

# *Een petabyte weinig? Spoorzoeken door data*



Nikhef

Weekend van de Wetenschap  
Amsterdam Science Park

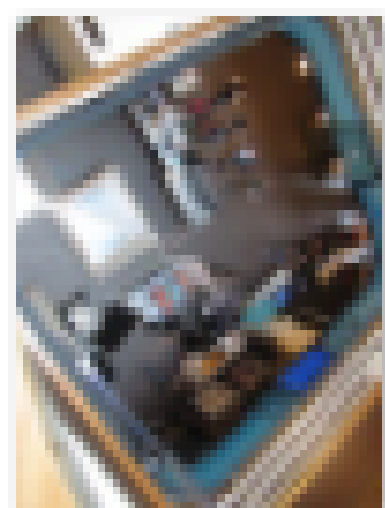
2018-10-06

David Groep  
[davidg@nikhef.nl](mailto:davidg@nikhef.nl)





# Botsingen 'fotograferen'?

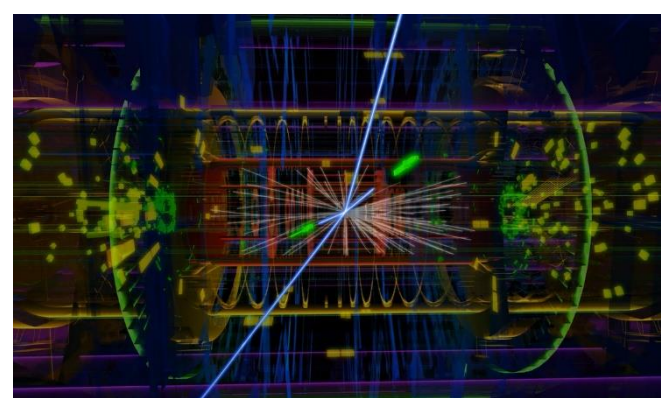


IMG\_6592.JPG

JPG File

1.51 MB

digitale compact-camera, 5MP



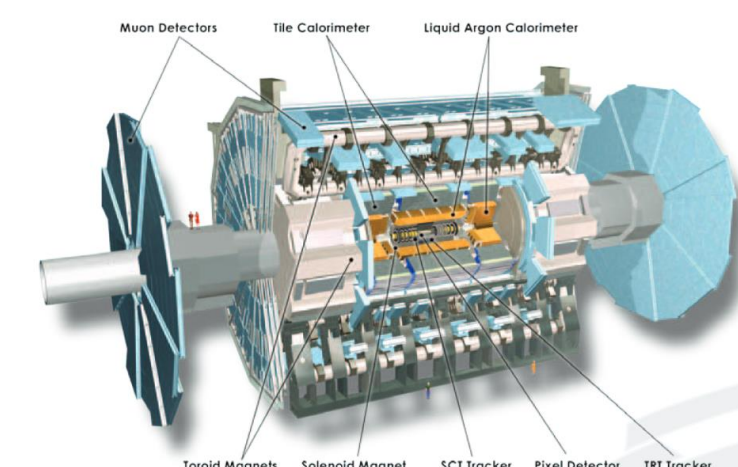
ATLAS\_RAW\_single\_event.data

ROD File

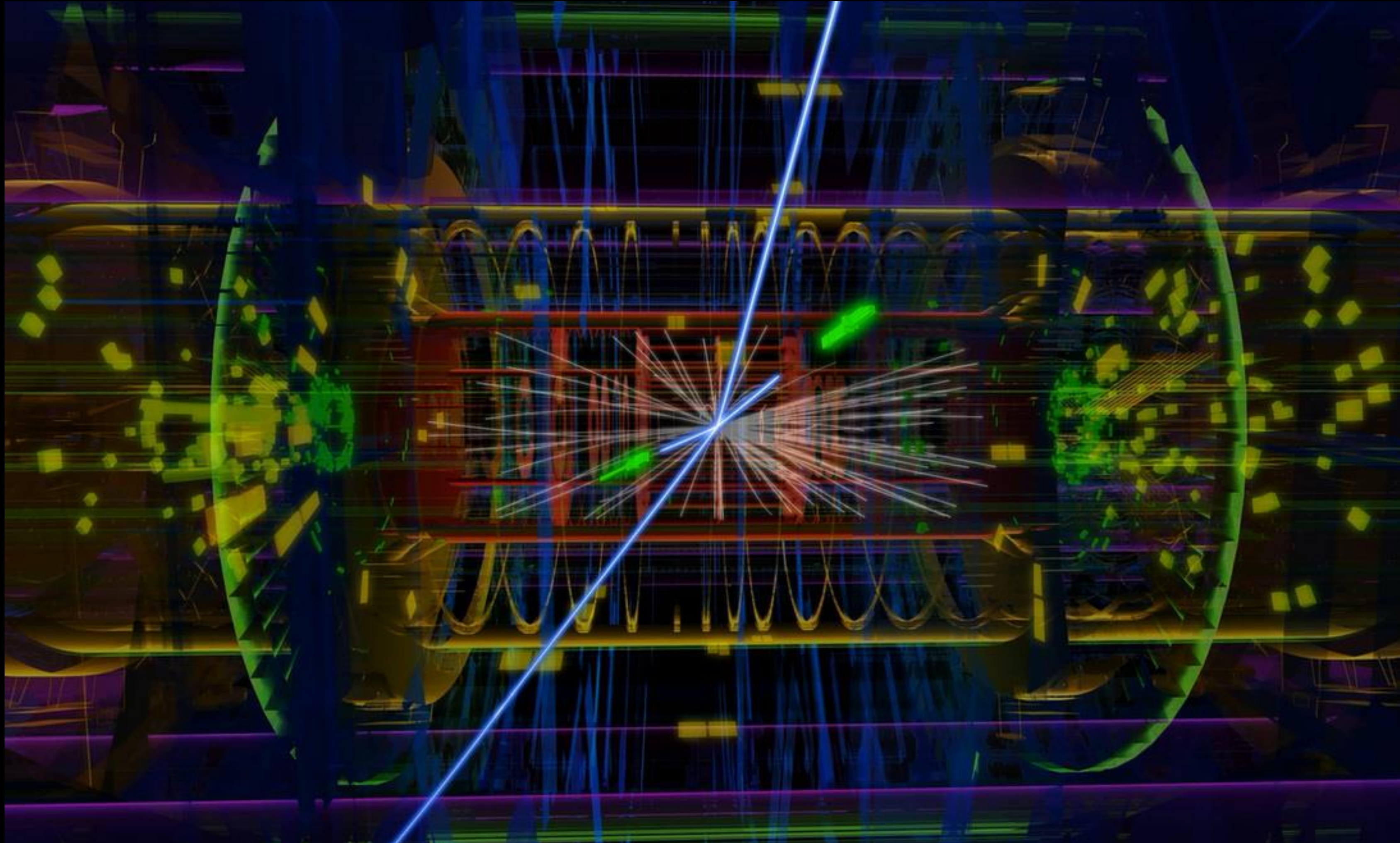
1.60 MB


ATLAS detector, ~88 MP

(maar niet alles licht  
elke keer op ...  
gelukkig!)



*~ 10 seconden rekenen per gebeurtenis voor een ATLAS event met 'jets' met daarin ~30 botsingen*





<i>Zachte botsingen</i>	$10^8$
$W^\pm \rightarrow e^\pm \nu$	15
$Z^0 \rightarrow e^+ e^-$	1
<i>Top-anti-top quarks</i>	1
$bb \rightarrow \mu + X$	$10^3$
<i>QCD jets, <math>p_T &gt; 150 \text{ GeV}</math></i>	$10^2$

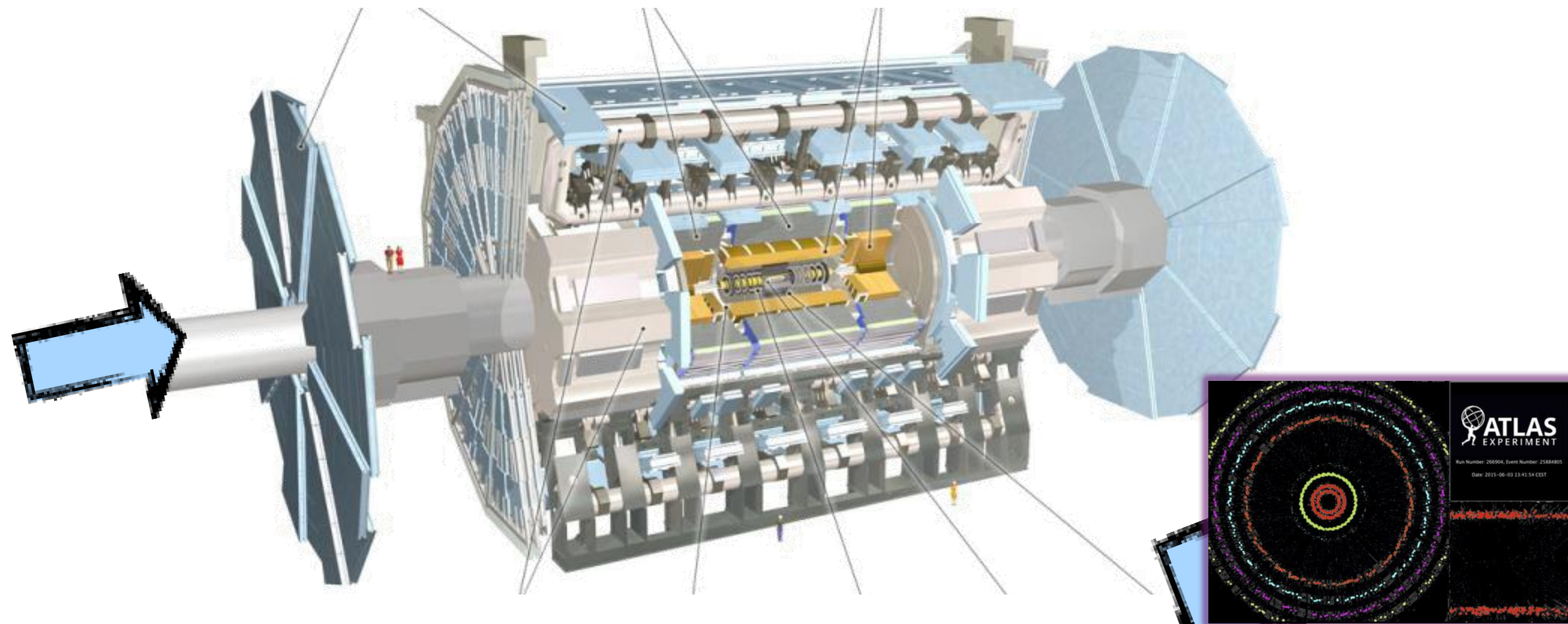
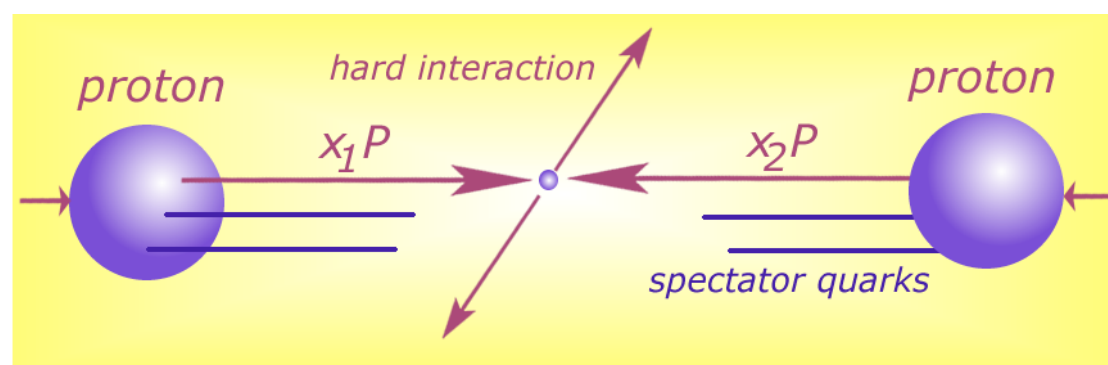
*Higgs deeltje:  $\sim 1$  per dag*

