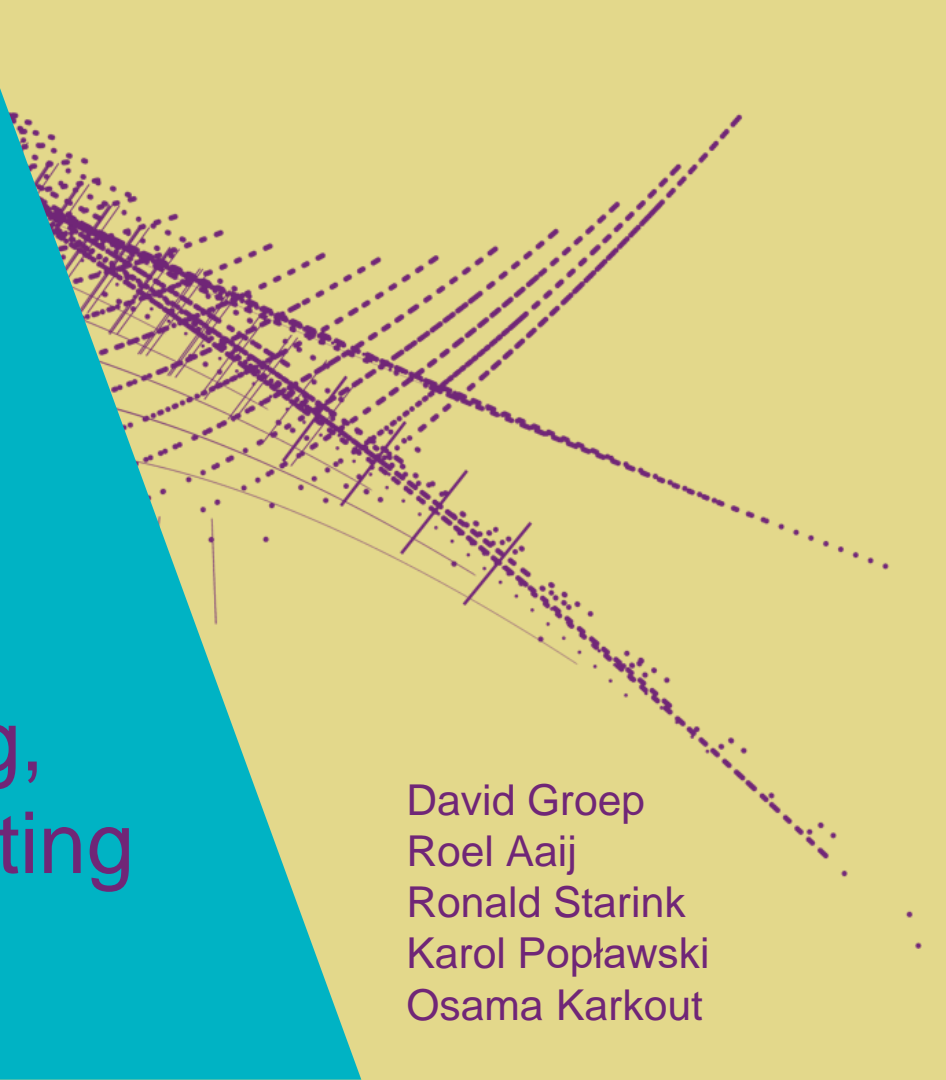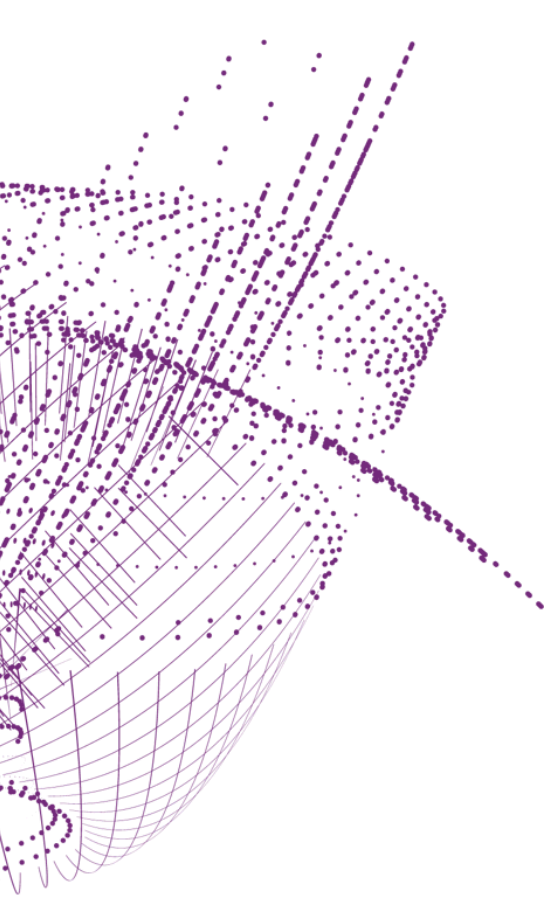Nikhef SEP panel visit November 2023

