Computing Infrastructures for Research & the Worldwide LHC Computing Grid

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



David Groep
DACS & Nikhef
4 November 2025
KEN3239/BCS3239 r1.2



Exploding data? the Large Hadron Collider at CERN

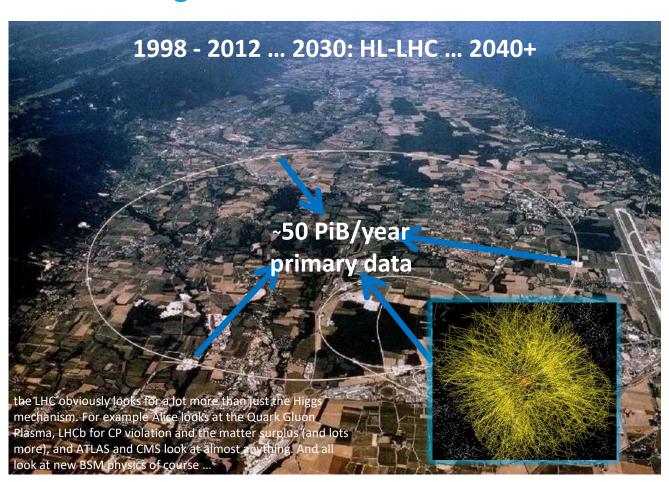
1964

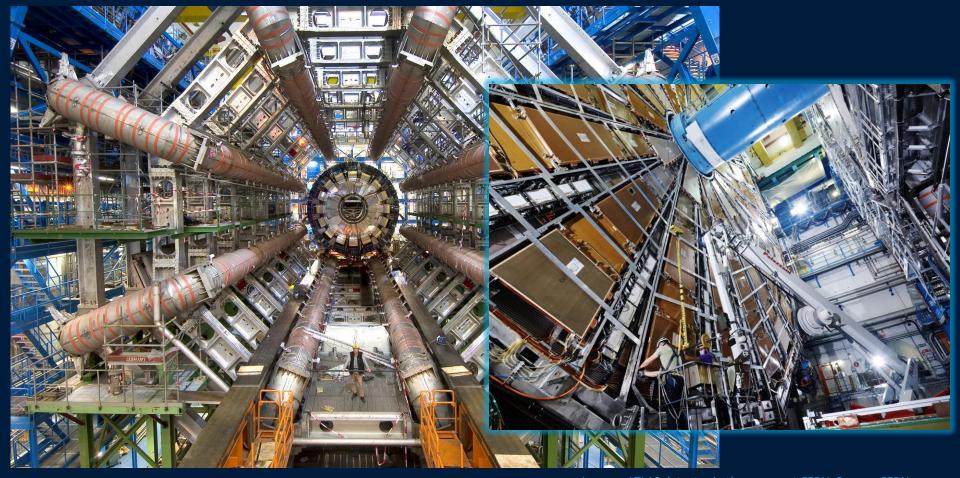


P. Higgs, Phys. Rev. Lett. 13, 508:

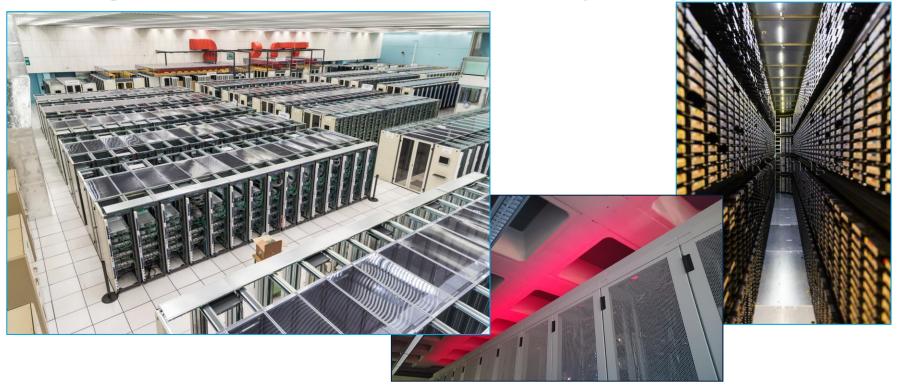
16823 characters, 165 kByte PDF

Maastricht University | DACS



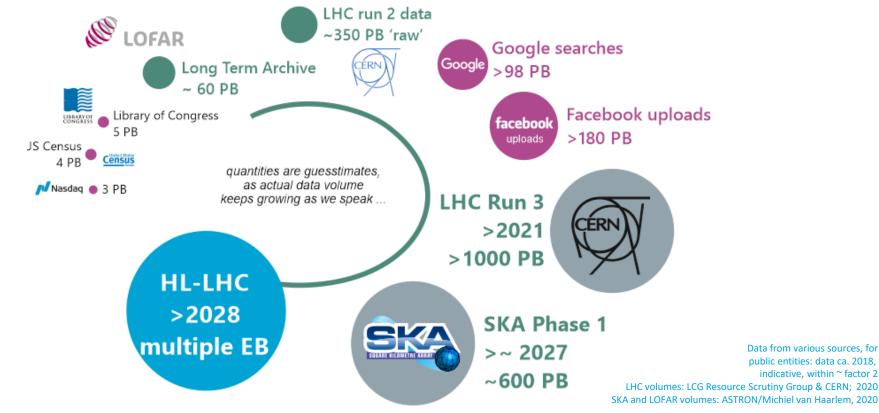


So 'big science' needs some computing ...



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

Processing at scale for data intensive science



Volume versus complexity



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

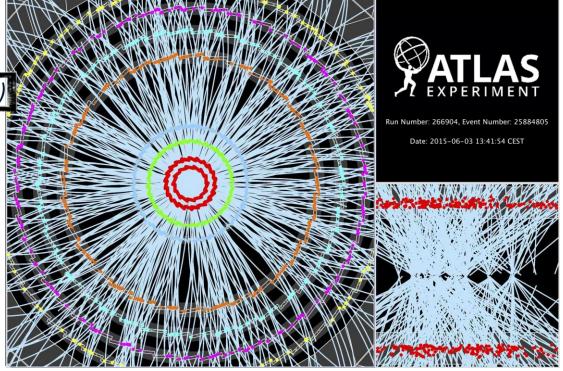
~60 TByte/s (compressed)

Trigger system selects 600 Hz ~ 1 GB/s data

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

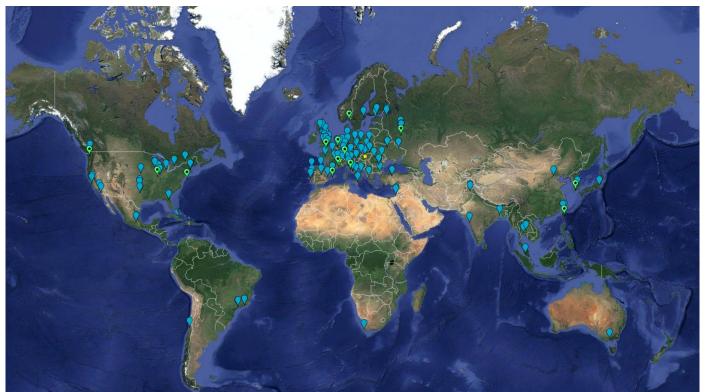
~10k researchers

CERN and ~170 institutes



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)

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

Earth background: Google Earth; Data and compute animation: STFC RAL for WLCG and EGI.eu; Data: https://home.cern/science/computing/grid ACCESS-CI
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

Scaling computing infrastructure – a common need



Sources: CERN https://wlcg.web.cern.ch/; HADDOCK, WeNMR, @Bonvinlab https://wenmr.science.uu.nl/; Virgo, Pisa, IT; SKAO: the SKA-Low observatory, Australia https://www.skatelescope.org/ - OpenMOLE simulation on EGI - https://cdn.egi.eu/app/uploads/2022/04/EGI_Use_Cases.pdf; agent-based modelling of ICAs: https://collective-action.info/research-on-icas/ Molood Dehkordi (TUDelft), Tine de Moor (EUR RSM)

This is a tour of {a,one} large-scale IT landscape

Today: a use case driven overview, looking at

- relationship of infrastructure and one application area: scientific computing
- scaling & interdependency of infrastructure components
 compute, storage, data, network, trust & identity, and security
 'how to identify scalability bottlenecks'
- no directly related questions or project tasks but don't fall asleep just yet

What you should be able to do by the end of today

"Build a globally distributed data processing infrastructure serving different customers and diverse workloads at the Exascale"

or at least

- where to look for relevant standards and reference architectures
- know what scalability issues can affect globally distributed large-scale IT
- realise why this LsIT course is essential to making things actually work ©

For a global large-scale IT system, we need to

1. Build 'sites' that scales to lots of systems

- Datacentre
- Power-efficient systems that can process lots of data
- Orchestrate these systems to function as a service
- sharing workload with brokering, data placement, and platforms

2. Connect sites, services, and users together. Across the world

- with high speed networks, that can deal with latency and severance of links
- for authenticated and authorised users
- that collaborative gain access to services and data
- and are secured and protected ... also from themselves

and then getting actual users ... is 'just' marketing and communications ...

Start with ... a green field approach

from field to IT 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,

You got the power!

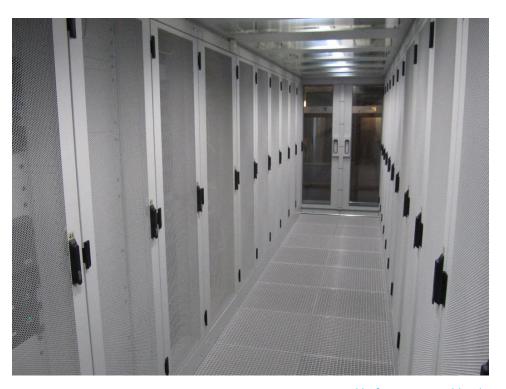




Images: Anton Mors, David Groep, Nikhef

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 ⊕)

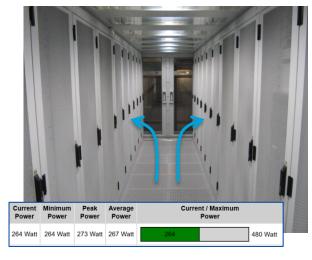
NikhefHousing: a cold aisle

Data centres: where we put large-scale IT

- 'tier-1' ... 'tier-4' datacenters increasingly redundant
- systems are 'lights out': the DC may be miles away
- small or large, in terms of power and cooling capacity:
 - smallish: Nikhef Housing Amsterdam is ~2.5 MW,
 - Dutch AI Factory to be built in Groningen: ~10 MW
 - Microsoft Middenmeer Venster West: 285 MW
- data centre efficiency metric: $PUE = \frac{E_{total}}{E_{IT_equipment}}$

Reducing cost and environmental impact:

- airflow engineering: prevent mixing of cold and hot air
- liquid cooling
- (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

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.nikhef.nl/housing/datacenter/floorplan/ SURF tape total: ~80 PByte by end 2022; image library at Schiphol Rijk from Sara Ramezani; NikhefHousing: https://www.nikhef.nl/housing/datacenter/floorplan/

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

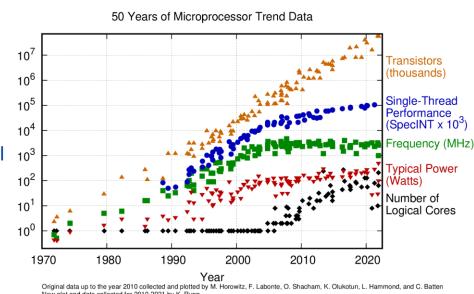
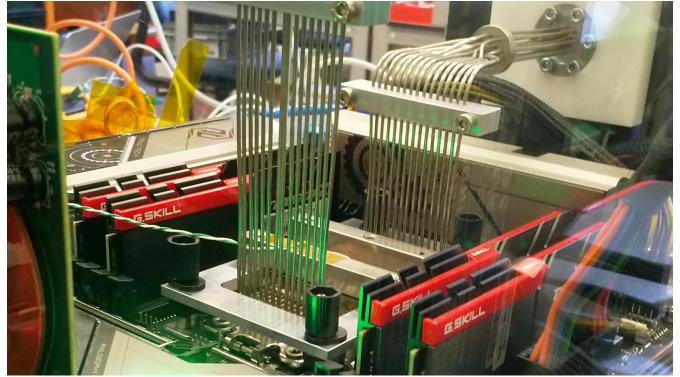


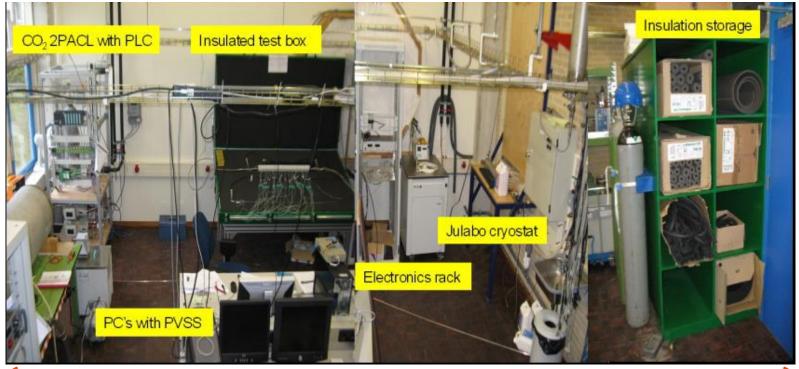
Image: K Rupp, https://github.com/karlrupp/microprocessor-trend-data

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

Step one: scale inside one system

- multiple cores and SMT on a single die
- '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 required between CPU cores and sockets
- 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 may limit single-thread performance (although AMD did better here than 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 HTC, limited by 'off-die memory'

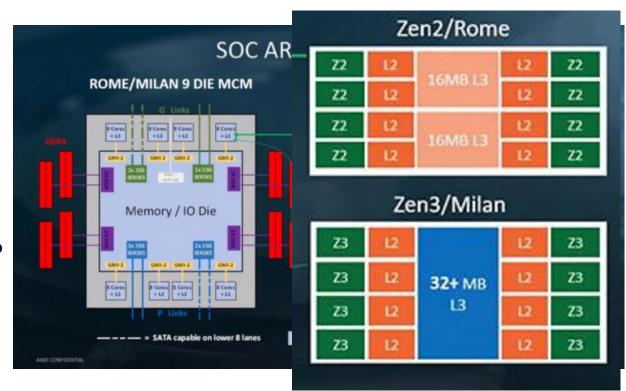
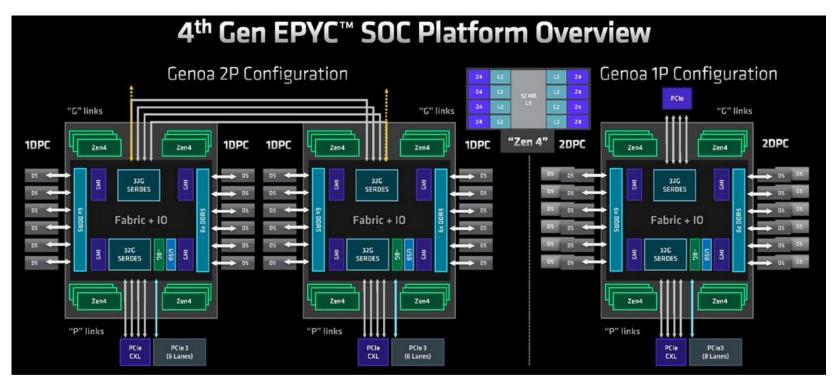


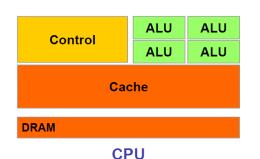
Image source: AMD, retrieved from https://m.hexus.net/tech/news/cpu/135479-amd-shares-details-zen-3-zen-4-architectures/

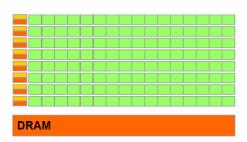
Scaling up, more examples



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

Accelerators – general purpose GPUs





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!