event rates from LHC Run 1

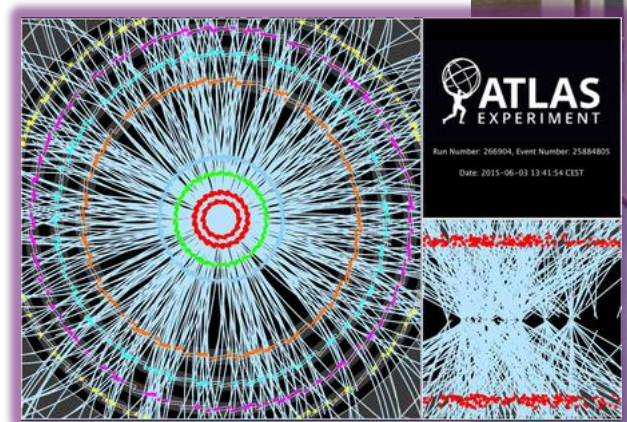
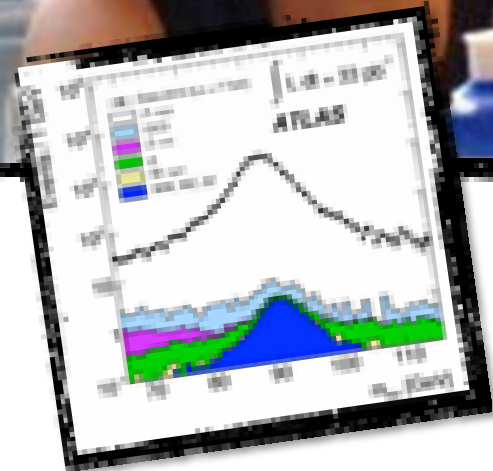
**Selectie interessante gebeurtenissen uit 'achtergrond' van 1 op  $10^{13}$  gebeurtenissen**  
- dit is equivalent met zoeken van  
1 persoon op 1000 wereldpopulaties  
- oftewel één speld in 20 miljoen hooibergen

# Hoe vind je 'interessante' events?

40 miljoen / seconde



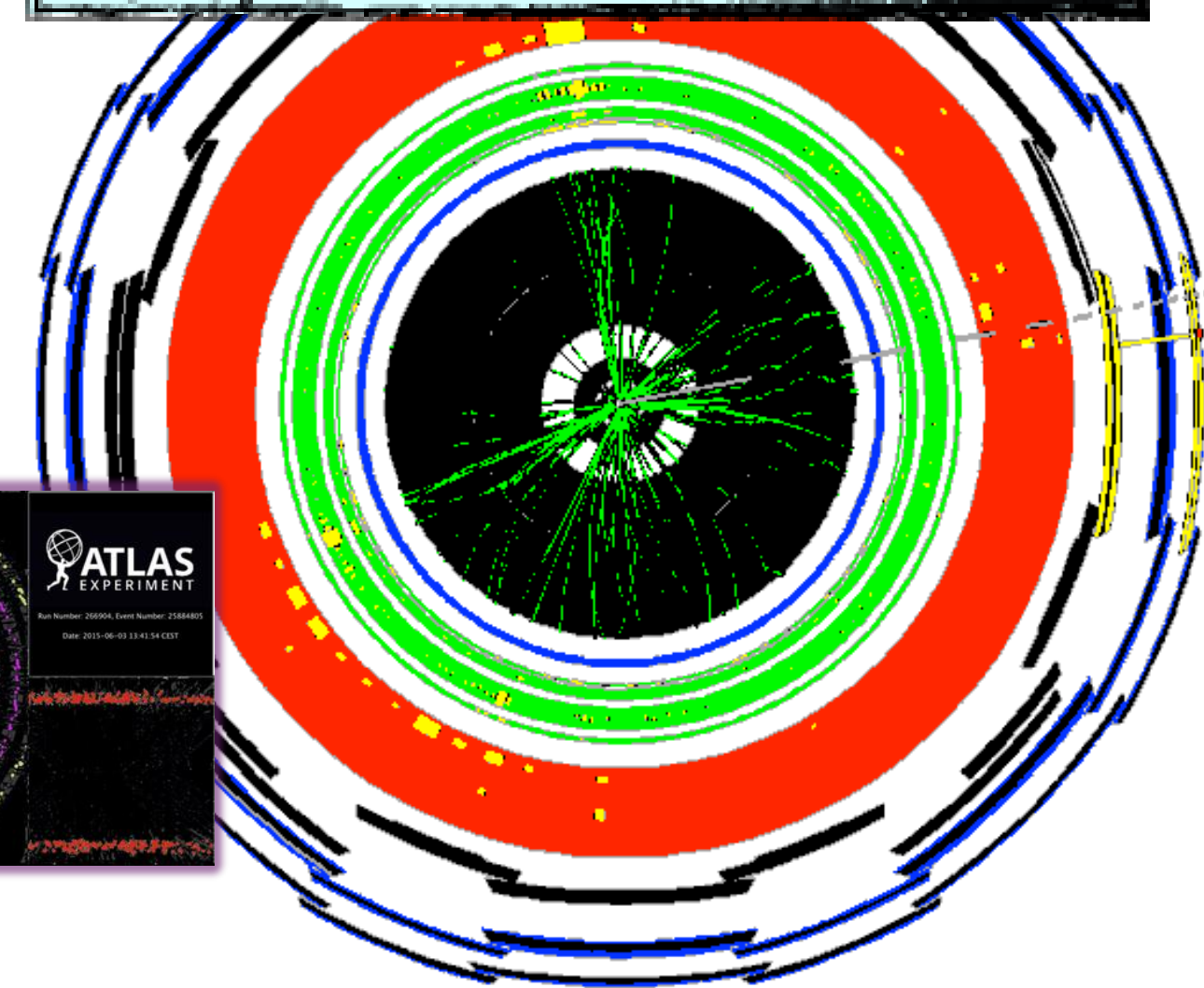
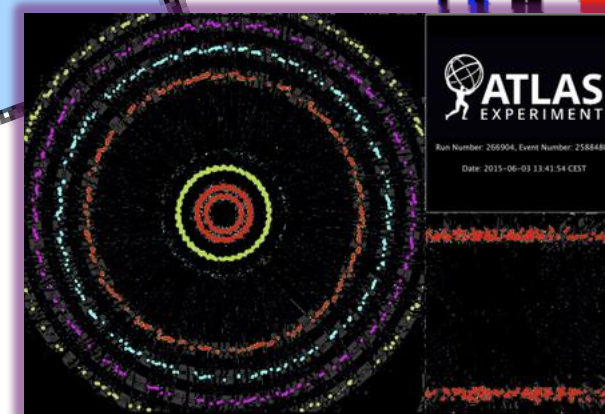
Analyse van botsingen door promovendi



Classificatie van deeltjes in de botsingen:

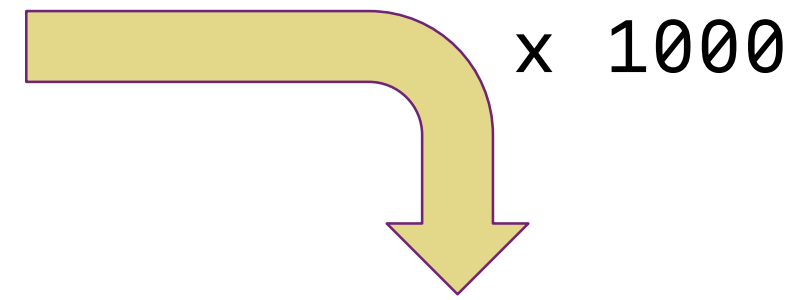
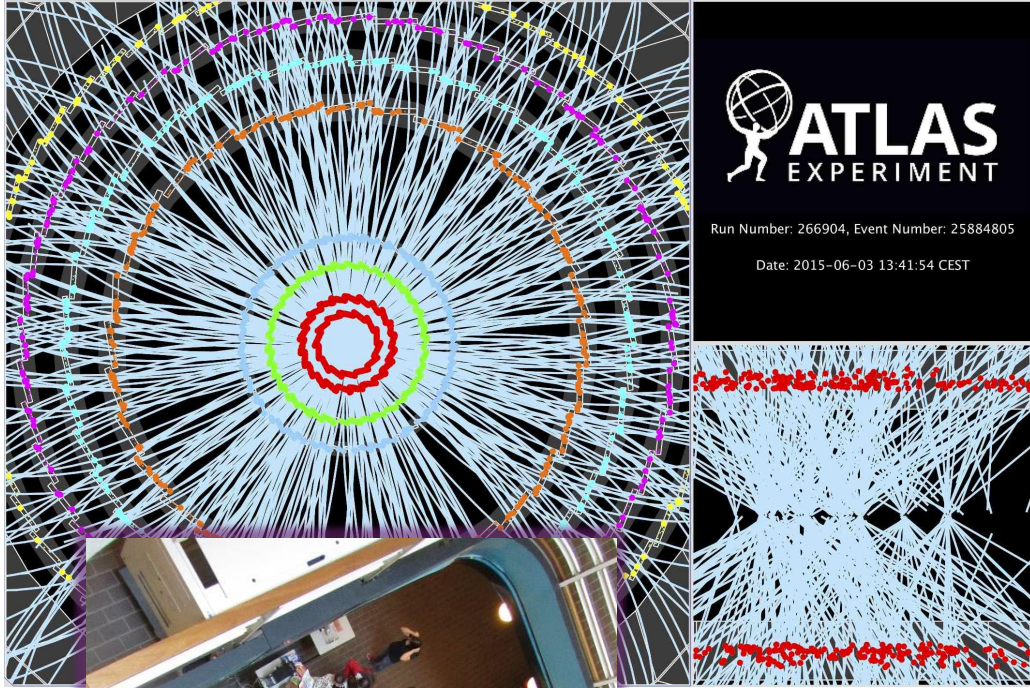
- *electronen*
- *muonen*
- *jets van hadronen*
- ...

Trigger system selecteert 600 Hz ~ 1 GB/s data

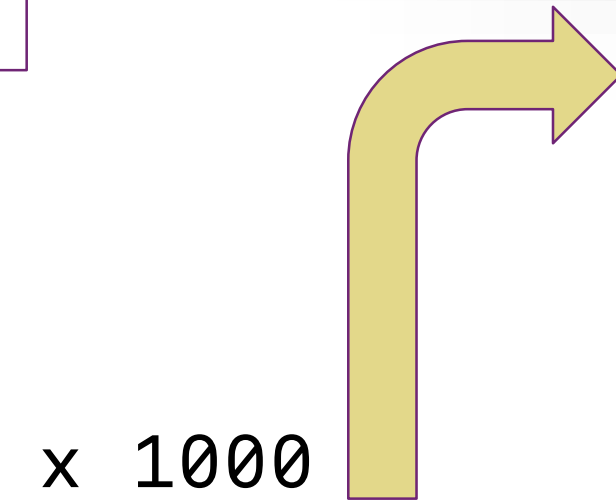


# Mega, Giga, Tera, ... en Peta byte

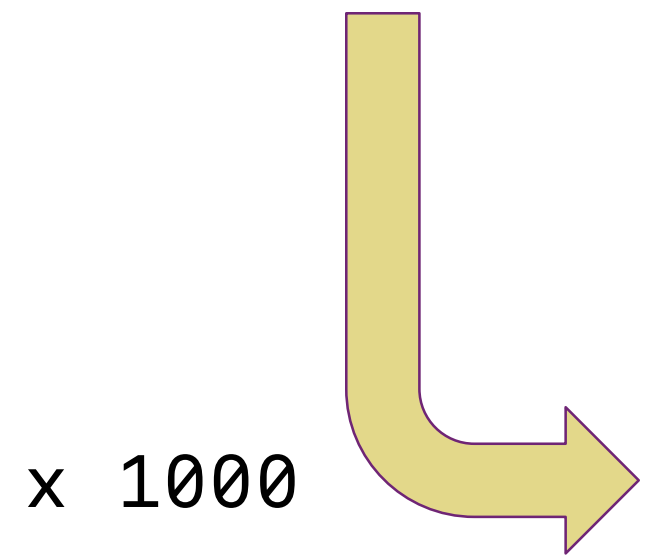
1 Megabyte, MB = 1 000 000 bytes  
 [en 1 Mibibyte, MiB = 1 048 576 bytes]



