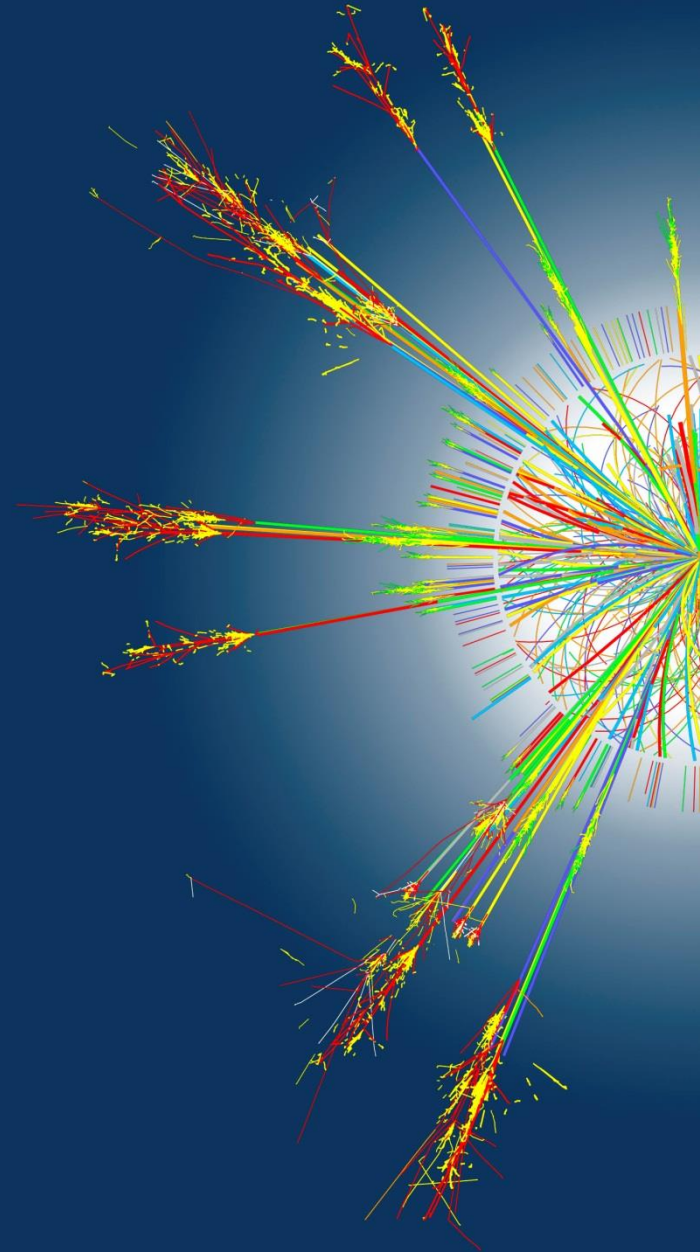
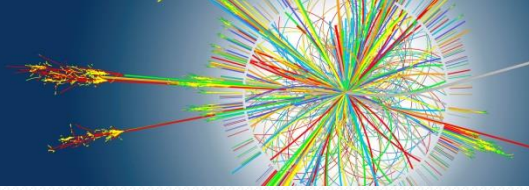


# Analysis of LHC data using globally distributed computing resources

**Oxana Smirnova**  
Div. of Particle Physics  
Lund University  
*3 May 2017*



# CERN physicists don't know everything



What are dark matter and dark energy that pervade our Galaxy?

What is gravity, really?

Why is the electrical charge of the proton equal and opposite to that on the electron?

Why are protons stable, or are they, really?

How many space-time dimensions do we live in?

Are elementary particles really elementary, or do they have structure?

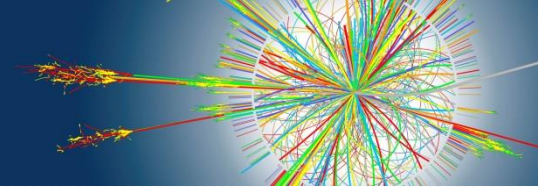
Why are there 3 generations of elementary particles? Why not 4 or 5?

Where did the anti-matter go?

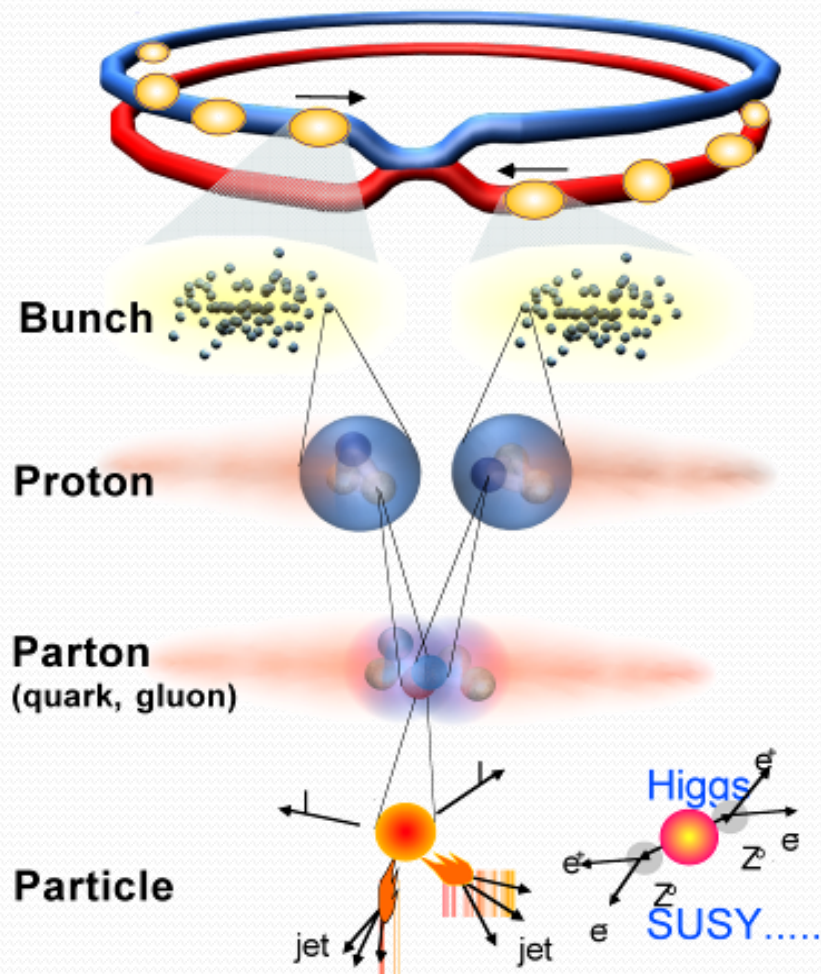
Why are the neutrinos so light, and which is the heaviest?

