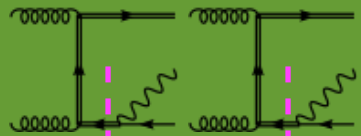FOM

pdp

# Nikhef – a Journey in Physics and Data Processing

David Groep
Nikhef
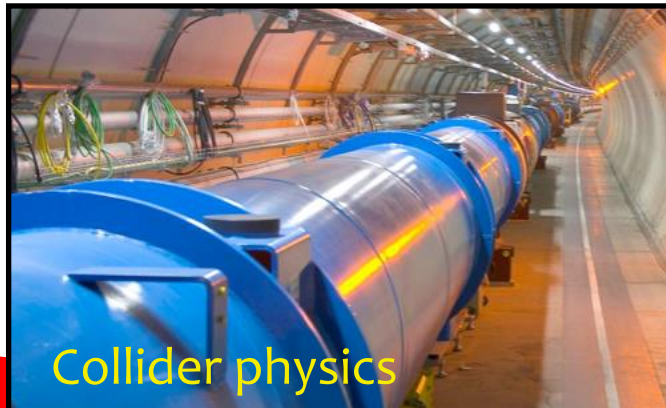*PDP - Advanced Computing for Research*

*Verleggen van de grenzen van onze kennis*

- **Accelerator-based particle physics**
  Experiments studying interactions in particle collision processes at particle accelerators, in particular at CERN;
- **Astroparticle physics**
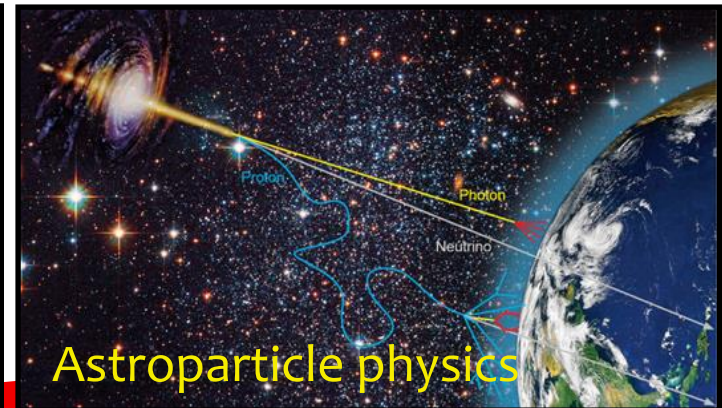  Experiments studying interactions of particles and radiation emanating from the Universe.

$$d\sigma^{(2)} + \sum_{\alpha\beta} \int \frac{dx_1 dx_2}{2x_1 x_2 S} \mathcal{L}_{\alpha\beta} \left( \hat{S}_{\alpha\beta} + \mathcal{I}_{\alpha\beta} + \mathcal{D}_{\alpha\beta} \cdot \right.$$
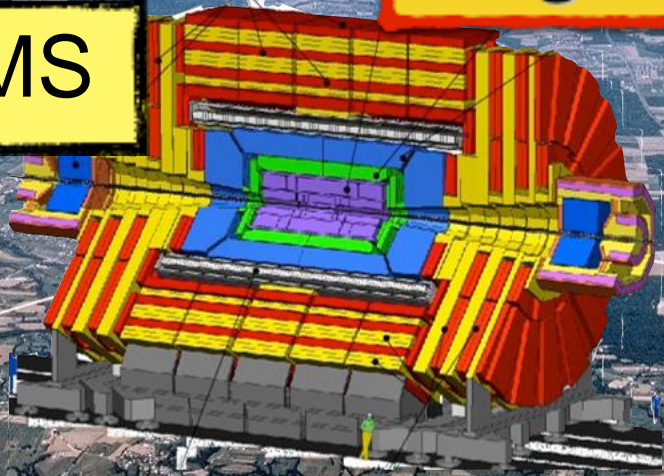
Phenomenology

