# Data life-cycle in particle physics

Alessandro Razeto – LNGS

Laboratori Nazionali del Gran Sasso - 30/11/23

# Physics vs Metaphysics



Physics focus on measurable quantities.

Large part of metaphysics by the Scholasticists is now modeled by physics

# Particle physics



Cosmogony

- **The matter is the result of the interaction of elementary particles**
- **Discovering the laws modeling the behavior of elementary particles**
- **The applications are all around us**
  - Semiconductors
  - Medicine
  - Cosmology
- **The models are point-to-point**
  - Not necessarily providing the tools to describe complex bodies

# Theory vs reality



- **Theoretical physicists provide plenty of models**
  - Or models with plenty of degree of freedom
- **Experimentalists have to pin the reality to the right model**
  - Building experiments to probe the reality
- **Experiments are getting bigger and bigger**
  - More and more expensive and complicated

# Water - Earth – Fire – Air



- **Elementary particles interact with matter producing measurable quantities**
  - Electrical Charge
  - Change of Temperature
  - Light Emission
- **Sensors are used to detect these quantities → into an electrical signal**
  - Photo-multiplier tube for light
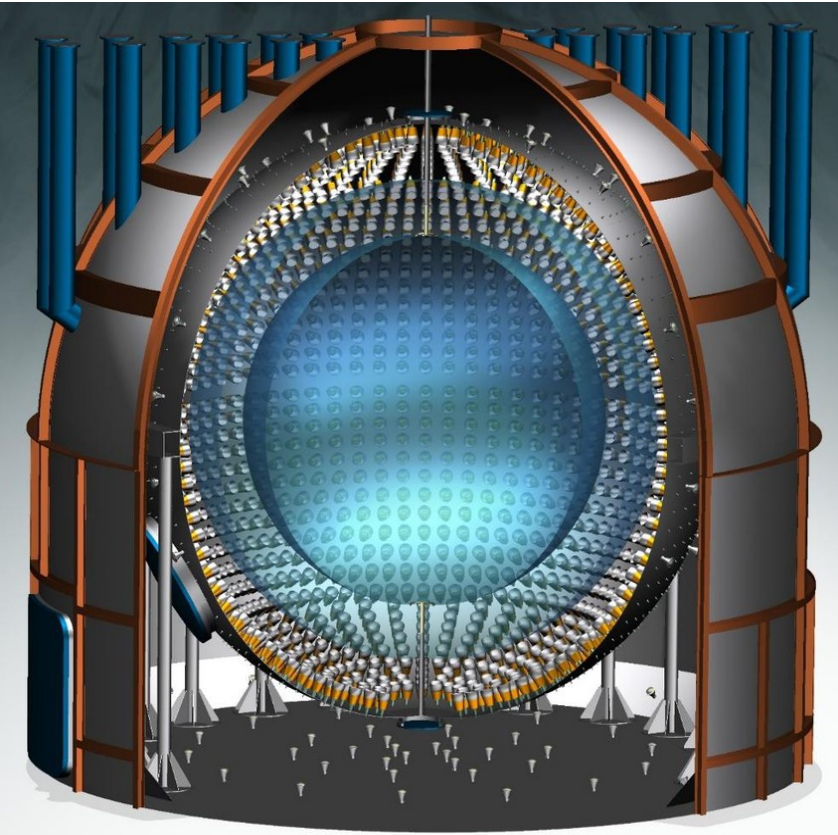  - Thermistors for temperature

# Detectors


Demiurge

- **The detectors are designed to maximize the signal**

- **Using elements having special properties to detect particles/radiation**

  – Scintillators are materials that emit faint light pulses during interaction with particles

- **Using active and passive shielding to reject unwanted events**
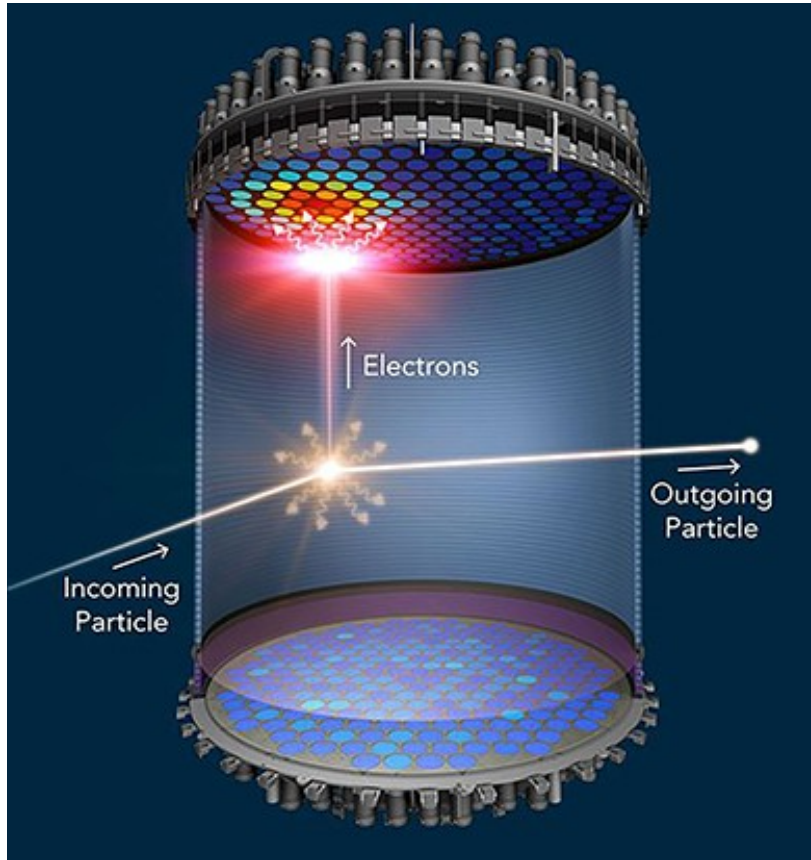
# LNGS Experiments