# Physics Data Processing, Engineering and Computing
*accelerating 'time to results'*
*through computing and collaboration*

David Groep
Roel Aaij
Ronald Starink
Karol Popławski
Osama Karkout

David Groep | PL Physics Data Processing

# Physics Data Processing
# & 'CT-PDP' engineering

Nikhef

# PDP: Computing as research and instrumentation
## validated through our real-life applications
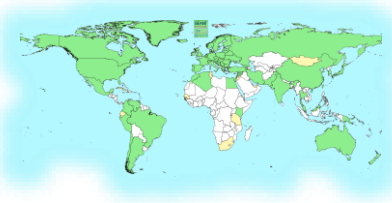
Physics Data Processing programme lines

1. **infrastructure, network & systems research**
   - **building** 'research IT facilities' through co-design & development
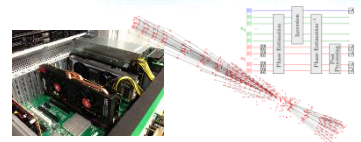   - big data science innovation: research **next gen IT infrastructure**

2. **infrastructure for trusted collaboration**
   - **trust and identity** for enabling communities
   - managing complexity of **collaboration mechanisms**
   - **securing** the infrastructure of our **open science** cloud

3. **algorithmic design patterns** - *more in Roel's introduction*
   - **GPU accelerated** computing, **Quantum Computing**, AI and **Machine Learning**

# Physics Data Processing and CT-PDP engineering effort

2.5 staff (**David Groep**, **Roel Aaij**, Jeff Templon)

1 postdoc (Maarten van Veghel)

~ 11 engineers: DevOps, research software, RDM, Collaboration / Trust & Identity
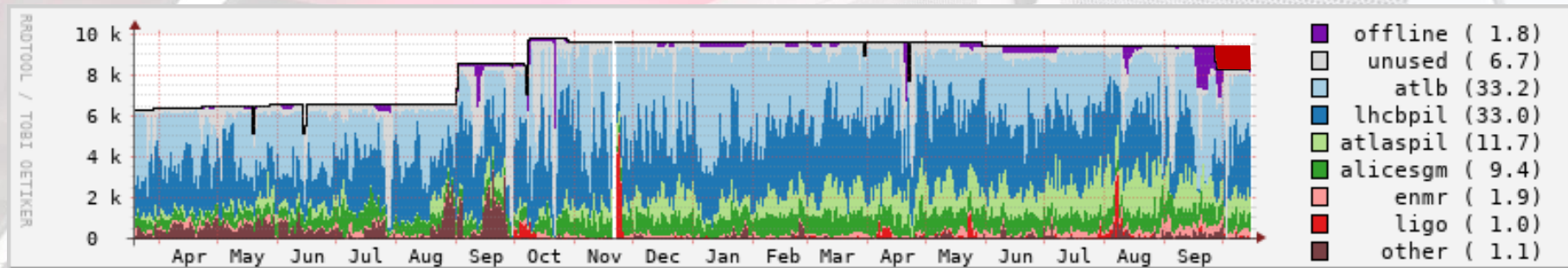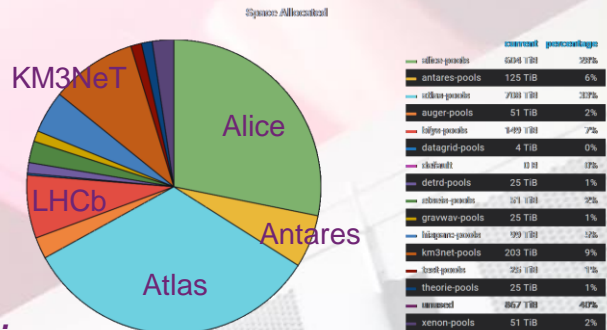*with organic embedding in the Computer Technology engineering group*

With wide range of activity leads, including for example **LHCb RTA Reco convener**, **SURF innovation expert group chair**, **EOSC** Security Coordinator, **AEGIS** Trust and Identity policy lead, **Interoperable Global Trust Federation** chair, Dutch **National Infrastructure Executive**, lead for the **Thematic DCC Natural and Engineering Sciences**, and members in board & committees for the (global) e-Infrastructure landscape: GEANT GCC, PC-GWI, CieDO…

# Infrastructure for Research

High-Through Compute (HTC) + HT Storage

- **National e-Infrastructure** coordinated by SURF
  *LHC NL-T1, IGWN, KM3NeT, Xenon, DUNE, WeNMR, MinE ...*

- **'Stoomboot' local analysis facility** + IGWN cluster & 'submit node'