Collider physics

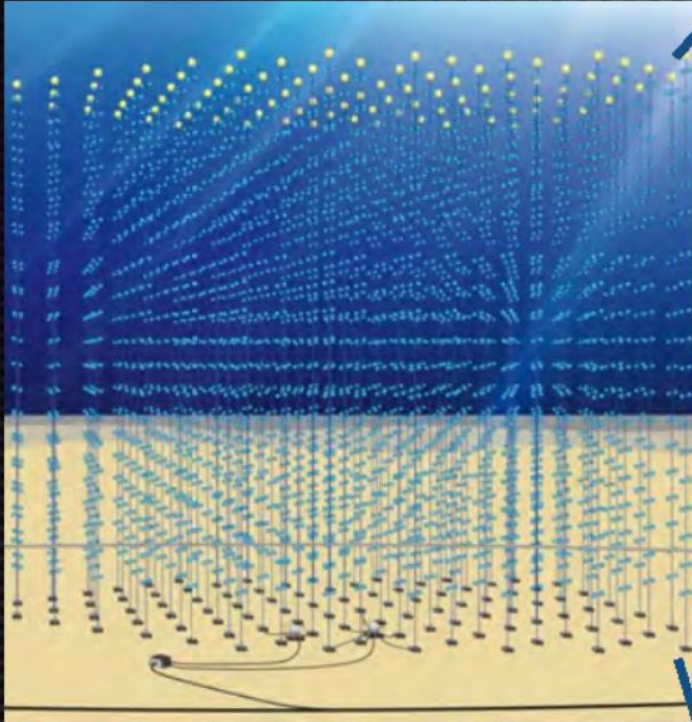Astroparticle physics

Large Hadron Collider

CMS

Nikhef

LHCb

ALICE
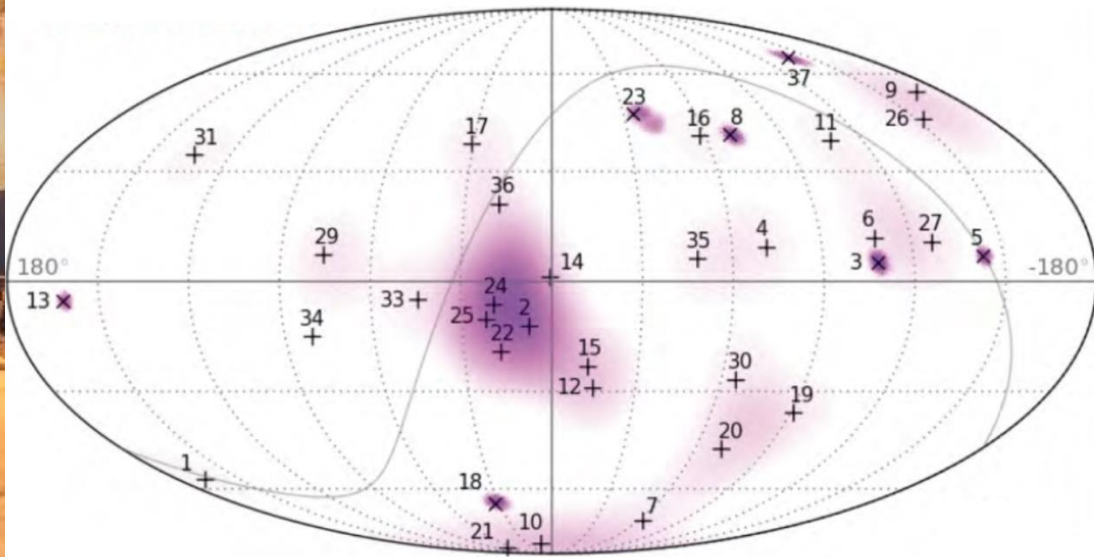
ATLAS

Imagery: CERN

# Nikhefs neutrino-detector: KM3NeT

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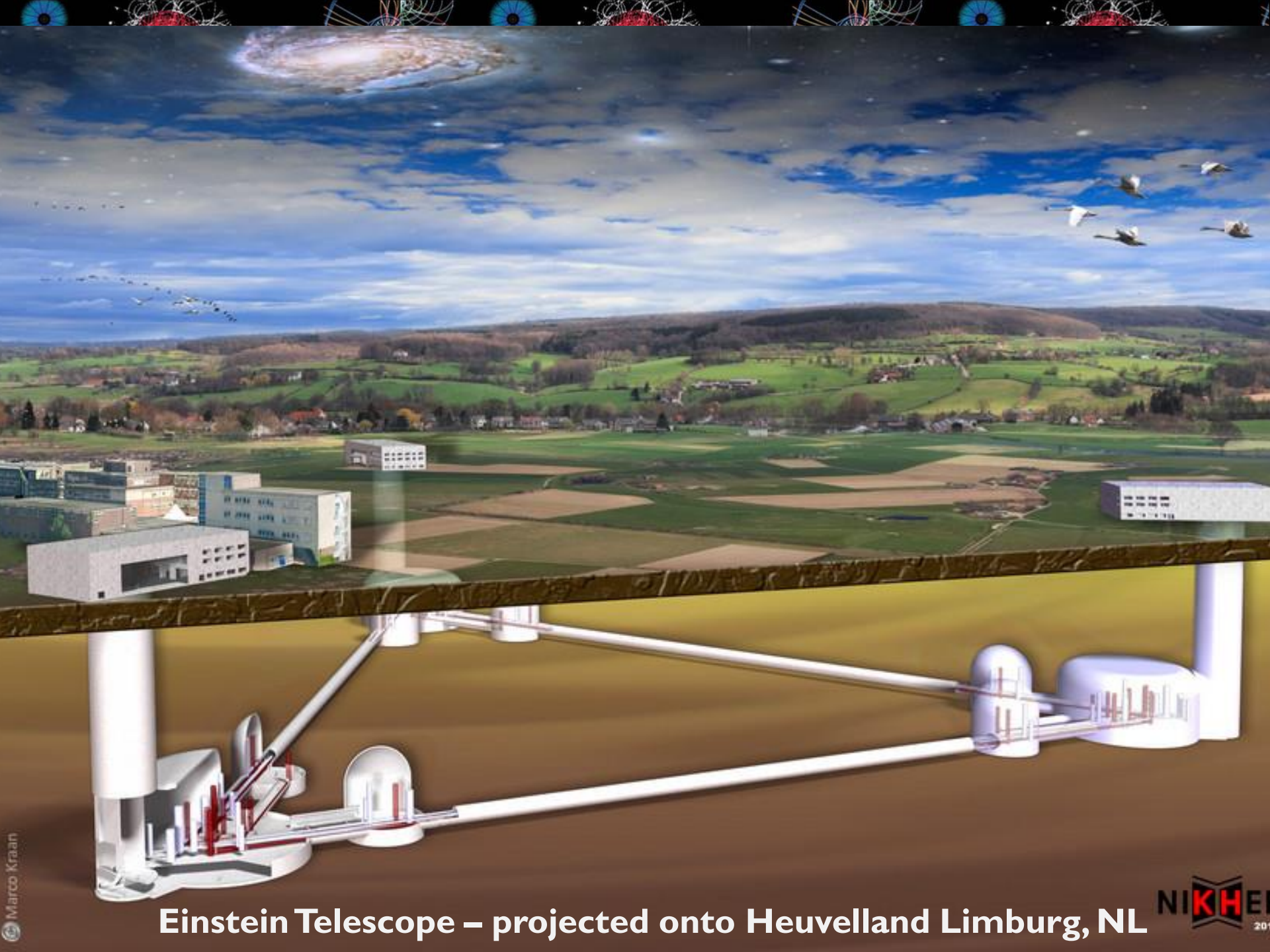
De Melkweg