- **Detect neutrinos from the Sun**
- **Low background design**
  - Radio-pure materials
  - Active and passive shielding
- **Reached unprecedented contamination levels**
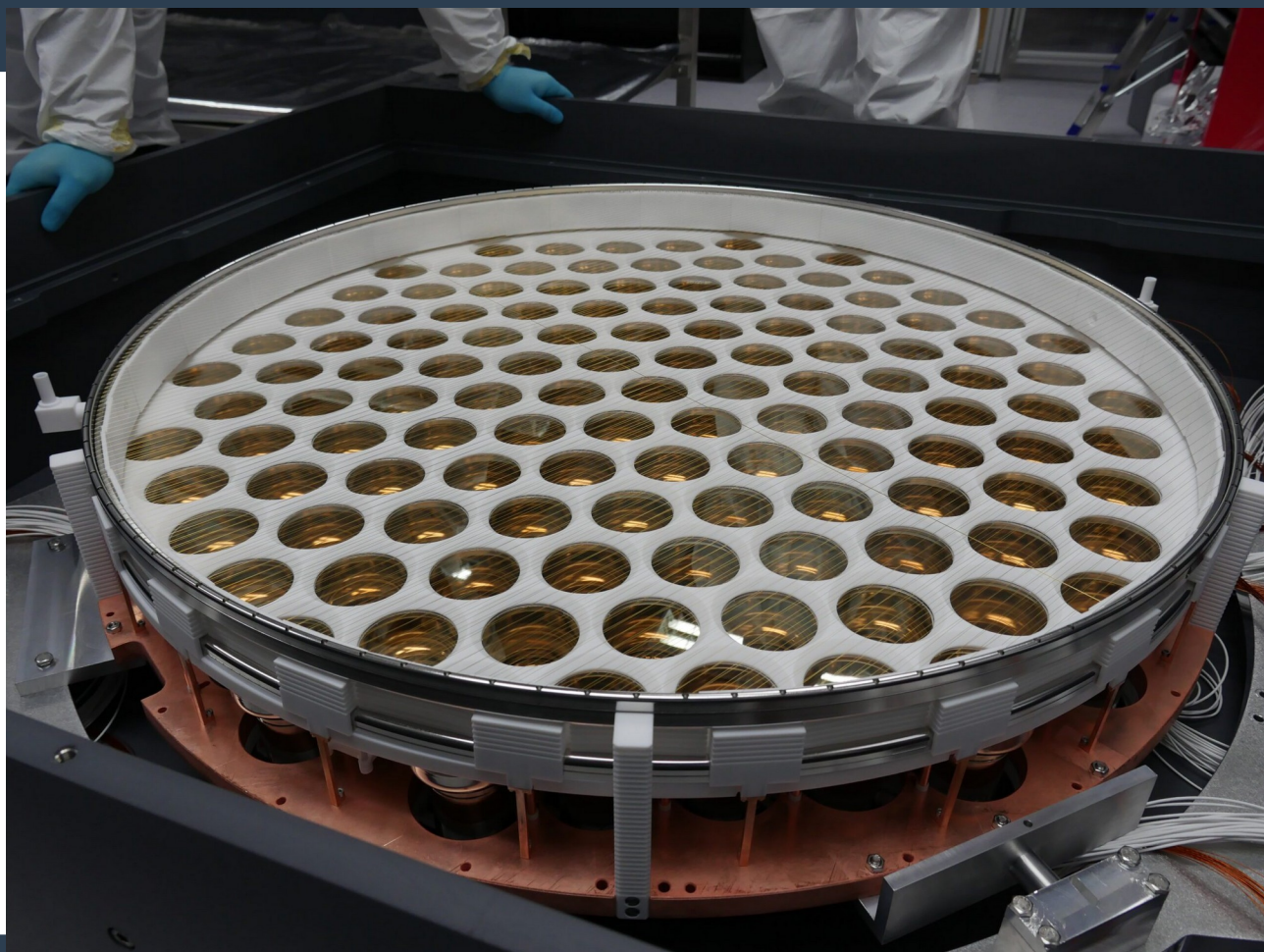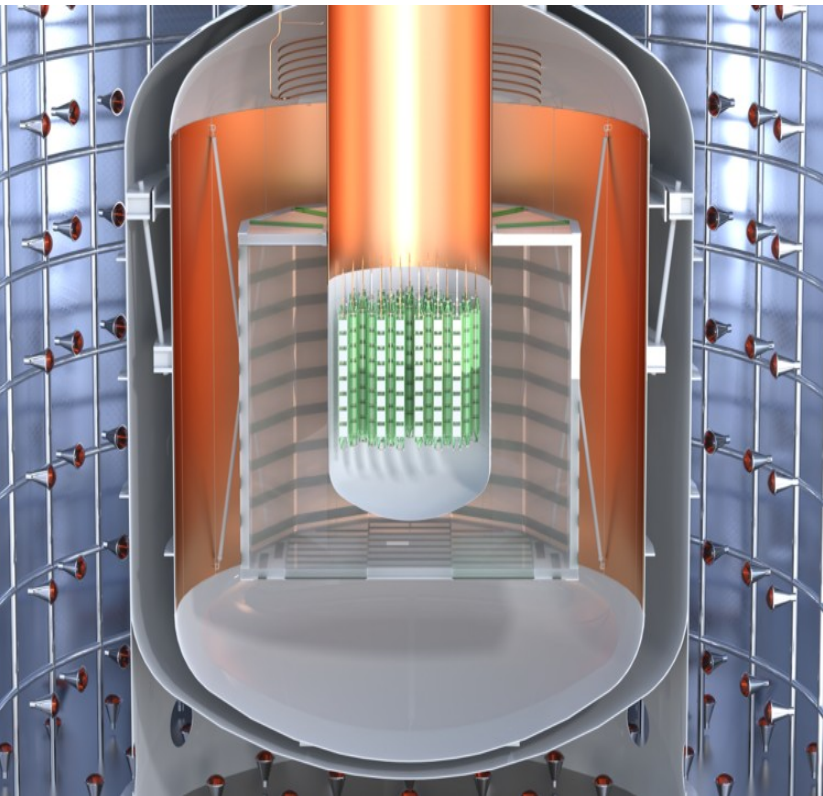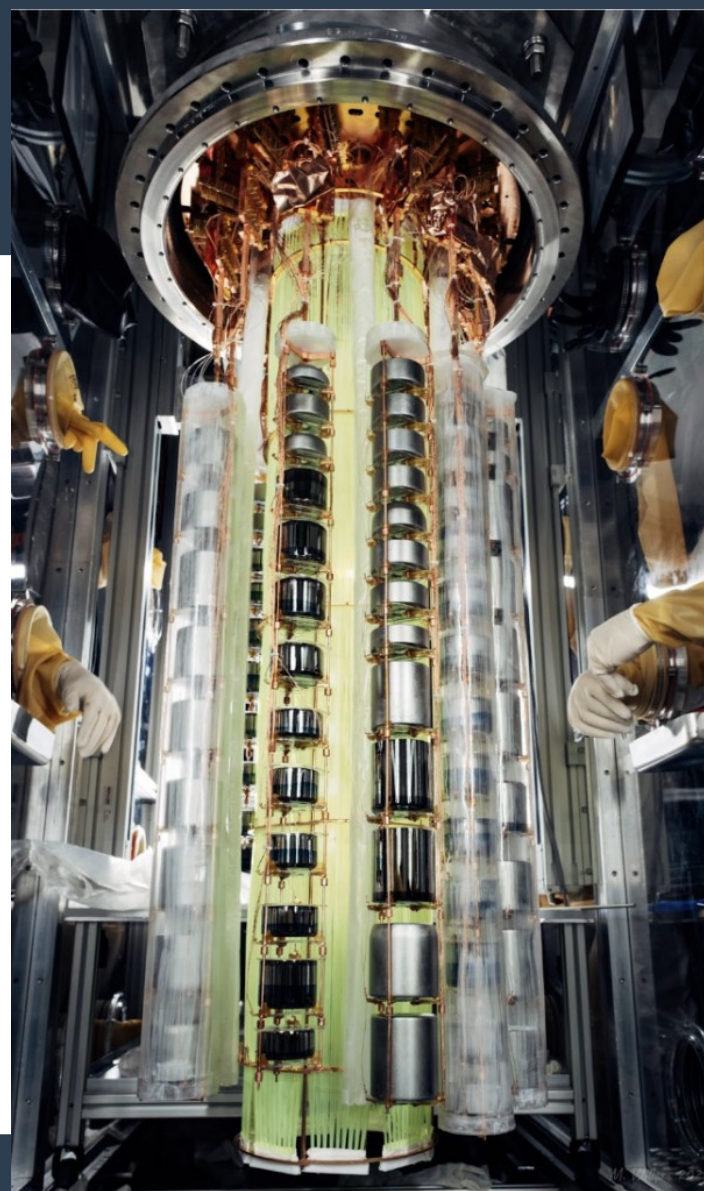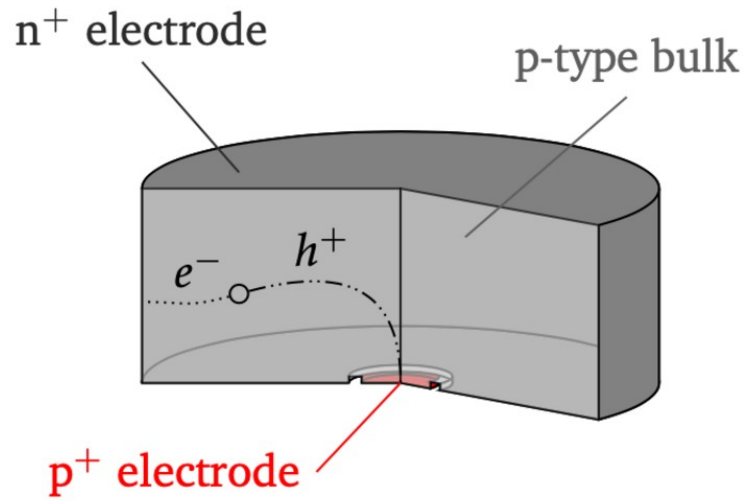- **300 tons of liquid scintillators + 2000 photomultiplier tubes**

