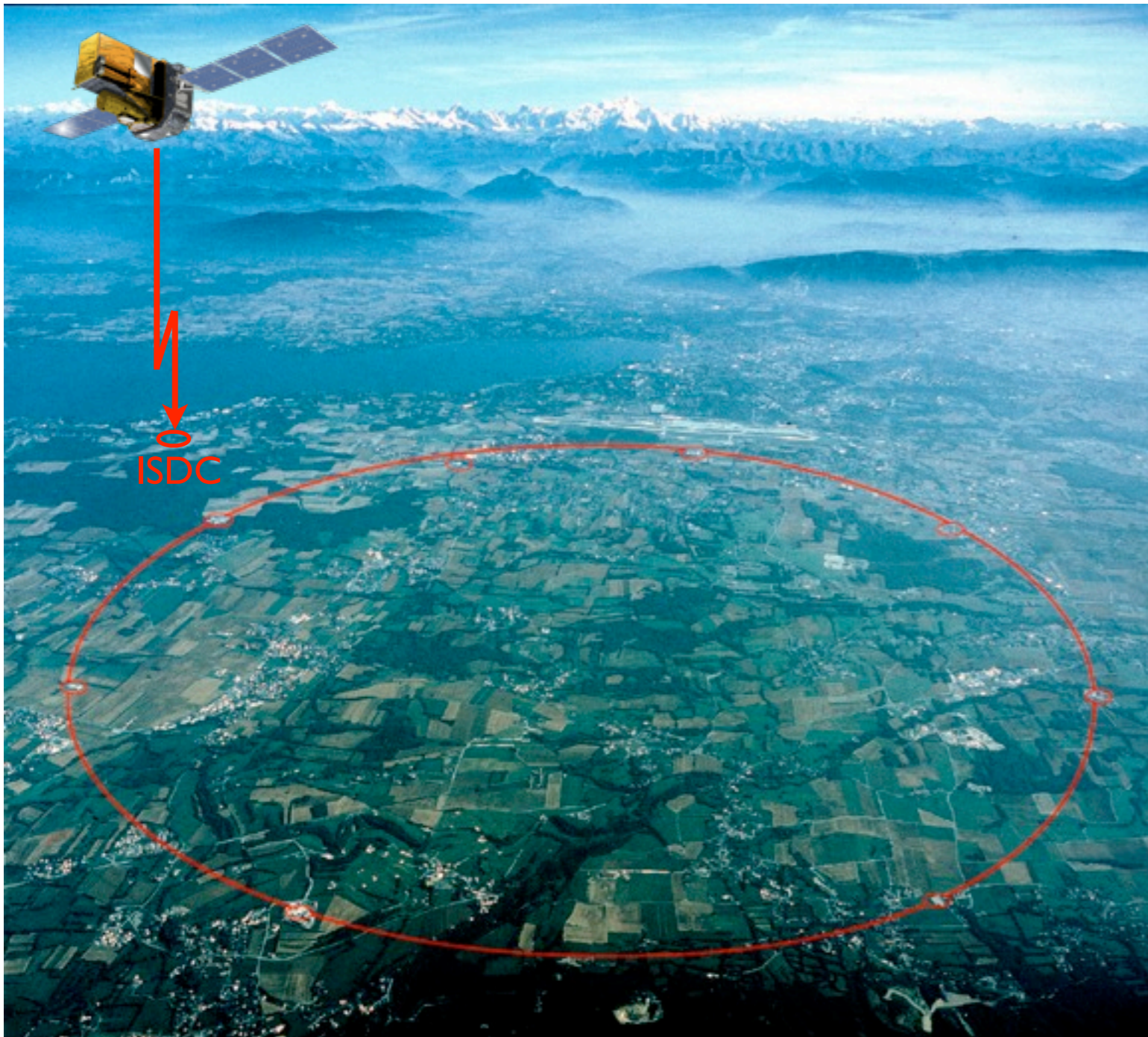


Models for open observatories and the multi-PB scale



Roland Walter
ISDC Data Centre For Astrophysics
University of Geneva

ASPERA computing workshop
Barcelona May 30, 2011

- High-energy observatories
 - specific requirements
 - computing model

- Multi PB age (CTA),
can we follow the same path ?
 - Petabytes, what does it means beyond 2014
 - do we need to change anything fundamental ?