Imagery:  Nikhef, NIOZ, KM3NET collaboration

Imagery: gw-astronomy collaborations, LSC

**Einstein Telescope – projected onto Heuvelland Limburg, NL**

Image sources: LNGS/INFN, Xenon collaboration;  Pierre Auger collaboration; Nikhef

# Deeltjes botsingen

# Deeltjes botsingen

# *Deeltjes botsingen*

$$E = mc^2$$

Higgs

muon

muon

muon

muon

Kans Higgs deeltje:
1 op de 1.000.000.000.000 bostingen
- Dit is equivalent met zoeken van 1 persoon op 1000 wereldpopulaties
- Oftewel één naald in 20 miljoen hooibergen

Higgs → ZZ*→ 4μ kandidaat,
M(4 leptonen)=125.1 GeV

Imagery: ATLAS experiment,, atlas.ch

ALICE

Run:244918
Timestamp:2015-11-25 11:25:36(UTC)
System: Pb-Pb
Energy: 5.02 TeV

12.5 MByte/event … 120 TByte/s … *and now what?*

proton — hard interaction — proton
$x_1 P$ — $x_2 P$
spectator quarks

**40 miljoen / seconde**

**Trigger systeem selecteert 600 Hz ~ 1 GB/s data**

**Analyse van botsingen door promovendi**

and processing

**Data distributie met GRID computers**

**50 PiB/year
primary data**

Image source: joint CERN (wLCG) and EGI

**Organisations participating in the global collaboration of e-Infrastructures**

_Even just for wLCG, supporting the CERN LHC programme_
**More than 200 independent institutes with end-users**
**More than 50 countries & regions**
**More than 300 service centres**
**One independent 'policy-bridge' identity service**
**Handful regional 'service coordination organisations'**
**500 000 CPU cores, 200+PByte storage**

# From SC04, CCRC08, STEP09, .. to today

Global transfer rates now > 40 GB/s –
acquisition:10 PB/mo (~x2 derived data)



David Groep
Nikhef
*PDP - Advanced
Computing for
Research*

# Atlas: ~50 TByte/day raw data to tape; 1000 TByte/day processed data transfers



Image source: CERN, Atlas

# … and tomorrow ?!

**Data estimates for 1st year of HL-LHC (PB)**

- ALICE
- ATLAS
- CMS
- LHCb

(Raw, Derived)

**CPU Needs for 1st Year of HL-LHC (kHS06)**

- ALICE
- ATLAS
- CMS
- LHCb

CPU (HS06)

Data:
- Raw 2016: 50 PB → 2027: 600 PB
- Derived (1 copy): 2016: 80 PB → 2027: 900 PB

CPU:
- x60 from 2016

Technology at ~20%/year will bring x6-10 in 10-11 years

David Groep
Nikhef
*PDP - Advanced Computing for Research*

# Infrastructure for research: balancing network, CPU, and disk

- CPU and disk both expensive, yet idling CPUs are 'even costlier'

- architecture and performance matching averts any single bottleneck

- but requires knowledge of application (data flow) behaviour data pre-placement (local access), mesh data federation (WAN access)

This is why e.g. your USB drive does not cut it – and neither does your 'home NAS box' … *however much I like my home system using just 15 Watt idle and offering 16TB for just € 915 …*

David Groep
Nikhef
*PDP - Advanced Computing for Research*

# Building the infrastructure for the LHC data



- From hierarchical data distribution to a full mesh and dynamic data placement

Amsterdam/NIKHEF-SARA

David Groep
Nikhef
*PDP - Advanced Computing for Research*

LHCOne graphic: lhcone.net

# Connecting Science through Lambdas



David Groep
Nikhef
*PDP - Advance
Computing for
Research*

# Network built around application data flow



Public Internet — 70 Gbps — SURFNET

SURFNET — 100 Gbps — Geant

100 Gbps

NL Light

LHCOPN LHCOne

2x10 Gbps

parkwachter LS Nikhef

20 Gbps

KIAE Atlas T1 via ReTN@Nikhef

100 Gbps

SURF SARA

40 Gbps

nx10 Gbps

> 100 Gbps

120Gbps

40 Gbps

4x10 Gbps

5x40Gbps

240Gbps
farmnet-core

10-40 Gbps per node

10 Gbps per node

roe
6x20 Gbps

SURFsara Compute & Storage

Tape Library

Desktops and servers 1Gbps

STBC dCache 500TB (kip)

NDPF DPM 2.5PByte (strijker, oliebol, haas)