Are there more states of matter beyond solid, liquid, gas and plasma?

# LHC holds answers



**LHC is the biggest ever machine and data generator**



|                      |  |
|----------------------|--|
| <b>Proton-Proton</b> | 2835 bunch/beam                          |
| Protons/bunch        | $10^{11}$                                |
| Beam energy          | 7 TeV ( $7 \times 10^{12}$ eV)           |
| Luminosity           | $10^{34} \text{ cm}^{-2} \text{ s}^{-1}$ |

Crossing rate 40 MHz

Collisions rate  $\approx 10^7 - 10^9 \text{ Hz}$

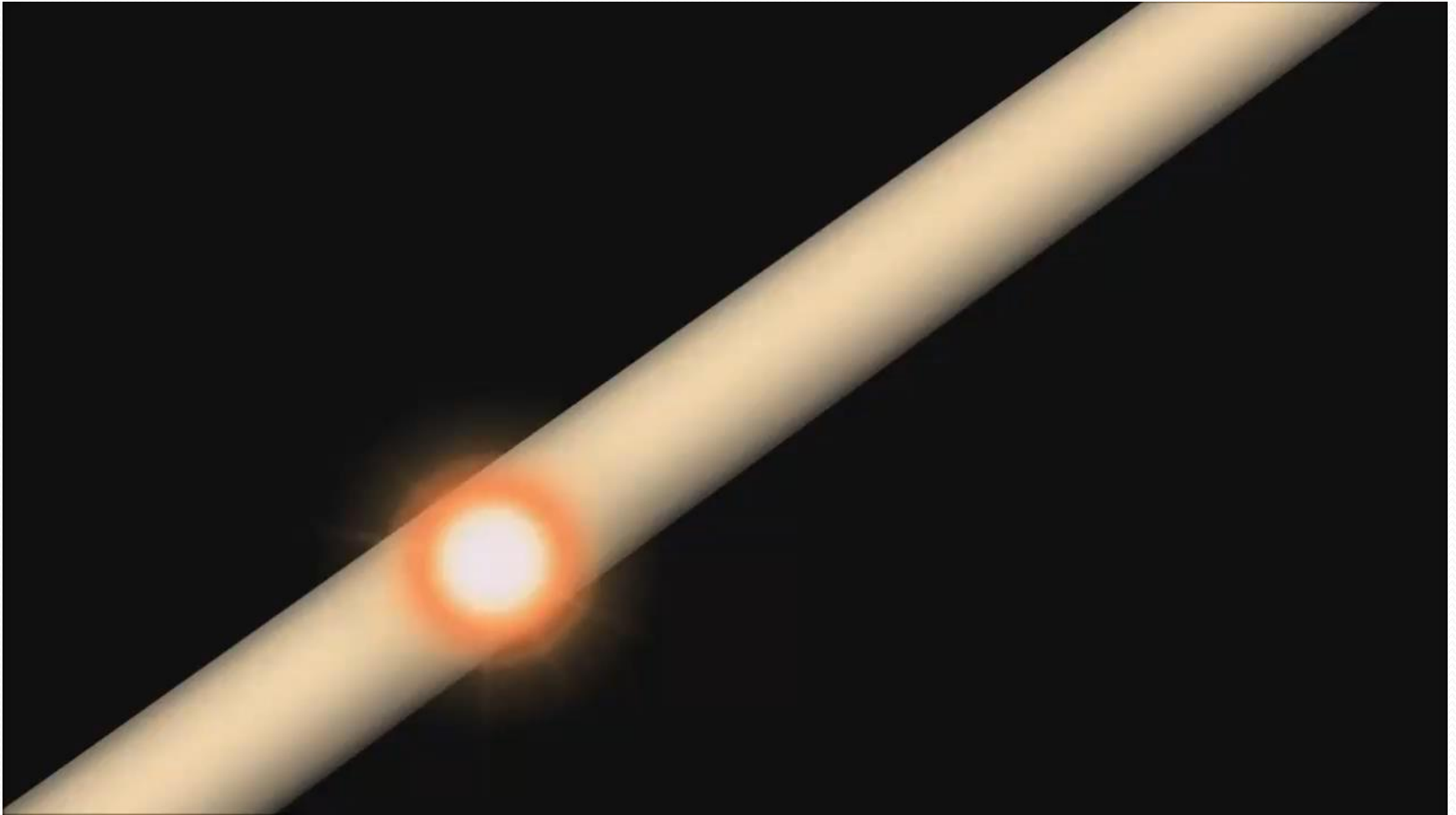
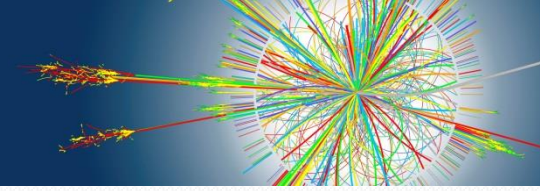
*There are also Pb-Pb and p-Pb runs*

New physics rate  $\approx .00001 \text{ Hz}$

**Event selection:**  
**1 in 10,000,000,000,000**

Graphics by CERN

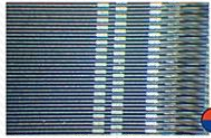
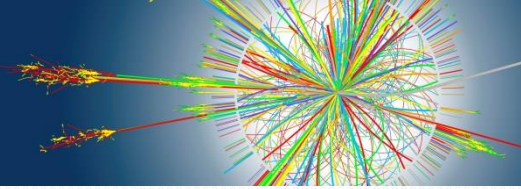
# A data sample: collision event at LHC



*ATLAS video*



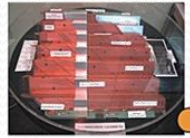
# An instrument at LHC: ATLAS



Stripsensor, en av de innersta halvledarsensorerna. Stripparna, 80  $\mu\text{m}$  breda baksända dioder med bonddrädar till höger.