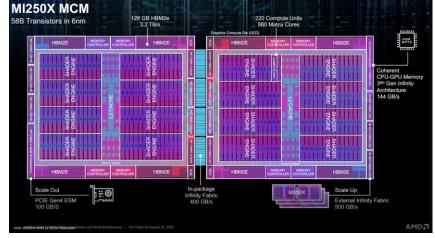
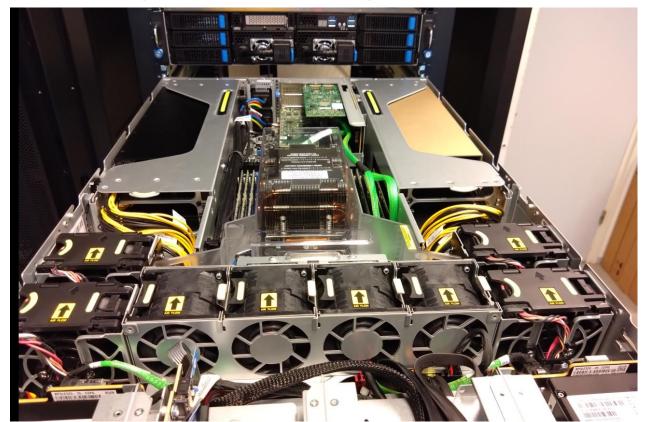


Image: 'Massively Parallel Computing with CUDA', Antonino Tumeo Politecnico di Milano, https://www.ogf.org/OGF25/materials/1605/CUDA_Programming.pdf Floorplan image of die: AMD MI250 GPU, slide source: AMD

GPU enhanced system, with 4 'partitionable' GPUs (L40)



DACS FSE CSLab

DACSGPU003, at Debeyeplein 1

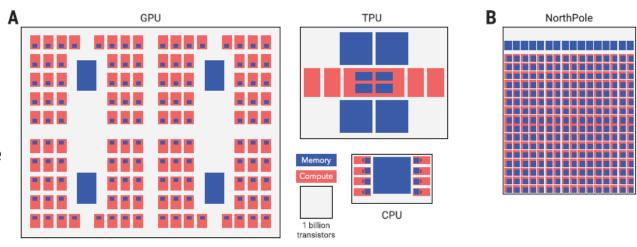
info@fse-cslab.nl

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



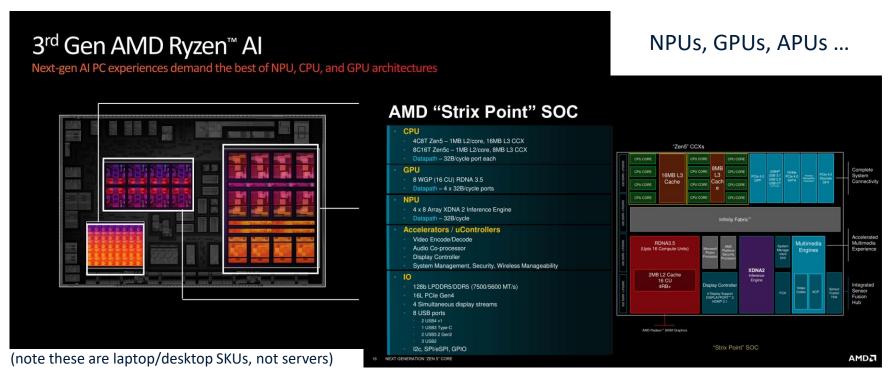


Physical organization of on-chip memory (blue) and compute (red) are diagrammed for representative processors, scaled to constant transistors per unit area. From Modha's paper Modha et al., *Science* **382**, 329–335 (2023)

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



Hybrid SOCs and heterogeneous architectures



Images: AMD Ryzen 9 HX 370 Al, Strix SOC – compare also Intel Lunar Lake architecture

but there is also a serious issue here ...

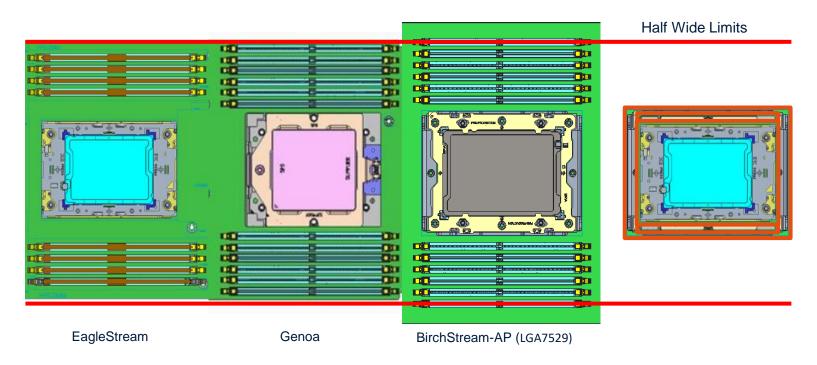
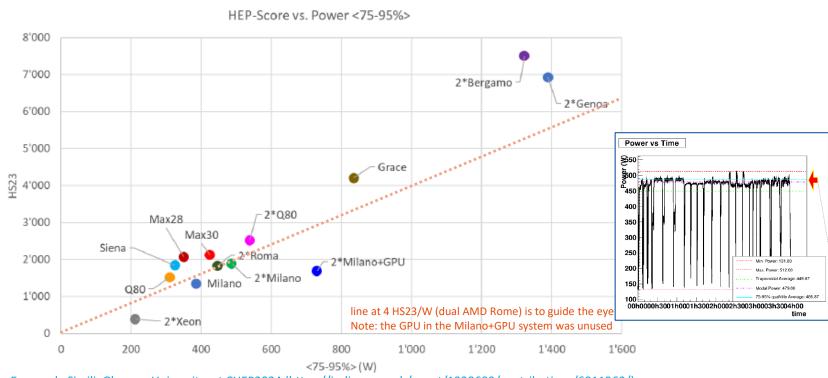


Image thanks go to Rick Koopman – Lenovo at the HTCondor Workshop 2024 https://indico.cern.ch/event/1386170/

The energy bottleneck: architecture figure of merit



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 port d'informació Strongly depends on systems engineering: when water inlet temperature can be >40

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

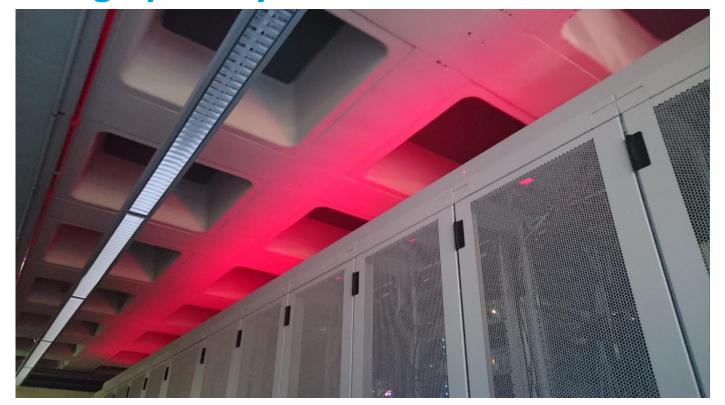
degC, you have almost always free cooling

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





Scaling up – beyond one lone motherboard



Several scalable computing models

- a-priori manual definition of scale: {podman,docker}-compose, cloud VM dashboards
- 'cloud-native' dynamic scaling 'horizontal' or 'vertical' works well for services that lend themselves to load balancing
 (used by kubernetes, docker swarm, many cloud APIs)
- 'infrastructure-as-a-service' tools to meta-manage these, like Hashicorp's Terraform
- Parameter sweeping and throughput computing management: batch systems
- Application-driven scaling: parallel computing (MPI) and distributed underlays (like parallel tasks in matlab or jupyter python notebooks)

← take the parallel computing course ☺

compute farms: 'milking' computer clusters

But workloads usually need more than *just* compute

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

For example for WLCG:

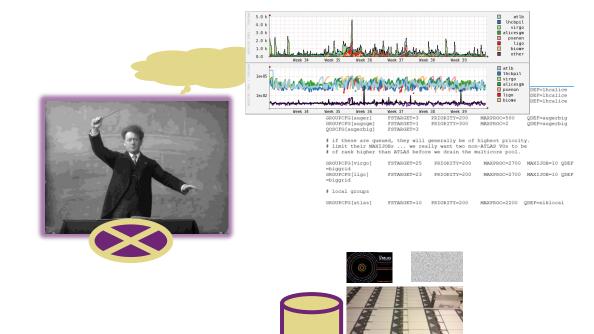
- single-socket multicore systems are fine, today typically 64-128 cores per system
- network: 2x25Gbps (matching #cores)
- memory: say ~ 8 GiB/core
- local disk: 8-16 TB NVME (~100GB/core)
- + space (physical + power) to add GPUs



Image: Cluster 'Lotenfeest' at the Nikhef NDPF, acquired March 2020. Lenovo SR655 with AMD EPYC 7702P 64-Core single-socket. Some with 4 L40s Nvidia GPUs

Cluster computing and 'conveniently parallel' HTC





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

Large-scale IT: worldwide LHC Computing and beyond (2025 ed)

Batch system platform

Many things in life are conveniently parallel

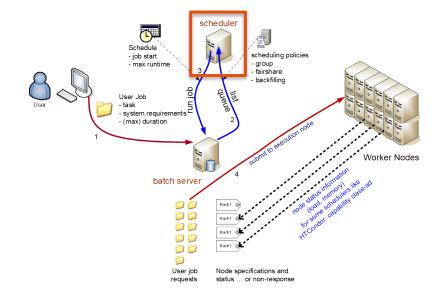
- HEP events & simulation
- ligand matching & drug discovery ——
- structural biochemistry
- ...

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



Job ID	Username	Queue	NDS	TSK	Req'd Memory	Req'd Time	S	Elap Time	
33134895.korf.nikhef.n	lhcbpi08	lhcb	1	1	5120m	41:59:57	R	37:46:21	wn-choc-023
33134901.korf.nikhef.n	lhcbpi08	lhcb	1	1	5120m	41:59:57	R	40:04:09	wn-smrt-12
33134908.korf.nikhef.n	lhcbpi08	lhcb	1	1	5120m	41:59:57	R	37:14:29	wn-choc-03
33134917.korf.nikhef.n	lhcbpi08	lhcb	1	1	5120m	41:59:57	R	14:23:42	wn-smrt-07
33135197.korf.nikhef.n	atlb019	atlasmc	1	4	16040	208:00:00	R	183:02:04	wn-mars-01
wn-mars-018+wn-mars-018	+wn-mars-018	3							
33135883.korf.nikhef.n	atlb019	atlasmc	1	4	16040	208:00:00	R	166:44:22	wn-mars-01
wn-mars-018+wn-mars-018	+wn-mars-018	3							
33142633.korf.nikhef.n	lhcbpi08	lhcb	1	1	5120m	41:59:57	R	37:30:47	wn-mars-04
33149106.korf.nikhef.n	lhcbpi08	lhcb	1	1	5120m	41:59:57	R	10:23:30	wn-car-027
33149132.korf.nikhef.n	lhcbpi08	lhcb	1	1	5120m	41:59:57	R	32:36:49	wn-mars-05
33149220.korf.nikhef.n	lhcbpi08	lhcb	1	1	5120m	41:59:57	R	32:50:19	wn-choc-04
33151669.korf.nikhef.n	lhcbpi08	lhcb	1	1	5120m	41:59:57	R	09:49:53	wn-choc-00
33152704.korf.nikhef.n	atlb019	atlasmc	1	4	16040	208:00:00	R	128:39:13	wn-mars-01

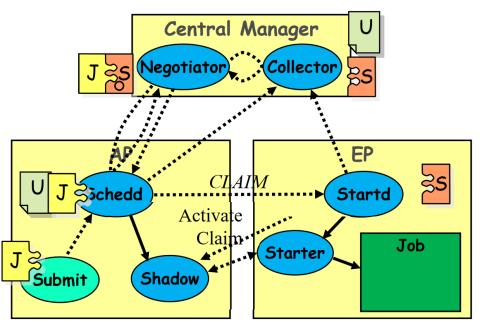
Scalable submission: HTCondor

Matchmaking based on 'ClassAds'

- both jobs and machines
 advertise their requirements
 and capabilities in
 'classified advertisements'
- matchmaking done
 by a negotiator
 execution nodes
 are autonomous (mostly)

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/

Estimated Response Time (and predicting it?)

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

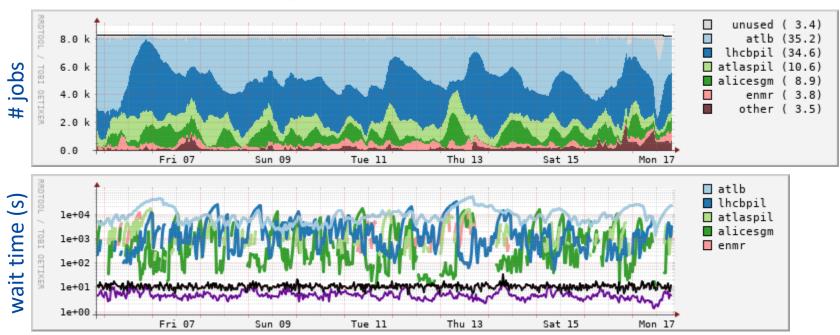


Image: Nikhef NDPF DNI "Grid" cluster. Period: October 6-17, 2022; top-5 communities; GRISview images: Jeff Templon
For work on run time prediction in high-occupancy clusters, see Hui Li Workload characterization, modeling, and prediction ... https://hdl.handle.net/1887/12574

Occupancy: service level differs per intended target audience

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

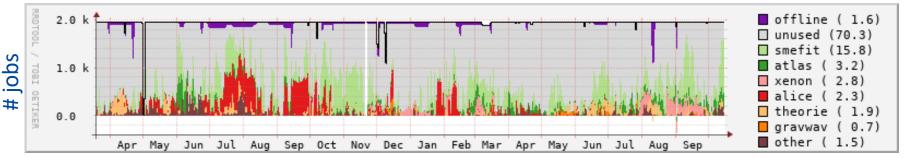
- response time is more important than efficiency
- fast turn-around/short waiting times by heterogeneous ('competing') user base
- data access must be parallelism-ready, but is 'always' local on-site
- local storage credentials and sharing with desktop and Jupyter environments

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

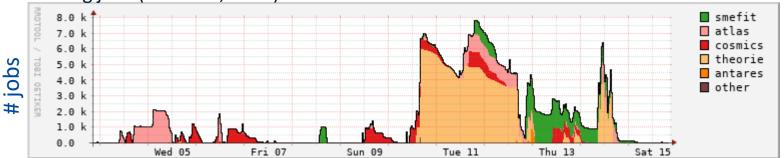
Local analysis clusters 'Tier-3 / Tier-2 computing'

period: March 2021 .. October 2022

Running jobs:





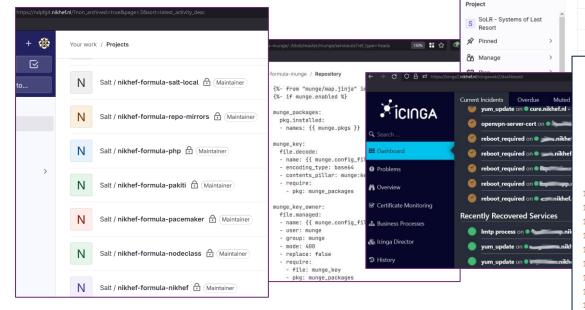


Source: NDPF Statistics overview, https://www.nikhef.nl/pdp/doc/stats/ - GRISview images: Jeff Templon for NDPF and STBC

But ... fancy an interactive (console) install?