- **Detect WIPMs as dark matter candidate**
- **XENON based experiments achieved the best sensitivity**
  - Negative result
- **Using the ultra-low background design (materials, vetoes)**
- **WIMPs interaction generate light and charge**
  - Light detected by ~500 PMTs
    - Between 3 to 100 photons
  - The charge is converted in a second light pulse

- **Searching for *forbidden* double beta decay in $^{76}$Ge**
  - Current limit is larger than $2 \cdot 10^{26}$ years
- **200 – 1000 kg of enriched Ge**
- **Using the ultra-low background design (materials, vetoes)**
- **Ge crystals are configured as large fully-depleted diodes**
  - Detect the ionized charge induced by the decay

12

n+ electrode

p-type bulk
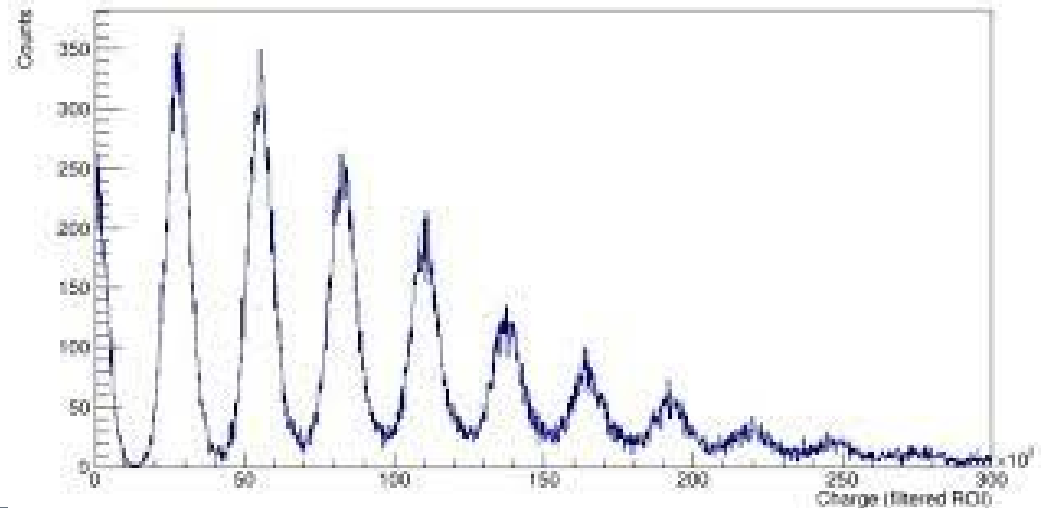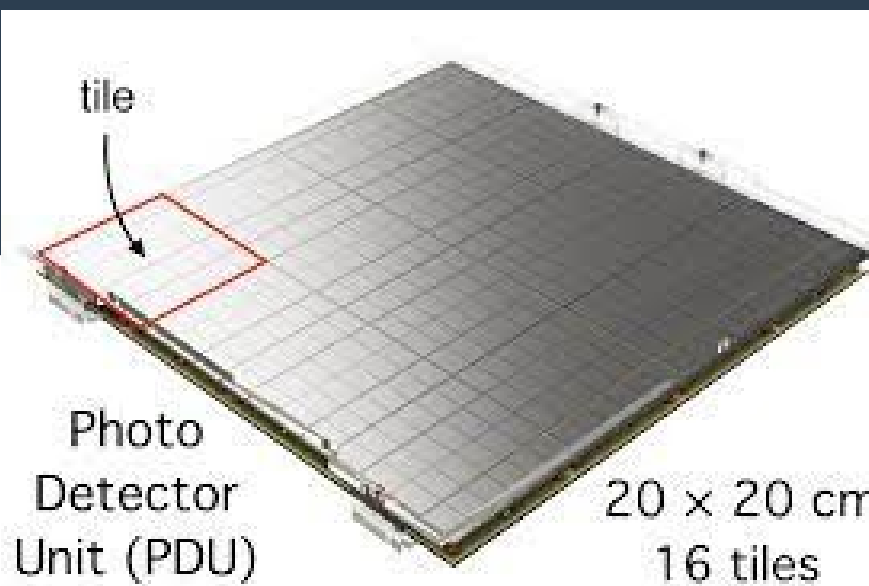
$e^-$ $h^+$