NDPF Compute 5600 cores smrt, car, knal, mars, choc

STBC 'gloei' cluster local analysis service for PhDs

**100G IP poort voor testen - interface traffic (5 minute average)**

| | In | Uit | Piek In | In | Piek Uit | Uit |
| Maximum: | 90.59G | 100.32G | | | | |
| Average: | 19.01G | 29.50G | | created on Sun Sep 7 21:36:37 2014 | | |

*Need to work together!*
without our SURFsara peering,
SURFnet gets flooded ☺
and: you really want many of your own peerings

# 100Gbit

Nikhef → SURFnet → RUG-CIT||UvA

T Suerink
Nikhef
Amsterdam
*PDP & Grid*

**Duration**

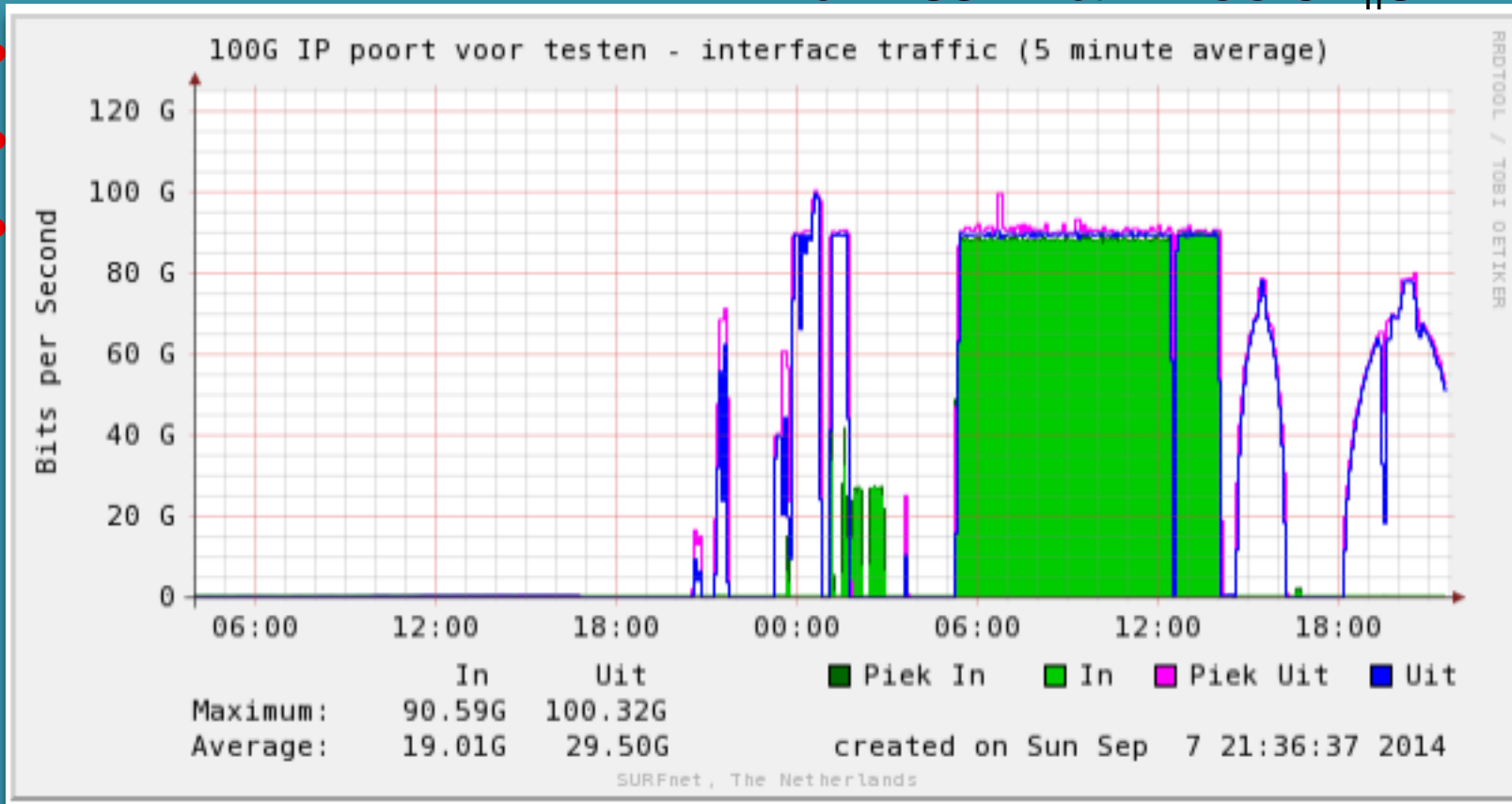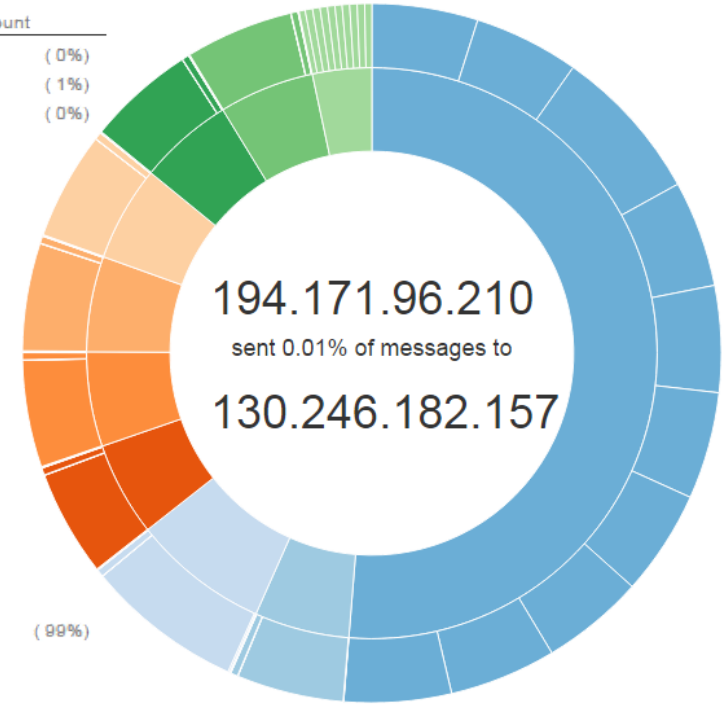| Source | Count | |
|---|---|---|
| 194.171.96.175 | 39 | ( 0% ) |
| 194.171.96.176 | 2507 | ( 7% ) |
| 194.171.96.177 | 2670 | ( 8% ) |
| 194.171.96.178 | 2677 | ( 8% ) |
| 194.171.96.181 | 2392 | ( 7% ) |
| 194.171.96.183 | 2390 | ( 7% ) |
| 194.171.96.185 | 2818 | ( 8% ) |
| 194.171.96.189 | 2516 | ( 7% ) |
| 194.171.96.190 | 2393 | ( 7% ) |
| 194.171.96.202 | 2555 | ( 8% ) |
| 194.171.96.203 | 0 | ( 0% ) |
| 194.171.96.205 | 2984 | ( 9% ) |
| 194.171.96.206 | 2373 | ( 7% ) |
| 194.171.96.207 | 2711 | ( 8% ) |
| 194.171.96.209 | 2607 | ( 8% ) |