Stripsensor med skiftregistren som håller signalerna i 2,5  $\mu\text{s}$



EM, elektromagnetiska kalorimetern, uppbyggd av veckat kretskort



Kryogeniska kärlet som innehåller inre solenoiden



Ändring till inre kryogeniska kärlet



Elektromagnetiska kalorimeterna monterade i kryogeniska kärlet  
Foto: CERN



Två bilder av änd-toroidens kryogeniska kär, dels från utsidan, dels från insidan. Utan lock.



Kryogenisk fabrik för tillverkning av Helium II till dipolen



Tryckkärl i den kryogeniska fabriken  
Foto: CERN



Provlinje för dipolrör  
Foto: CERN

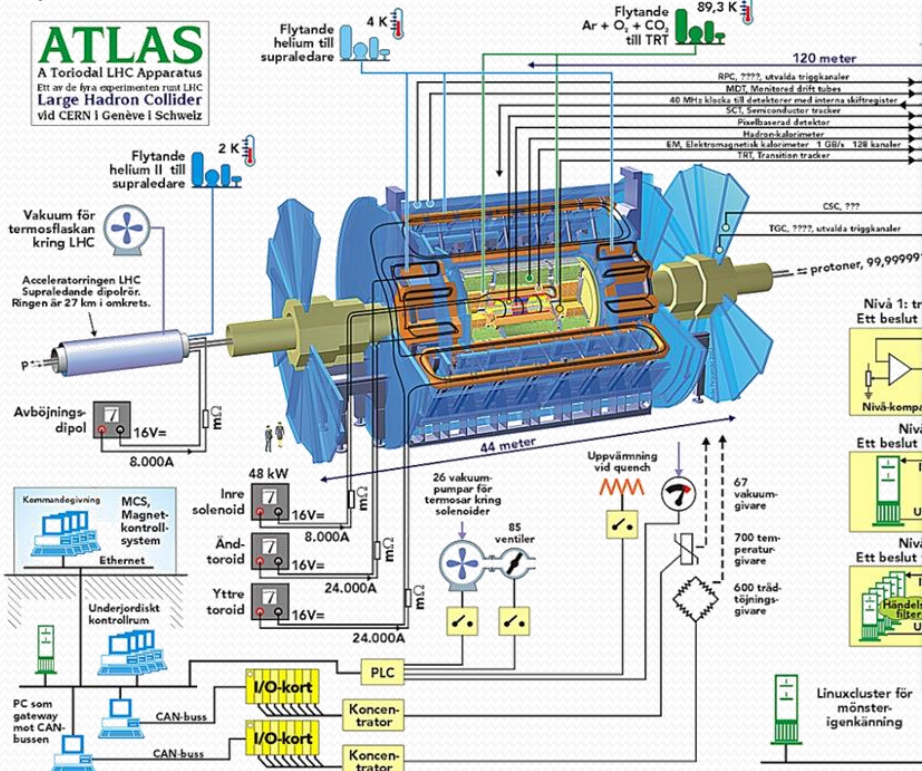


Två dipolrör kopplas ihop



Dipolrörets ände. Vakuumkärlet ytterst, sedan det heliumfyllda kärlet. Den massiva biten i mitten är magneten, med de två strålrörerna sida vid sida.

**ATLAS**  
A Toroidal LHC Apparatus  
Ett av de fyra experimenten runt LHC  
Large Hadron Collider  
vid CERN i Genève i Schweiz



- Think of a high-resolution photo camera 44 m long
- Built and operated by scientists from 38 countries
- Produces so much data, no one single data center can accommodate it
- There are 4 such instruments at LHC!



Supraleadarna till yttre toroid ingjutna i glasfiber, men utan aluminiumhölje



Vakuumkärl i titan till yttre toroid. Supraleadarna ska läggas i och en övre halva svetsas på.



Gigabitswitch



En array av linuxmaskiner

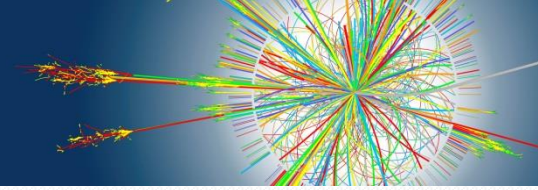


Forskarnas arbetsstationer



Det som alla väntar på: Higgs-sönderfallet (gula linjer)  
Foto: CERN

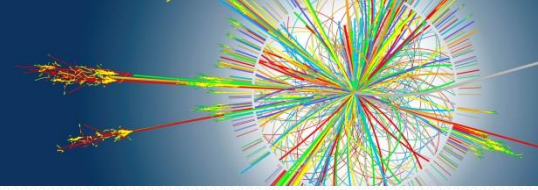
# LHC produces Real Big Data



- Number of read-out channels:  $\sim 15$  million
- Event rate: 40 MHz
- At 1 byte per channel per event, we'd need to write at the speed of  $\sim 600 \text{ TB/s}$ 
  - And process it all...
  - Luckily, not all channels get fired at once
- **But** most events are “normal” and thus not interesting
  - Much like mowing a field hoping to get an orchid in a haystack of daisies
  - E.g., the Higgs production rate is only  $\sim 0.1 \text{ Hz}$ , detection rate is even less
- Most of our data is “background noise”!
  - This is true for other sciences, and is characteristic of **Big Data**