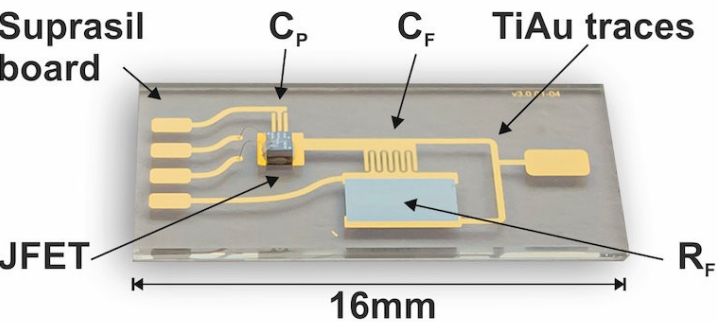p+ electrode
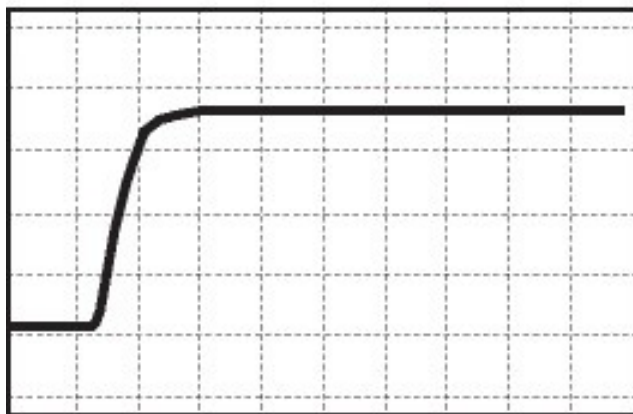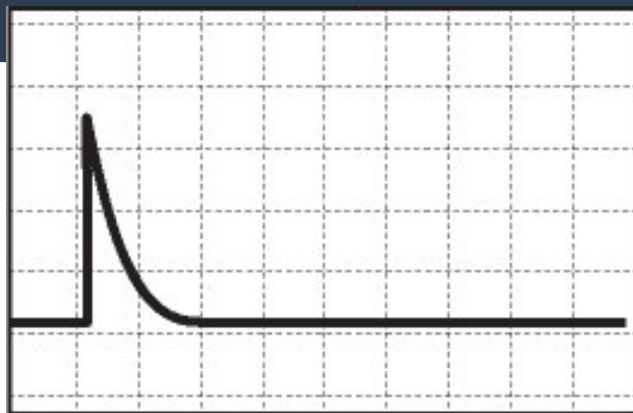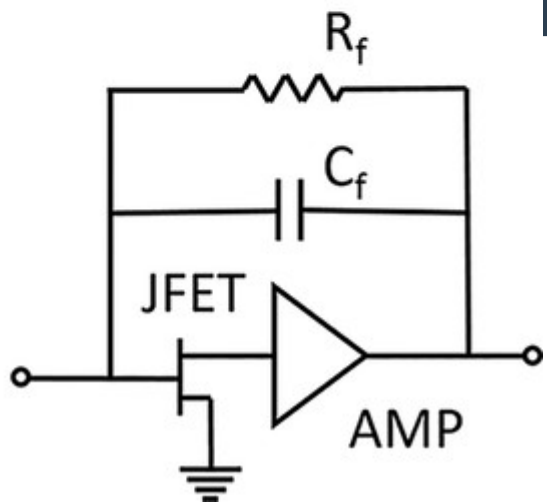


13

# Signal extraction

- **Typically the output of the sensors is very small**
  - For the signal acquisition
  - For the transmission on cables
- **Very low noise amplifiers are developed**
  - To optimize the signal integrity
- **The front-end can be at room temperature or in cryogenic environment**
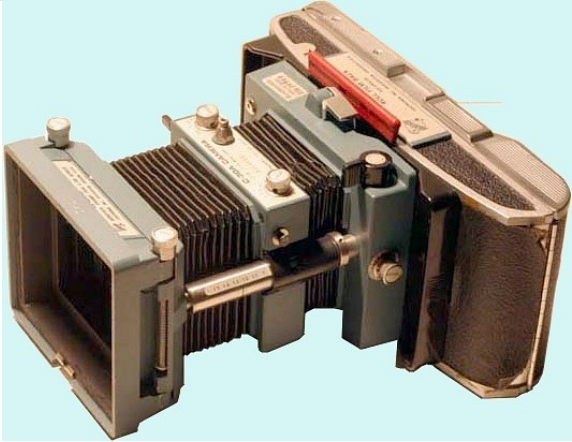
# LNGS Examples