1 Gigabyte, GB = 1 000 MB  
 = 1 000 000 000 bytes



1 Terabyte, TB = 1 000 GB  
 = 1 000 000 MB



1 Petabyte, PB = 1 000 TB  
 ofwel 1 000 000 000 MB



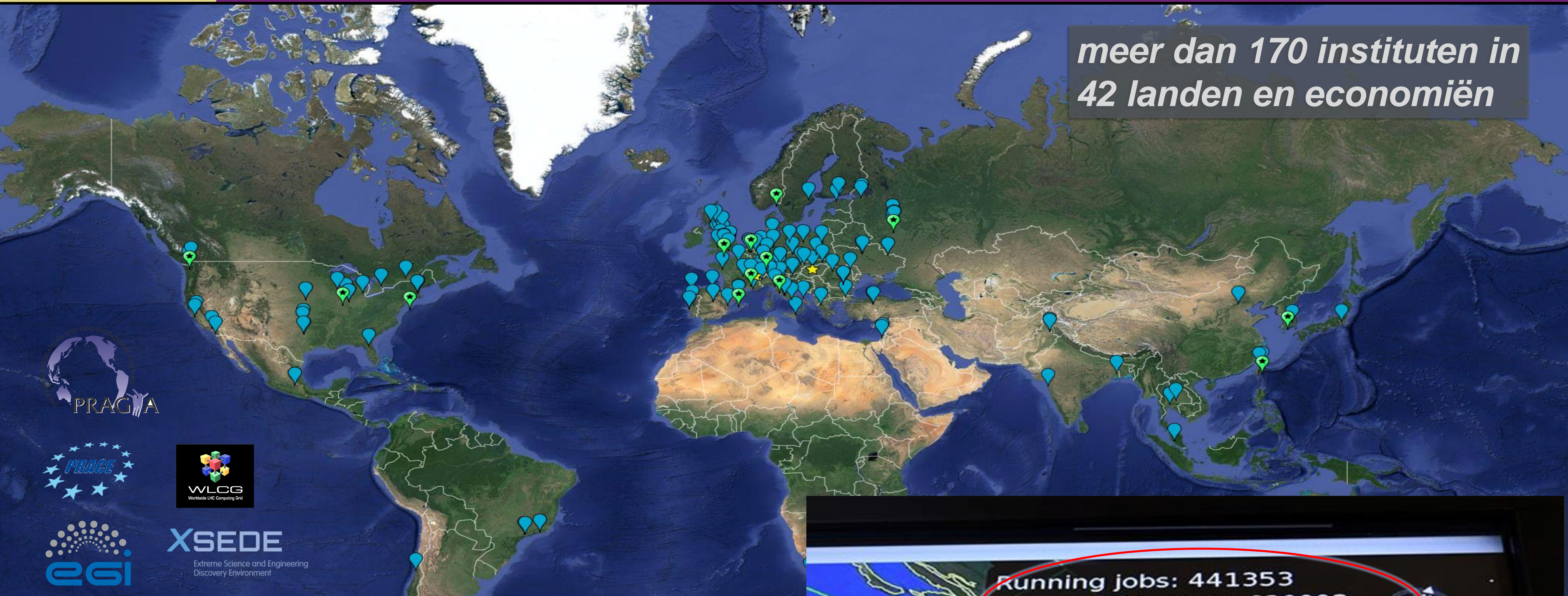
x 1000



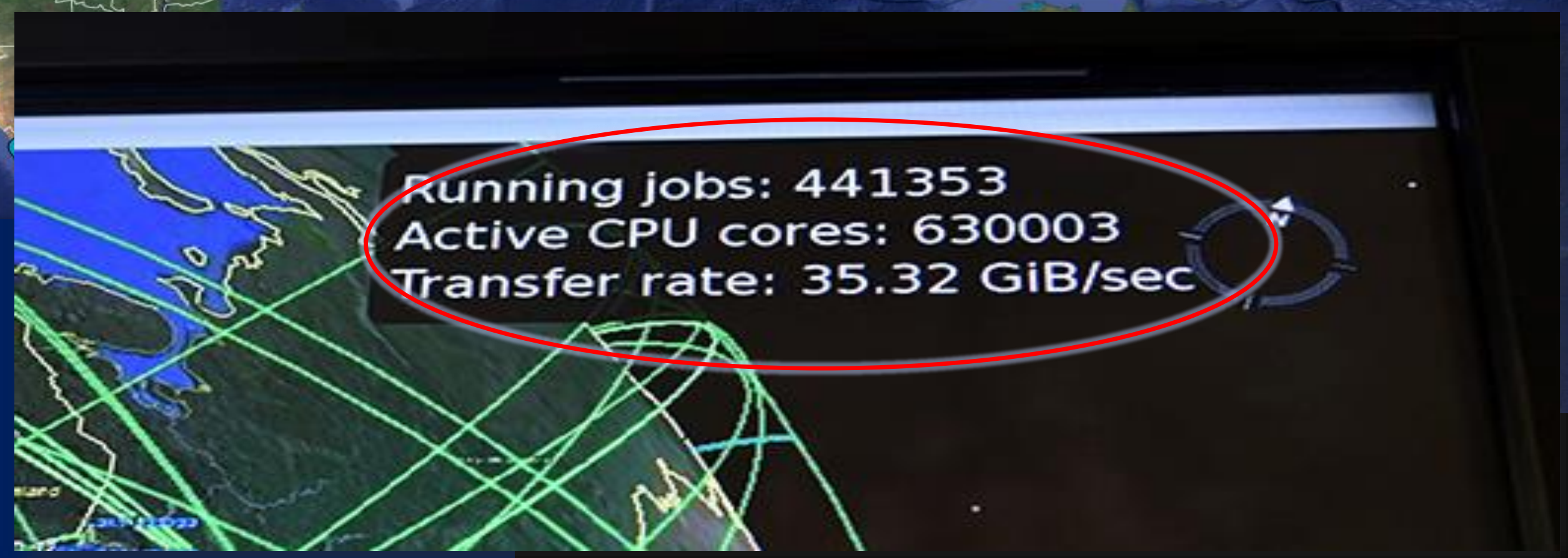


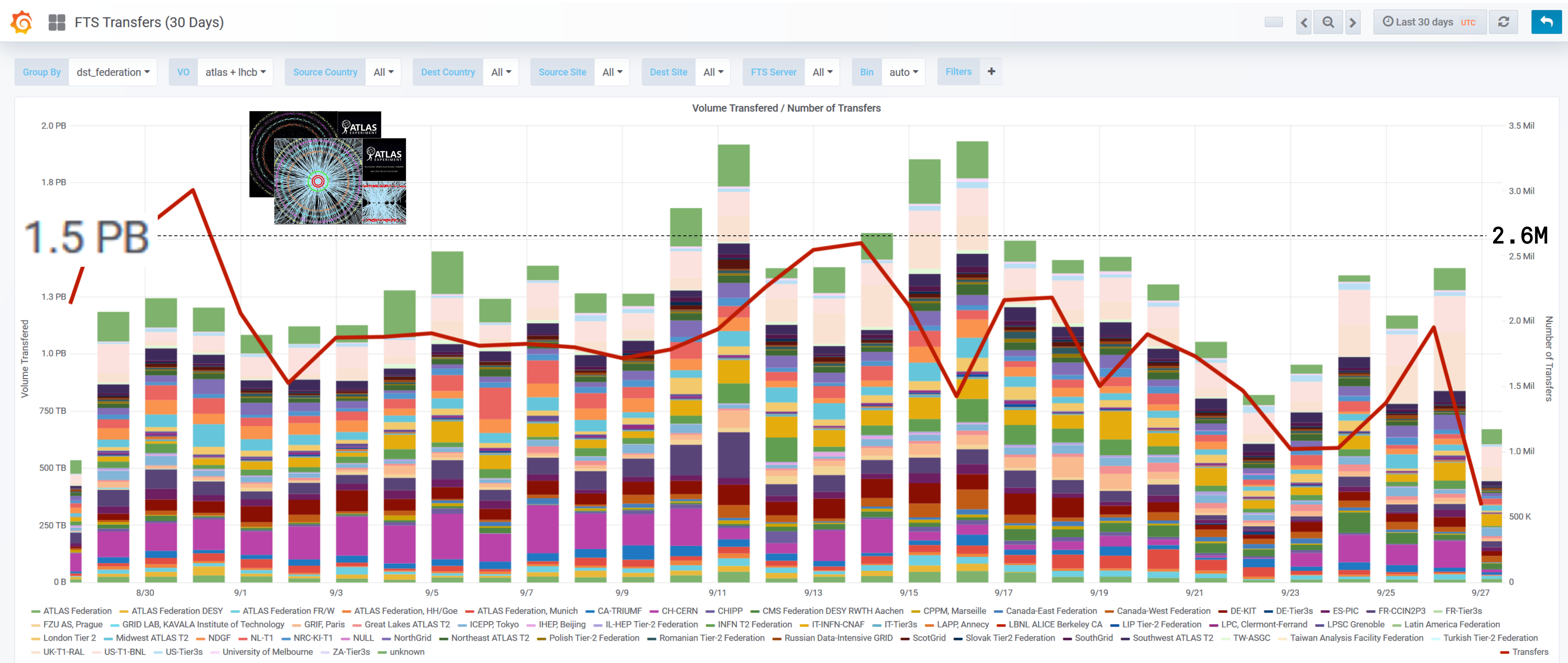
**50 PiB/year gegevens  
– alleen al vanuit CERN**