Images: Nikhef Housing H234b NDPF science processing data centre

On the way to 'Software Defined Infrastructure' ...



services: datahaga. image: postgres:13.4-alpine environment: - POSTGRES USER=hedgedoc - POSTGRES PASSWORD-SECRET David Groep / SoLR - Systems of Last Res - POSTGRES DB=hedgedoc - database:/var/lib/postgresql/data fabman_centralsyslog deploy: resources: limits. fabman_control_hosts memory: 1G healthcheck: disable: true fabman_core # Make sure to use the latest release from https://hedgedoc.org/latest-release image: quay.io/hedgedoc/hedgedoc:1.9.9 environment: fabman install pxe - CMD DB URL=postgres://hedgedoc:SECRET@database:5432/hedgedoc - CMD URL ADDPORT=true - CMD DOMAIN-sharemd.nikhef.nl - CMD USECDN=false - CMD URL ADDPORT=false # Copyright B - CMD PROTOCOL USESSL=true - CMD ALLOW ORIGIN=['localhost', 'sharemd.nikhef.nl', 'nikhef.nl'] # SPDX-License - CMD HSTS ENABLE=false - CMD CSP ENABLE=false annotations: category: CMS licenses: Apache-2.0 images: - name: apache-exporter image: docker.io/bitnami/apache-exporter:1.0.9-debian-12-r2 10 - name: os-shell image: docker.io/bitnami/os-shell:12-debian-12-r32 - name: wordpress image: docker.io/bitnami/wordpress:6.6.2-debian-12-r15 14 apiVersion: v2 15 appVersion: 6.6.2 16 dependencies: - condition: memcached.enabled 17 18 name: memcached repository: oci://registry-1.docker.io/bitnamicharts 20 version: 7.x.x 21 - condition: mariadh enabled 22 name: mariadb repository: oci://registry-1 docker io/hitnamichart

oot@protosaurus ~1# cat docker-compose-clean.vml

ersion: '3'

រា

Q Search or go to...

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/



Scaling things '... as a service'

Systems or storage today are usually not physical

- 'predefined' virtualization
 of systems, network and (block) storage
- can be single site, distributed, outsourced, or federated
- laaS: Infrastructure as a Service ('EC2-like' VM hosting, VPS, S3 storage)
- PaaS: Platform as a Service (container hosting, batch systems, database platforms like DynamoDB ...)
- SaaS: Software as a Service (science application portals, ERP systems, ...)

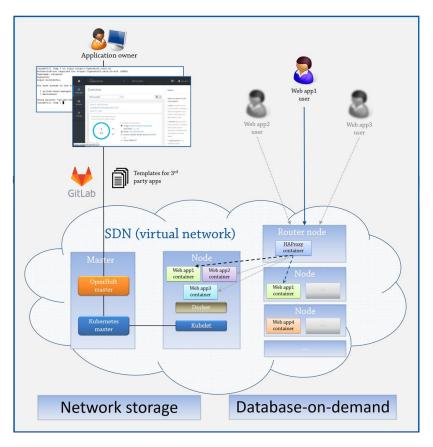


Image from CERN's OpenShift, A Lossent et al 2017 J. Phys.: Conf. Ser. 898 082037 https://doi.org/10.1088/1742-6596/898/8/082037

Moving the management boundary

Infrastructure-as-a-Service **Application** Data Runtime environment Guest Middleware Operating system Hyper Virtualisation layer visor Physical server Storage devices Host Network

Platform-as-a-Service **Application** Data Runtime environment Middleware Operating system Virtualisation layer Physical server Storage devices Network

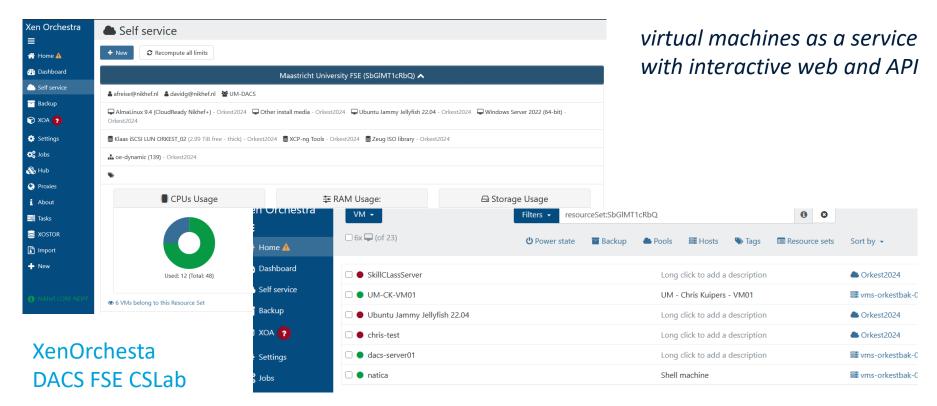
Application Data Runtime environment Middleware Operating system Virtualisation layer Physical server Storage devices Network

Software-as-a-Service

Maastricht University

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

Infrastructure as a Service example: DACS FSE CSLab



What is a 'service': on Service Management Systems

Structuring service management

- ISO 20000 standard
- ITIL (now at ITIL v3)
- https://www.fitsm.eu/

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

FitSM: ITSM process framework



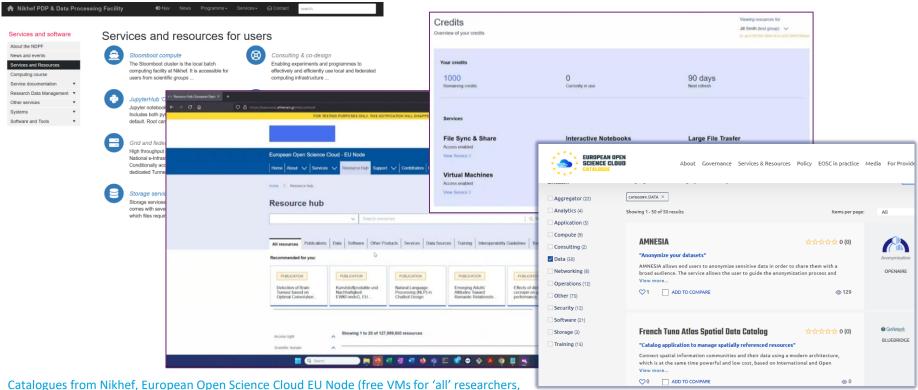
- 1. Service portfolio management (SPM)
- 2. Service level management (SLM)
- 3. Service reporting management (SRM)
- 4. Service availability & continuity management (SACM)
- 5. Capacity management (CAPM)
- Information security management (ISM)
- Customer relationship management (CRM)
- 8. Supplier relationship 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

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

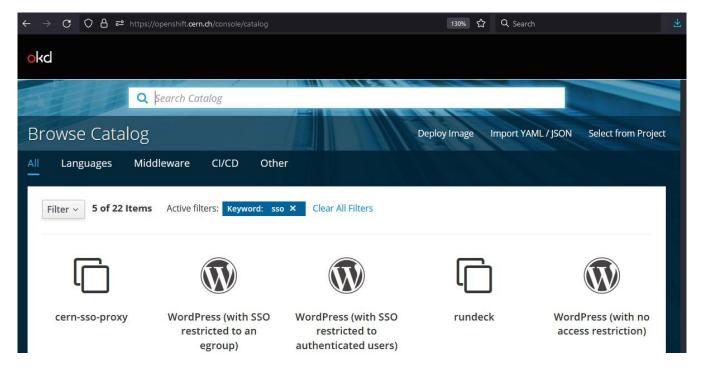
Maastricht University

Service portfolios – what do you offer, and to whom



subject to https://open-science-cloud.ec.europa.eu/system/files?file=2024-10/EOSC-EU-Node-User-Access-Policy-v1.0.pdf)

A service catalog is still built on systems & networks



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?

Tools and applications USER APPLICATIONS Directory brokering, diagnostics, and COLLECTIVE SERVICES Secure RESOURCE AND access CONNECTIVITY PROTOCOLS to resources and services Diverse resources such as FABRIC computers, storage media, networks, and sensors

'protocol hourglass'

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

Standard interfaces for compute and data?

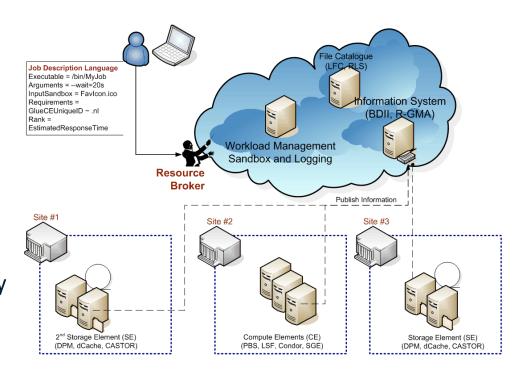
hourglass model 'kind-of' worked for IP and web with http as common standard

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

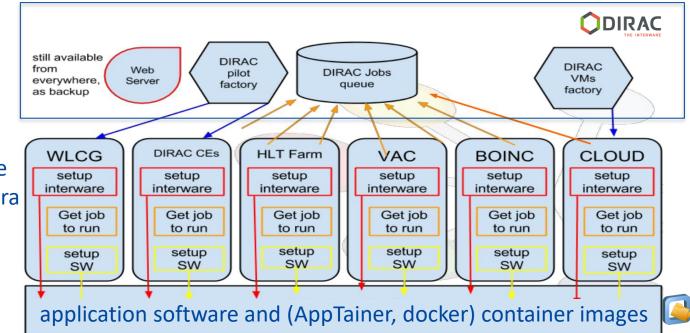
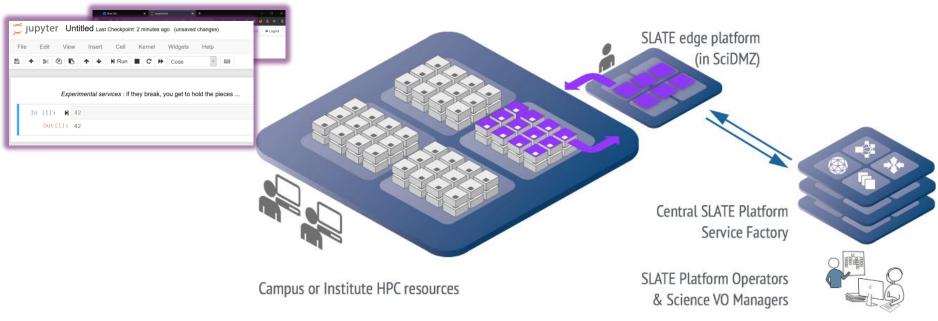


Image: DIRAC project, A. Tsaregorodtsev et al. CPPM Marseille, from https://dirac.readthedocs.io/; CVMFS (CERN VM File System) is a common software distribution platform using distributed signed data objects in a cached hierarchy using CDN techniques, see https://cernvm.cern.ch/fs/

An overlay network of containers

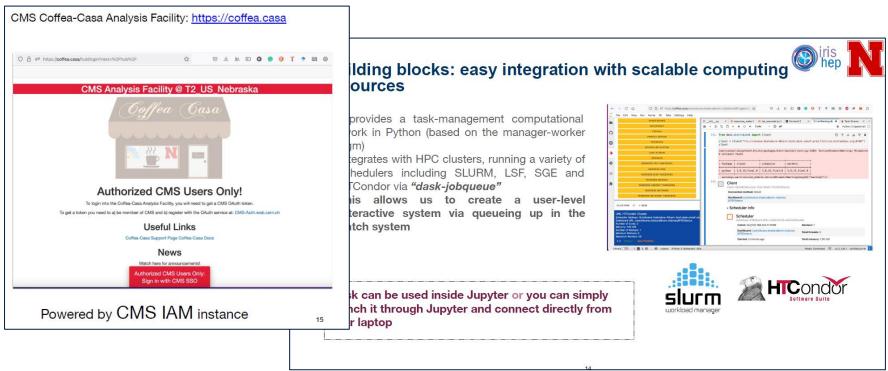
Nobody wants a cloud per-se ... what folk want is a solution ...



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

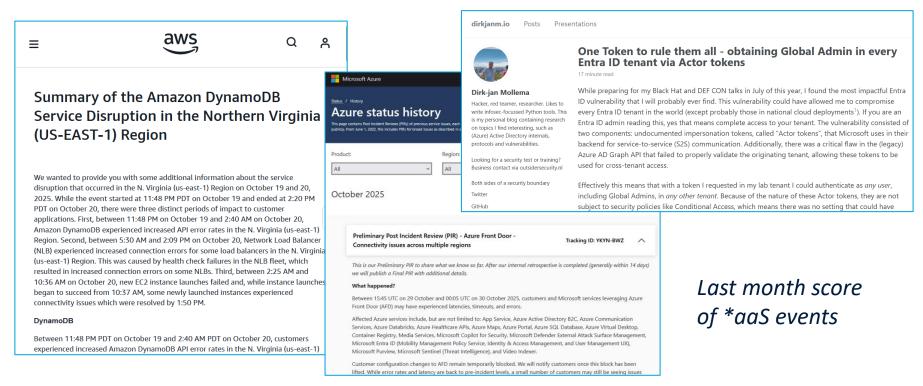
Image sources: NDPF JupyterHub service "Callysto"; SLATE: Service Layer At The Edge – Rob Gartner (UChicago), Shawn KcMee (UMich) et al. – slateci.io

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/

On leaky abstractions and circular dependencies



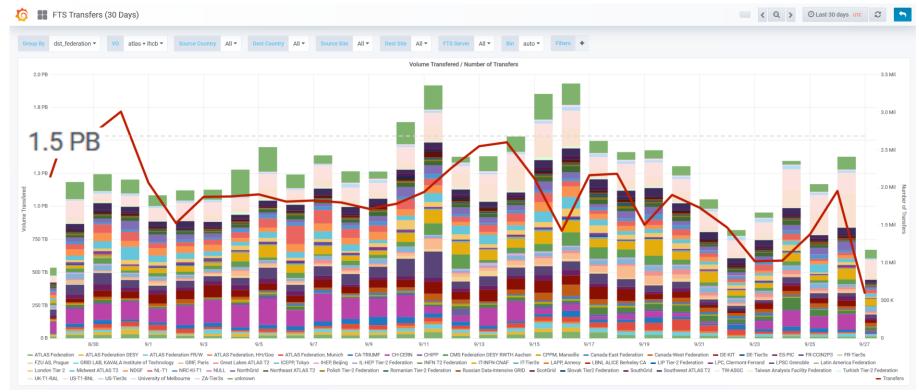
https://aws.amazon.com/message/101925/, https://azure.status.microsoft/en-us/status/history/, https://dirkjanm.io/obtaining-global-admin-in-every-entra-id-tenant-with-actor-tokens/

Putting more than one together

The Internet Is Not Enough! More than one user ... Perspectives on global federated access



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

'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

Traceroute measurement to linuxsoft.cern.ch (multihomed)

General Information Probes Map TraceMON IPmap (beta) Results



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

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

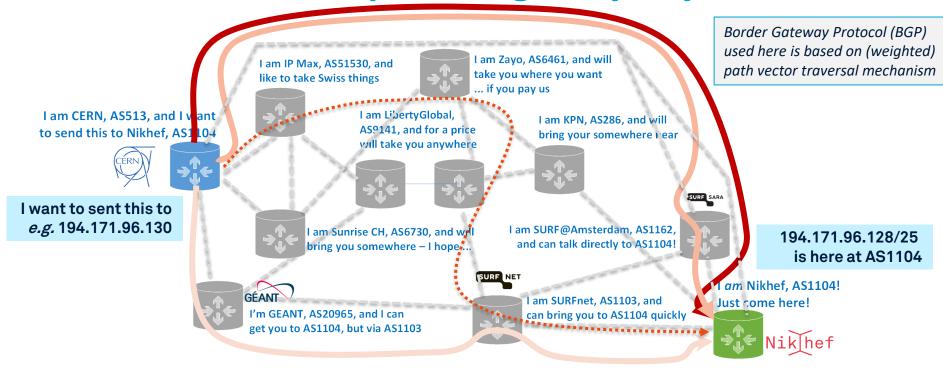
5 irb-4.asd002a-jnx-06.surf.net (2001:610:f00:1120::121) [AS1103] 2.041 ms 2.138 ms 2.138 ms

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?