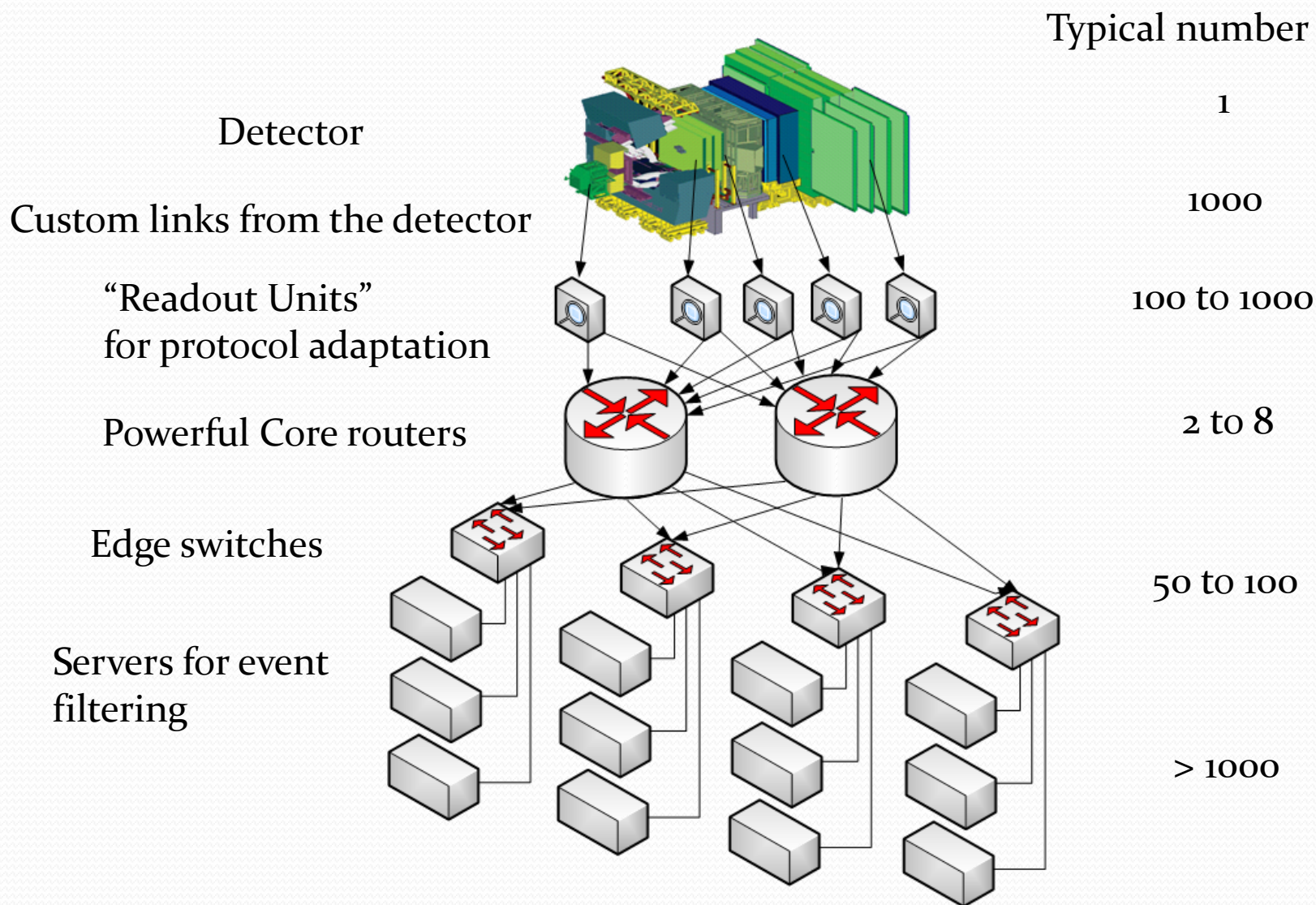
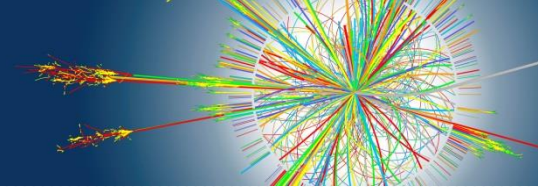
# Do we need to store all this data?



- Luckily, “interesting” events have some obvious signatures
  - For example, particles produced with high energy at large angles w.r.t. the beam direction indicate that something unusual might have happened
- We don’t need high granularity to detect such signatures, thus we can quickly distinguish them
  - So we **trigger** data taking only if such signatures occur
  - Trigger involves software making a yes/no decision
  - To do it fast (real fast), serious hardware is needed
    - Mind it, this is needed before even recording data!

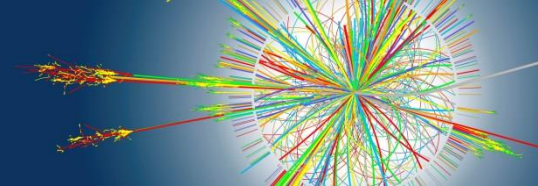


# Trigger is a powerful system itself





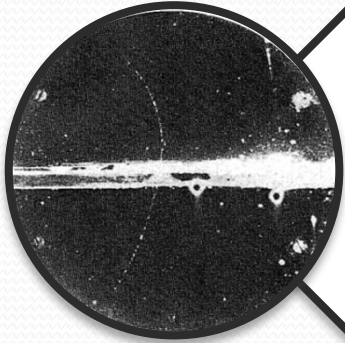
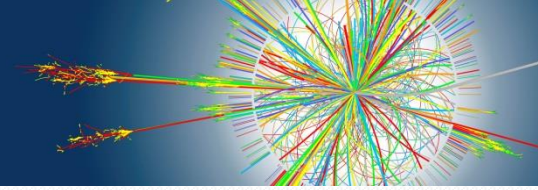
# LHC experiments and their triggers



|       | #<br>Trigger | Level-0,1,2<br>Rate (kHz) | Event<br>Size (Byte)               | Network<br>Bandw.(GB/s) | Storage<br>MB/s(Event/s)          |
|-------|--------------|---------------------------|------------------------------------|-------------------------|-----------------------------------|
| ALICE | 4            | Pb-Pb 5<br>p-p 1          | $5 \times 10^7$<br>$2 \times 10^6$ | 25                      | 4000 ( $10^2$ )<br>200 ( $10^2$ ) |
| ATLAS | 3            | LV-1 75<br>LV-2 6         | $1.5 \times 10^6$                  | 10                      | ~1000 (>400)                      |
| CMS   | 2            | LV-1 75                   | $10^6$                             | 200                     | ~1000 ( $10^3$ )                  |
| LHCb  | 2            | LV-0 1000                 | $6 \times 10^4$                    | 55                      | >600 ( $1.2 \times 10^4$ )        |