meer dan 170 instituten in  
42 landen en economiën



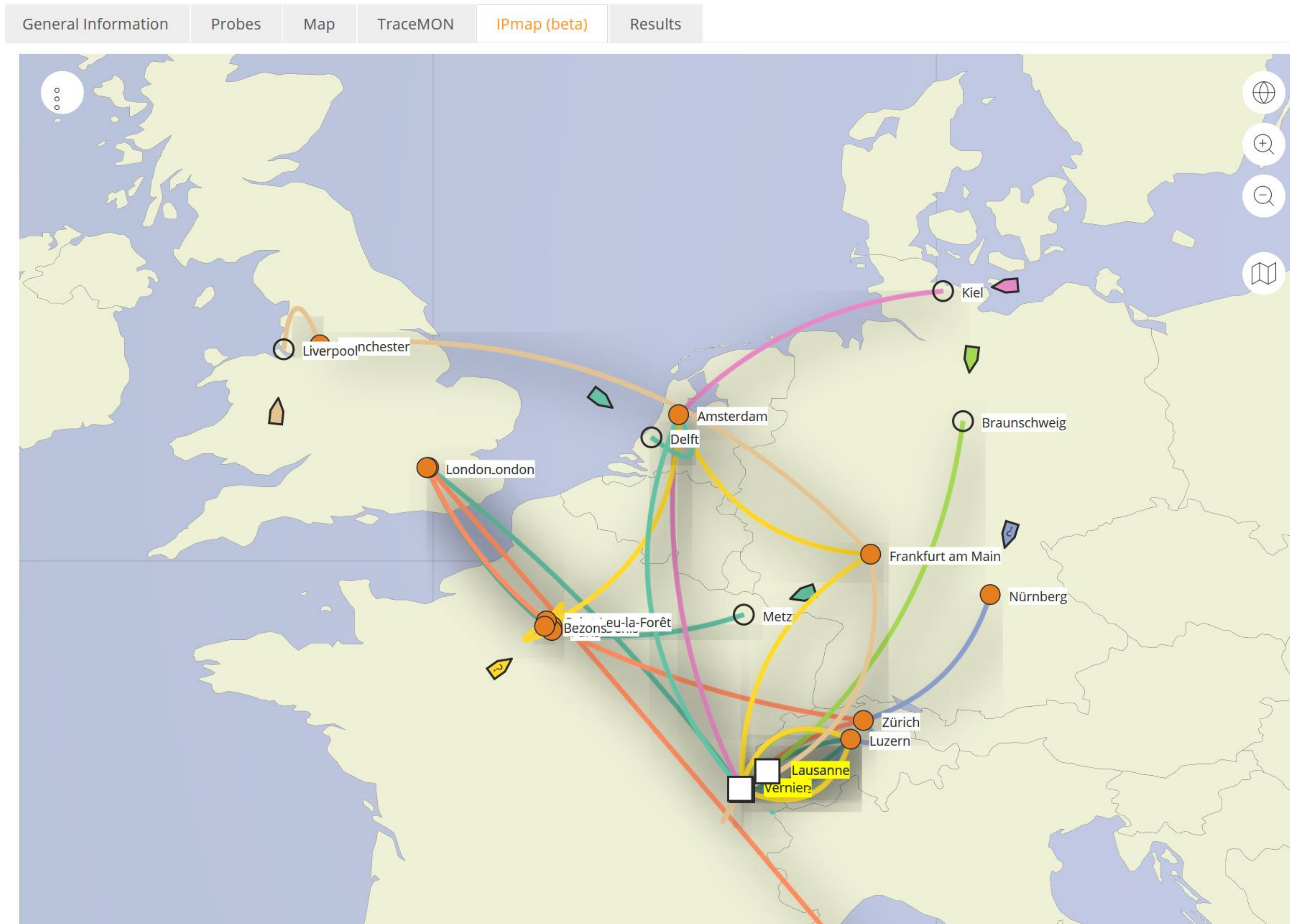
- ❑ **CPU: ~ 350,000 modern rekenkernen**
- ❑ **Disk 310 PB**
- ❑ **Tape 390 PB**



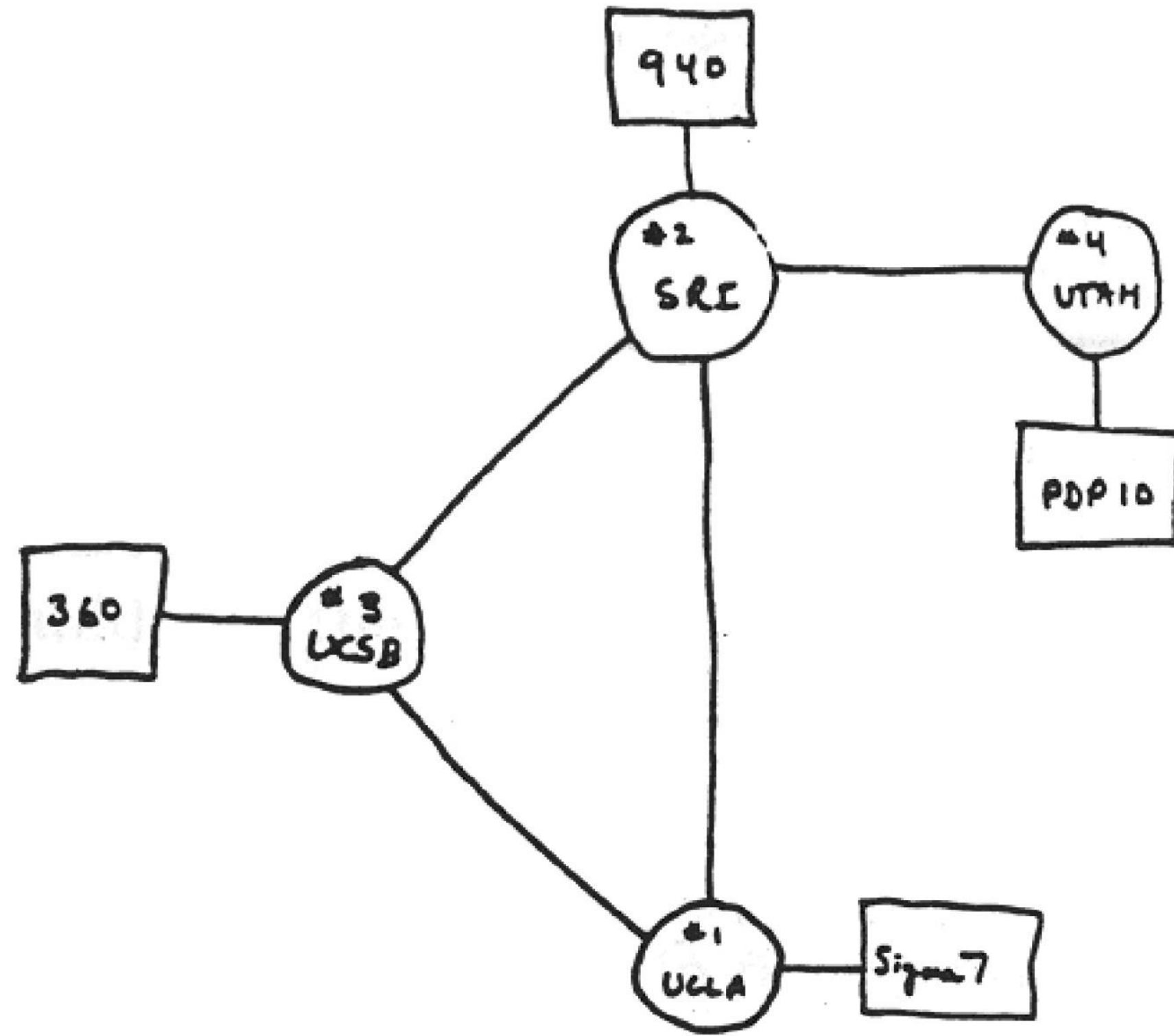


# hoe gaan je thuis naar CERN?

## ⚡ Traceroute measurement to linuxsoft.cern.ch (multihomed)



Data: TraceMON IPmap  
 from RIPE NCC Atlas  
[atlas.ripe.net](https://atlas.ripe.net)  
 measurement 9249079



THE ARPA NETWORK

DEC 1969

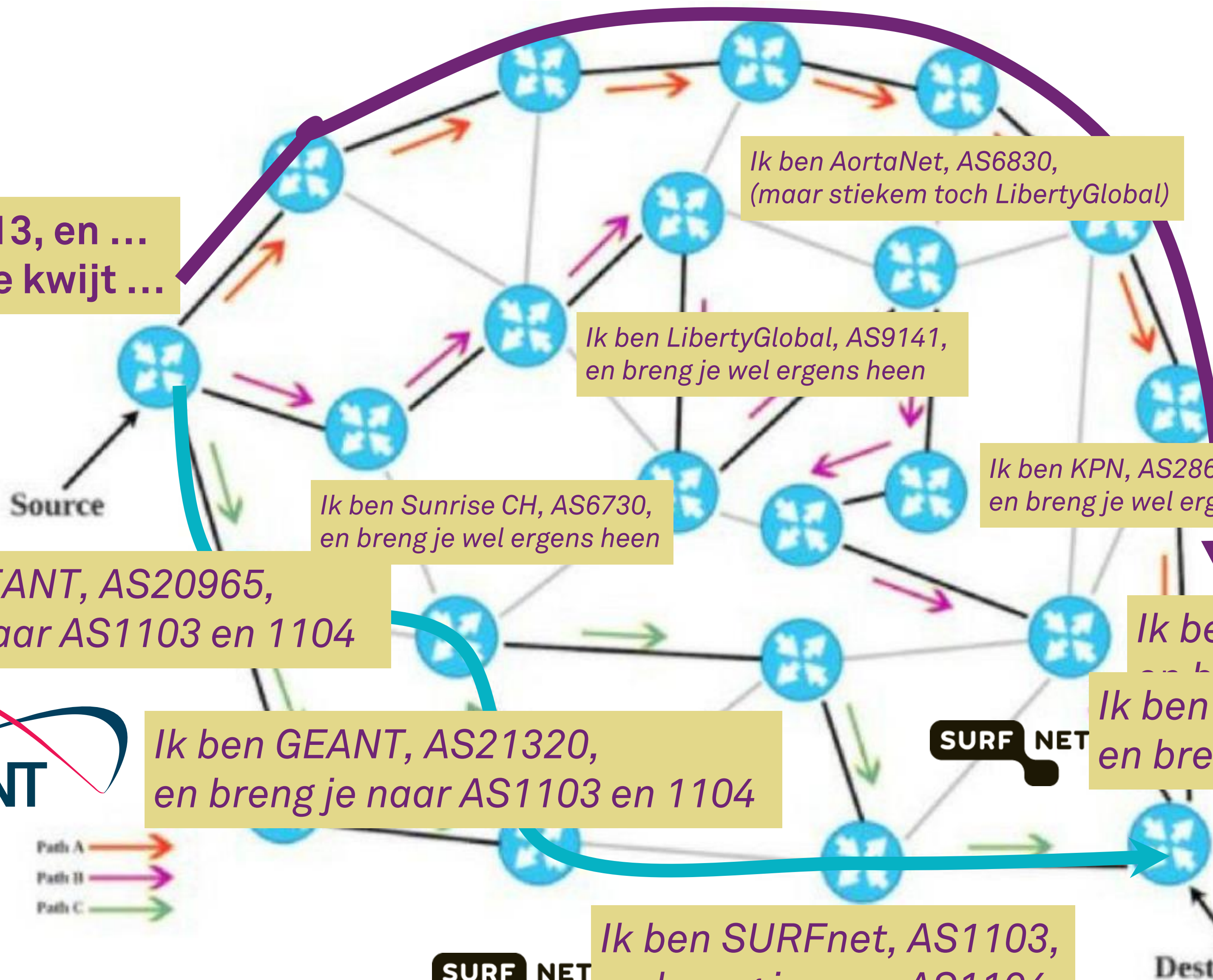
4 NODES





**Ik ben CERN, AS513, en ...  
ik wil m'n pakketje kwijt ...**

188.184.38.9



*Ik ben AortaNet, AS6830,  
(maar stiekem toch LibertyGlobal)*

*Ik ben LibertyGlobal, AS9141,  
en breng je wel ergens heen*

*Ik ben KPN, AS286,  
en breng je wel ergens heen*

*Ik ben Sunrise CH, AS6730,  
en breng je wel ergens heen*

*Ik ben ook GEANT, AS20965,  
en breng je naar AS1103 en 1104*

*Ik ben GEANT, AS21320,  
en breng je naar AS1103 en 1104*

*Ik ben SURFsara, AS1162,  
en breng je direct naar AS1104!  
Ik ben SURFnet, AS1103,  
en breng je naar AS1104*

