(m)) = E(D(m))$
if a.o. if $\quad de = 1 \bmod(\phi(p,q))$
where $\quad \phi(p,q) = (p\text{-}1)(q\text{-}1)$

**(3,55)**

## If you just have (3,55), it's hard to get the 27…

*but also: the maximum plaintext is limited by the modulus length*

# The most used asymmetric crypto application

Asymmetric crypto underpins
the transport layer security
of all of the web today

- ASN.1 syntax data with
  X.509 (RFC5280) structure
- mostly RSA or Elliptic Curves (EC)
- used to negotiate a
  (symmetric) bulk cipher (typically AES)

then used to protect channel to usually
*unauthenticated* client application (browser)

# Other ancillary materials

these generic slides do not form part of
the module, but are just general
background knowledge and example

# Open Systems Interconnection model (OSI model)

| Layer | | | Function |
|---|---|---|---|
| Host layers | 7 | Application | High-level protocols (resource sharing, remote file access) |
| | 6 | Presentation | Translation of data between a networking service and an application |
| | 5 | Session | Managing communication sessions, i.e., continuous exchange of information in the form of multiple back-and-forth transmissions between two nodes |
| | 4 | Transport | Reliable transmission of data segments between points on a network |
| Media layers | 3 | Network | Addressing, routing and traffic control |
| | 2 | Data link | Transmission of data frames between two nodes connected by a physical layer |
| | 1 | Physical | Transmission and reception of raw bit streams over a physical medium |

OSI X.200 layering model, ITU-T (CCITT), https://www.itu.int/rec/T-REC-X.200; image adapted from https://en.wikipedia.org/wiki/OSI_model

# OSI vs Internet Protocol Architecture model

# Private (direct) peerings to distribute traffic load



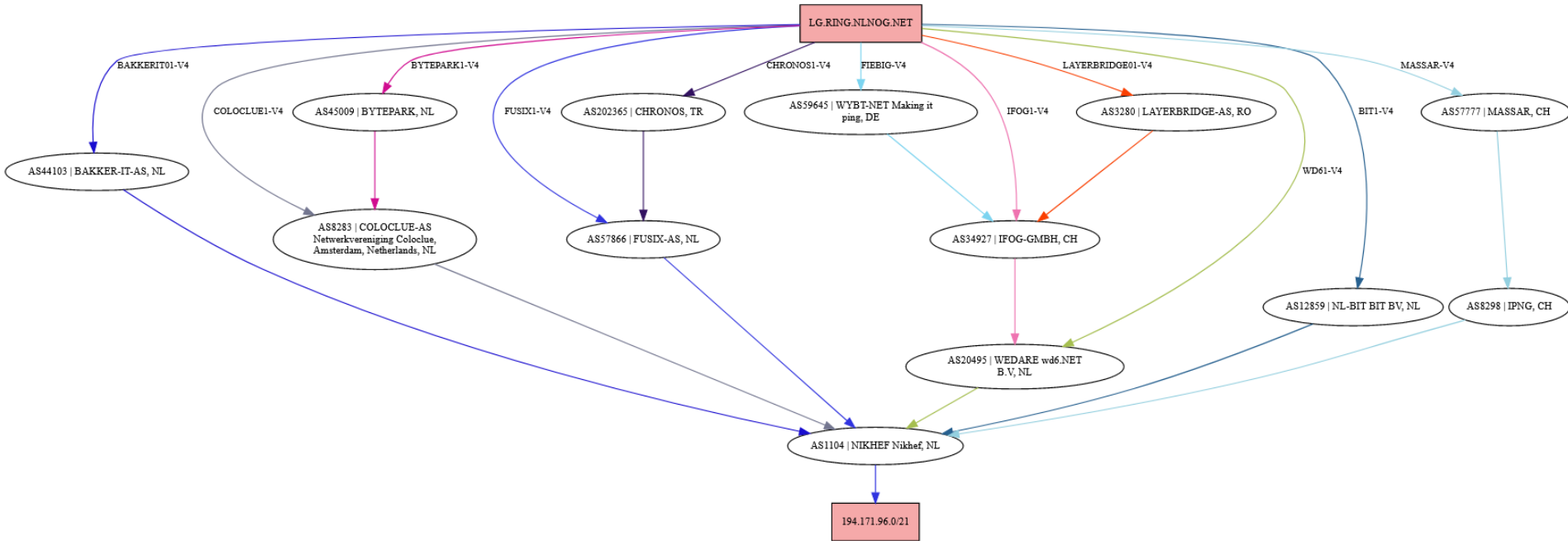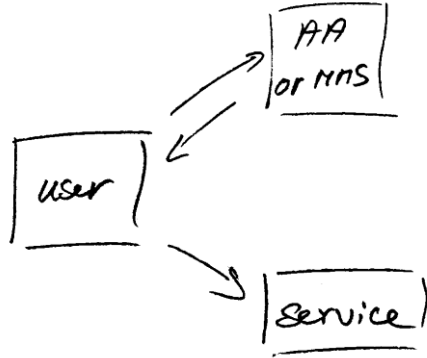Image sources: NLNOG RING map https://lg.ring.nlnog.net/

# RFC2904 authorization models: three AuthZ flows



'push'                    'pull'                    'agent'

Authorization models: AAA Authorization Framework, RFC2904, Vollbrecht et al.