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@deelgfx-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
                                            10
                192.16.166.21
                                                              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
                                                              513 43475 T
```

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

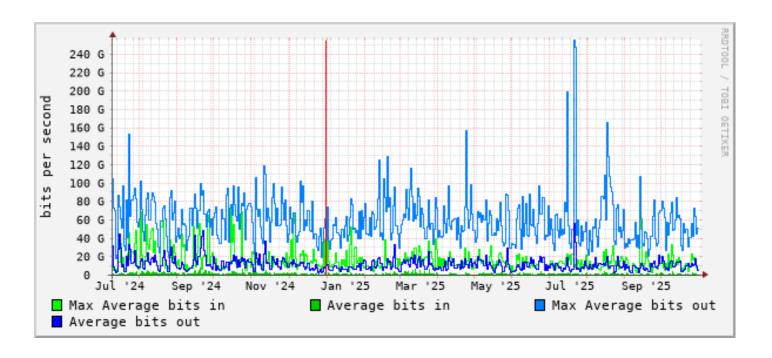
Lclpref

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



Source: Nikhef cricket graph for compute clusters only on deelqfx – https://cricket.nikhef.nl/

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

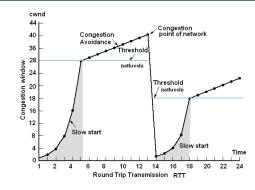
Theoretical network limit

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://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 ...

latency AMS-GVA 17 ms congestion event @20ms: 2 ms of UDP traffic to GVA

- 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

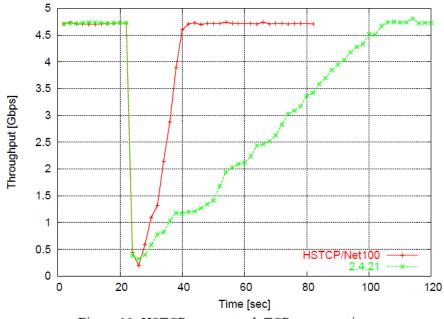
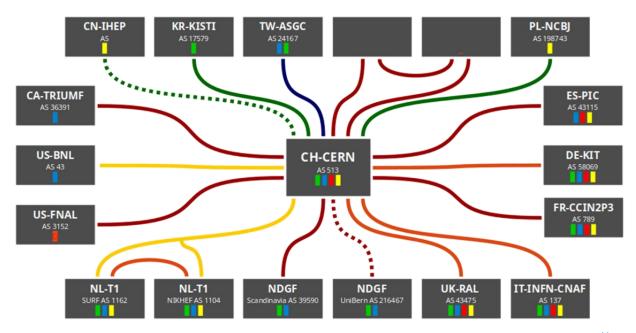


Figure 10: HSTCP versus stock TCP recovery time

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

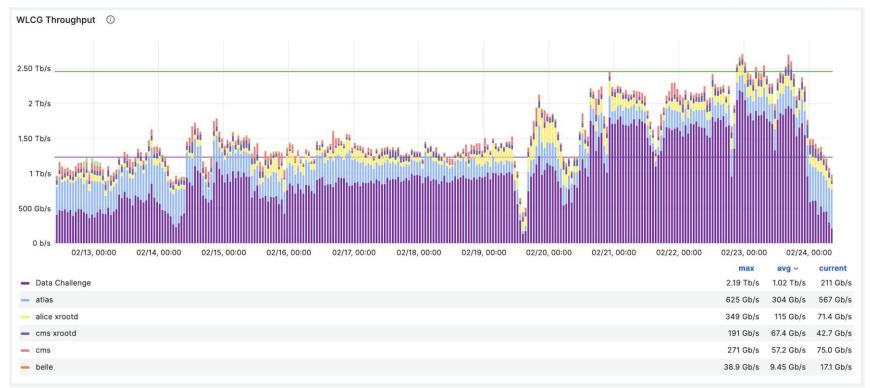
LHCOPN – distributing raw data LHC PN





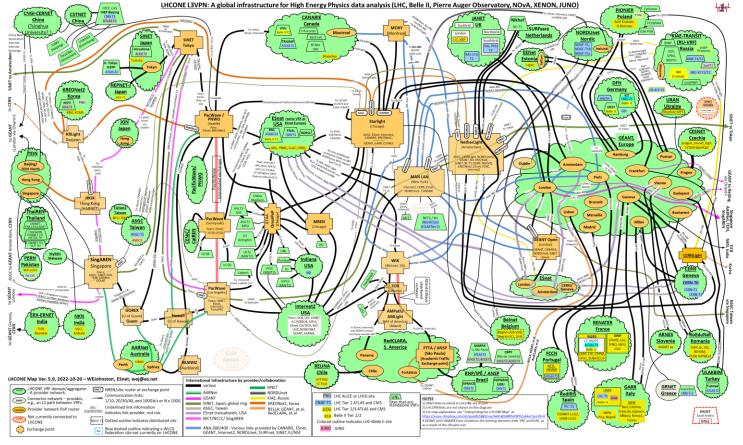


LHCOPN – traffic levels for data transfer (DC24)



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

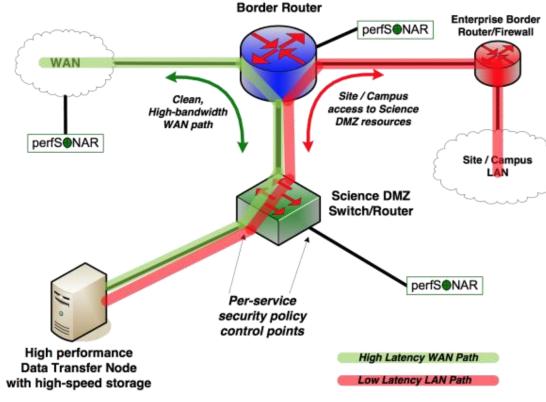
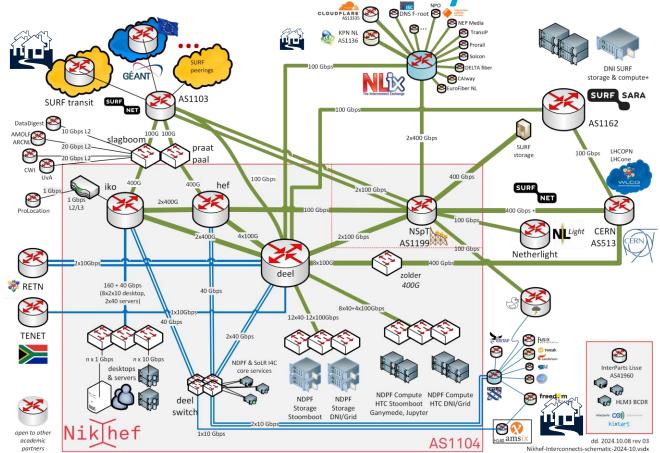


Image and 'ScienceDMZ' concept promulgated by ESnet (see fasterdata.es.net)

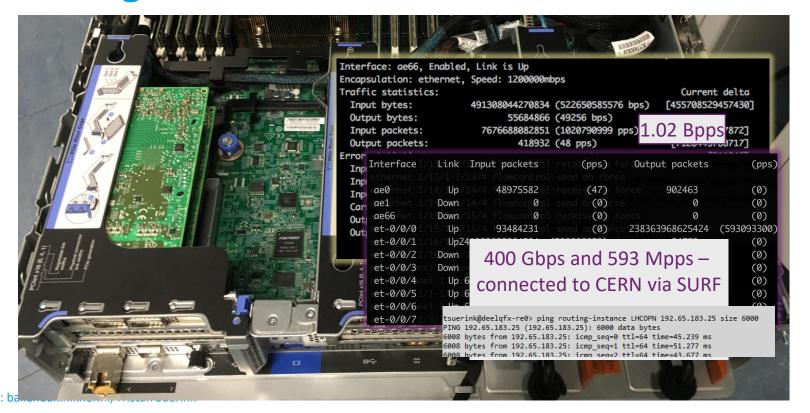
Just one random autonomous system: AS1104



state as of Oct 2024

AS1104

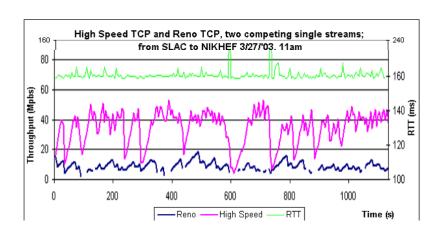
Exercising the network – sensor data and events



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



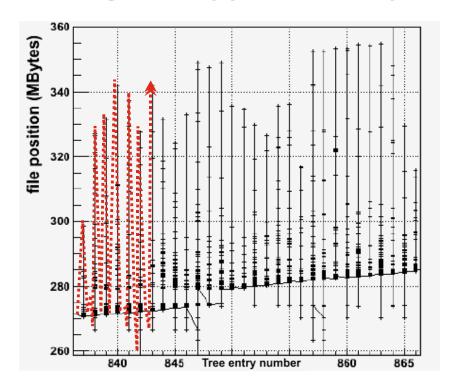


Image of TCP slow-start and packet loss impact (in Mpps): Antony Antony et al., Nikhef, for DataTAG, 2003(!)
Right: base graphic: Philippe Canal "Root I/O: the fast and the furious", CHEP2010 Access pattern reflects Root versions < 5.28, before Ttree caching and 'baskets'

And some traffic is triggered by researchers scaling up 'accidentally' from a laptop to a cluster without too much thought

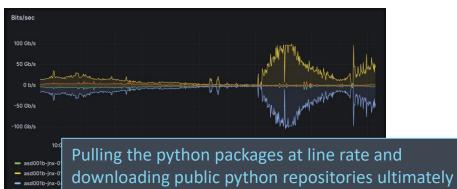
Min = 194.6 Gbps

Copyright (c) 2023 AMS-IX B.V.

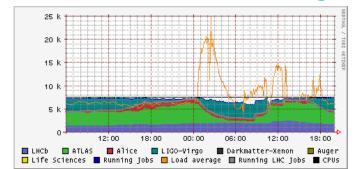
A researcher doing mass creation of containers, rebuilding their python 'virtual env' for each job, running on >> 4000 cores

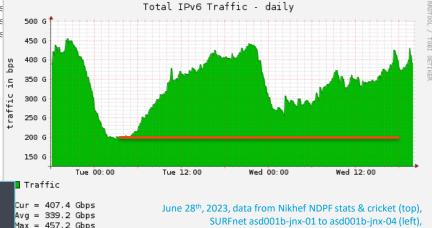
```
[root@wn-pep-002 ~]# top
top - 09:40:47 up 71 days, 12:17, 2 users, load average: 110.38, 101.43, 106.3
Tasks: 700 total, 7 running, 666 sleeping, 0 stopped, 27 zombie
%Cpu(s): 17.0 us, 2.0 sy, 0.0 ni, 81.0 id, 0.0 wa, 0.0 hi, 0.0 si, 0.0 st
KiB Mem: 39462902+total, 23514457+free, 10406320 used,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



will trigger Cloudflare and flood SURFnet





AMS-IX SFlow https://stats.ams-ix.net/sflow/index.html (bottom)

Updated: 28-Jun-2023 19:55:02 +0200



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 \rightarrow 1600G AMS-GVA \odot







Home BTG BTG Services INTUG Innovatielab Activiteiten Lobby & Opinie Publicaties

Minister Adriaansens opent testomgeving voor volgende generatie netwerktechnologieën

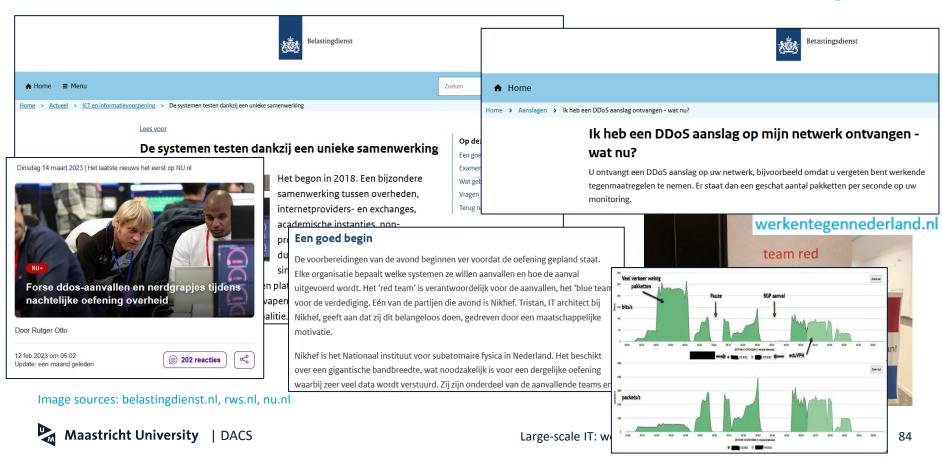


in Amsterdam is door minister Micky Adriaansens van Economische Zaken en Klimati ierotonde is een testomgeving waar SURF en Nikhef gaan experimenteren met Rieuave ng beschikt over een internetsnelheid van 800 Gbit/s, wat meer dan 1000 keer sneller ngemiddeld huishouden in Nederland. De innovatierotonde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

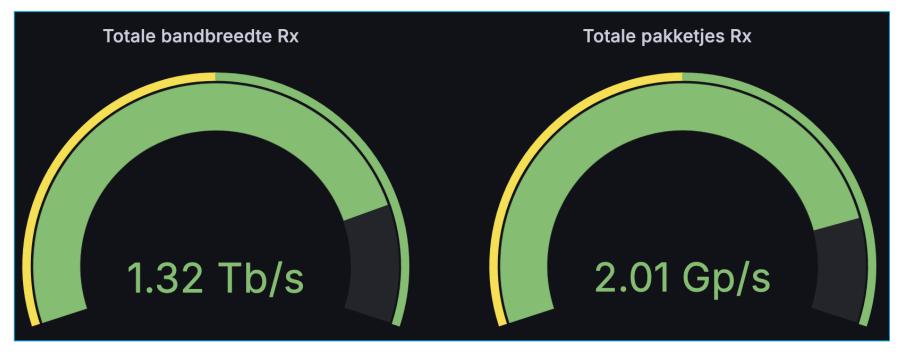
an onderzoek naar bandbreedte op het internet groeit. Onderzoekers willen steeds meer pver de landsgrenzen heen met elkaar delen. De bandbreedte van het netwerk speelt ote hoeveelheden data snelt te kunnen verwerken, is de verwachting dat 8000bit/s J. De innovalterotonde maakt het mooellik om te experimenteren met nieuwe

Research data traffic looks like ... a DDoS to others ©





with packets being more destructive than bandwidth ...



https://wiki.nikhef.nl/grid/2Bpps Machine - in preparation for the 2025 Resilience Exercise -

Access: Trust & Identity

More than one user, from more than one organization, in more than one country



WLCG: when we met a global trust scaling issue



- 170 sites
- ~50 countries & regions
- ~12000 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