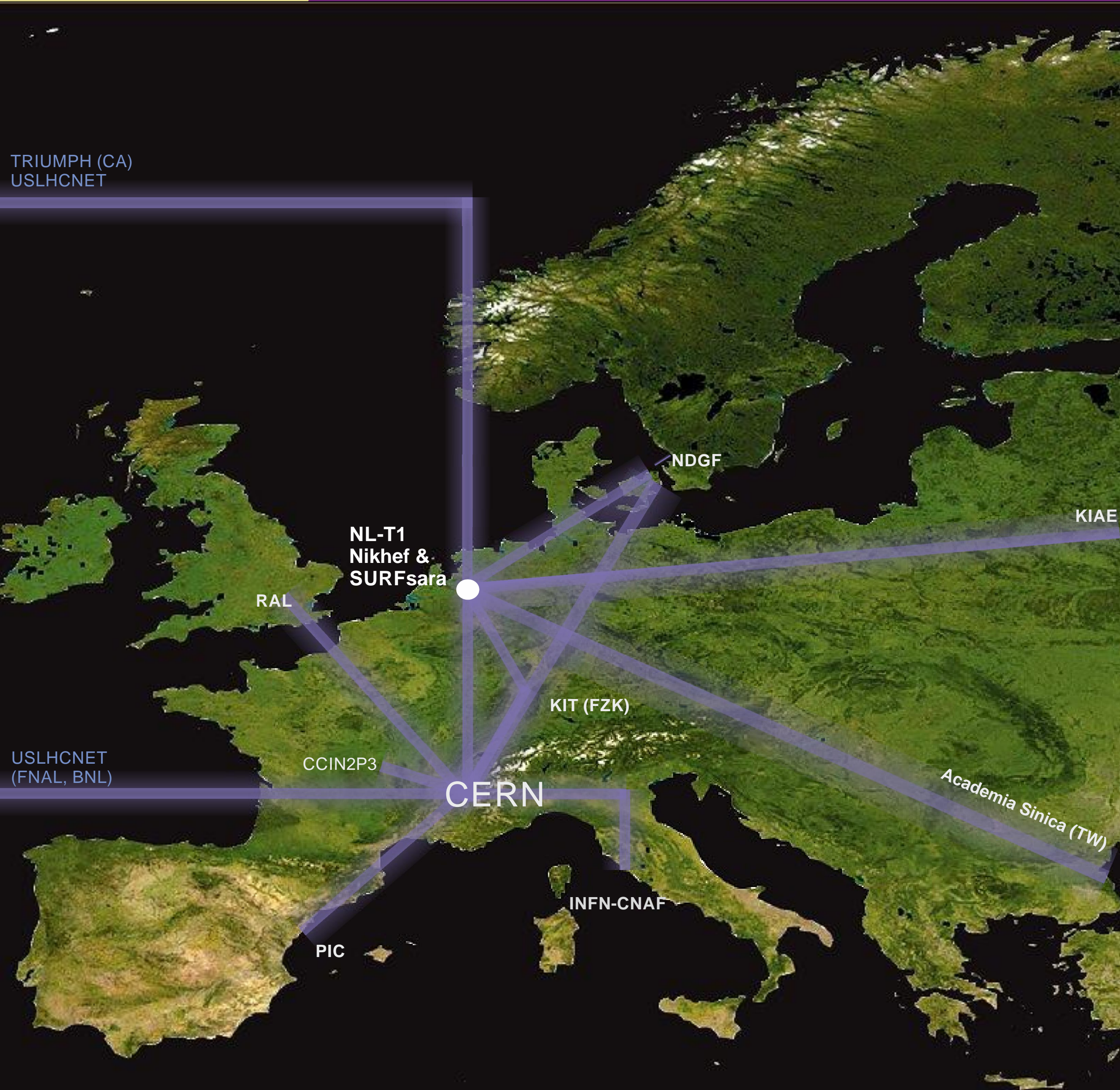
**Ik ben Nikhef, AS1104  
en ik heb schijfruimte voor je**

*Ik ben SURFnet, AS1103,  
en breng je naar AS1104*

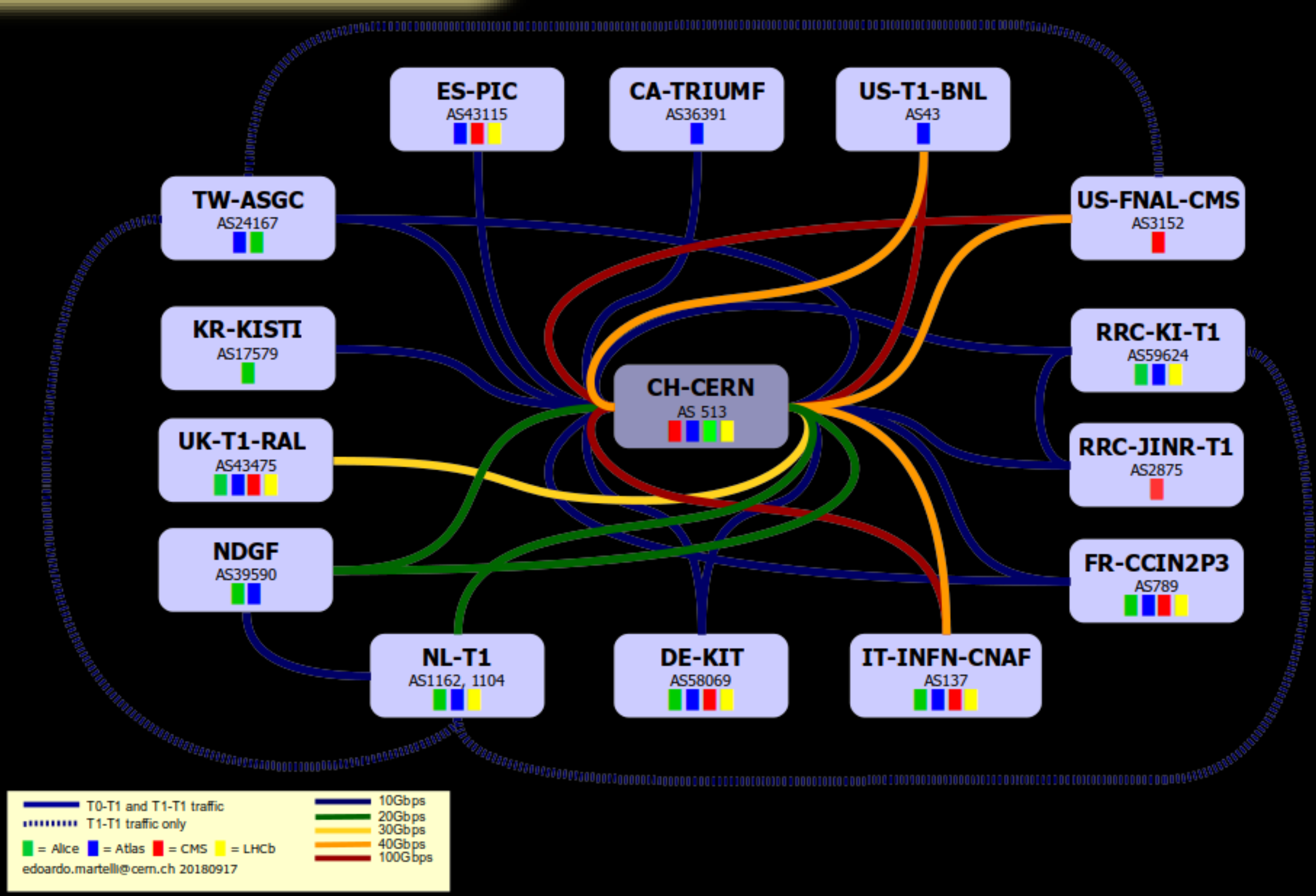


194.171.96.128/25

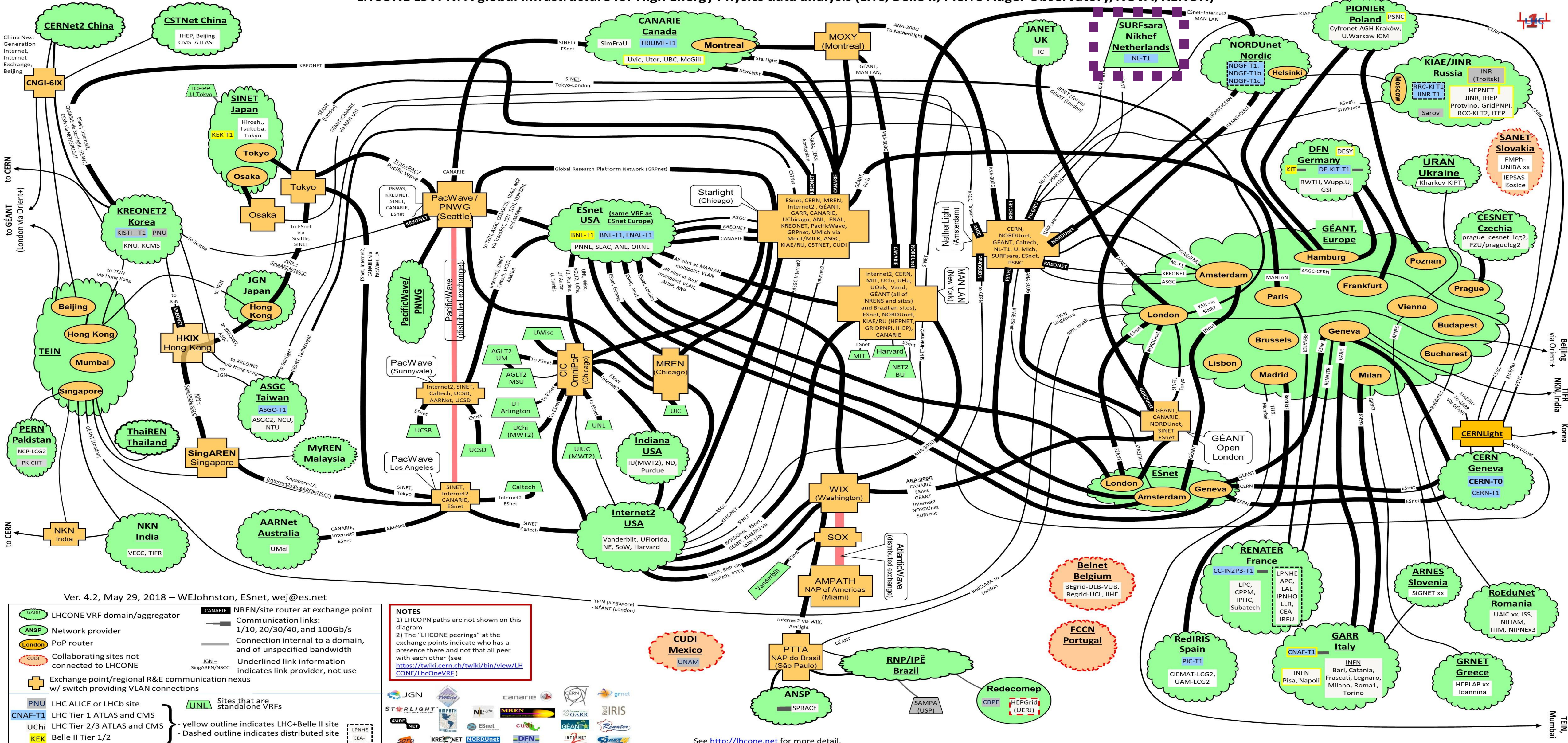
"segment routing"  
image by David Penaloza, CISCO