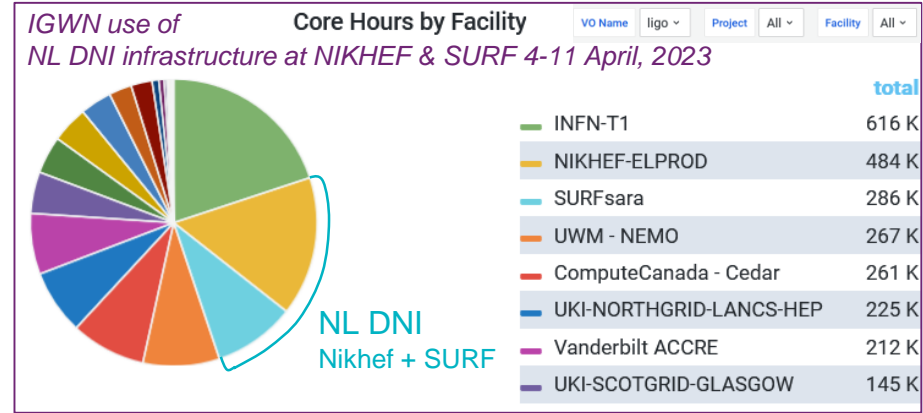~ 12 000 cores (total), 13 PByte storage – very competitive *w.r.t.* commercial cloud



Occupancy: NDPF DNI processing facility in the period March 2021 .. October 2022. Top-right: storage capacity allocated in DNI Nikhef segment

# Common solutions are essential for our 'national' facility

**Alignment of common e-Infrastructure**
and shared use by experiments
(LHC, GW, KM3NeT, DUNE, Xenon, ...)

- *common solutions,* since bespoke systems
  for each experiment do not scale for Nikhef (or NL)

- efficient sharing of both hardware and DevOps effort

- synergy with other domains helps sustainable funding



*IGWN use of NL DNI infrastructure at NIKHEF & SURF 4-11 April, 2023*

**Core Hours by Facility**

VO Name `ligo` ⌄  Project `All` ⌄  Facility `All` ⌄

| | total |
|---|---|
| INFN-T1 | 616 K |
| NIKHEF-ELPROD | 484 K |
| SURFsara | 286 K |
| UWM - NEMO | 267 K |
| ComputeCanada - Cedar | 261 K |
| UKI-NORTHGRID-LANCS-HEP | 225 K |
| Vanderbilt ACCRE | 212 K |
| UKI-SCOTGRID-GLASGOW | 145 K |

NL DNI
Nikhef + SURF

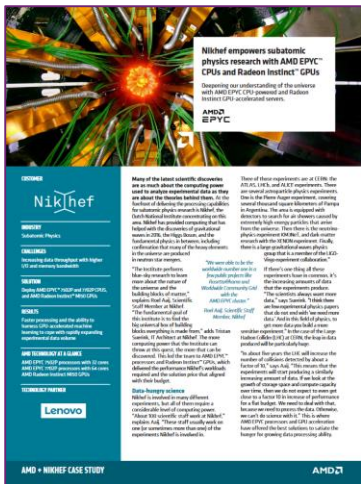Continuation of our long-term strategy - from EU DataGrid in 2000 onwards:

- drives collaborative efforts we work with for identity management and **common protocols globally**:
  ACCESS-CI and CILogon (US) – key players for LIGO, DUNE, and US-ATLAS
- common **processing framework** development (e.g. together with KM3NeT in its INFRADEV project)

Data on Tue April 11th 2023 from the OSG accounting for the LIGO VO for past week (with also SURFsara fully in production)
https://gracc.opensciencegrid.org/d/9u1-Q3vVz/cpu-payload-jobs?orgId=1&var-ReportableVOName=ligo&var-Project=All&var-Facility=All&var-Probe=All&var-interval=1d&from=1680566400000&to=1681257600000

# Innovation *on* infrastructure

- Network-to-systems integration
- Storage throughput & parallelism
- Systems integration design and tuning



798.49 Gb/s



FUNGIBLE

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

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



- early **engineering engagement with vendors** to build us suitable systems
- **co-design** of our national HPC systems ('Snellius')
- **data-intensive compute** with DPUs, or on-NIC FPGAs?
- networks > 800Gbps, >1 Bpps (today: 400G to CERN)

Image: Minister of Economic Affairs M. Adriaansens launched the Innovation Hub with Nikhef, SURF, Nokia and NL-ix, January 2023. Composite image from https://www.surf.nl/nieuws/minister-adriaansens-lanceert-testomgeving-voor-supersnelle-netwerktechnologie; Bluefield Hackathon by Nvidia/Mellanox; *abbreviations*: **DPU**: Data Processing Unit; **on-NIC FPGAs:** on-network interface card field programmable gate arrays; **pps**: packets per second