tile

Photo Detector Unit (PDU)
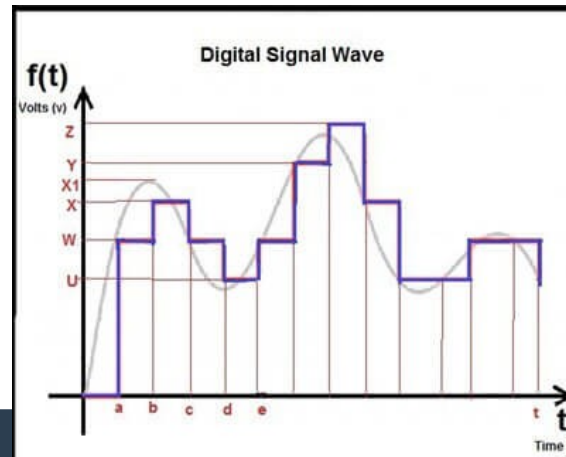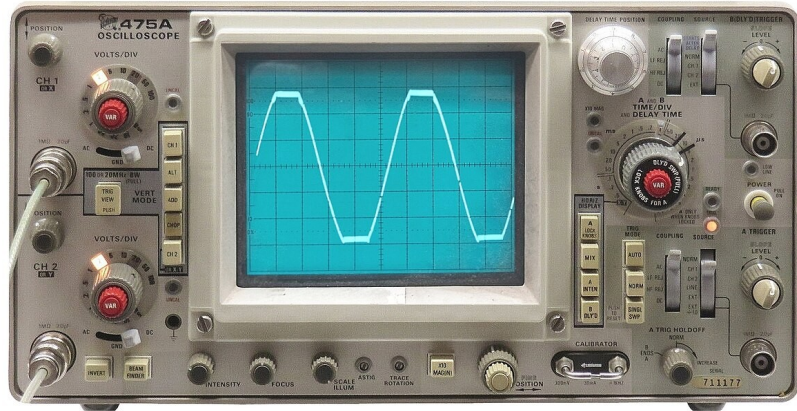
20 × 20 cm, 16 tiles

# Charge sensitive amplifiers



$R_f$

$C_f$

JFET

AMP

Suprasil board    $C_P$    $C_F$    TiAu traces

JFET    $R_F$

16mm

$i_d(t)$: detector current pulse: 10ns/div

$V_{out}(t)$: CSP output pulse: 10ns/div

# Waveform Digitization

- **Before 1971 photos on oscilloscope waveforms**
- **Modern experiments are based of fast digitizers (>100 MHz) with 12-16 bits**
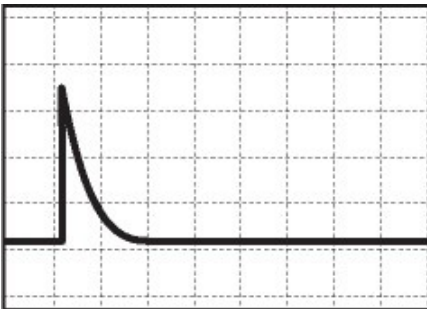


100 GS/s
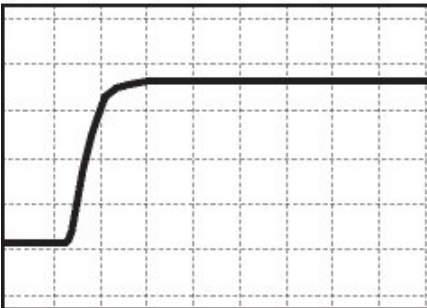
18

# Digital Signal Processing



- **Digitized waveform can be further processed to maximize the signal to noise ratio**
  - Typically filtering is used to abate noise
    - At the expense of losing part of the signal
- **DSP use is ubiquitous in modern devices**
  - Multi-rate filtering ← iPOD (mp3)
  - RADAR

$i_d(t)$: detector current pulse: 10ns/div

$V_{out}(t)$: CSP output pulse: 10ns/div

$\tau = R_i C_i$

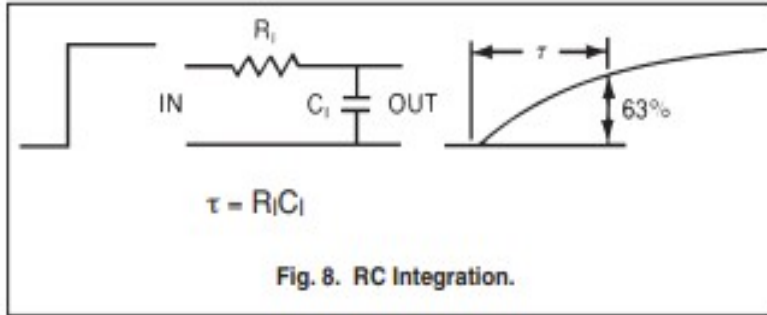Fig. 8. RC Integration.

$\tau = R_D C_D = R_i C_i$

Fig. 9. CR-RC Pulse Shaping.