NIKHEF, NATIONAAL INSTITUUT VOOR KERN	state of Grid a	and	the L	.HC c	omp	outing i	n 20
Guest/students form (pleas	Fermilab				Office Use O		
This form is completed in work experien			ID: Insurance:	Action		ID Exp: Safety:	- 1
connection with: otherwise, viz.			Computer: Stkrm:			Family:	1
CERN/User Registration			NON-473:	Sensitive:	Verifier:	Date:	1
CERN COMPUTER CENTRE - US		_				-	1
http://cern.ch/it/documents/ComputerUsage/CompA	Name: SWIETZER	JOHN			JAMES		- I
To be returned to the User Registration box at the en	Last	First		Midd	le		- I
completed by a user who requires a computer accou Department, and is not yet registered in another gro	University or Institution Name: FLORIDA STATE UNIVERSITY	7	Telephone: 850-644-XXXX				- I
	Experiment/Department:						'
treated confidentially and only be used for ensuring Supply name as registered by the Users' Off	Exp. / Dept. Spokesperson		Home Institution Contact SHARON HAGOPIAN			act Telephone	⊒ I
FAMILY NAME(S):	D0 WOMERSLEY/V	VEERIS	SHAKON HA	AGOPIAN	850-6	44-4777	
FIRST NAME(S):							_
	Month Year ERVISOR MBER (as on CERN card)			CAUPTOCure	G2 1 GH 4 PRS 7	ABC 05F 06F 06F 06F 06F 06F 06F 06F 06F 06F 06	

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
balt:x:7:0:balt:/sbin/sbin/halt
```

root:\$6\$s8ciAG5gLuv2bPQS\$6EcskgtKvQ.rHbif davidg:\$6\$nDYcIez2Uaufbtlg\$R1hS/Qjn0gYQZk

marianne:\$6\$p3CeevG6jfNDqZj1\$HKHqUTnt2fEqQfkA/m5J3oAOA0zSvgLCKOSQhPS

Passport image: cropped from original by Jon Tyson on Unsplash https://unsplash.com/photos/Hid-yhommOg



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

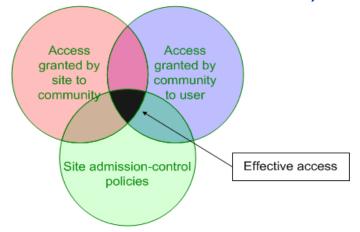
```
ldap_access_order = filter,authorized_service
ldap_access_filter = (|(memberOf=cn=gridSrvAdministrators,ou=DirectoryGroups,dc=farmnet,dc=nikhef,dc=nl)(memberOf=cn=gridMWSecurityGroup,ou=DirectoryGroups,dc=farmnet,dc=nikhef,dc=nl)(memberOf=cn=nDPFPrivilegedUsers,ou=DirectoryGroups,dc=farmnet,dc=nikhef,dc=nl))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

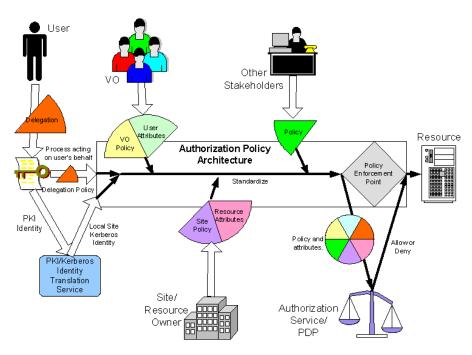
Policy example: Argus system, https://argus-documentation.readthedocs.io/en/stable/misc/examples.html; service-specific: sssd.conf ldap auth provider

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: Open Grid 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

Scaling credentials: per service per user

Many start with *credentials* dedicated to each service where you need access

Number of to-be-protected credentials:

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

usually creates a strong link to authorization:

different accounts for different roles, multiplying the number of credentials per user

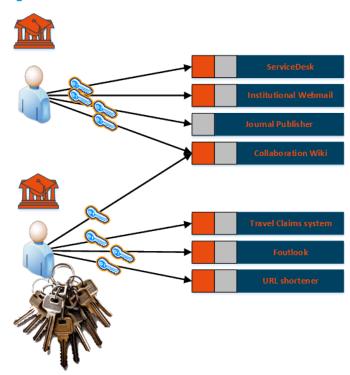


Image imspired by AARC NA2 training module "Authentication and Authorisation 101" – keychain image created by generative AI

bilateral SSO: per service single identity source per 'home' organisation

#credentials required \downarrow from previously

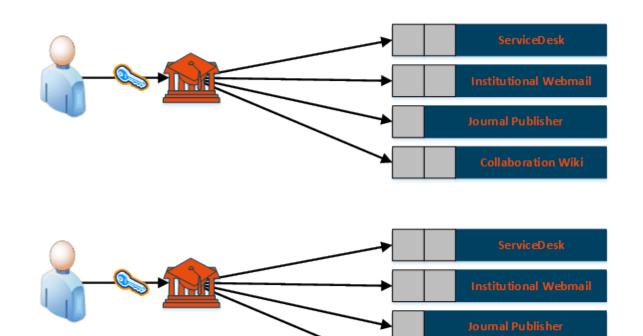
$$\mathcal{O}(n_{\text{services}}) * \mathcal{O}(n_{\text{users}})$$

to

$$\mathcal{O}(n_{users})$$

+ $\mathcal{O}(n_{services}^*n_{home-orgs})$

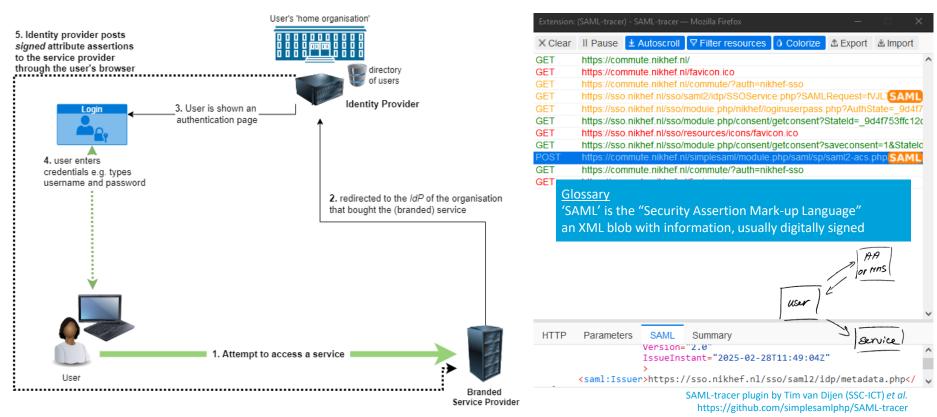
in first order at least





Collaboration Wiki

Enterprise (SAML) Single sign-on: browser redirect as login flow

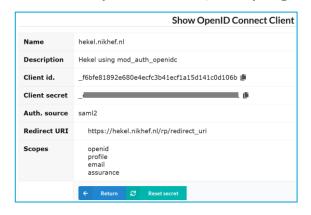


Under the hood, sends a (signed) XML document

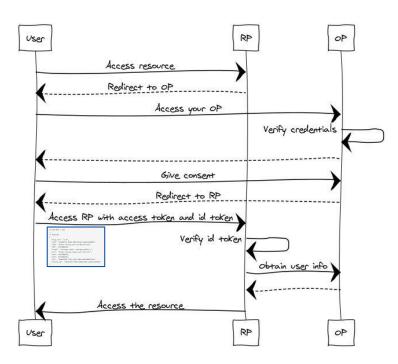
```
<saml:Subject>
    <saml:SubjectConfirmation Method="urn:oasis:names:tc:SAML:2.0:cm:bearer">
      <saml:SubjectConfirmationData NotOnOrAfter="2022-10-21T18:16:40Z"</pre>
        Recipient="https://attribute-viewer.aai.switch.ch/Shibboleth.sso/SAML2/POST"
        InResponseTo=" 64c10a60c382bdaeb328653d9d25951c" /></saml:SubjectConfirmation>
  </saml:Subject>
   <saml:Conditions NotBefore="2022-10-21T18:11:39Z"</pre>
                   NotOnOrAfter="2022-10-21T18:16:402">
    <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 | <saml:AttributeStatement>
                                                            <saml:Attribute Name="urn:mace:dir:attribute-def:cn"</pre>
                       SessionNotOnOrAfter="2022-10-22T0
                                                                           NameFormat="urn:oasis:names:tc:SAML:2.0:attrname-format:uri">
                       SessionIndex=" 90f745f18f712b6a56
                                                              <saml:AttributeValue xsi:type="xs:string">David Groep</saml:AttributeValue>
   <saml:AuthnContext>
                                                            </saml:Attribute>
       <saml:AuthnContextClassRef>urn:oasis:names:tc:SAM
                                                            <saml:Attribute Name="urn:oid:2.5.4.3"</pre>
       <saml:AuthenticatingAuthority>https://sso.nikhef
                                                                            NameFormat="urn:oasis:names:tc:SAML:2.0:attrname-format:uri">
   </saml:AuthnContext>
                                                              <saml:AttributeValue xsi:tvpe="xs:string">David Groep</saml:AttributeValue>
   </saml:AuthnStatement>
                                                            </saml:Attribute>
                                                            <saml:Attribute Name="urn:mace:dir:attribute-def:eduPersonAffiliation"</pre>
                                                                            NameFormat="urn:oasis:names:tc:SAML:2.0:attrname-format:uri">
                                                              <saml:AttributeValue xsi:type="xs:string">employee</saml:AttributeValue>
                                                              <saml:AttributeValue xsi:tvpe="xs:string">member</saml:AttributeValue>
                                                              <saml:AttributeValue xsi:tvpe="xs:string">facultv</saml:AttributeValue>
                                                            </saml:Attribute>
                                                            <saml:Attribute Name="urn:oid:1.3.6.1.4.1.5923.1.1.1.1"</pre>
```

OpenID Connect and OAuth2: the 'modern' way

- Quite .well-known (used by lots modern 'non-enterprise' SSO)
- shows signs of its initial design objective: one source of identity (Openid Provider, 'OP'), and many services (Relaying Parties, 'RP')





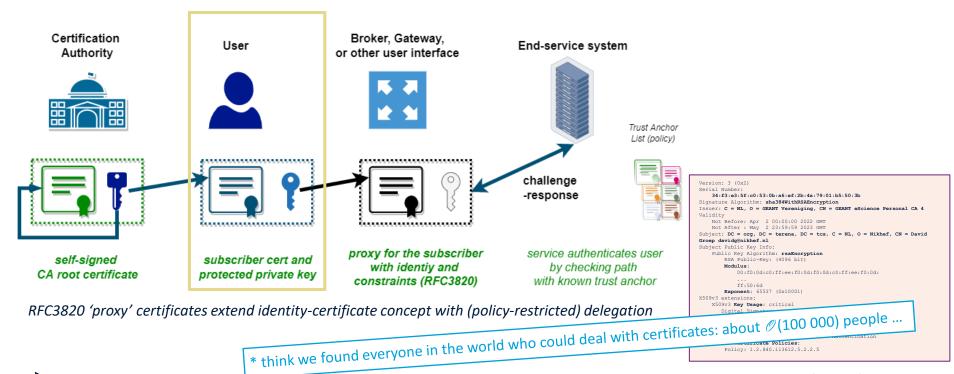


Shown is the 'implicit flow', other flows possible. Image source: AARC NA2 training on AAI 101

See https://openid.net/ for protocols and standardization work

PKI client certificates – user* client held credentials

You have seen https, but the same PKI Certificates can be used for clients, not servers ...