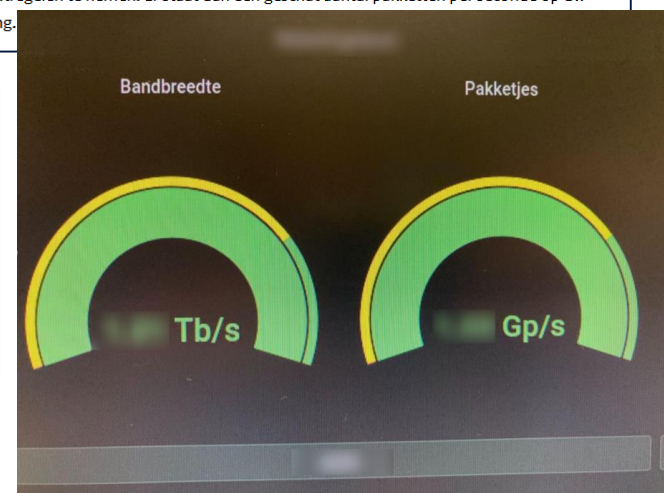# Our science data flows are somebody else's DDoS attack



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

# Infrastructure for Collaboration

Target impactful areas in architectures for 'AAI' & 'OpSec'
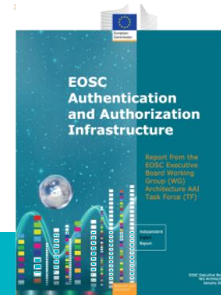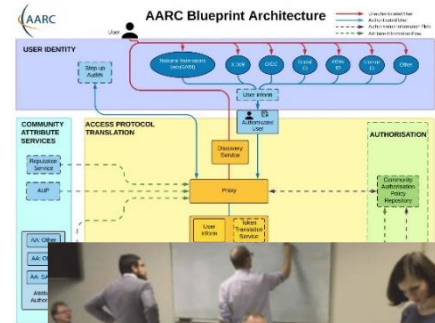


- **authentication & authorization for research collaboration**
  - **AARC** project & community: GEANT Framework projects, R&E federation, identity and credentialing services, EOSC Future, …
  - recently awarded: **EOSC Core** security and **AARC-TREE**
  - **policy frameworks for interoperability** for data protection and global seamless service access
  - continuous **technical evolution** driving IGWN, WLCG in line with AARC and global AAI architecture

- embedding data **processing needs** of our experiments in the (EOSC) landscape
  - **EOSC** Interoperability Framework: Security Baseline, AAI Architecture
  - **eduGAIN** Operational Security for the global R&E inter-federation service
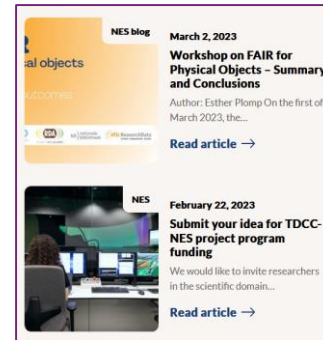  - **EGI** Advanced computing for research federation**, GEANT community**

AAI: Authentication and Authorization Infrastructure; **OpSec**: Operational Security (incident response); **AARC**: Authentication and Authorisation for Research Collaboration community and projects**; EOSC**: European Open Science Cloud; **EGI** and **GEANT**: pan-European e-Infrastructure and network collaborations; **IGWN**: International Gravitational Waves Observatory Network;

# Collaboration: Research Data Management beyond 'FA'

> FAIR for **live data**, in large volumes, from 'FA' towards the 'I' and 'R'

- not *that* many disciplines with really **voluminous data**
  - so nationally join forces with those who do: ASTRON (SRCnet), KNMI (earth observation, seismology), *&c*

- work with those who care about ***software*** to bring data to life on ***infrastructure***
  - NLeSC, 4TU.RD/TUDelft, CWI, and with those who ensure the *infrastructure*: SURF

- for our own analyses and the local (R&D) experiments we work towards *continuous deposition* of **re-usable** data and software:
  - **Thematic Digital Competence Centre** for the Natural and Engineering Sciences
  - **co-develop Djehuty RDM** repository software link 'Stoomboot' analysis cluster storage



**NES blog** March 2, 2023
**Workshop on FAIR for Physical Objects – Summary and Conclusions**
Author: Esther Plomp On the first of March 2023, the...
Read article →

**NES** February 22, 2023
**Submit your idea for TDCC-NES project program funding**
We would like to invite researchers in the scientific domain...
Read article →