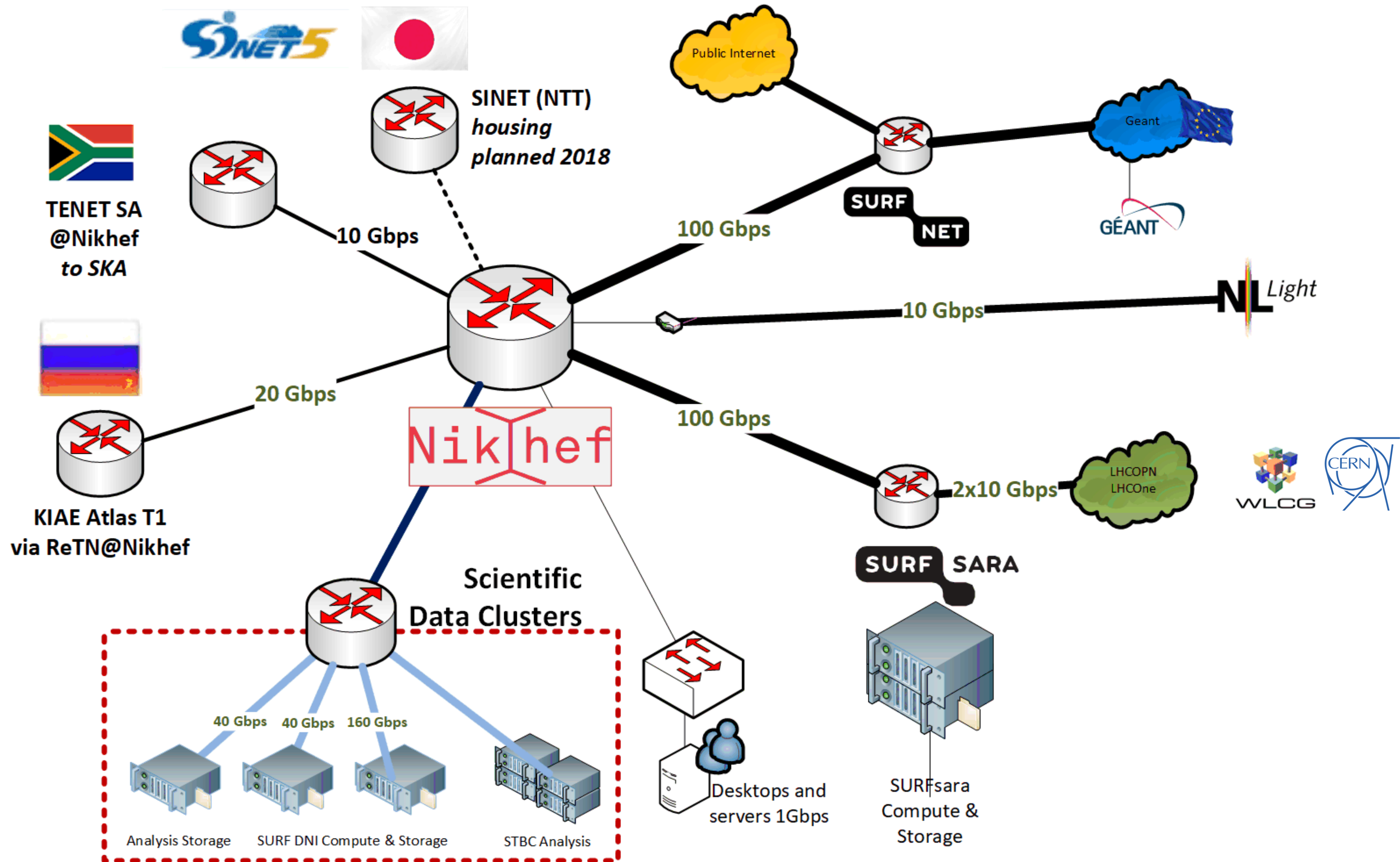
## LHCOPN



LHCONE L3VPN: A global infrastructure for High Energy Physics data analysis (LHC, Belle II, Pierre Auger Observatory, NOvA, XENON)



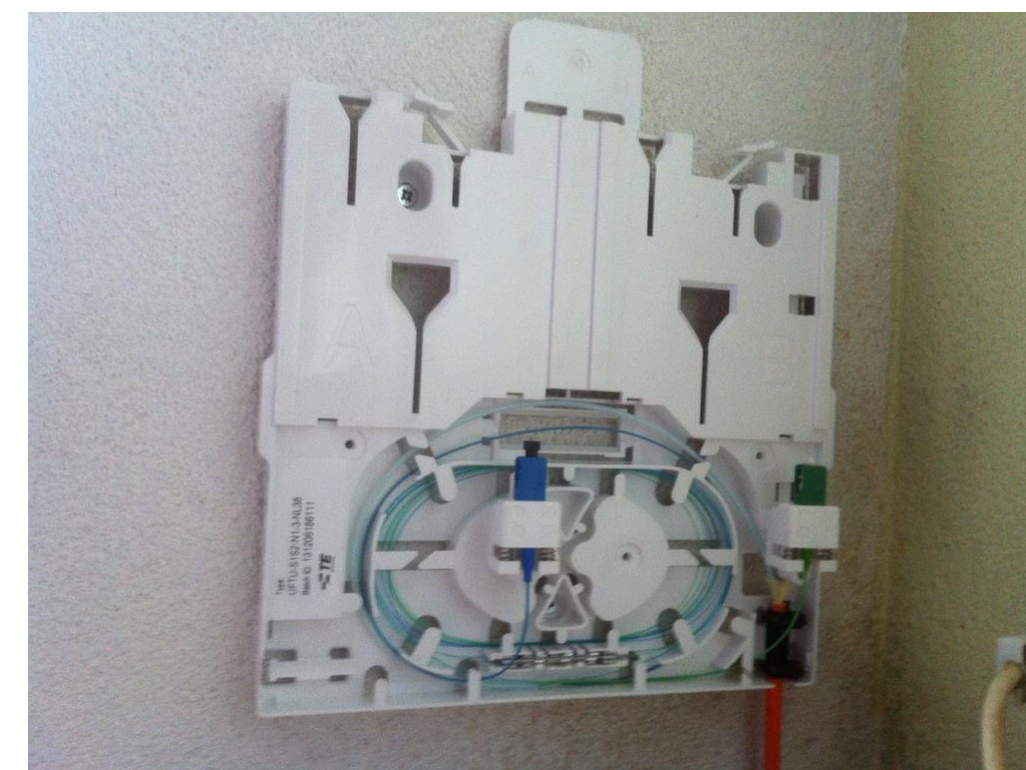




# Hoe ziet 100Gbps eruit?

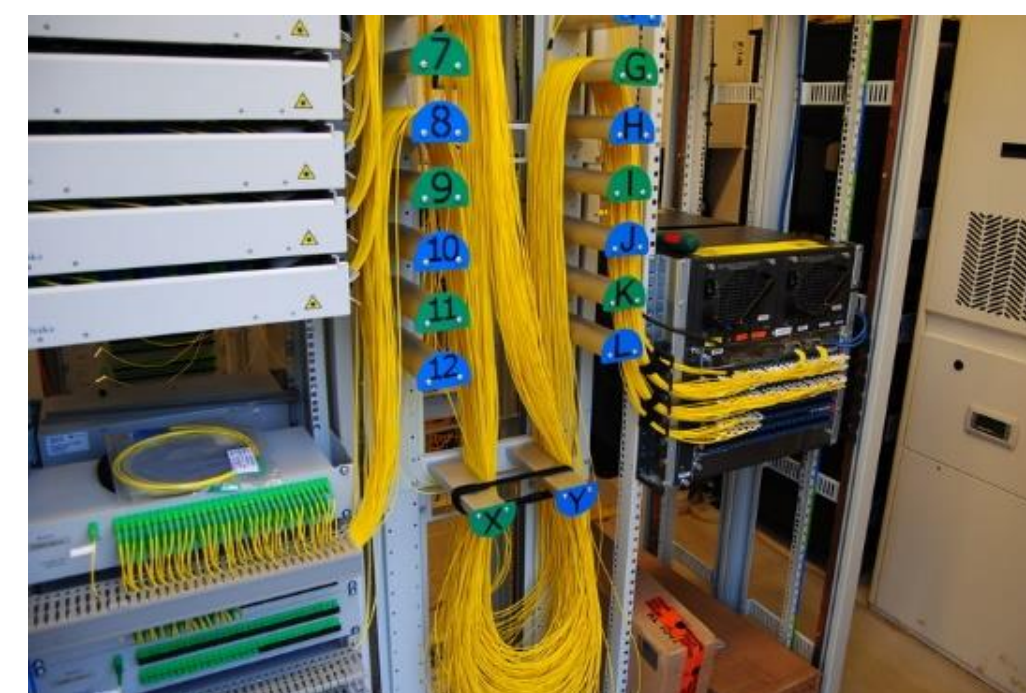


Thuis 'FttH'  
~1Gbps BX  
single strand, SC



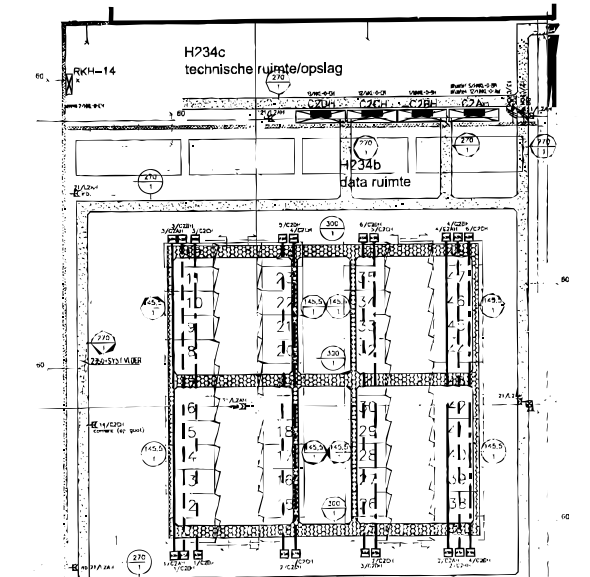
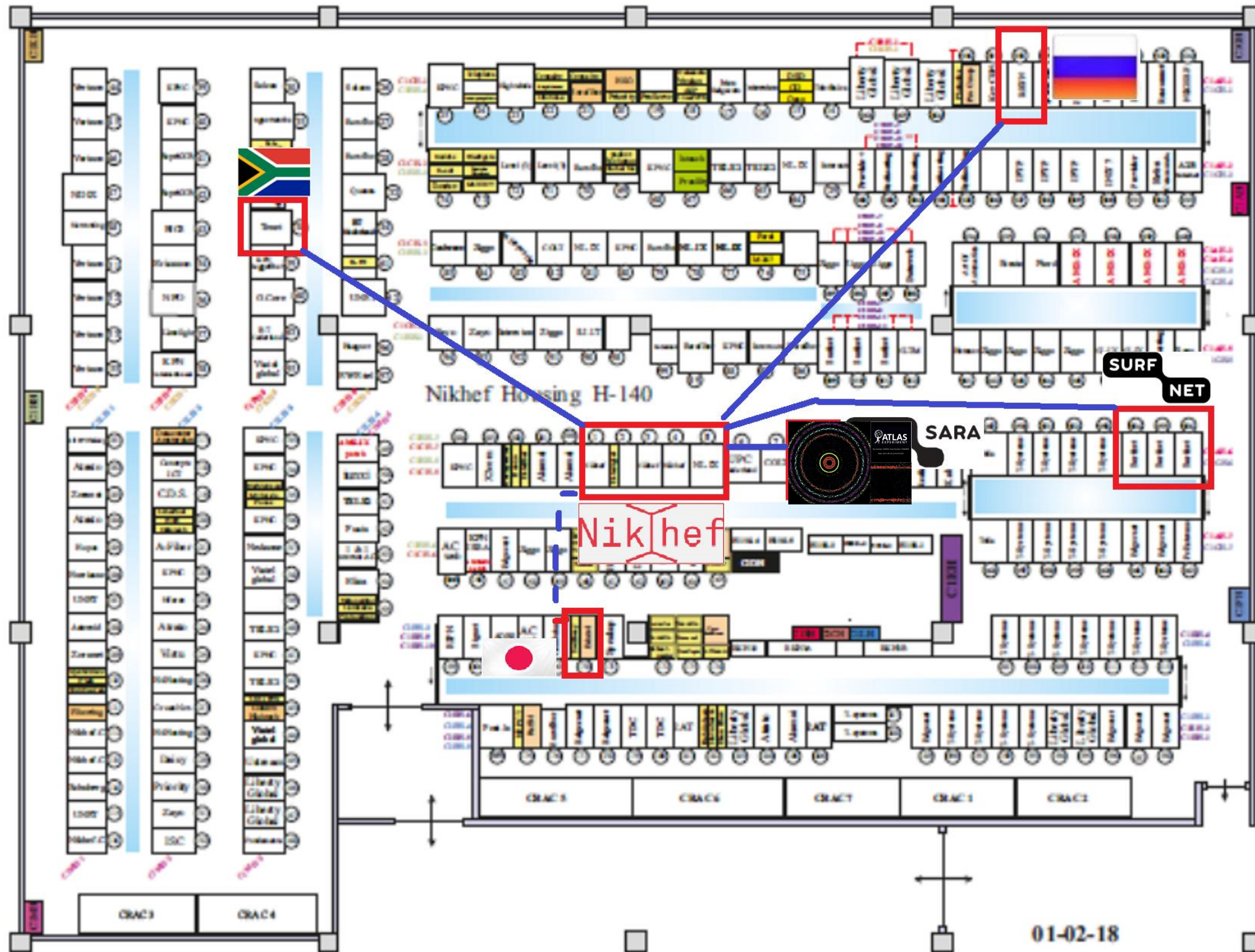
Nikhef  
Data Processing  
Facility  
router 'deel'