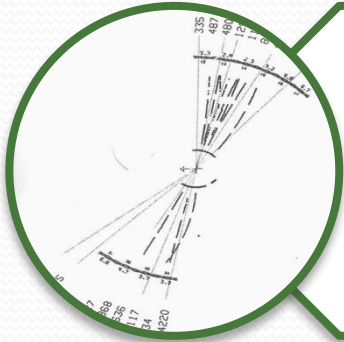
Slide from ISOTDAQ16 TDAQ for LHC - Niko Neufeld, CERN

# We still have a lot of data to analyse



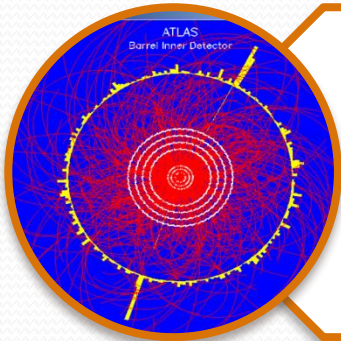
## A discovery in 1930-ies

- ~2 scientists in 1 country
- pen-and-paper



## A discovery in 1970-ies

- ~200 scientists in ~10 countries
- mainframes

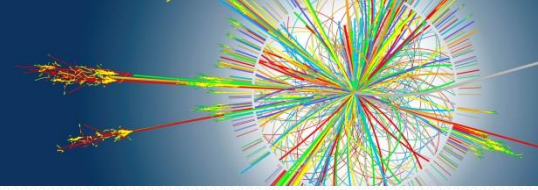


## A discovery today

- ~2000 scientists in ~100 countries
- **Grids**

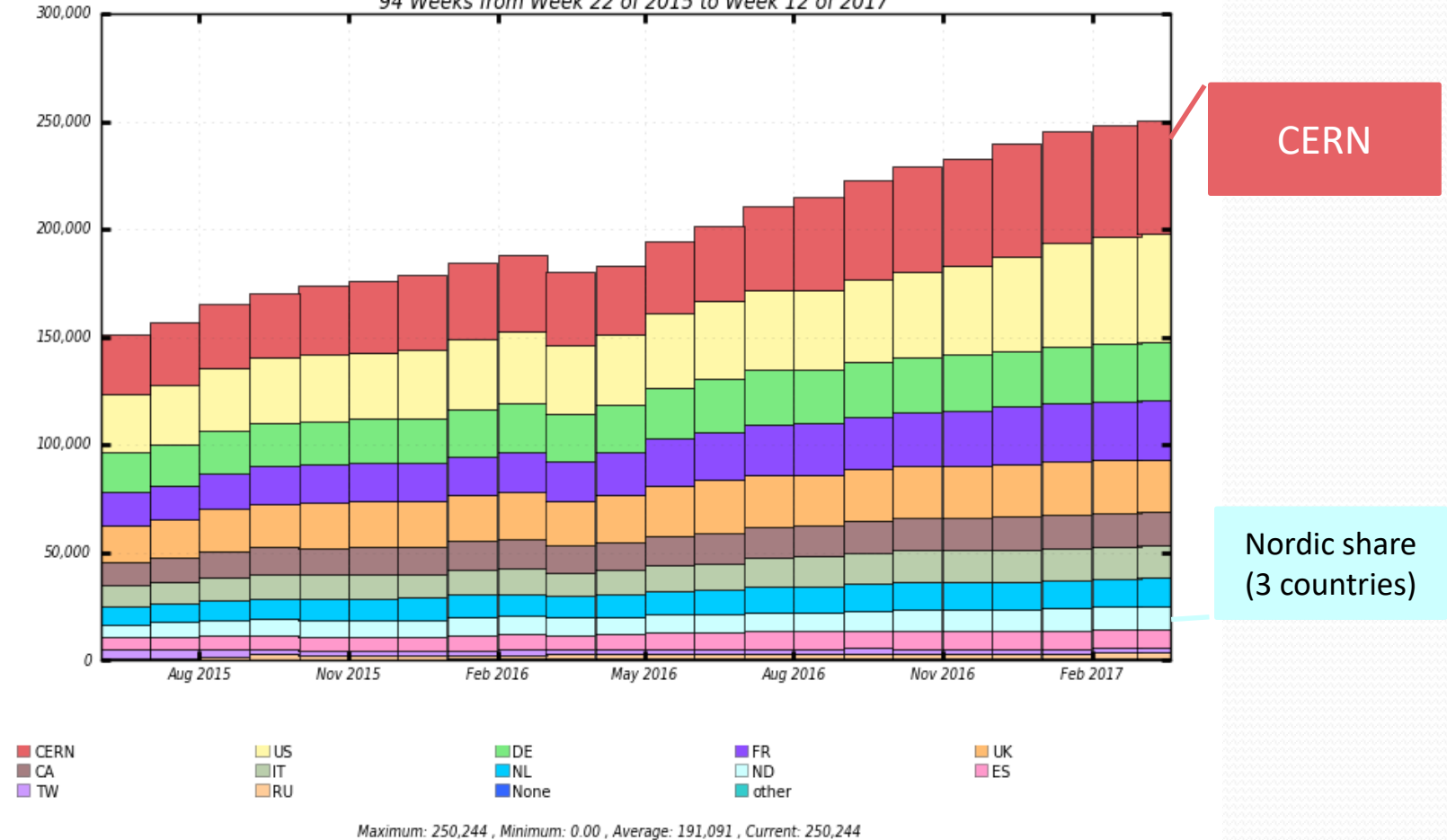
- Relatively simple algorithms
- Very large samples
- Analysis is easily split in very many computing jobs
- Easy to distribute data and jobs
- **Distributed computing:**  
**Grid of data centres** (as in power grid)