# PDP strategic projects mechanism

Join initiatives and projects that

- *strengthen* the strategic areas
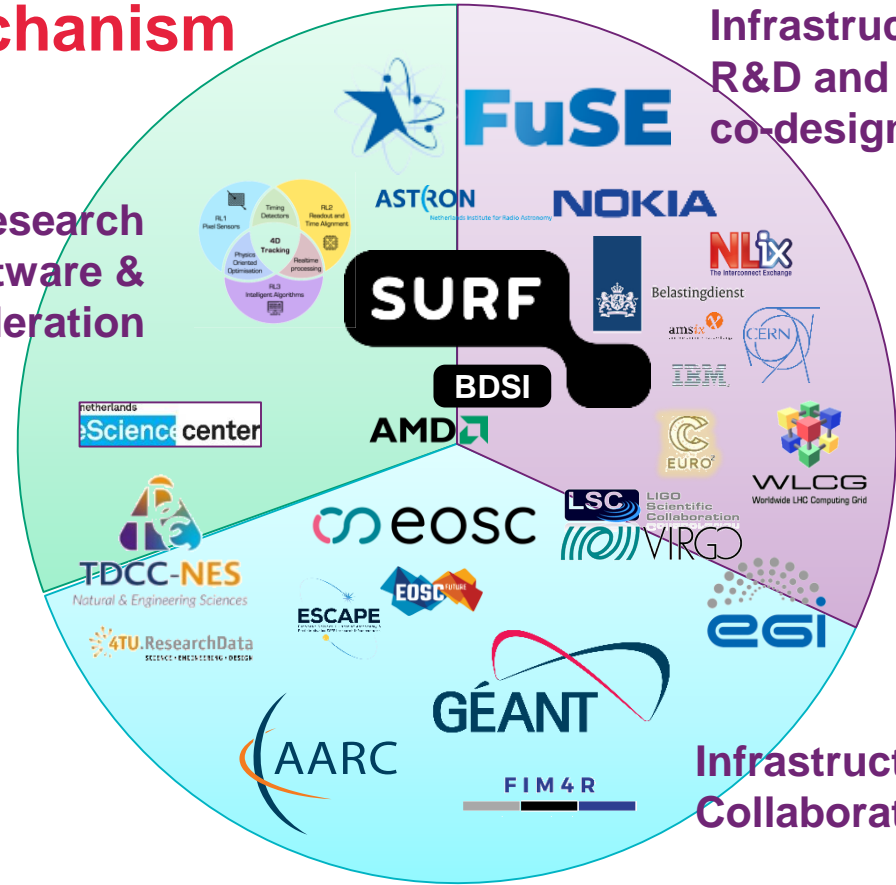- ensure *continuity* of research and infrastructure

*project pathways include*
SURF innovation, GN5-*, AARC-TREE, EOSC Core, LHC4D (planned), …

**Public partner R&D engagement**
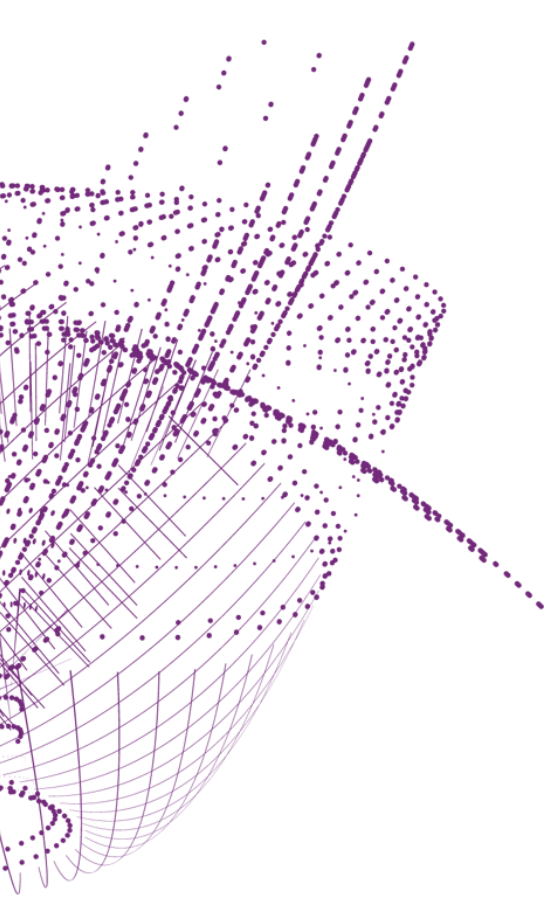AMD, Nokia, Nvidia/MLNX, NL-ix, …
Dutch national government

# Sustained infrastructure for advanced computing?

Data processing: a **persistent need** for all our experiments. And 'we' are not alone!

- many disciplines: ESCAPE (our experiments plus astronomy), bio-informatics, health, SSH, …

  need infrastructure to exploit the collected data with **long-term ICT capabilities**

- project-based funding for 'upgrades' of infrastructure not the appropriate way of funding **persistent requirements** … but only thing we have at the moment:



vl·e

**BiG** *Grid*
the dutch e-science grid

SURF

FuSE …?

2003    2007    2011    2015    2020

Roel Aaij | Senior Scientific Researcher PDP

# Algorithms and Acceleration

# Efficient and Scalable Computing For HEP

- CPU-only is not going to be affordable;
  In other words: compute accelerators offer
  more physics for less money

- Most mature compute accelerators: GPUs

- How to program?