een KPN FttH  
PoP in de wijk

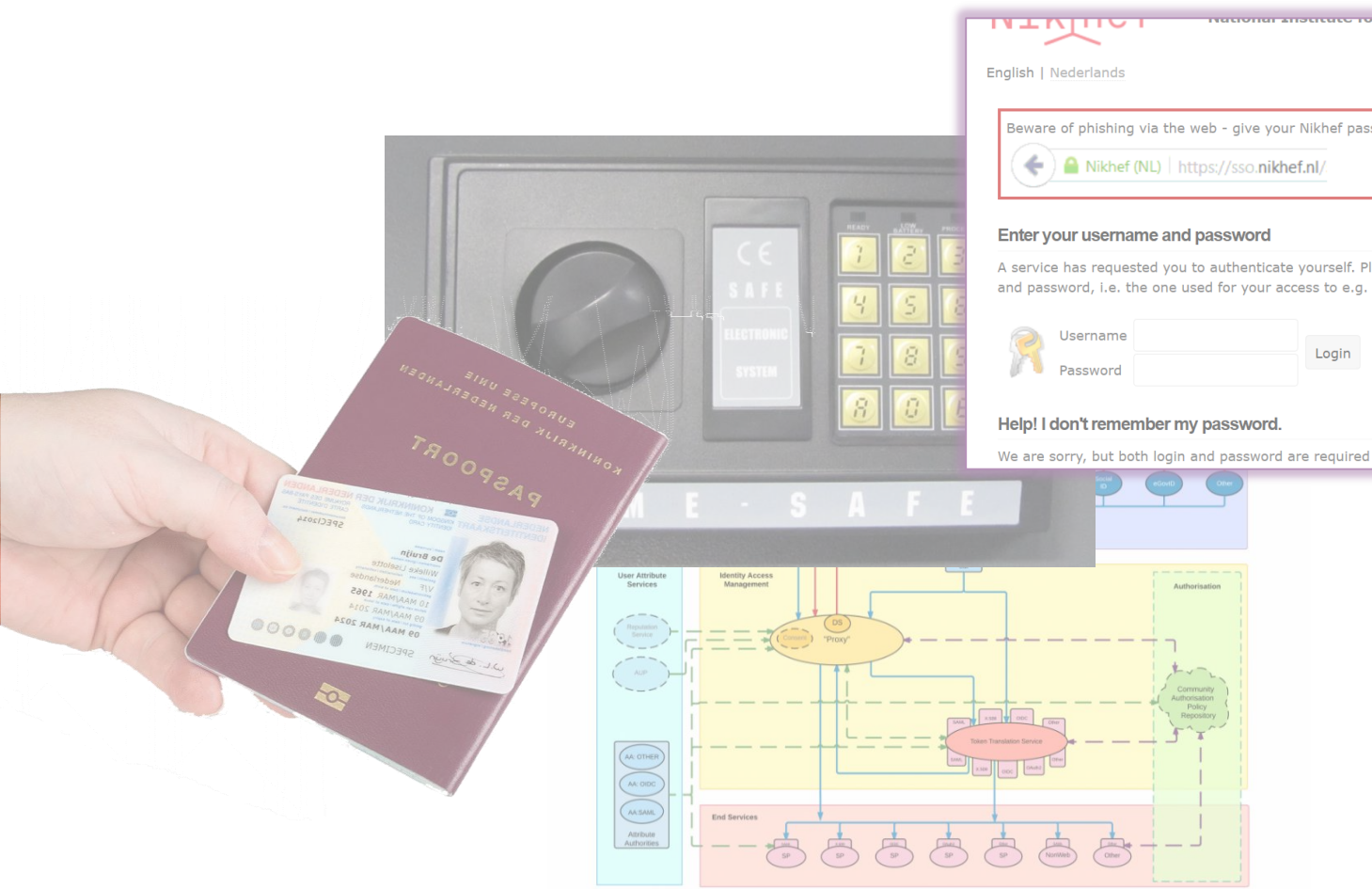


vergelijk:  
VDSL BR straatkast  
voor als je nog  
op xDSL koper zit



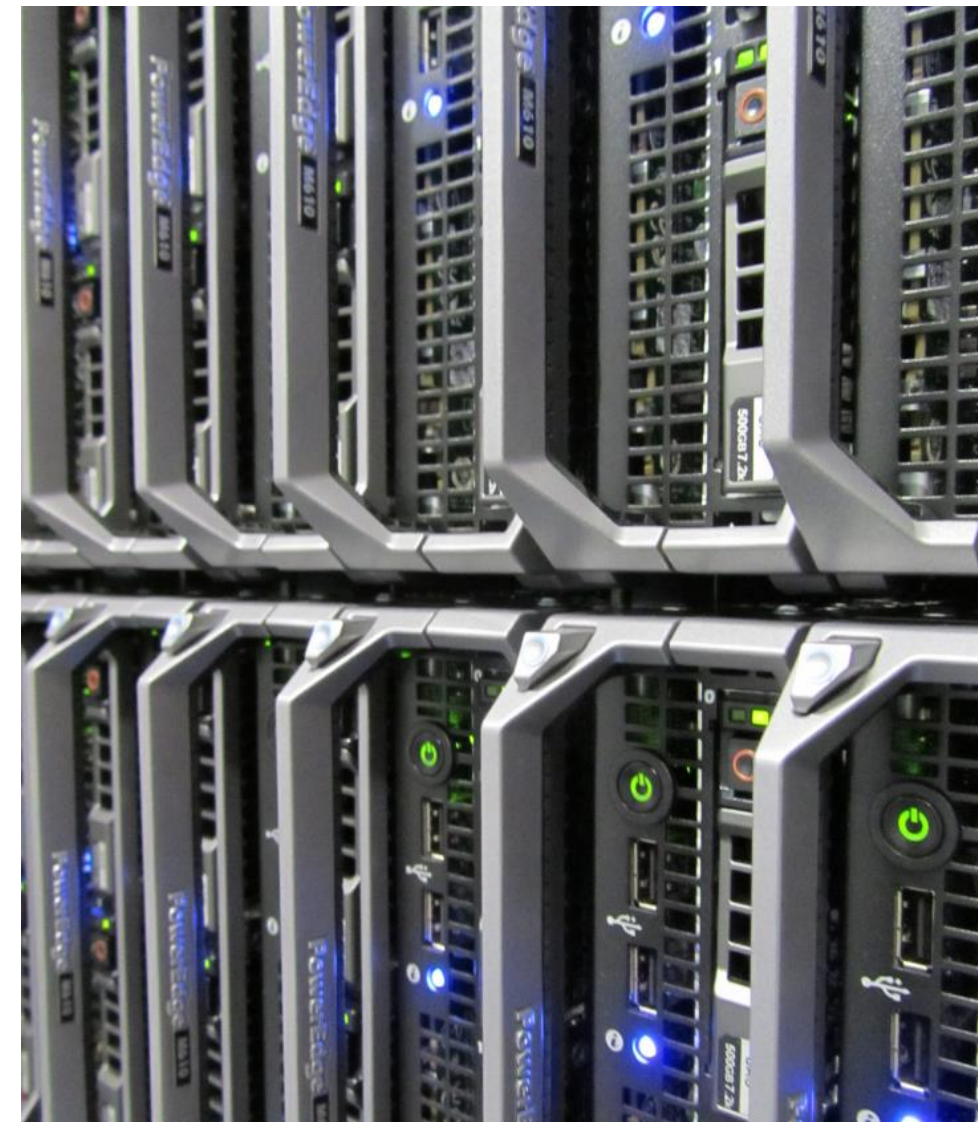


## Samenwerking en beveiliging



we houden bij wie er mag rekenen, welk experiment wat gebruikt, en of onze rekenkracht goed gebruikt wordt ...  
 ... en we niet misbruikt worden om het internet aan te vallen!

## Parallel rekenen



7000 rekenkernen (is ~300 servers) staan dag en nacht te rekenen aan onderzoeksdata

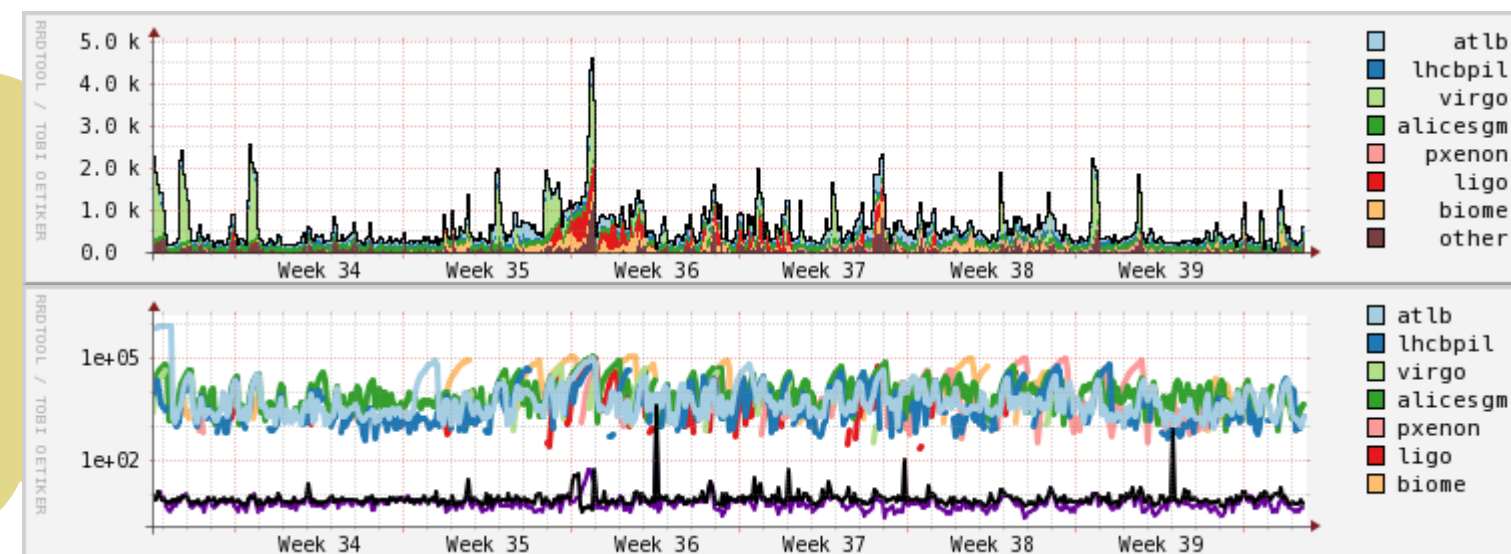
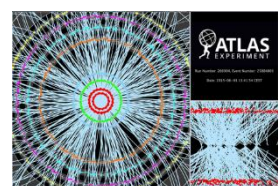
## en parallele data opslag



de ~5 petabyte gegevens bij ons moeten we ook in één dag kunnen lezen, anders staan die 7000 rekenkernen te wachten ... die lezen ieder ~10MB per seconde!



?