# Some real data volumes



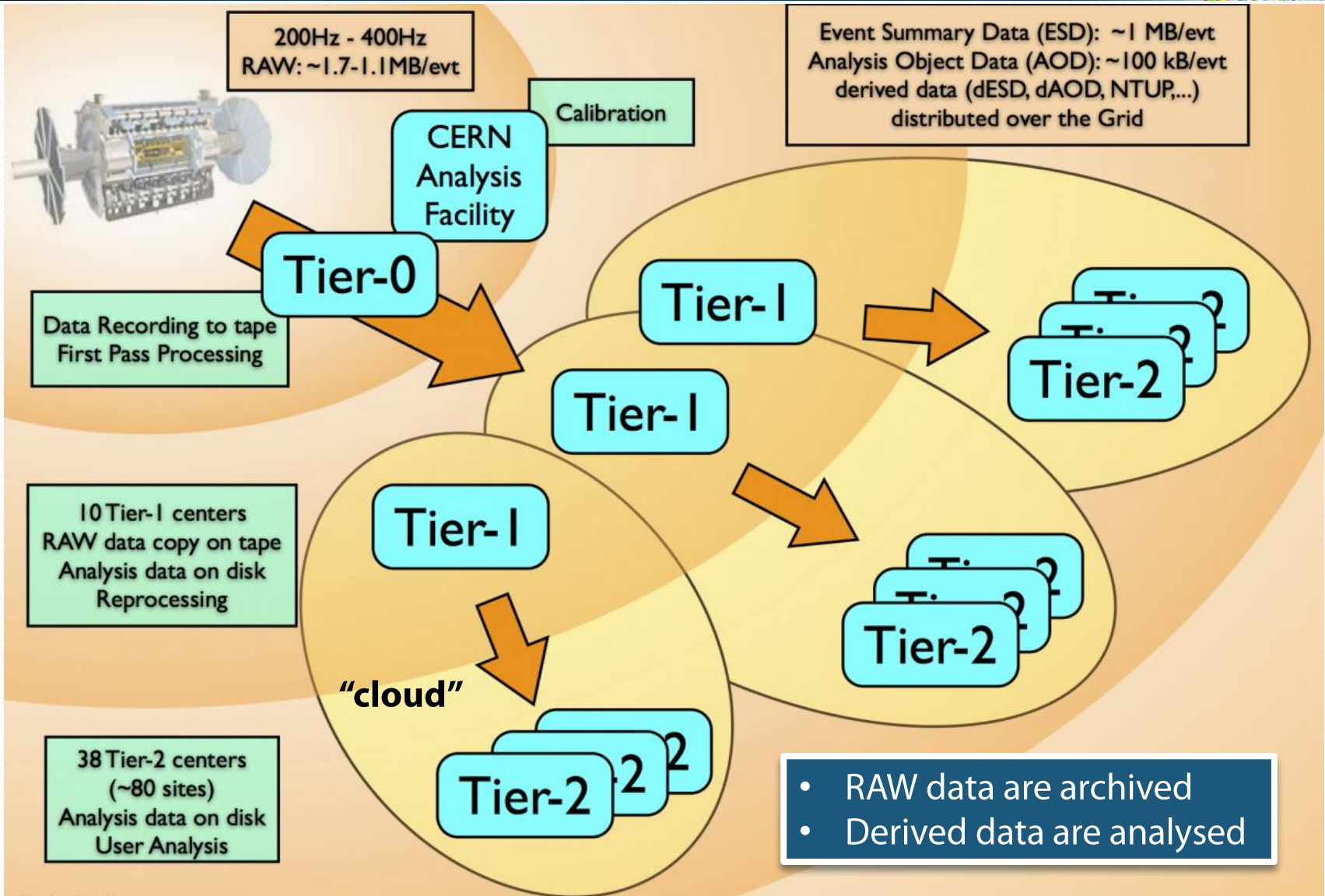
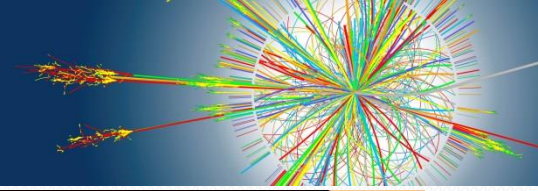
## Number of Physical Bytes (in TBs)