Different tech also an AAA push concept: X.509 and a trust PKI

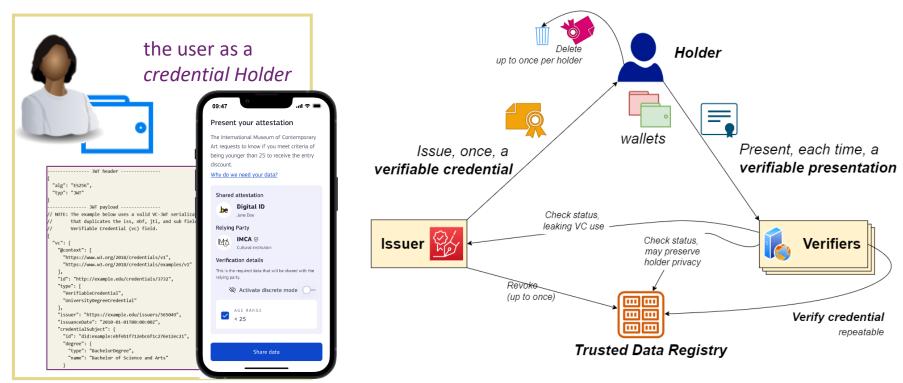
```
Version: 3(0x2)
Serial Number:
    34:f3:e3:5f:c0:53:0b:a6:ef:2b:4a:79:01:b5:50:3b
Signature Algorithm: sha384WithRSAEncryption
Issuer: C = NL, O = GEANT Vereniging, CN = GEANT eScience Personal CA 4
Validity
    Not Before: Apr 2 00:00:00 2022 GMT
    Not After: May 2 23:59:59 2023 GMT
Subject: 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



Identity wallets with VCs held by the user, are another



Flow diagram inspired by: Lifecycle Details (5.1), Verifiable Credentials Data Model V1.1, W3C Recommendation 03 March 2022, https://www.w3.org/TR/vc-data-model/EU eID Wallet from https://commission.europa.eu/strategy-and-policy/priorities-2019-2024/europe-fit-digital-age/european-digital-identity_en
Appimage: European Commission, at https://ec.europa.eu/digital-building-blocks/sites/display/EUDIGITALIDENTITYWALLET/Security+and+Privacy



Interoperable wallet federation example: life-long learning

time to think less institution-centric?

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

Study

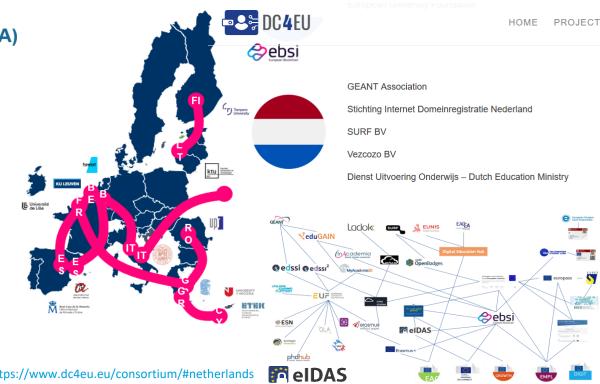
- 01 A student gets a diploma with a list of course units validated from Erasmus (<u>Transcript of Records Credential</u>) (<u>ES/BE/IT</u>)
- 02 A student applies for a PhD with a Bachelor / Master degree from a foreign country (<u>Bachelor/Master Diploma Credential</u>) (<u>RO/GR/FR</u>)
- 03 A student gets access to local discounts using student credential (European Student IDentity) (BE/ES)
- 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

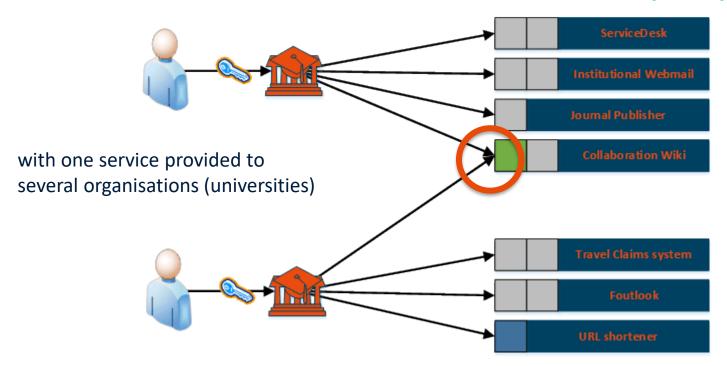
06 A PhD student applies for specific courses in a foreign country (Cross-border Micro-credentials) (FI/LT)



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

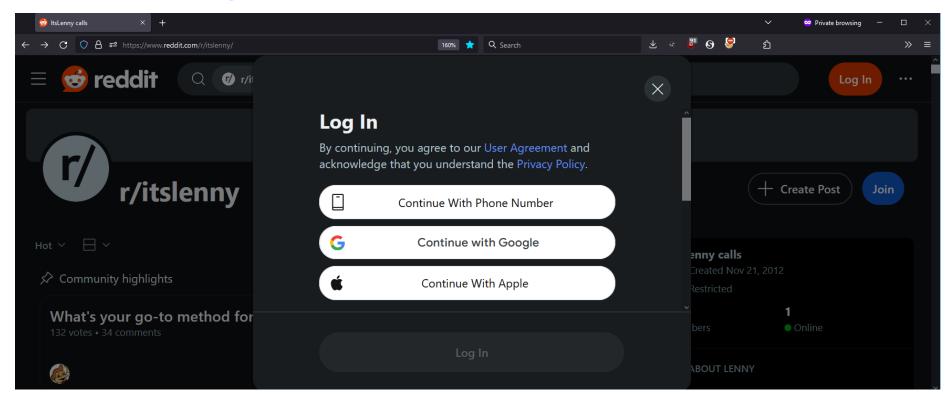


Can we scale better with a 'federated' Authentication and Authorisation Infrastructure ('AAI')



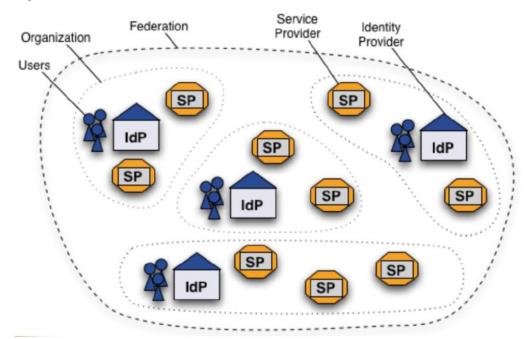
Looking at authentication first ...

Where are you from??



Federation

portability of identity information across otherwise autonomous administrative domains

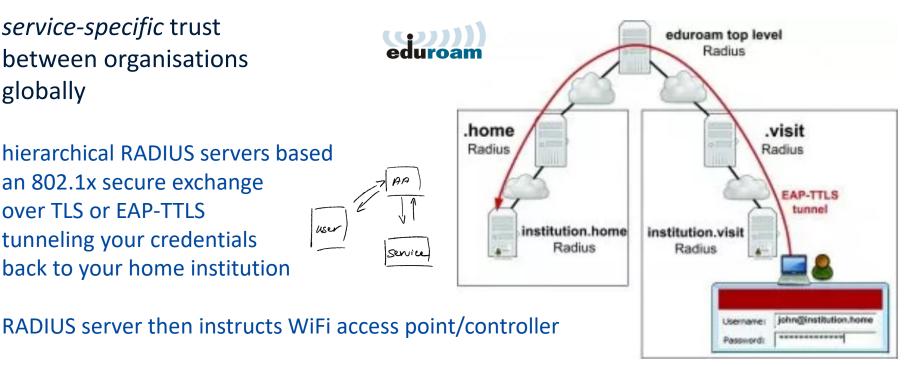


Shibboleth IdP image and SAML2 auth flow by SWITCH (CH) – see also https://refeds.org/ on federation structure and (assurance and security) guidelines

One simple federation you know: eduroam

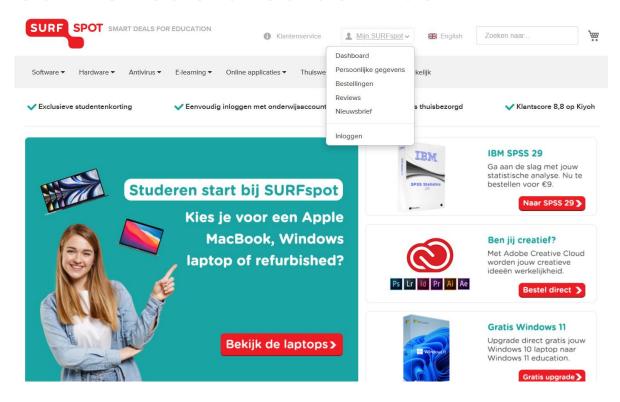
service-specific trust between organisations globally

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



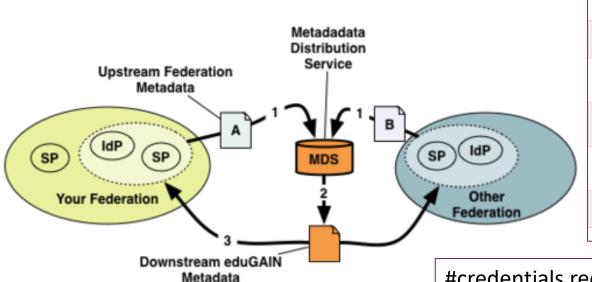
eduroam: Klaas Wieringa et al., image from https://eduroam.org/how/, GEANT; RADIUS: RC2865 https://www.rfc-editor.org/rfc/rfc2865; see also freeradius.org

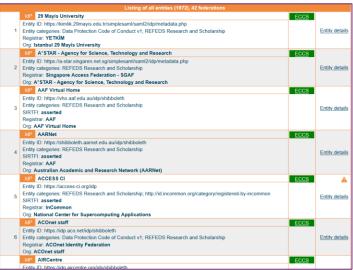
Your favourite federated service?



https://surfspot.nl/

Multilateral federation and entity meta-data



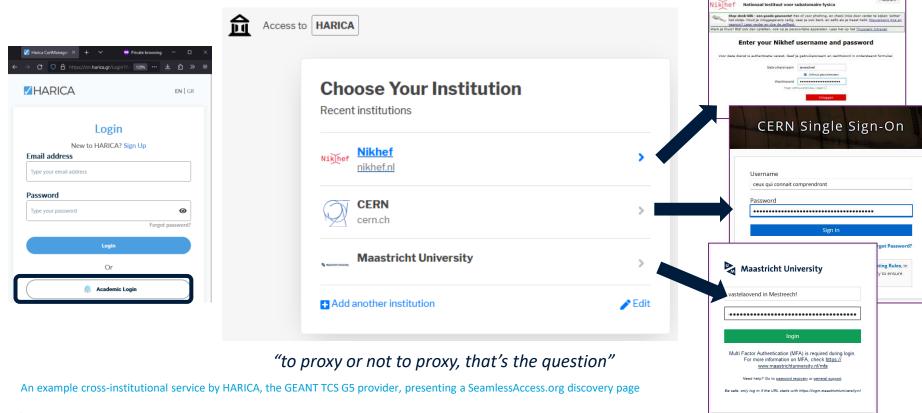


#credentials required?

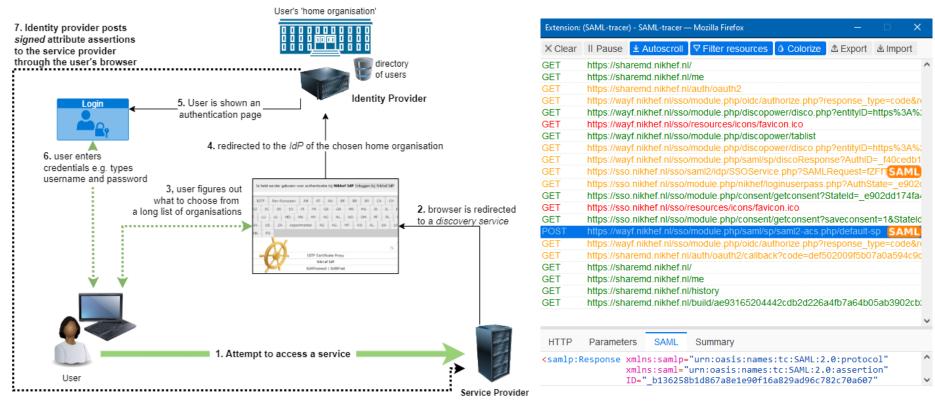
 $\mathcal{O}(n_{users}) + \mathcal{O}(n_{services} * n_{home-orgs})$ from $\mathcal{O}(n_{users}) + \mathcal{O}(n_{home-orgs}) + \mathcal{O}(n_{services})$ to

MDS meta-data flow: https://wiki.geant.org/display/eduGAIN/Metadata+Flow+in+eduGAIN eduGAIN meta-data https://mds.edugain.org/edugain-v2.xml; table excerpt from https://technical.edugain.org/entities showing only R&S IdPs, i.e. those supporting research ...

Federation and discovery can be different aspects



Full mesh or hub-n-spoke, with SAML meta-data



Access to the ShareMD service through token translation service SAML->OIDC with the simpleSAMLphp suite – simplesamlphp.org, sharemd.nikhef.nl

AAI: different technologies, same federation idea

SAML - Security Assertion Markup Language WebSSO ('SAML2Int' federation)

- XML-formatted 'attribute statements' over web transport (usually POST)
- **SAML-Metadata**: list of entities with description of bindings with entityAttributes

PKI - Public Key Infrastructure

- 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 ...

See www.oasis.org for SAML; RFC5280 (tech) & RFC3247 (policy) for PKIX, https://igtf.net/ and https://cabforum.org; OpenID Connect Federation: https://openid.net/specs/openid-connect-federation-1 0.html



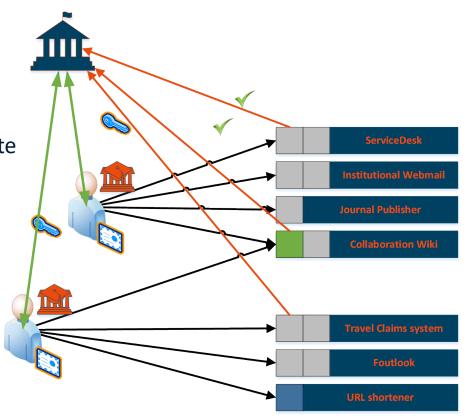
Federation with user-centric AAI

A trusted authority giving the user a 'self-managed' credential, like a passport

a personal authentication digital certificate

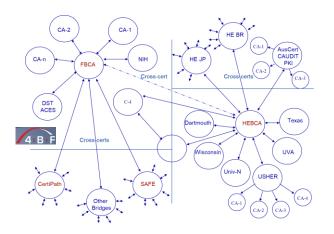
- a verifiable credential in a wallet
- ...

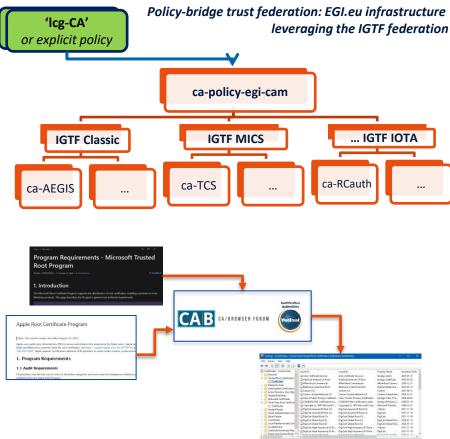
verified (on-line and also offline) at the original trusted issuer or at an independent trusted verifier



Federation: technological or policy bridge

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

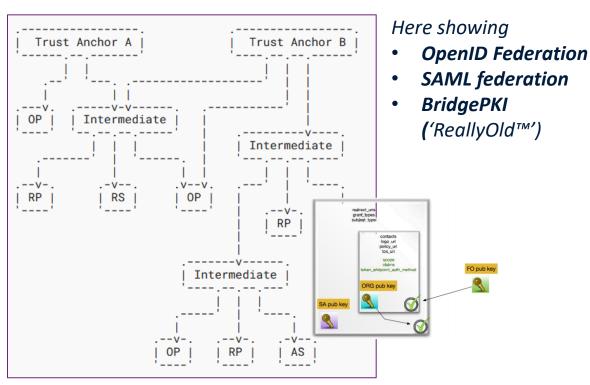




Left-hand image: 4 Bridges Forum, source: Scott Rea (then: Dartmouth University)

Images: cabforum.org, WebTrust logo: from DigiCert.com; image MS root store, https://learn.microsoft.com/en-us/security/trusted-root/program-requirements

Federation concept is technology agnostic



Metadadata Distribution Upstream Federation Metadata Your Federation Downstream eduGAIN Metadata

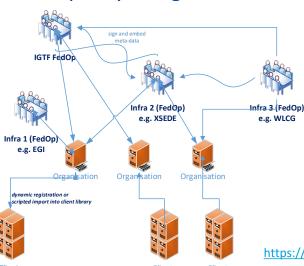
OpenID Federation signing image by Roland Hedberg; text diagram from OpenID Federation spec https://openid.net/specs/openid-federation-1_0.html

MDS meta-data flow: https://wiki.geant.org/display/eduGAIN/Metadata+Flow+in+eduGAIN

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



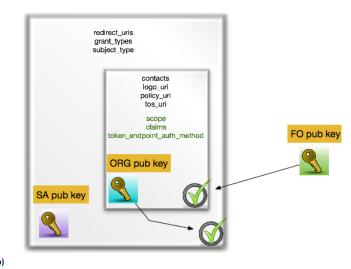
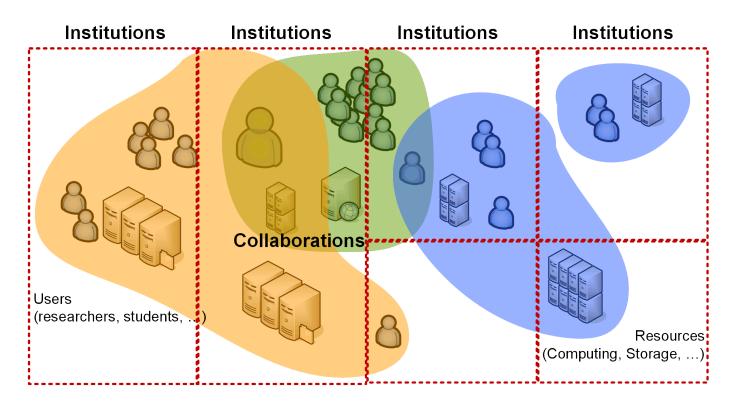


Image: Roland Hedberg, University of Umeå
OpenID Connect Fedrration:

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



'Orthogonal' sources of Authority





Multiple sources of authority: the community

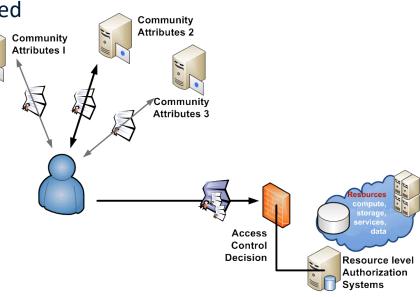
• authorization assertion providers (attribute authorities) use the identifier(s) from authentication in their membership services

source of authority for attributes is distributed

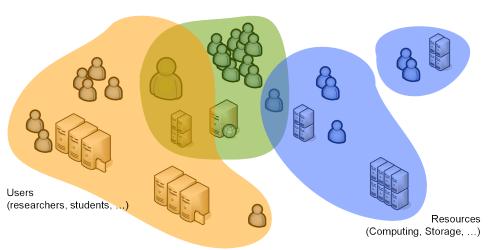
for example:

- community membership from an experiment
- affiliation status from home organisation

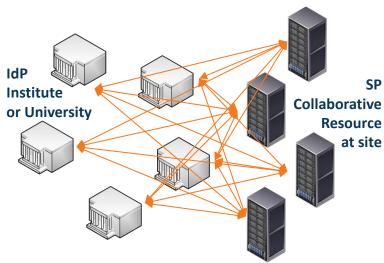
may be jointly needed to access sensitive data that is subject to medical-ethical clearance



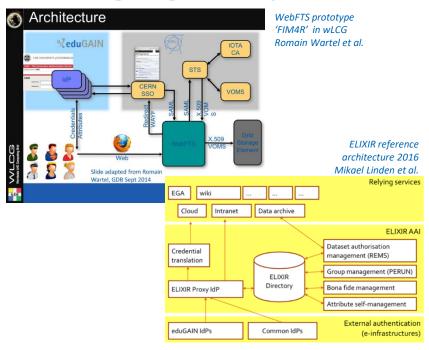
A fundamental scaling issue remained unique to research (and to joint venture business ...)



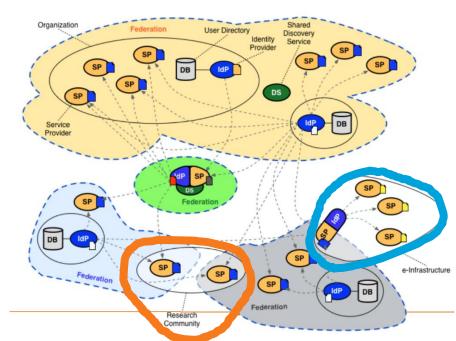
for identity and user data 'n x m' agreements remain(ed)



Managing complexities of federation & identity



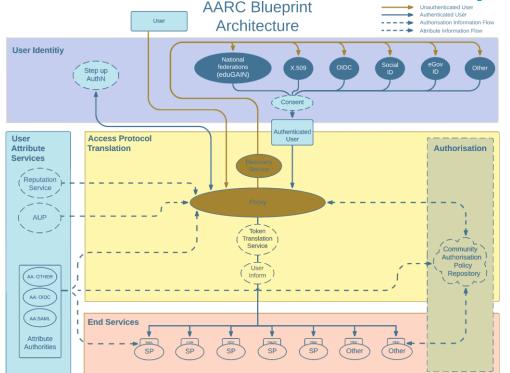
communities had either invented their own 'proxy' model to abstract complexity

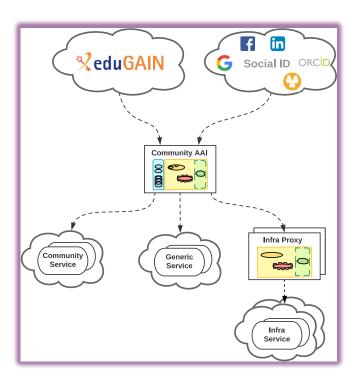


or they were composed of many services each of which had to manage federation complexity