# OAuth2 & JWTs: assertions can be quite detailed

```
$ echo $AT | jwt
...
❉ Payload
{
  "wlcg.ver": "1.0",
  "sub": "a1b98335-9649-4fb0-961d-5a49ce108d49",
  "aud": "https://wlcg.cern.ch/jwt/v1/any",
  "nbf": 1593004542,
  "scope": "storage.read:/ storage.modify:/",
  "iss": "https://wlcg.cloud.cnaf.infn.it/",
  "exp": 1593008142,
  "iat": 1593004542,
  "jti": "da0a2f89-3cbf-42a7-9403-0b43d814551d",
  "client_id": "edfacfb1-f59d-44d0-9eb6-a745ac52f462"
}
```

OAuth2 Access Token following the WLCG AuthZ WG Profile, from: https://wlcg-authz-wg.github.io/wlcg-authz-docs/token-based-authorization/

# Example flow in the European Open Science Cloud



EOSC Portal & Marketplace Amnesia service by the OpenAIRE e-infrastructure, EOSC Helpdesk: Zammad hosted by KIT https://eosc-helpdesk.eosc-portal.eu

# RCauth demonstrator

RCauth is an AARC BPA token translation service that forges X.509 end-user certificates that are managed in a central portal for you (the portal is 'elevator.nikhef.nl')

Qualified users are all those in eduGAIN with basic assurance (Sirtfi version 1 + Research & Scholarship entity categories), and everyone in a Dutch SURF 'Annex IX' institution – such as UM

Your end-entity certificate is globally IGTF trusted under the 'Identifier Only Trust Assurance' (IOTA) profile

# RCauth do-it-yourself demo

- Go to https://rcdemo.nikhef.nl/
- select "Basic Demo"

- Enable browser Inspector (F12) on the network tab, (and start the SAML tracer extension if you have it)



- Run the "non-VOMS demo"

- From eduGAIN, select "Maastricht University"



Maastricht University  | DACS

Large-scale IT: workflows and computing clusters, FNA (2023-24)

# RCauth: SAML to UM, but OIDC for your credential management service

- Approve transfer in OIDC flow & see your PKI X.509 user cert!

- Review the network interactions with
  - engine.surfconext.nl
  - login.maastrichtuniversity.nl
  - pilot-ca1.rcauth.eu
  - elevator.nikhef.nl (this is the credential management service where your long-term private key is)

- What is the difference in the POSTs?
- Can you see the difference in the SAML and the OIDC flow?

**RCauth.eu Online CA consent page**

The Master Portal below is requesting access to your personal
The white-label Research and Collaboration Au

If you approve, please accept, otherwise, cancel.

Details on which attributes are released, why, to whom, and h
For further information on the CA see the RCauth.eu homepag

☑ Remember

[Yes, continue]  [No, cancel]

**Master Portal Information:**

| | |
|---|---|
| *Name:* | Nikhef MasterPortal |
| *Description:* | Nikhef MasterPortal |
| *URL:* | https://www.nikhef.nl/ |

**Information that will be sent to the Master Portal:**

*sub :*  P70081609@unimaas.nl

aBMALkQAo5smbNx+PW7fOoNbfzReSJfGt7DaYqekEO/yvvxH0OO8xf20w+rmPCEA
-----END CERTIFICATE-----

**Proxy information:**

```
subject    : /DC=eu/DC=rcauth/DC=rcauth-clients/O=maastrichtuniversity.nl/CN=Groep, David (DACS) KWwWAnhI4psmiGTw 1/CN=1435128199/CN=318189164
issuer     : /DC=eu/DC=rcauth/DC=rcauth-clients/O=maastrichtuniversity.nl/CN=Groep, David (DACS) KWwWAnhI4psmiGTw 1/CN=1435128199
identity   : /DC=eu/DC=rcauth/DC=rcauth-clients/O=maastrichtuniversity.nl/CN=Groep, David (DACS) KWwWAnhI4psmiGTw 1/CN=1435128199
type       : RFC compliant proxy
strength   : 2048 bits
path       : /tmp/x509up_uLK91hN
timeleft   : 12:00:00
key usage : Digital Signature, Key Encipherment, Data Encipherment
Certificate:
    Data:
        Version: 3 (0x2)
        Serial Number: 318189164 (0x12f72e6c)
    Signature Algorithm: sha256WithRSAEncryption
        Issuer: DC=eu, DC=rcauth, DC=rcauth-clients, O=maastrichtuniversity.nl, CN=Groep, David (DACS) KWwWAnhI4psmiGTw 1, CN=1435128199
        Validity
            Not Before: Nov  2 12:04:47 2024 GMT
            Not After : Nov  3 00:09:47 2024 GMT
        Subject: DC=eu, DC=rcauth, DC=rcauth-clients, O=maastrichtuniversity.nl, CN=Groep, David (DACS) KWwWAnhI4psmiGTw 1, CN=1435128199, CN=3
```

**Maastricht University** | DACS

# Nulla folia post hoc sunt

Thanks for watching!

*"En daarmee, geachte luisteraars, laat ik u over aan
de verpozing die uw babbelklant u gemeenlijk pleegt te bieden."*