- How to optimize at algorithm and system level?

- How to integrate in frameworks and infrastructure?

- How to do this efficiently in terms of people's time?

- How to maintain?

- What will future (compute) accelerators look like?

- How to provide career perspective to people who focus on (accelerated) software?



empirical GPU FLOP/s per dollar

- Our data (2x every 2.46 years)
- Moore's law slope (2x every 2.00 years)
- Huang's law slope (2x every 1.08 years)
- Bio anchors report slope (2x every 2.50 years)
- empirical CPU slope (2x every 2.32 years)
- Top FLOPs/dollar GPUs (2x every 2.95 years)
- ML GPUs (2x every 2.07 years)

# Efficient and Scalable Computing For HEP

- LHCb's first-stage GPU trigger (Allen)
  - Bespoke application with all-custom kernels
  - 4 TB/s of detector data on 400 GPUs
  - Nominal luminosity next year
  - Focused on integration (DAQ, LHCb stack, etc.)
- - From idea to R&D to production in 5 years
- With NLeSC: fast ML inference
  - Using standard format (ONNX) and libraries
- User (software) support for GPUs (AMD and NVIDIA)
- FASTER: computing for HL-LHC & '4D' reconstruction

https://github.com/LHC-NLeSC/run-allen-run

Ronald Starink | TGL Computing Technology

# Computing group at Nikhef

# CT organisation

## CT-B: system administration & service desk

- General support for ICT, end user support
- Staff: 7

## CT-PDP: dedicated support for the PDP programme

- Software & infrastructure innovation
- Staff: 11
- (→ presentations DG and RA)

## CT-PO: support for projects by experiments

- Software engineering: slow controls, data acquisition, analysis framework support
- Staff: 6
- (→ presentation KP)

# Challenges

- General: recruitment in a competitive labour market
- CT-B: balance flexibility ↔ standardization
- CT-PO: matching resource with experiments' needs (time, expertise)

Karol Popławski | Software Engineer CT-PO

# Project engineering & support

Nik|hef

# KP – Software, Controls and more

| | | |
|---|---|---|
| **ATLAS** | Barrel Alignment | DAQ of ~5800 channels |
| | MDT DCS Module | 1200 chambers: T, B-field, electronics monitoring & configuration |
| **LHCb** | SciFi tracker | FEE configuration with DB, granularity up to 1.5mln thresholds |
| **Nikhef** | Scrum Master | MT, ET, R&D |

# CT-PO

| | | |
|---|---|---|
| Henk | FELIX | |
| | MDM firmware | |
| Ton | SciFi FEE calibration | |
| | PTOLEMY | |
| Kostis | Future 4D tracking with White Rabbit | |

Osama Karkout | happy ATLAS & stoomboot user

# How I use stoomboot

Nikhef

# Data Analysis

designed for data analysis:
economic & efficient

Hi Osama,

How's it going? Could i bug you with sth?

As part of the unblinding approval checks we were asked to run toy studies. A while ago I wrote a code who can create these toys, but I don't have the infrastructure at CERN to run it myself.

Would you have some time to get this stuff running at Nikhef?

If we want fast results i can also make 5 ws and run 1k toys each

Brian, 18:53

mhm, but 5ws means effectively 5 x 70 = 350 jobs that run for > 10h each. Will this be ok or will you get death threats from Jeff Templon? ☺

I can also run the hadhad btw that's not an issue

apparently 350 jobs for 10 hours is not much at all

Brian, 11:29

cool, in that case it would actually be better if you could run the fully combined ☺

- **/project/atlas (3Tb total space, 4Tb since this week)**
  - Backed up daily: reliable, but expensive storage
  - NFS disk: slow (max speed 30% of network speed)
  - *Usage: only for code and sensitive information, not for bulk ntuples*

- **/data/atlas (40Tb total space)**
  - No backup: cheap larger-volume storage
  - NFS disk: slow (max speed 30% of network speed)
  - *Usage: Intended for bulk data that is not intensively analyzed*

- **/dcache/atlas (350Tb total space (? TBC)**
  - High Performance File system (masquarading as network file system)
  - Only accessible from stoomboot and Tier-1 computing facilities
  - Files can *not* be modified once written (but can be deleted, and recreated)
  - Not suitable, nor efficient for small files
  - *Usage: for storage of and intensive usages of ntuples, dAODs etc*