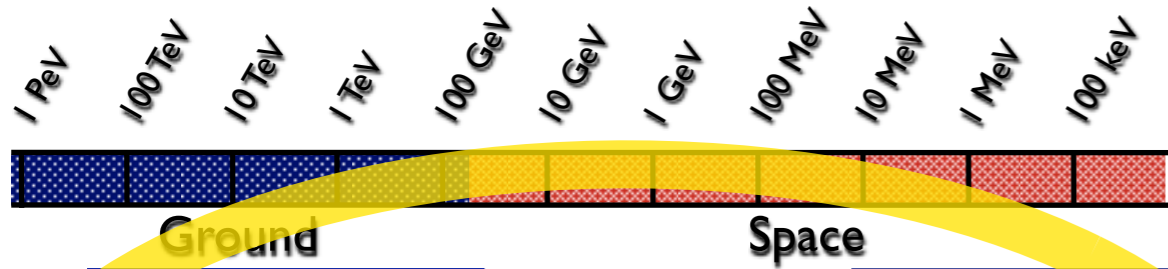
Open observatories

high energies: XMM-Newton, INTEGRAL, RXTE, Chandra, Fermi, Swift, Suzaku, ...,
other wave-bands: VLT, ..., HST, ..., Herschel, ALMA, VLA, ATCA, ...

Observatories are somewhat different from experiments:

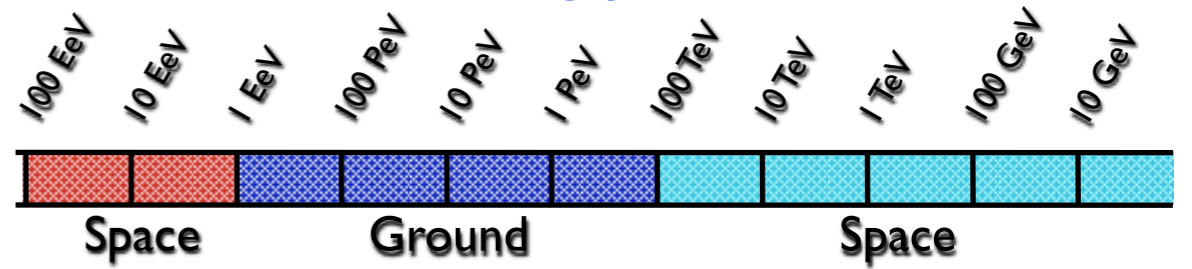
- they detect many (≈ 1000) sources
- they are used by **observers** with many various science goals and analysis requirements
- most observers use a small fraction of the data
- **archive users** are important, for a long time, also post mission