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

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, https://doi.org/10.2777/8702

Composite AAIs: proxies beyond just the research infrastructures

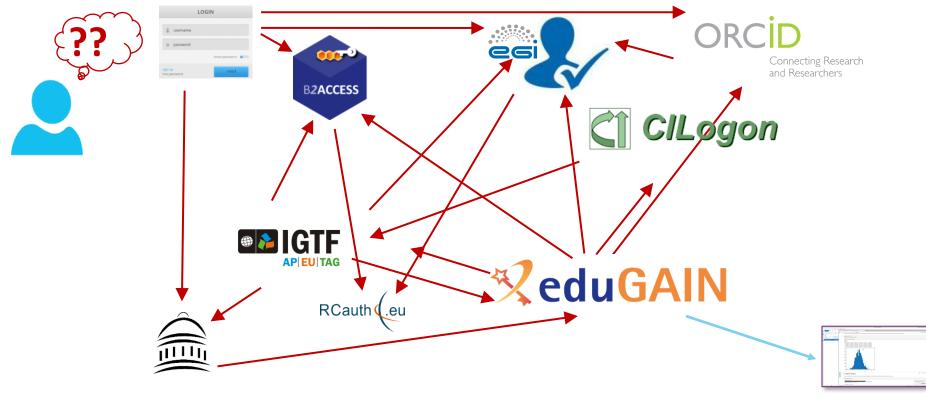
Proxy model harmonizes IdPs from many sources

- 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 carreer

₹eduGAIN Community AAI

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

But institutional identity sources may make loops ⁽²⁾

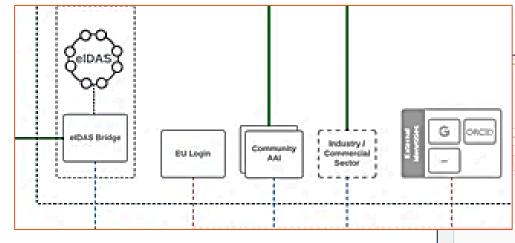




European Open Science Cloud (EOSC)

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)



EOSC – it's a federation yet again

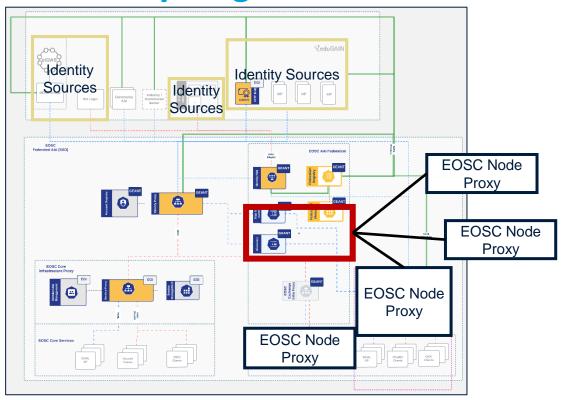


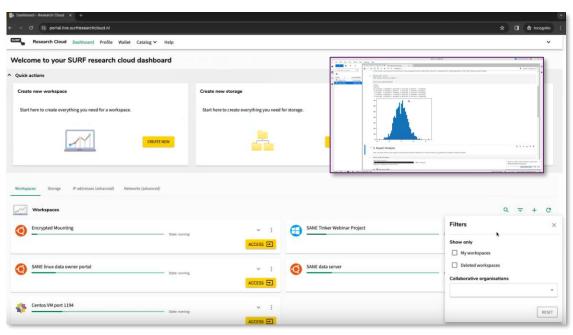
Image: EOSC AAI for the EOSC Core and Exchange Federation

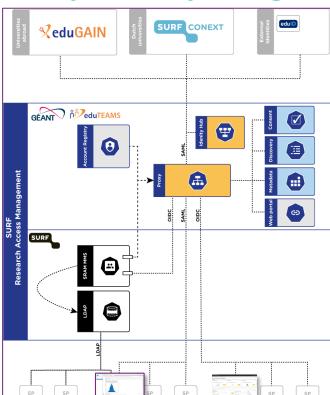
for the EOSC European Node

Christos Kanellopoulos, Nicolas Liampotis, David Groep (June 2023)



Example: SURF Research Cloud Secure Supercomputing





SURF SRAM architecture. Raoul Teeuwen et al. from

https://servicedesk.surf.nl/wiki/display/IAM/Dienstbeschrijving+SURF+Research+Access+Management

SURF Research Cloud capture: from Introduction to SANE (Secure ANalysis Environment)

webinar February 2024, by Martin Brandt et al., SURF

https://www.surf.nl/themas/onderzoeksinfrastructuur/sane-veilige-omgeving-voor-analyse-van-gevoelige-data



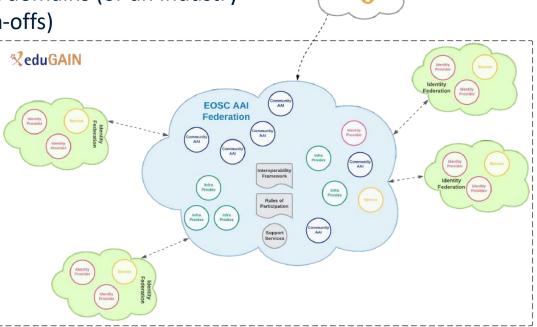
When many proxies from different groups come together

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

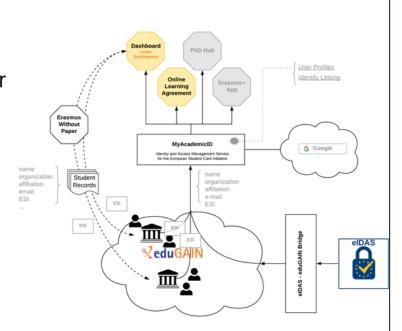


Christos Kanellopoulos (GEANT) for the EOSC AAI Federation in "The EOSC Core", https://eoscfuture.eu/wp-content/uploads/2022/04/EOSC-Core.pdf

Also building blocks for your student identity, Erasmus+, and EWP

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





The AARC Blueprint – a very digestible architecture ... so



Putting it back together again

Common patterns in scalability



Make and treat computing as the research instrument it is today

- institutionally and globally



Institutional: e.g. FSE CSLab, DSRI,



National Infrastructure SURF Snellius HPC



There are today as much part of science as detectors are to physics and: users should move seamlessly between tiers

Photos: Nikhef NDPF, DelftBlue/TUDelft, SURF Data Repository, Snellius, SURF @ DigitalRealty; EuroHPC images: EuroHPC, LUMI Consortium, Jules Verne consortium

Education labs are much like ad-hoc research as well

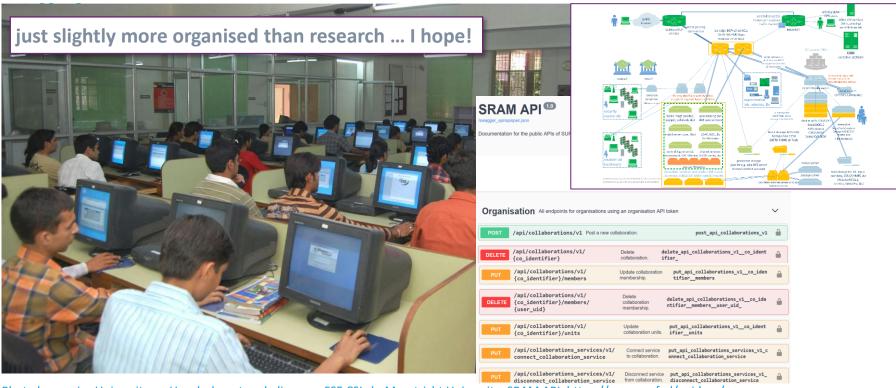
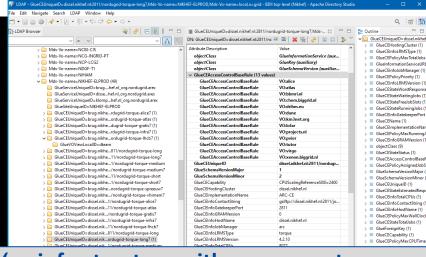


Photo by sunrise University on Unsplash; network diagram: FSE CSLab, Maastricht University; SRAM API: https://sram.surf.nl/apidocs/

A global infrastructure



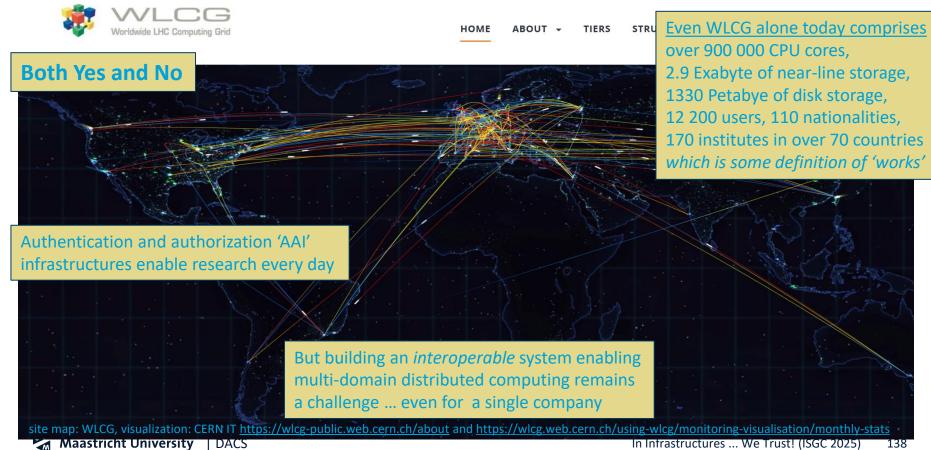


'an infrastructure with components matched to application need'

- 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 Idap://bdii03.nikhef.nl:2170/o=grid; Earth visualization: https://dashb-earth.cern.ch/, Google Earth

So: did we solve this inherently-cross-domain issue ...?



Look for the 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 $1 \rightarrow 2$
- 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

And: you can move problems around, but it's hard to actually solve them!





or enjoy the remained of today ...



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

(c) BY





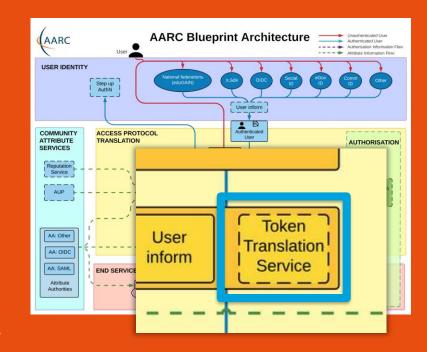






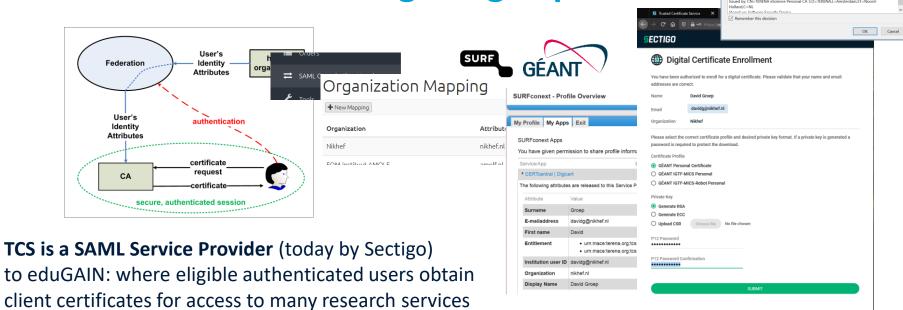
Distributed collaborative services a more technical example with RCauth.eu

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



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



www.eugridpma.org:443

Seamless in-line token translation services from

'SAML' to PKIX









Infrastructure Master Portal Credential

Store



Policy Filtering WAYF to eduGAIN

Certificate Authority

(Myproxy Server)

User Home Org or Infrastructure IdP

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

REFEDS R&S Sirtfi Trust

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

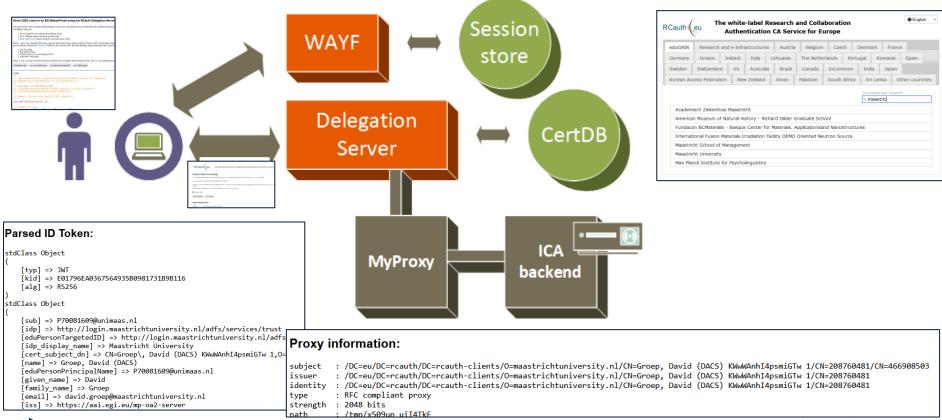
Plus the (domain-name based) schacHomeOrganisation and a 'representation of the real name'

/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



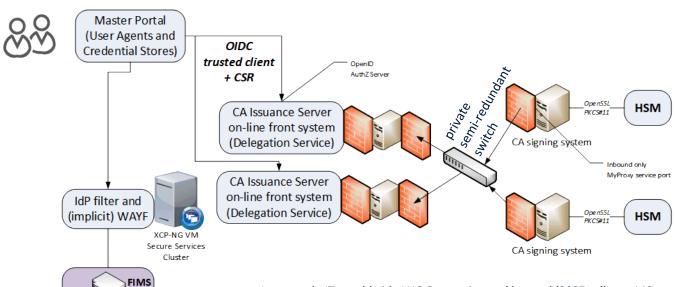
With a single, yet fully compliant, 'Heath Robinson' CA



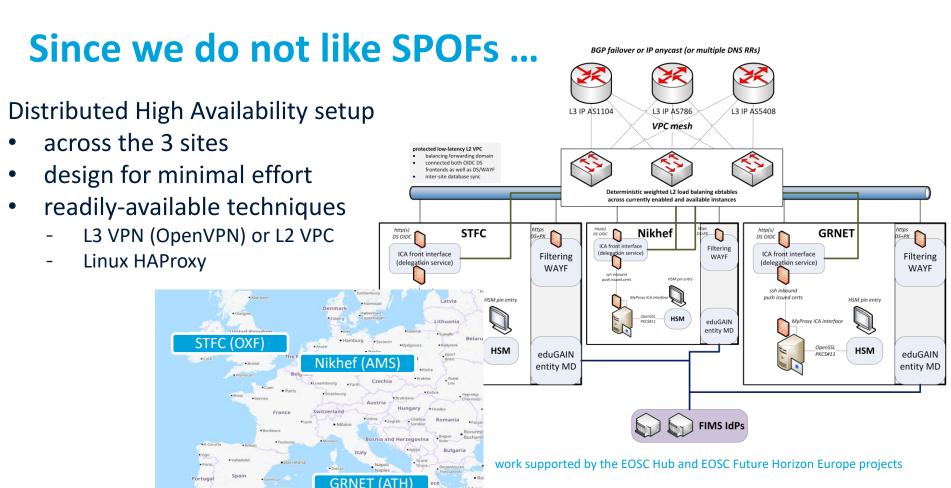
A single-site locally-highly-available RCauth at Nikhef Amsterdam

- 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'







Maastricht Univ

A transparent multi-site setup is needed for the user

User

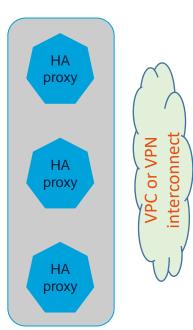
- connects to HA proxy at {wayf,pilot-ica-g1}.rcauth.eu
- HA proxy sends users to "closest" working service
- primarily forward to its own DS when available



Straightforward proven solution is IP anycast

wherever the user is, the service is at

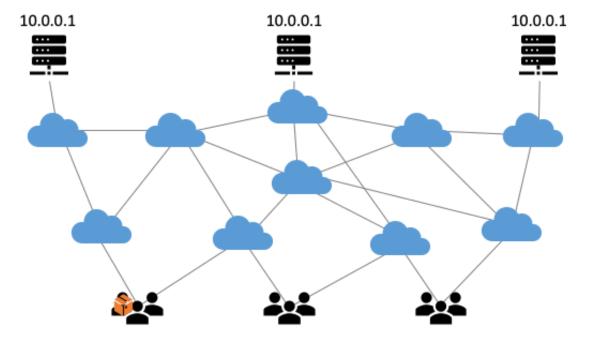
- 2a07:8504:01a0::1
- or for legacy IP users at 145.116.216.1



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

Anycast: when the same place exists many times



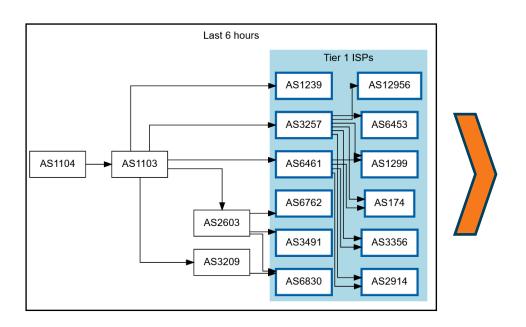
So we used

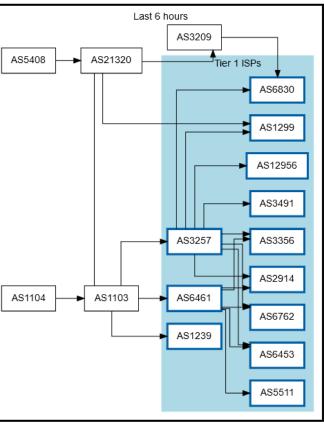
- 3 (for now: 2) sites
- one VM at each site exposing 2a07:8504:01a0::1
- smallest v6 subnet (/48)
- bird + a service probe
- each site's own ASN
- some IRR DB editing
- IPv4 is similar, with a /24

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



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

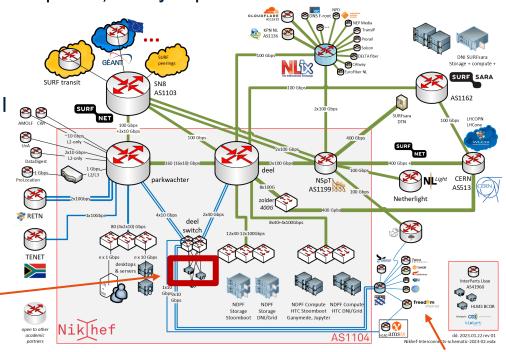
Shortest path, also when mixing with the default-free zone

[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. reauth. eu (145. 116. 216. 1) [AS786/AS5408/AS1104]

rcauth.eu HA proxy

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



at home

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')

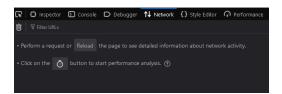
Qualified users are all those in eduGAIN with basic assurance (Sirtfi version 1 + Research & Scholarship entity categories), and everyone in a Dutch SURF 'Annex IX' institution – such as UM

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



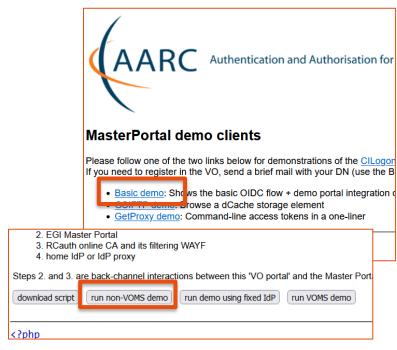
RCauth do-it-yourself demo

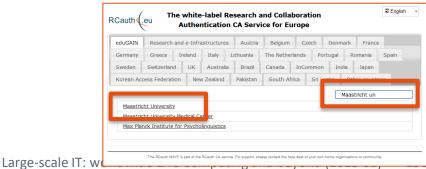
- Go to https://rcdemo.nikhef.nl/
- select "Basic Demo"
- Enable browser Inspector (F12) on the network tab,
 (and start the SAML tracer extension if you have it)



Run the "non-VOMS demo"

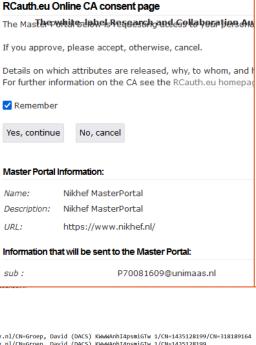
From eduGAIN, select "Maastricht University"





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?



aBMALkQAo5smbNx+PW7fOoNbfzReSJfGt7DaYqekEO/yvvxH0O08xf20w+rmPCEA

Not After: Nov 3 00:09:47 2024 GMT

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/X599up_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
```