| Destination | Count | |
|---|---|---|
| 188.184.66.250 | 1 | ( 0% ) |
| 192.108.46.8 | 189 | ( 1% ) |
| 192.108.46.89 | 39 | ( 0% ) |
| 192.5.19.10 | 33403 | ( 99% ) |

194.171.96.210
sent 0.01% of messages to
130.246.182.157

*Data flows: a user at UC Irvine reading a data set from Nikhef, with some background from CERN and KIT Karlsruhe*

*Graphics courtesy Jouke Roorda and Olivier Verbeek*

# Getting more bytes through?

- Power vs x64: more PCI lanes & higher clock should give more throughput – *if all the bits fit together*

- Only way to find out is … by trying it!
  *joint experiment with Nikhef and SURFsara*
  *on comparing IO throughput between x86 & P8*
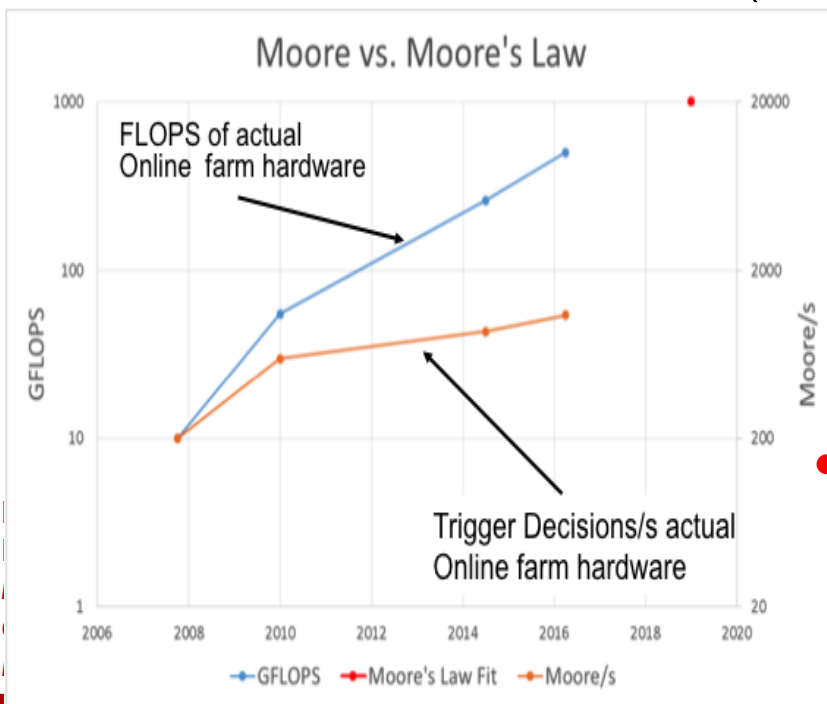
*HGST: 480 TByte gross capacity/4RU*

## yet more is needed

- RAID card are now a performance bottleneck
- JBOD changes CPU-disk ratio
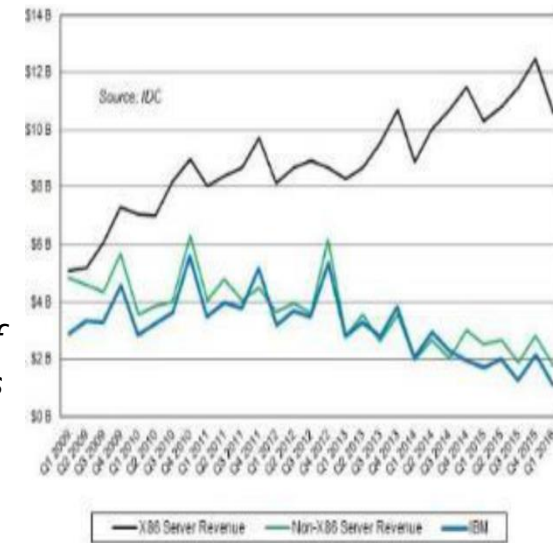- closer integration of networking to get >100Gbps