Detecting gamma-rays

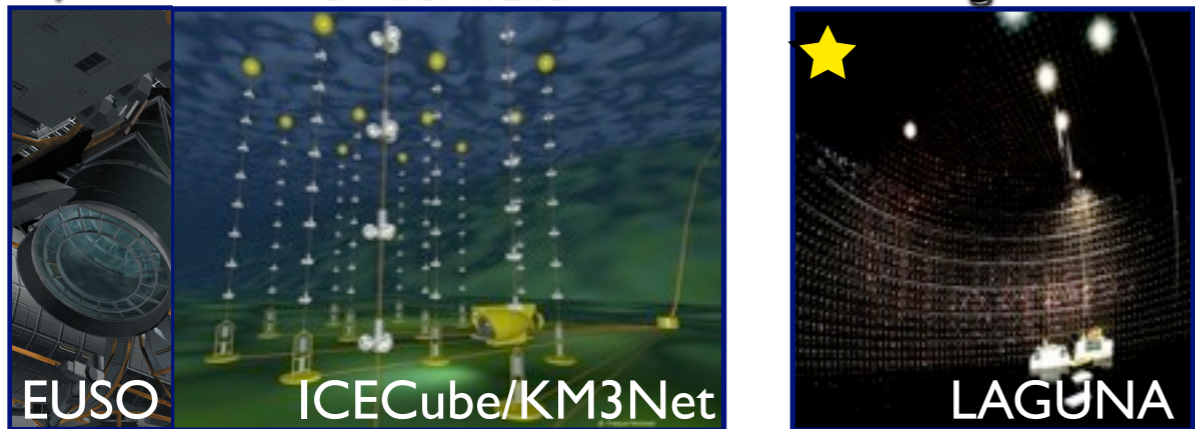
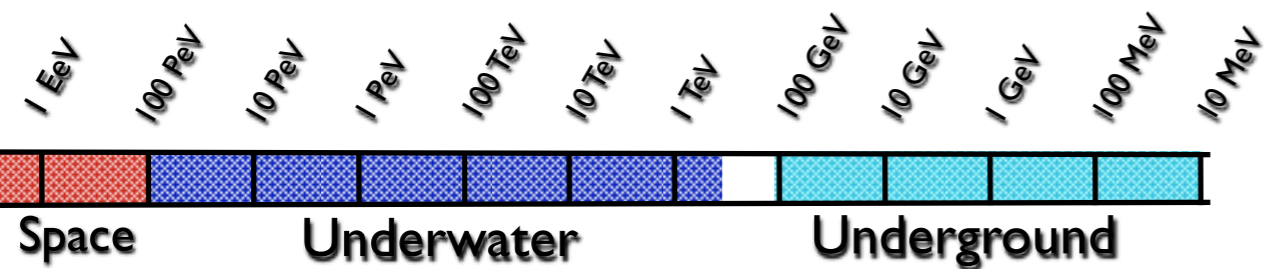