94 Weeks from Week 22 of 2015 to Week 12 of 2017



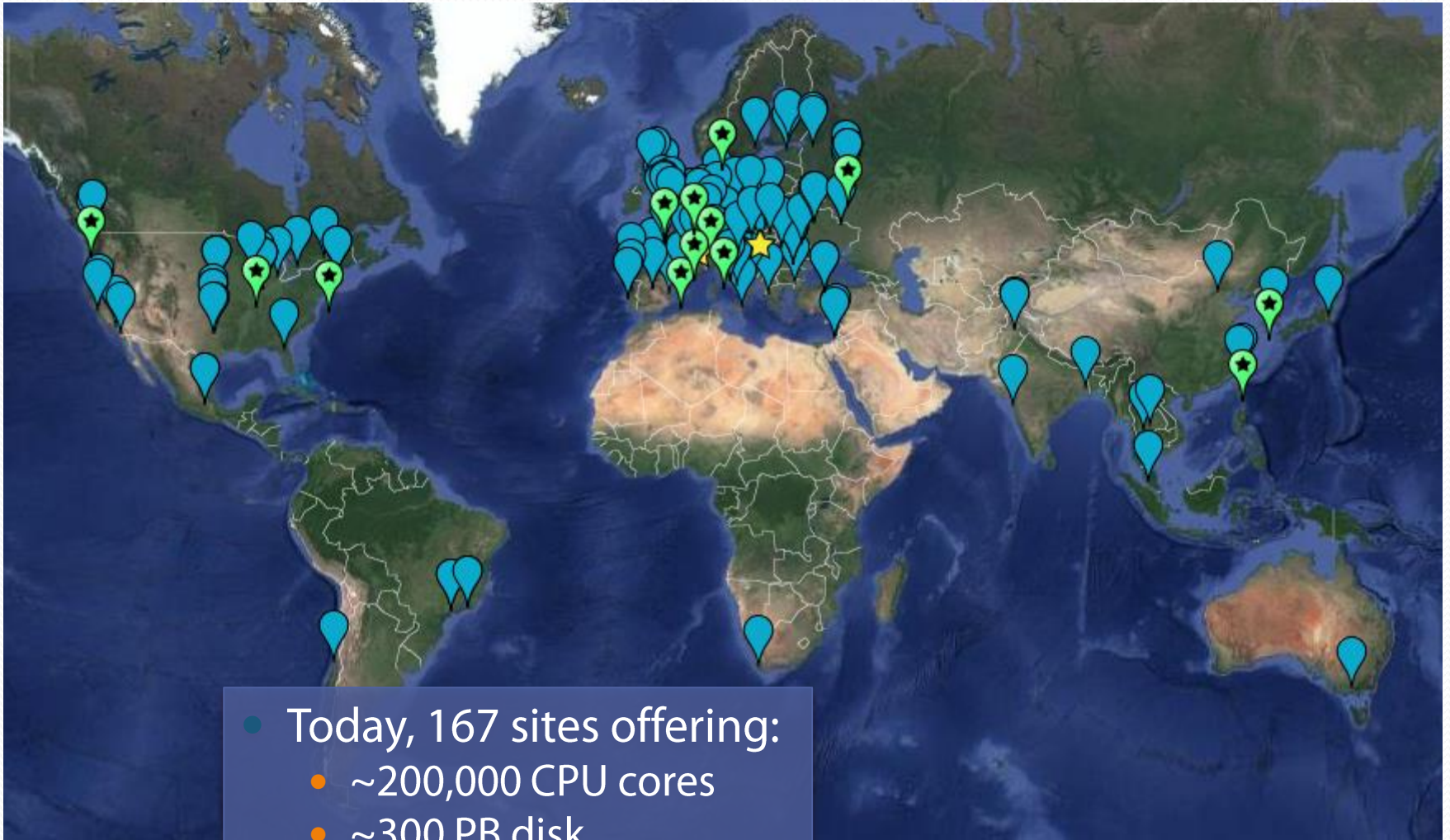
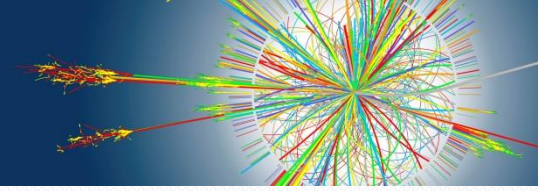
Data stored by the ATLAS experiment alone so far, by region  
(comparable with annual Facebook uploads)

# Stored and analysed around the World



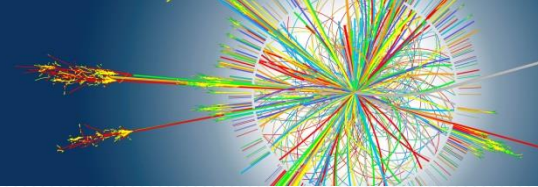


# LHC Computing Grid sites



- Today, 167 sites offering:
  - ~200,000 CPU cores
  - ~300 PB disk
  - ~380 PB tape

# For comparison: Hazel Hen



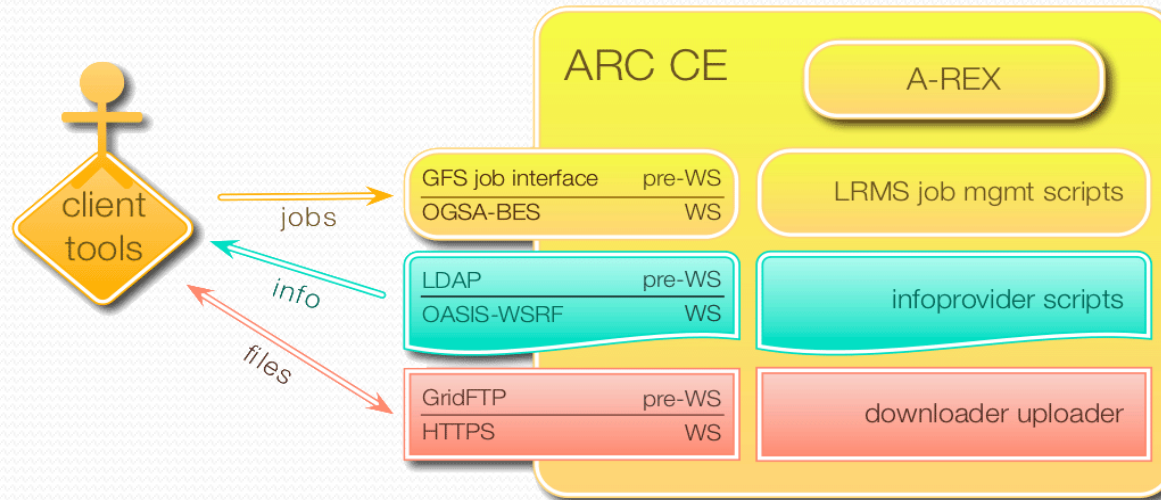
## Europe's fastest supercomputer

- Cray XC40
- 7.42 Petaflops
- **185 088** Intel Haswell E5-2680 v3 compute cores, 128 GB memory/node
- **10 PB** disk capacity
- 3.2 MW power consumption