# Matching systems architecture

- Most applications using x86 today, and probably will for a long time
  - alternatives (GPGPU or Power) not quite viable
    … although for 'dedicated farms' FPGAs help,
    and KNH works better (we need the memory)

*sales volume of different architectures*

- Yet change must be:
  most gain to be had from SIMD
  vectorization and improved
  memory access patterns
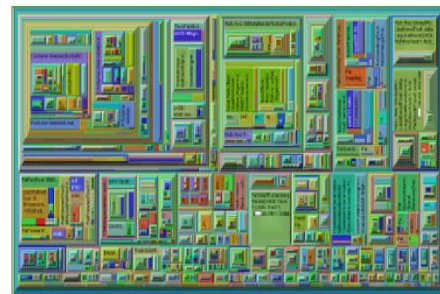
# Improvements at the application layer

- 'traditional' (1990's) style HEP applications were 'lean', and fail to scale even in pipelining

- let alone vector instructions or multicore



High level C++ code → `if (abs(point[0] - origin[0]) > xhalfsz) return FALSE;`

Assembler instructions →
```
movsd 16(%rsi), %xmm0
subsd 48(%rdi), %xmm0   // load & subtract
andpd _2il0floatpacket.1(%rip), %xmm0 // and with a mask
comisd 24(%rdi), %xmm0 // load and compare
jbe ..B5.3    # Prob 43% // jump if FALSE
```

Same instructions laid out according to **latencies** on the Core 2 processor →

NB: Out-of-order scheduling not taken into account.

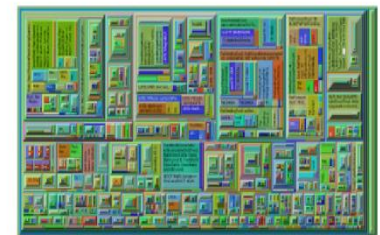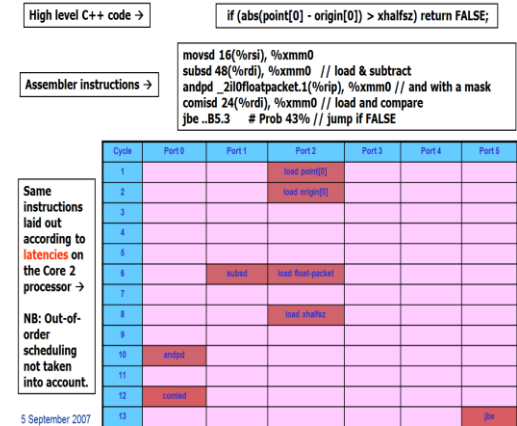5 September 2007

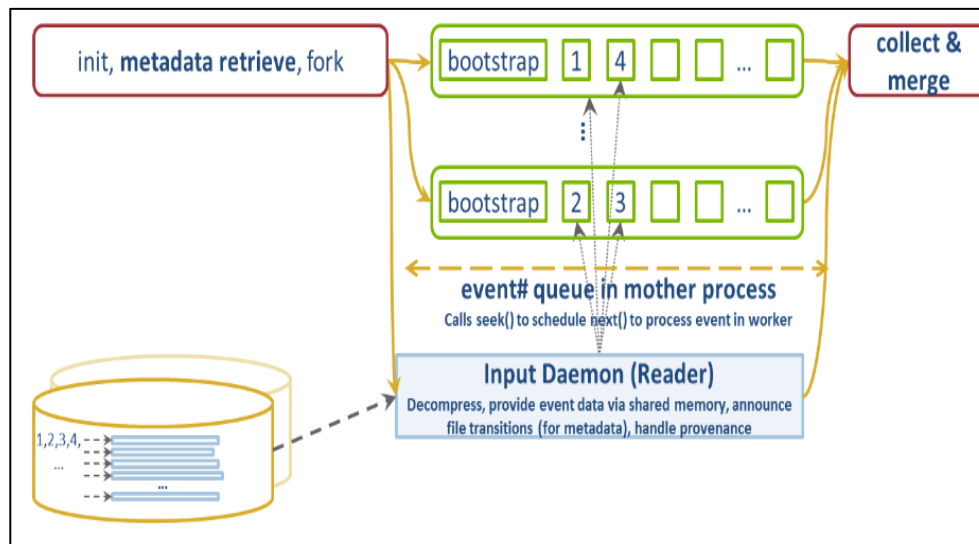2012    v45r1          v48r1          v48r1 (2015 reco)

*review of algorithms gave overall +34% in LHCb – memory layout still to be done …*

David Groep
Nikhef
*PDP - Advanced Computing for Research*

**To use current processor generations, you need better – machine-aware! – code**

# Systems architecture and your application

Many things you only find in production ...

- When you're 'embarrassingly parallel' with a memory challenge
  *why not try 'priming' of memory for the first few events and then fork?*

David Groep
Nikhef
*PDP - Advanced Computing for Research*

- Towards single-socket systems: cache coherence limits performance – there's a penalty to pay for massive multi-socket-big-memory hosts!