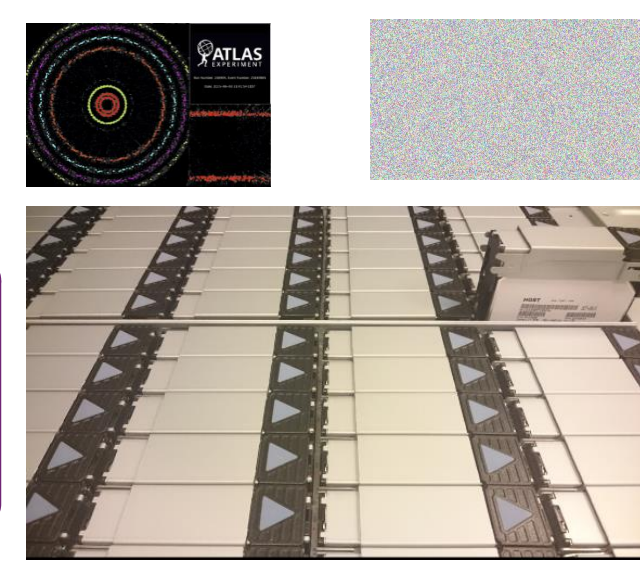
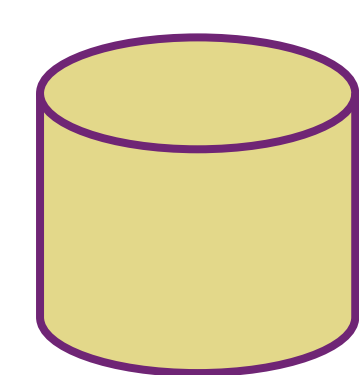
DEF=lhcalice  
DEF=lhcalice  
DEF=lhcalice

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

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

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

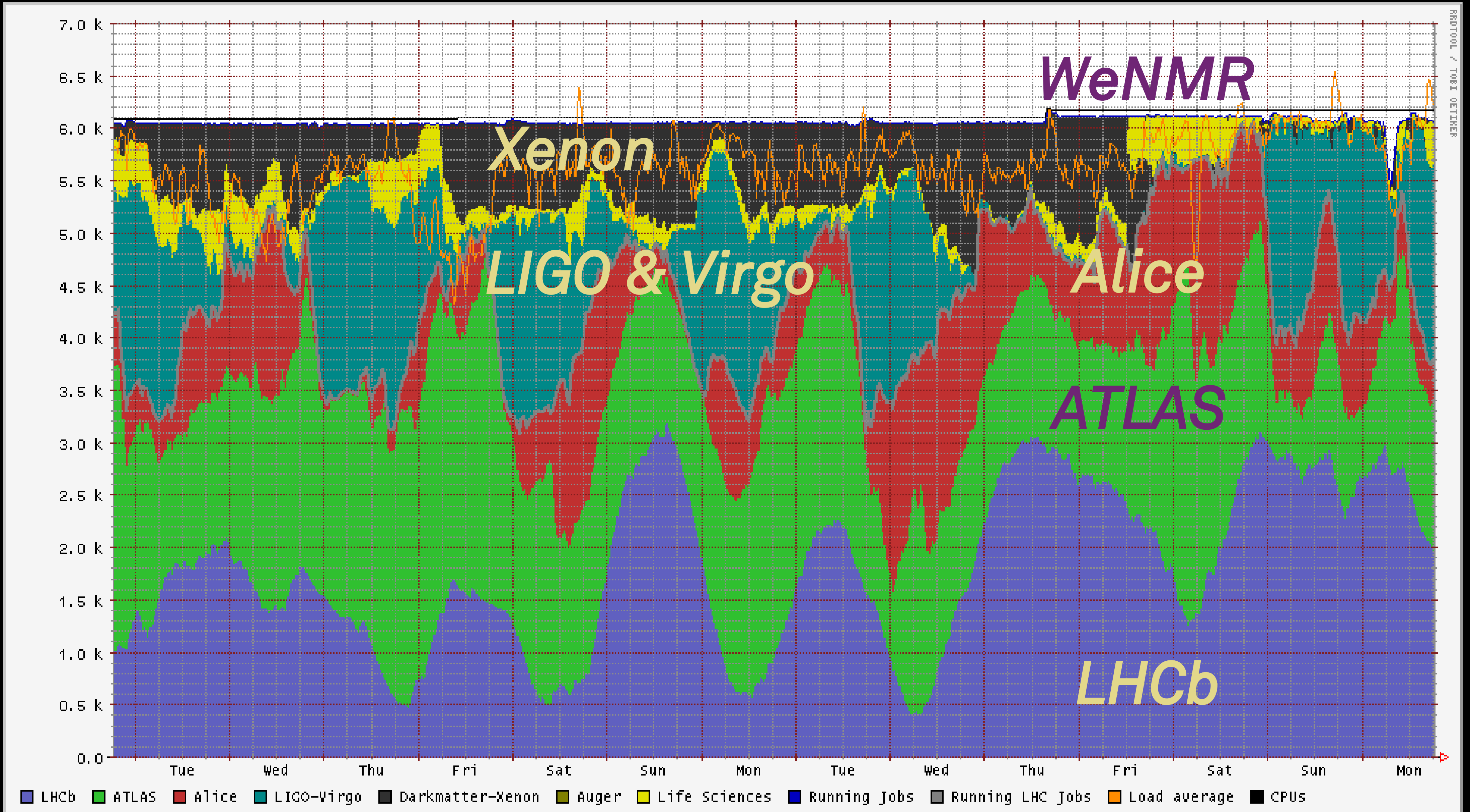
# local groups
GROUPCFG[atlas]      FSTARGET=10   PRIORITY=200  MAXPROC=2200 QDEF=niklocal
```



korf.nikhef.nl:

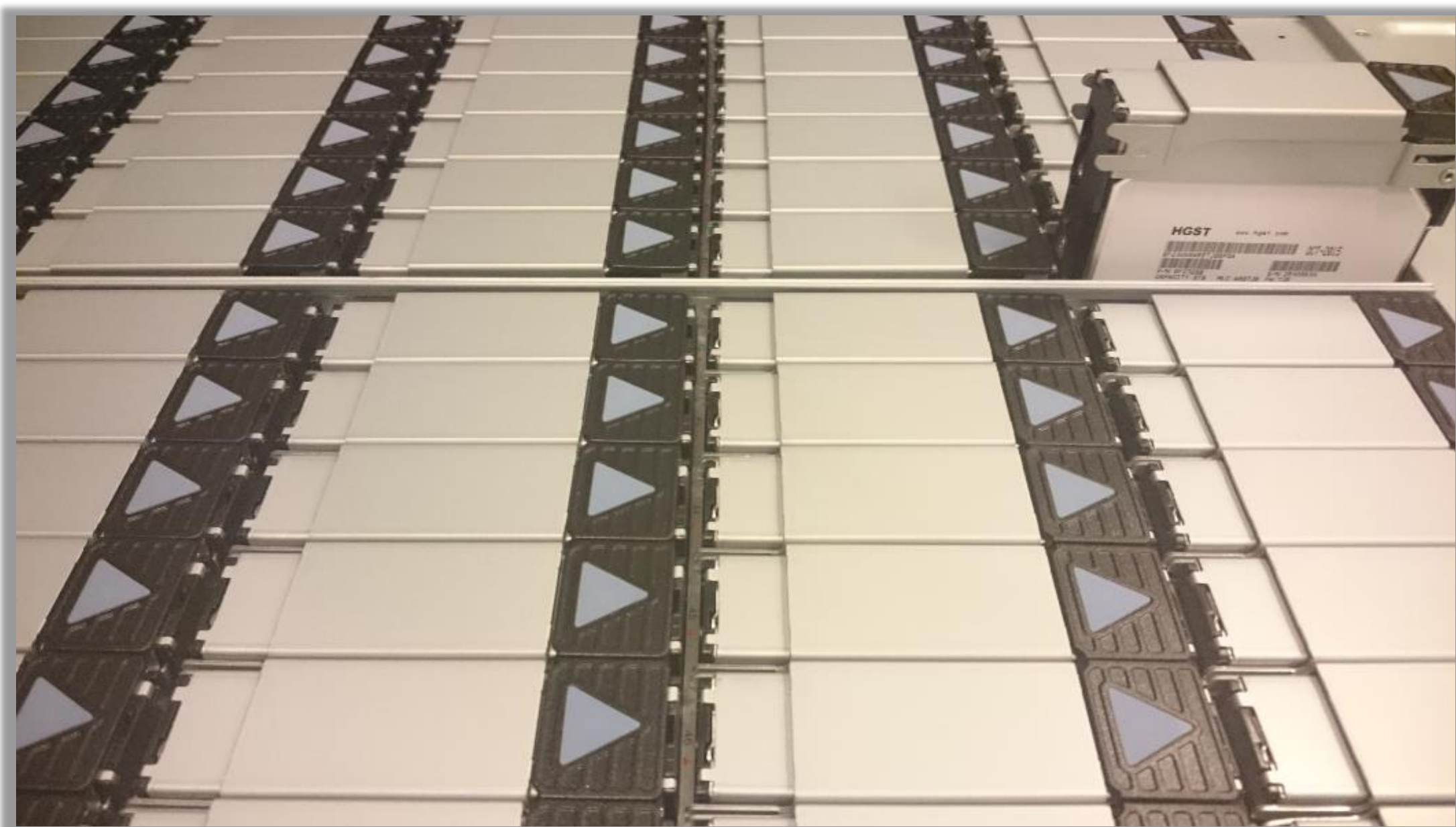
Job ID	Username	Queue	NDS	TSK	Req'd Memory	Req'd Time	S	Elap Time	
33134895.korf.nikhef.n	lhcbpi08	lhcb	1	1	5120m	41:59:57	R	37:46:21	wn-choc-023
33134901.korf.nikhef.n	lhcbpi08	lhcb	1	1	5120m	41:59:57	R	40:04:09	wn-smrt-128
33134908.korf.nikhef.n	lhcbpi08	lhcb	1	1	5120m	41:59:57	R	37:14:29	wn-choc-030
33134917.korf.nikhef.n	lhcbpi08	lhcb	1	1	5120m	41:59:57	R	14:23:42	wn-smrt-072
33135197.korf.nikhef.n	atlb019	atlasmc	1	4	16040	208:00:00	R	183:02:04	wn-mars-018+
wn-mars-018+wn-mars-018+wn-mars-018									
33135883.korf.nikhef.n	atlb019	atlasmc	1	4	16040	208:00:00	R	166:44:22	wn-mars-018+
wn-mars-018+wn-mars-018+wn-mars-018									
33142633.korf.nikhef.n	lhcbpi08	lhcb	1	1	5120m	41:59:57	R	37:30:47	wn-mars-043
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-057
33149220.korf.nikhef.n	lhcbpi08	lhcb	1	1	5120m	41:59:57	R	32:50:19	wn-choc-044
33151669.korf.nikhef.n	lhcbpi08	lhcb	1	1	5120m	41:59:57	R	09:49:53	wn-choc-009
33152704.korf.nikhef.n	atlb019	atlasmc	1	4	16040	208:00:00	R	128:39:13	wn-mars-018+
wn-mars-018+wn-mars-018+wn-mars-018									

# Zo hou je een rekencluster vol:



PROTOCOL / TOBI OETIKER

**DotHill (HGST): 480 TByte gross capacity/4RU**



### Kengetallen van opslag

- capaciteit ('terabytes')
- bandbreedte ('megabyte per seconde')
- aantal 'opdrachten' per seconds ("IOPS")



Daarom kan onze data niet op een USB stick  
– en doet je 'thuis NAS' oplossing het ook niet  
*... hoe leuk ik mijn eigen opslagdoosje ook vind  
van slechts 15 Watt met 16 Terabyte voor € 915 ...*

**5 Petabyte lezen in 1 dag?  
Dat is dus 61 GByte per seconde!  
(en dus ~500 Gbps)**



# En ... hoeveel gebruikt dat dan?



Eén server gebruikt zo'n 260W!

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

en het onderzoeksdatacentrum Nikhef (de 'glazen doos') kan 400kW aan – waar blijft dat dan?

De snelste CPU is voor ons niet altijd de beste (*sorry gamers!*). Want 5 jaar energie en beheer zijn even kostbaar als de server zelf!

**WKO: Warmte Koude Opslag**

*21% van het vermogen is nodig om te koelen*

*maar: we mogen 3500GJoule/jaar (~112 kW-jaar) aan de studenten leveren!*



... want dit is wel leuk, maar niet echt sustainable ...

# Take a T-Byte!



Nikhef

David Groep

davidg@nikhef.nl

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

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