Si(Li) X-RAY DETECTOR (4 mm DIA × 3 mm) and 117B PREAMP

$\tau_i = \tau_0 = \tau$

TOTAL NOISE

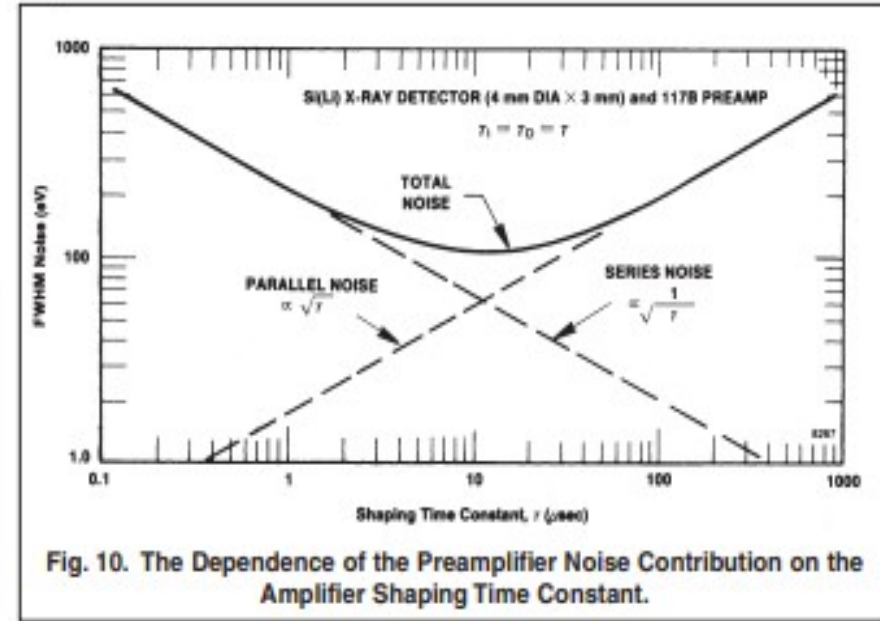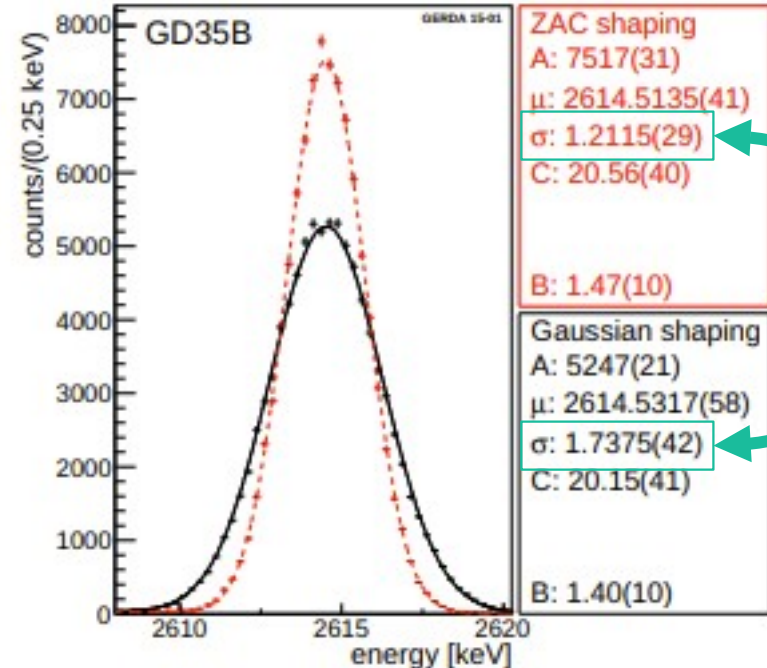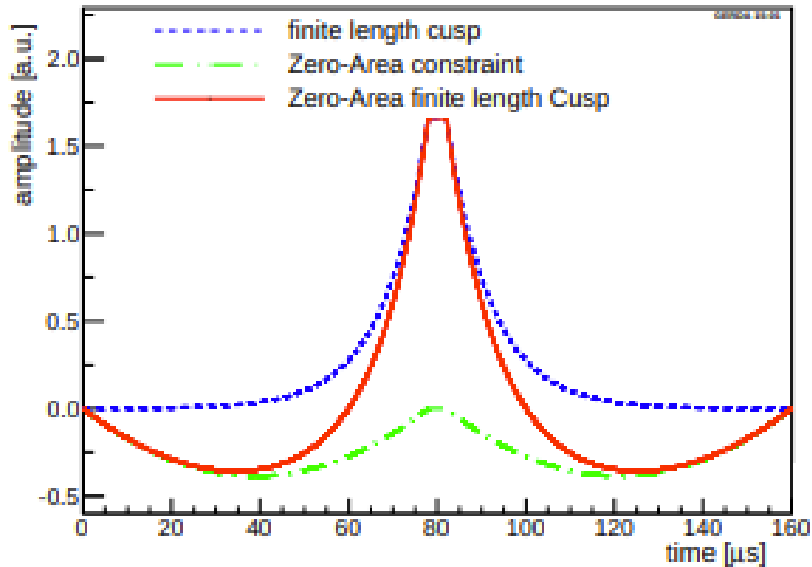PARALLEL NOISE $\propto \sqrt{\tau}$

SERIES NOISE $\propto \sqrt{\frac{1}{\tau}}$

Fig. 10. The Dependence of the Preamplifier Noise Contribution on the Amplifier Shaping Time Constant.