Observatories

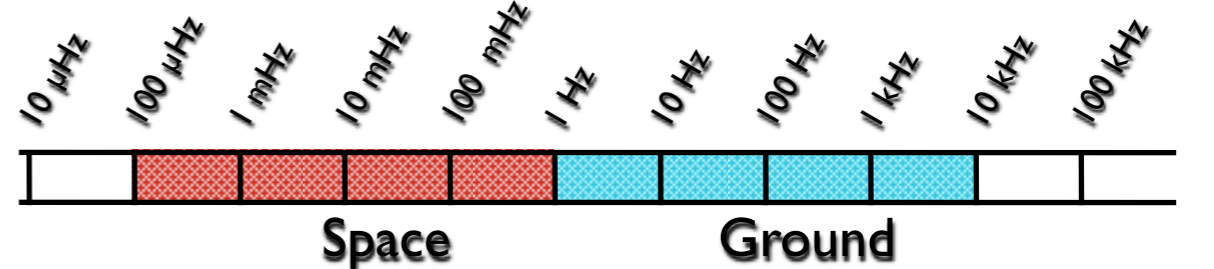
Detecting particles



Detecting neutrinos



Detecting gravitational waves



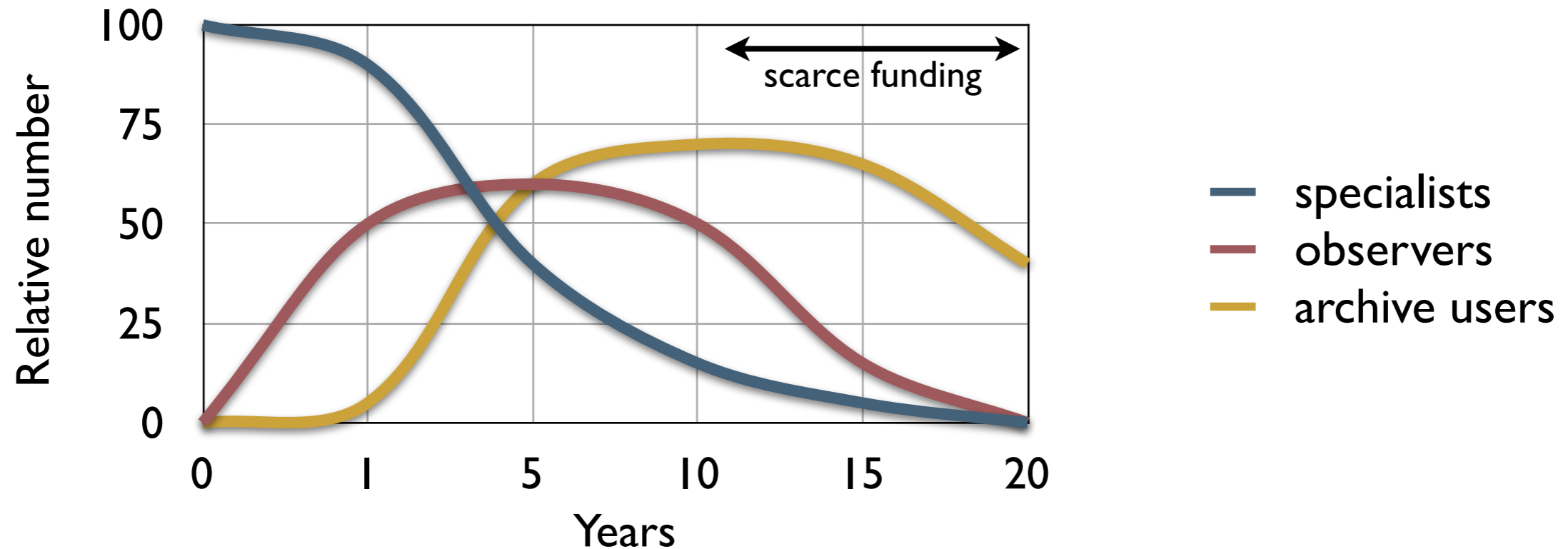
Observatory users

Various types of users

- **archive users:**
need to obtain high-level products made for their own needs
- **observers:**
need to generate high-level products themselves and to understand the limitations of their analysis
- **specialists:**
start from raw data, often members of the observatory team

On the long run the user community is important to support the observatory operations. One should make all of them happy !

Observatory life-cycle



Needs:

on-the fly analysis

corrected data, analysis software, simulation results

everything

High-energy observatory model



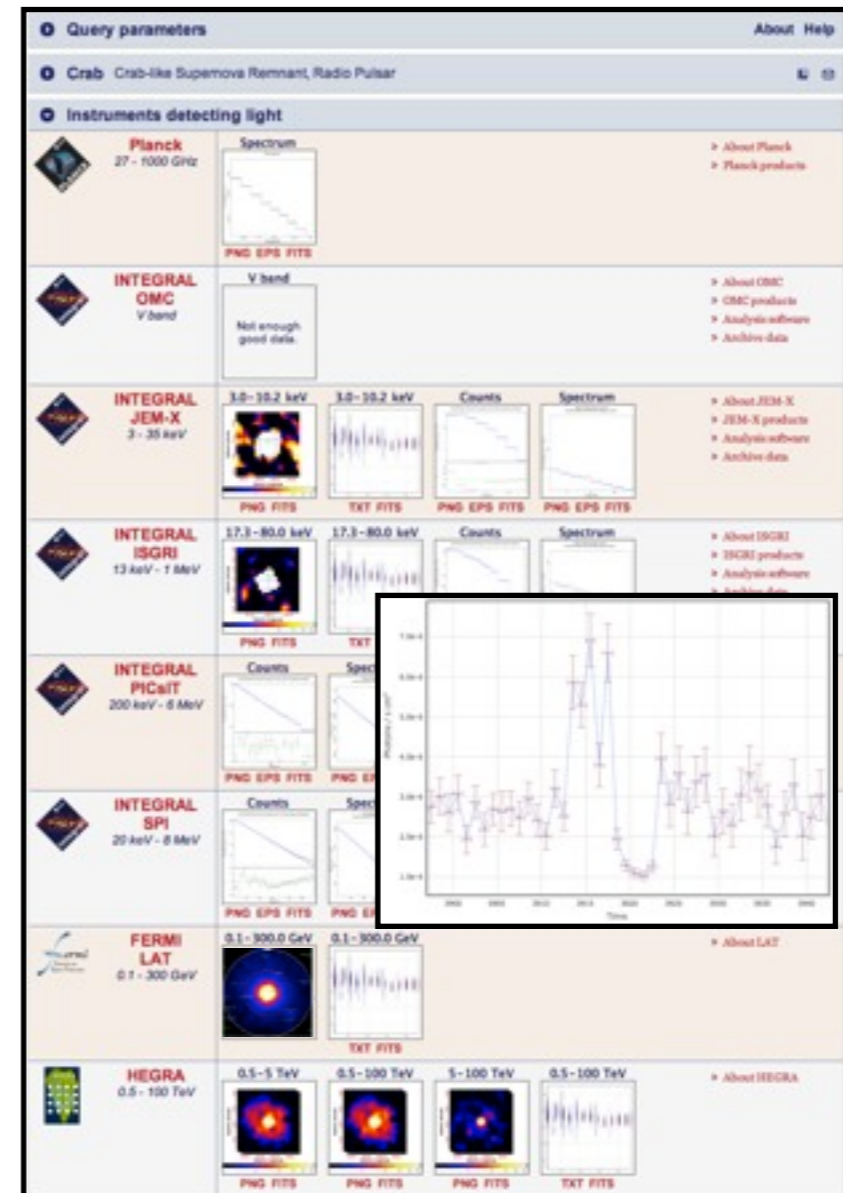
FTOOLS

A General Package of Software to Manipulate FITS Files

Portability, maintainability, simplicity

corrected data, analysis
software, simulation results

Current typical data size: 100TB



on-the fly analysis

Observatory in the PB age (CTA)

Can we follow the same path ?

- Standard data formats ✓
- Simple software technologies ✓
- Data from all levels available ?

Storage market

We just performed a study of the PB storage market (for the FACT telescope, and beyond):
hierarchical/disk-based/backups/bandwidth/security/...
We contacted Oracle/sgi/Hitachi/EMC(Isilon)/...



Our conclusions:

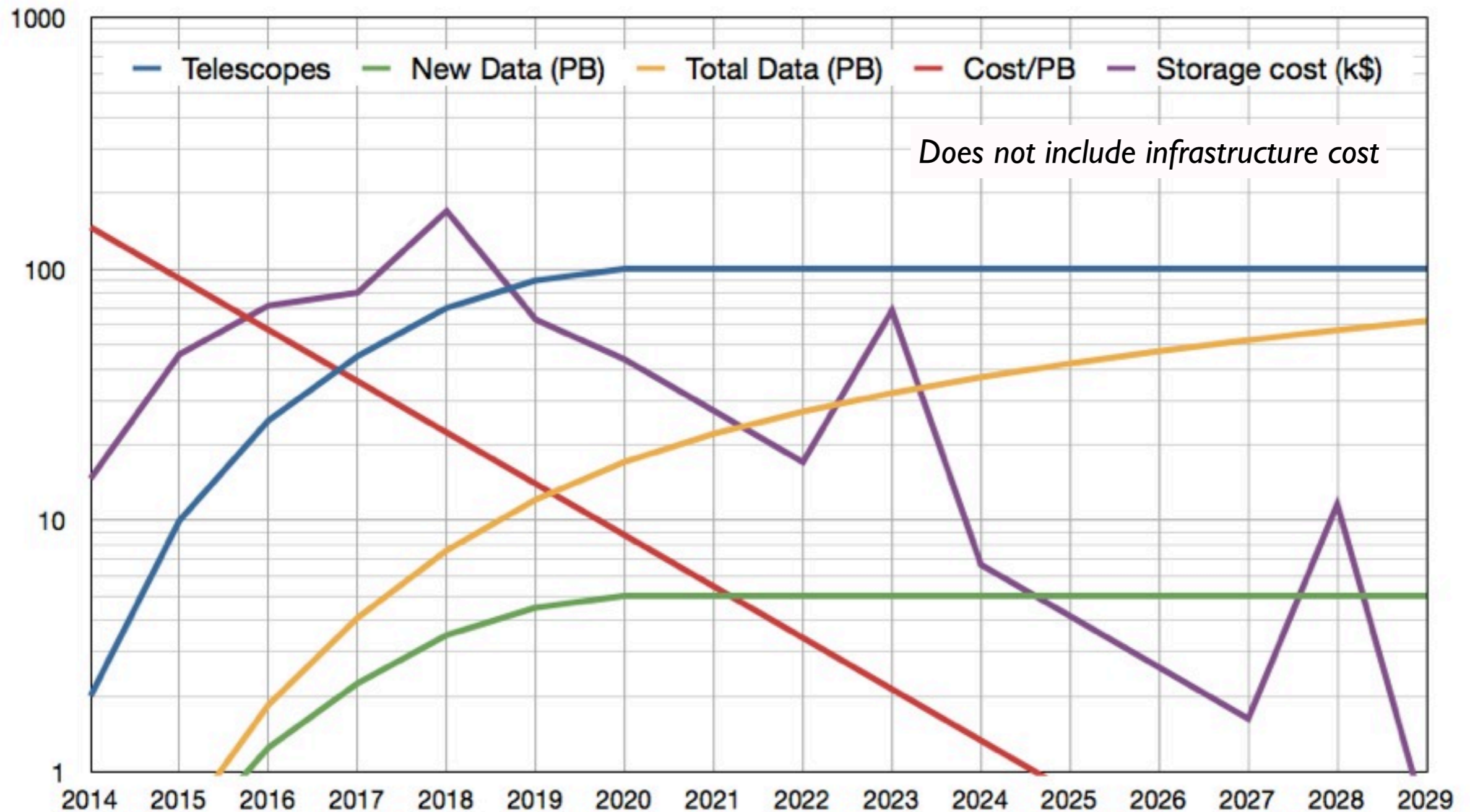
- 1) CTA should go for a disk (MAID) based archive
- 2) 2.5 PB (usable) per rack, now
- 3) Current hardware (50Gb/s) allows to read >0.5PB/day
- 4) Moore's law will continue for several years
- 5) 1PB (usable) in 2011: 400-600 k\$