# Neural Network training: ATLAS GNNs and Transformers

Arjen van Rijn
David Groep

# Data centre visit and networks

# Our datacenter is all about connectivity



#7 on the worldwide peering list

# Summary data of our datacentre

Maximum power: 2,7 MW

Power Usage Efficiency (PUE): 1,3

Residual heat used for:
- Nikhef building
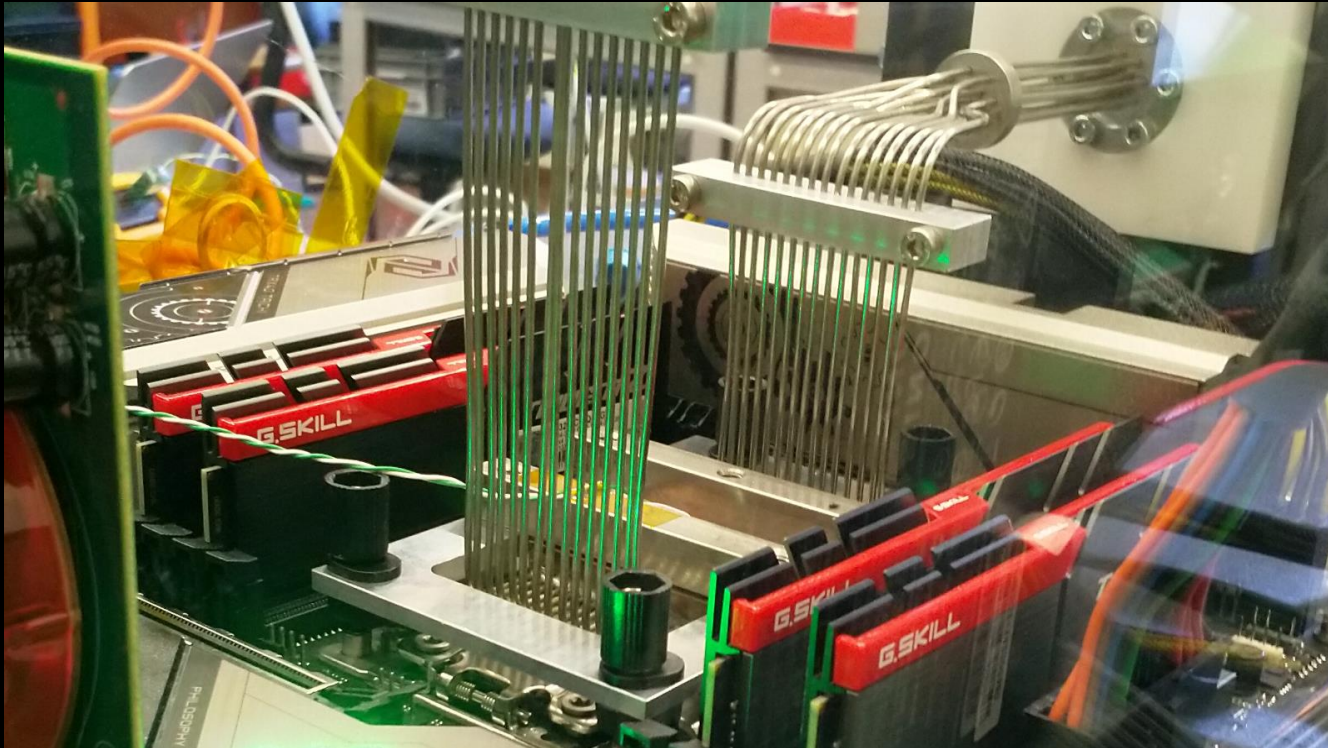- Amsterdam University College
- Student housing



| Room | Purpose | # racks max | IT power average per rack (kW) | IT power total (kW) | IT power currently (kW) |
|------|---------|-------------|-------------------------------|---------------------|-------------------------|
| H234b | Scientific computing | 47 | 6,25 | 300 | 180 |
| H140 | Nikhef Housing | 282 | 1,9 | 660 | 528 |
| H142 | Nikhef Housing (extension) | 112 | 4,0 | 540 | 10 |

Financial figures (annually):

- Turnover currently ~5 M€; will grow to about 6,5 M€;
- Running costs will grow to 3 – 3,5 M€ (energy costs!);
- Depreciation of extension: 1 M€ annually until end 2028;
- Net result: 2 – 3 M€

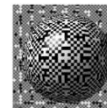# Because we can … does not mean it's the scalable way ☺