*Supercomputers: focus on data generation,  
less on data processing*

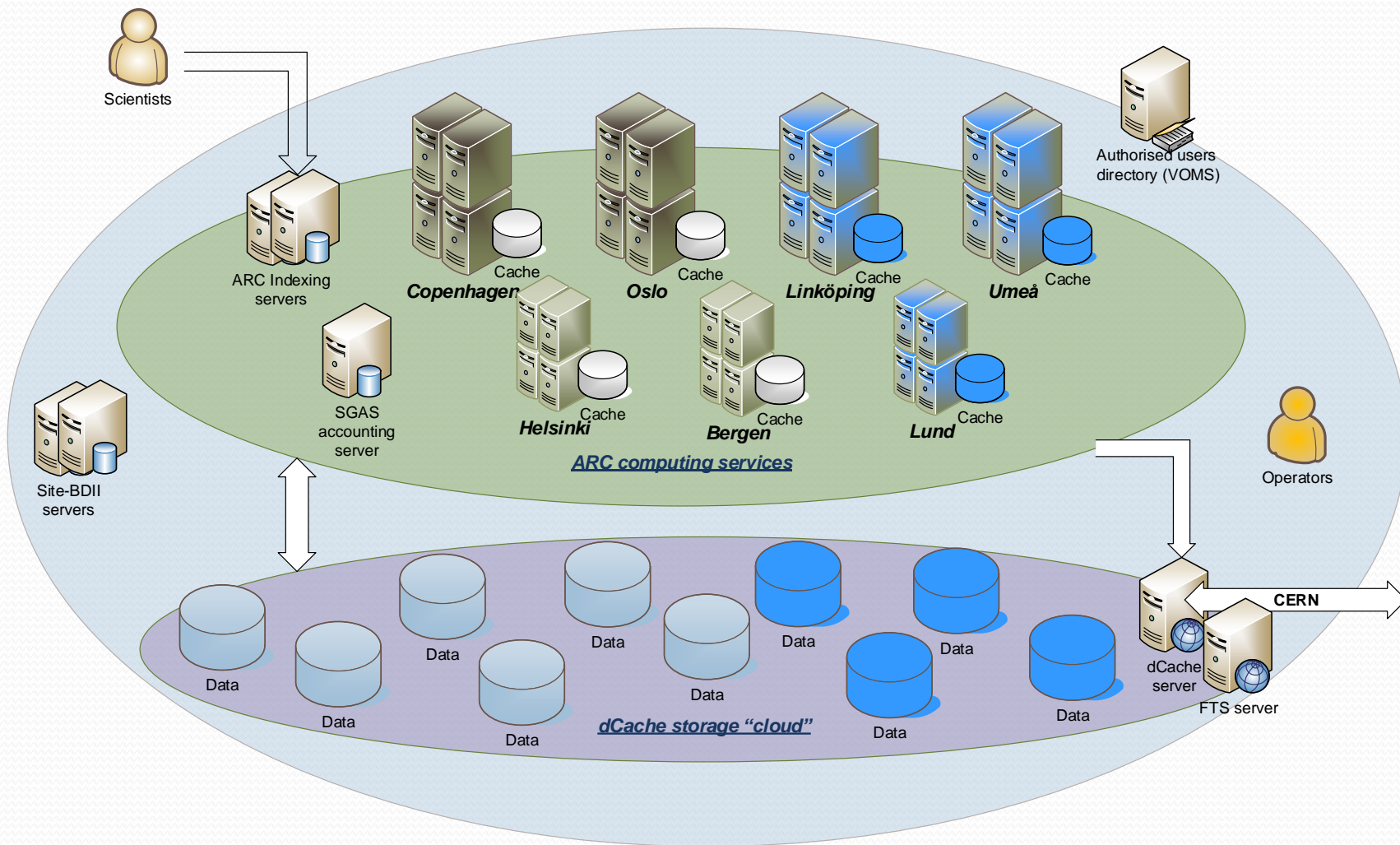
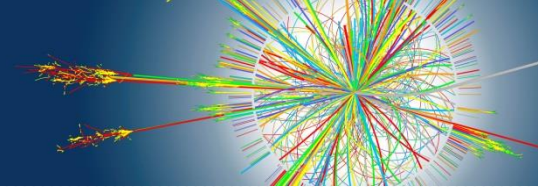
# We develop Grid solutions ourselves



- Big businesses are all doing Clouds, so we do things ourselves
- ARC– the key component of the Nordic Grid infrastructure
  - File handling on behalf of the user
  - Universal front-end for different batch systems
  - Status information publishing
  - Standard and custom interfaces
  - Developed in Nordic countries (in Lund, too), used by ~20% of LHC Grid
  - Is used by many other scientists (not just LHC)

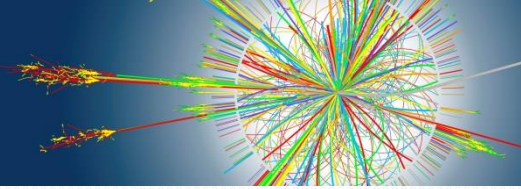


# Nordic Grid infrastructure: NDGF-T1





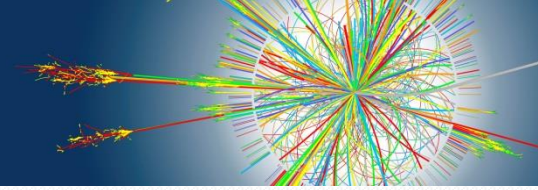
# How is Grid different from Clouds?



- Technical view:
  - In Clouds, environment is tailored for your tasks
  - In Grids, your tasks have to be tweaked for different environments
- The economy, stupid:
  - Clouds are capitalism
  - Grids are communism
- The reality:
  - We plug Cloud resources into our Grid



# Summary



- All sciences face rapid increase in digital data volumes
- LHC developed a working solution for very large volumes of data
  - Do not try to store and process data in one place: **share and distribute** instead
  - Use own open source software
  - This allows LHC to achieve scientific results almost instantaneously, despite huge data sets
- LHC solutions can be generalized to other data-intensive sciences and even industry