EMC (Isilon) has a smart solution



Big science is not driving IT (if it ever did)

CTA storage needs



- the CTA archive will fit in one rack
- the storage hardware investment is of the order of 50-100k\$/year
- GEANT2 allows to transfer data between well connected centers

CTA analysis

One observation (100hrs):

- 1) 5TB raw data
- 2) 5TB simulations



for specialists

- 3) 150(GB) simulated showers par.

necessary for
specific analysis

- 4) 50 GB shower parameters
- 5) 0.5 GB event lists



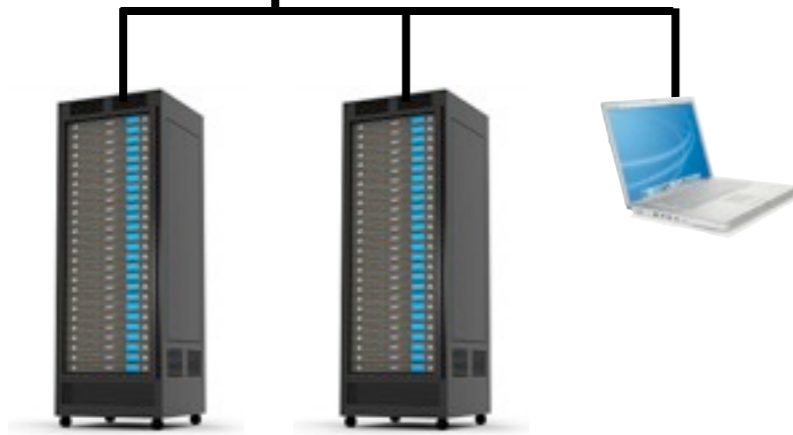
enough for
many observers

Observers will be able to get their full dataset.
Many will be happy with smaller datasets.

Possible CTA data analysis infrastructure



CTA operations, simulations, archives & replicates in several centers



Specialists (with duplicate of the data or access to operational machines)



Individual observers

Conclusions

- High-energy observatories are using a reasonably standardized computing model, defined by long term (>10 year) maintenance of data and analysis software and legacy
- The CTA archive will fit on one rack, the data (even raw) of one observation will fit in your laptop
- The CTA observatory can use a similar model than current high-energy observatories, despite the increase of data size