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

Supplementary materials

# Backup – Balign



ATLAS Barrel Alignment

*Hardware Setup (1)*

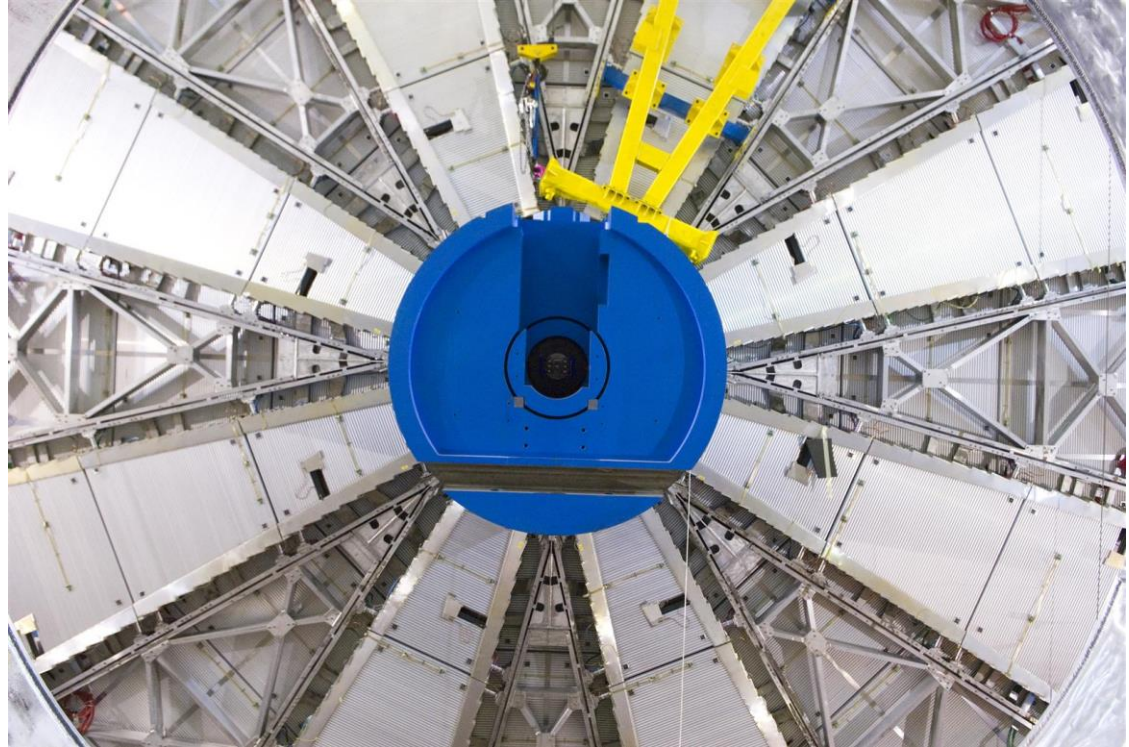Cavern — RasCam, ± 5800, RasLed → RasMux ± 600 → MasterMux 48 → USA15 TopMux 8

Source: Future Barrel Alignment by Robert Hart

# Backup – MDM

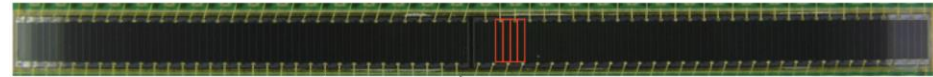1,171 chambers with total 354,240 tubes (3 cm diameter, 0.85-6.5 m long)

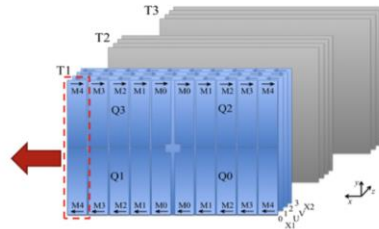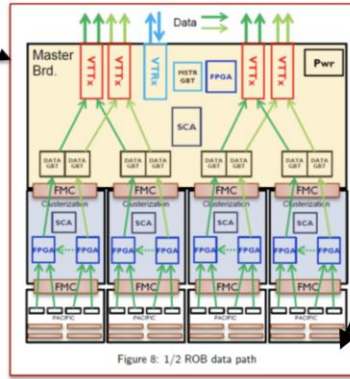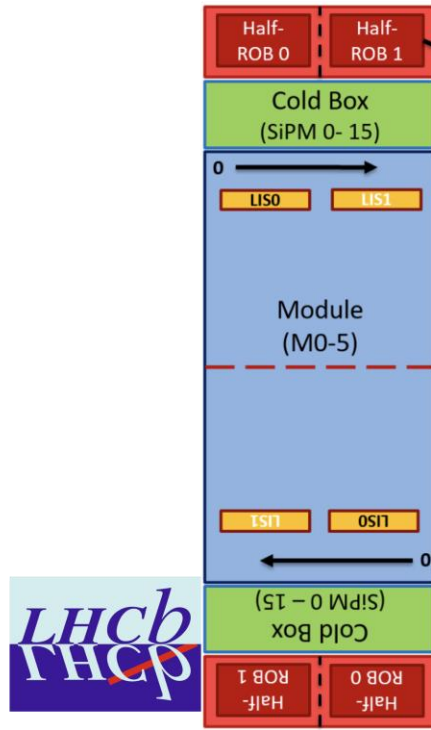Tube resolution 80 μm



Source: https://atlas.cern/Discover/Detector/Muon-Spectrometer

# Backup – SciFi

## SciFi - Channels



128-channel SiPM (Hamamatsu)

512 HalfROBs:

- 512 MBs

- 2048 CBs

- 2048 PBs

- 8192 ASIC chips on PBs

- e.g. 524288 SiPM channels with 3 thresholds

Source: SciFi Collaboration

Figure 8: 1/2 ROB data path