Subject: DC=eu, DC=rcauth, DC=rcauth-clients, O=maastrichtuniversity.nl, CN=Groep, David (DACS) KWWNAnhI4psmiGTw 1, CN=1435128199, CN=

RSA Crypto

Just in case ... you cannot factor '55'



Establishing trust at a distance

Remote trust needs cryptography in some way

Client authentication

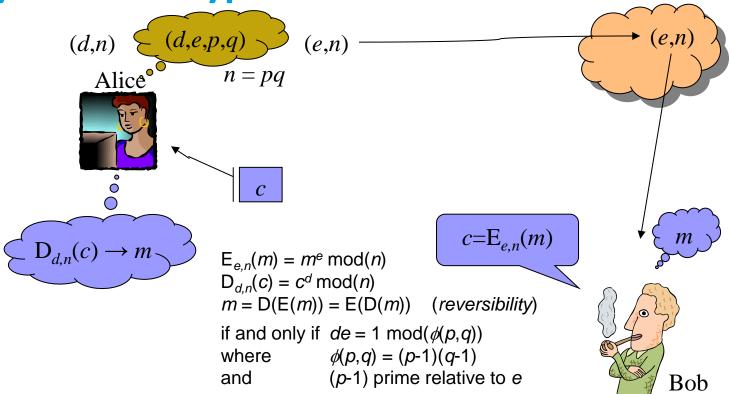
- pre-shared secrets, may be salted hashed on service side
- required: secure one-way hash function
- need a **protected channel** between identifiable end-points

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?



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)

```
\begin{aligned} & \mathsf{E}_{e,n}(m) = m^e \, \mathsf{mod}(n) \\ & \mathsf{D}_{d,n}(c) = c^d \, \mathsf{mod}(n) \\ & m = \mathsf{D}(\mathsf{E}(m)) = \mathsf{E}(\mathsf{D}(m)) \\ & \mathsf{if a.o. if} \quad de = 1 \, \mathsf{mod}(\phi(p,q)) \\ & \mathsf{where} \quad \phi(p,q) = (p\text{-}1)(q\text{-}1) \end{aligned}
```

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)
- $E_{27:55}(2) = 2^{27} \mod(55) = 18!$



$$E_{e,n}(m) = m^e \mod(n)$$

$$D_{d,n}(c) = c^d \mod(n)$$

$$m = D(E(m)) = E(D(m))$$
if a.o. if $de = 1 \mod(\phi(p,q))$
where $\phi(p,q) = (p-1)(q-1)$

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

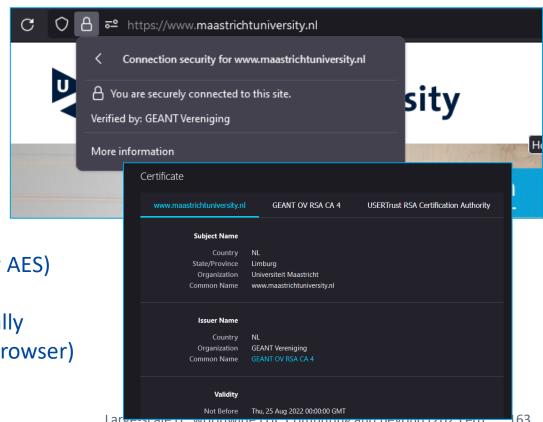
but also: the maximum plaintext is limited by the modulus length

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

Following slides are provided here merely as generic background

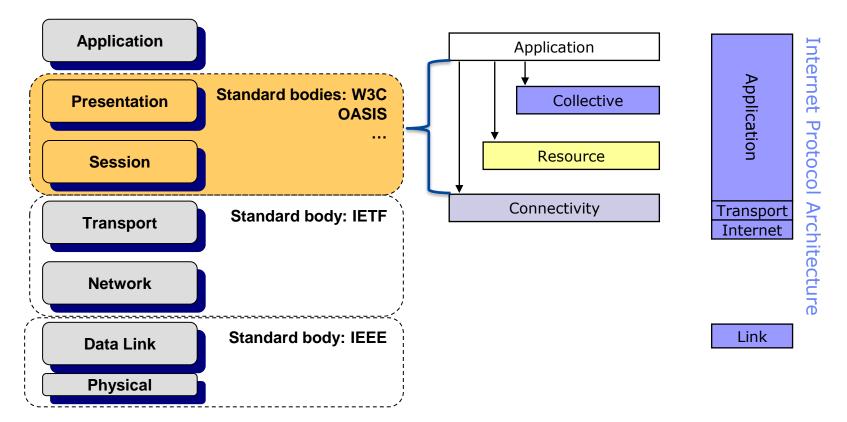


Open Systems Interconnection model (OSI model)

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	<u>Session</u>	Managing communication sessions, i.e., continuous exchange of information in the form of multiple back-and-forth transmissions between two nodes
	4	<u>Transport</u>	Reliable transmission of data segments between points on a network
Media layers	3	<u>Network</u>	Addressing, routing and traffic control
	2	<u>Data link</u>	Transmission of data frames between two nodes connected by a physical layer
	1	<u>Physical</u>	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

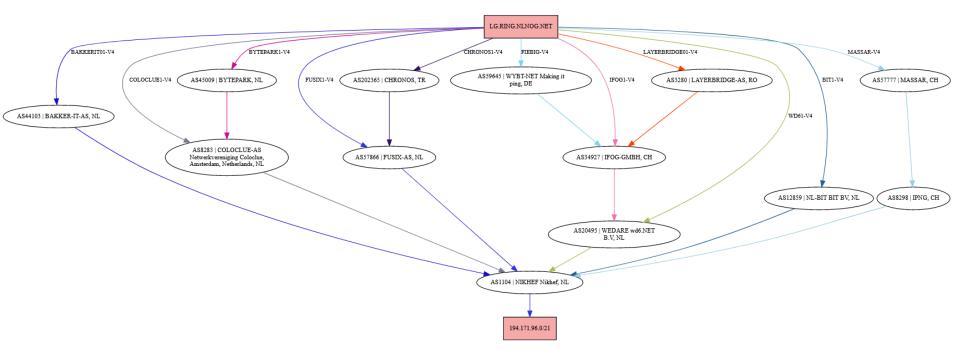
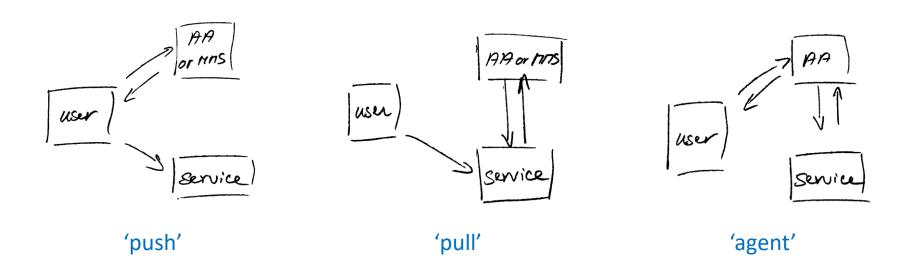


Image sources: NLNOG RING map https://lg.ring.nlnog.net/

RFC2904 authorization models: three AuthZ flows



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

OAuth2 & JWTs: assertions can be quite detailed

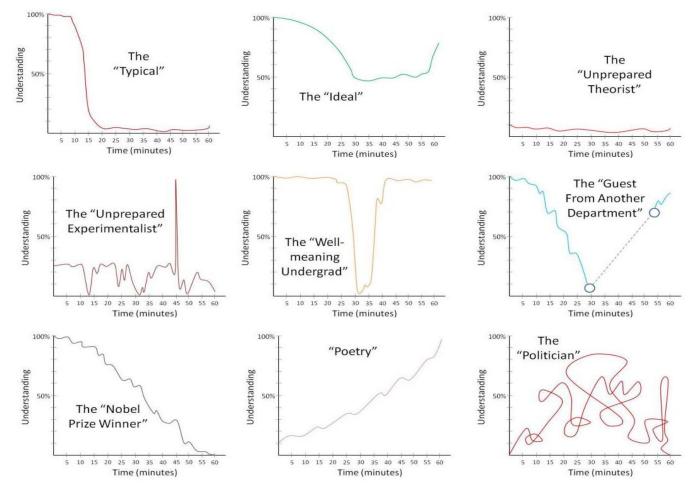
```
$ echo $AT | jwt
* Payload
  "wlcg.ver": "1.0",
  "sub": "a1b98335-9649-4fb0-961d-5a49ce108d49",
  "aud": "https://wlcg.cern.ch/jwt/v1/any",
  "nbf": 1593004542,
  "scope": "storage.read:/ storage.modify:/",
  "iss": "https://wlcg.cloud.cnaf.infn.it/",
  "exp": 1593008142,
  "iat": 1593004542,
  "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



N W

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

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."