# Systems for Research @ Nikhef

David Groep
Nikhef
*PDP - Advanced Computing for Research*

# Statistics

Dutch National e-Infrastructure coordinated by **SURF**

*"BiG Grid" HTC and storage platform services*

- 3 core operational sites: SURFsara, Nikhef, RUG-CIT
- 25+ PiB tape, 10+ PiB disk, 12000+ CPU cores
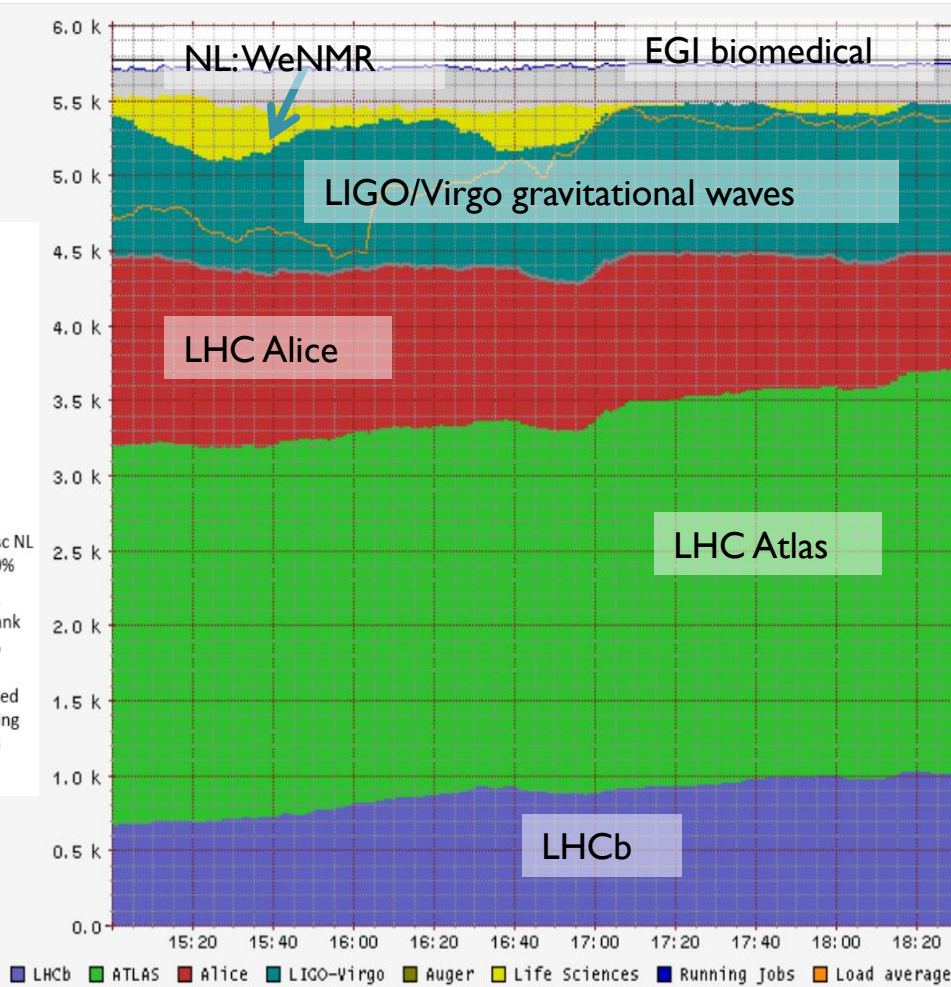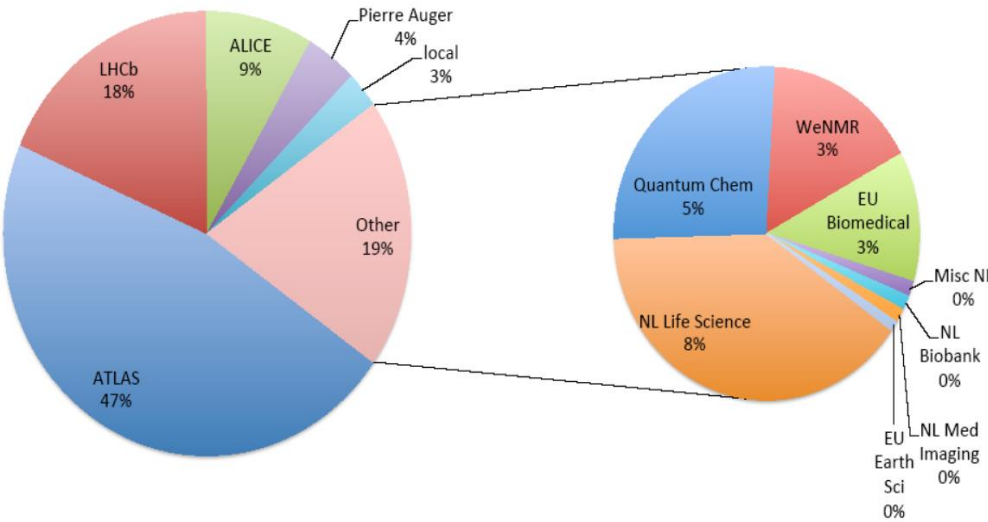
**@Nikhef**

~ 5500 cores and 3.5 PiB

focus on large/many-core systems

> 45 install flavours (service types)

*and a bunch of one-off systems*

David Groep
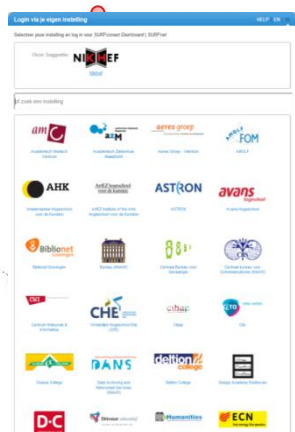Nikhef
*PDP - Advanced Computing for Research*

# Shared infrastructure, efficient infrastructure!

- >98% utilisation, >90% efficiency

David Groep
Nikhef
*PDP - Advanced Computing for Research*

# Federation of high-throughput services



**wLCG FIM4R pilot**

# 'cloud' is a means, not an end-all solution

**/cvmfs/softdrive.nl**

Nikhef  SURF SARA

E. Tejedor et al., CERN, SWAN Service for Web-based Analtsis, CHEP 2016

- Docker: single thin image, not managed by the user!
- CVMFS: configurable environment via "views"
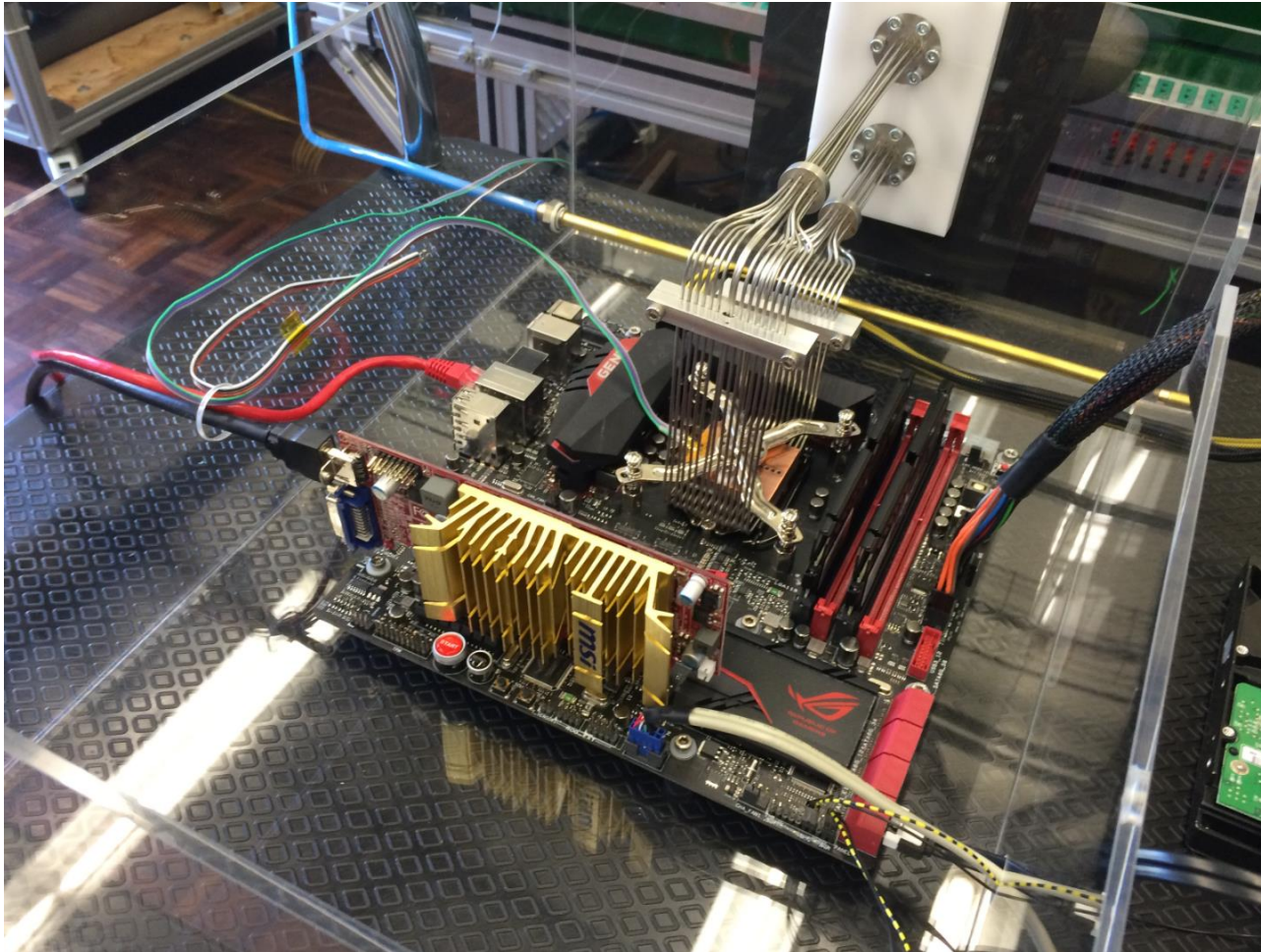- CERNBox: custom user environment

**LCG releases**

docker

C

CernVM File system

**CERN software**

CERNBox    **User software**

**@cern.ch - CHEP 2016 - Experiences with the ALICE Mesos infrastructure**

**Collaborative advantage:**
*joint effort of infrastructure and users*

# For (informed) fun & testing – some random one-off systems …



David Groep
Nikhef
*PDP - Advanced Computing for Research*

CO2-cooled Intel CPUs @6.2GHz

# For (informed) fun & testing – some random one-off systems …



David Groep
Nikhef
*PDP - Advanced
Computing for
Research*

Fun, but not the solution to single-core performance …