20

In the digital domain the filtering can be much more advanced
CRC filters can be replaced by gaussian or other symmetric shapes
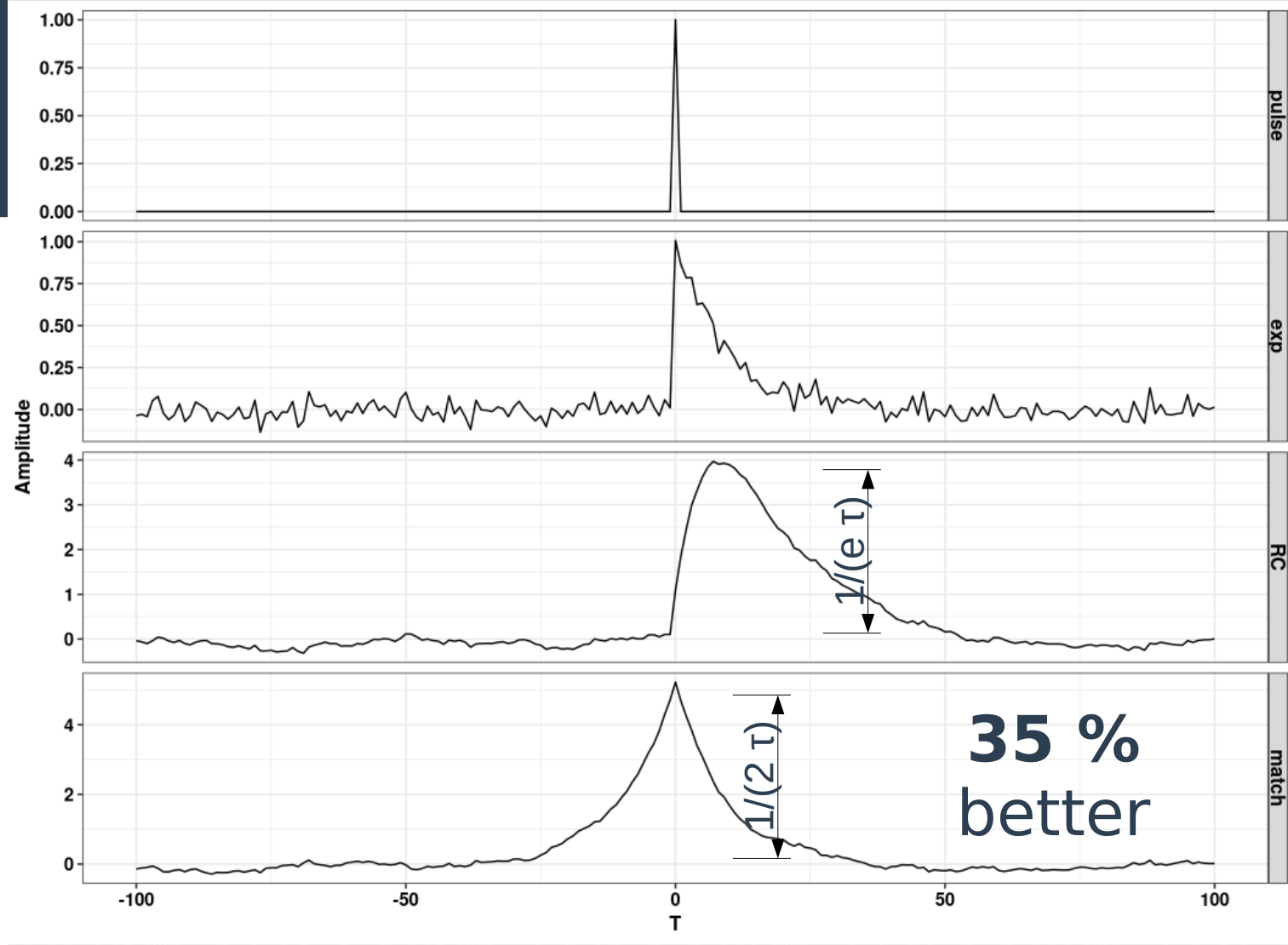
# Matched filter

Turin 1960
Maximize SNR
(amplitude)

Introduced for RADAR

Correlate the signal
with the known ref

Minimize the phase
Dispersion

Can be implemented
only in digital
(anti-causal)



$\frac{1}{(e\,\tau)}$

$\frac{1}{(2\,\tau)}$

**35 %**
better

# Low Level Analysis

- **Reconstruct the waveforms in high level data**

- **Include calibrations**

- **Position reconstruction**

- **In general the most complicated part of the analysis chain**

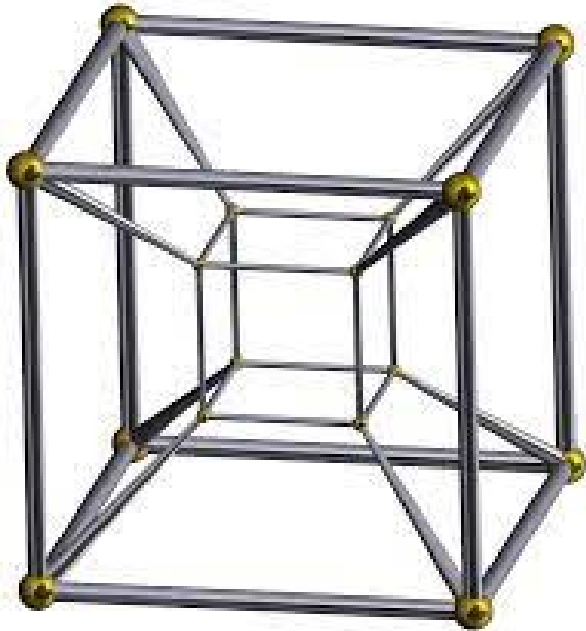  – Knowledge of the detector details

# Borexino



- **Interaction of solar neutrino generates about 500 photons**
  - ~ hundreds per day over a background of 1 million
- **Detected by the PMTs**

- **We record the time, amplitude and position of each signal**
- **We reconstruct**
  - Interaction energy
  - Pulse shape
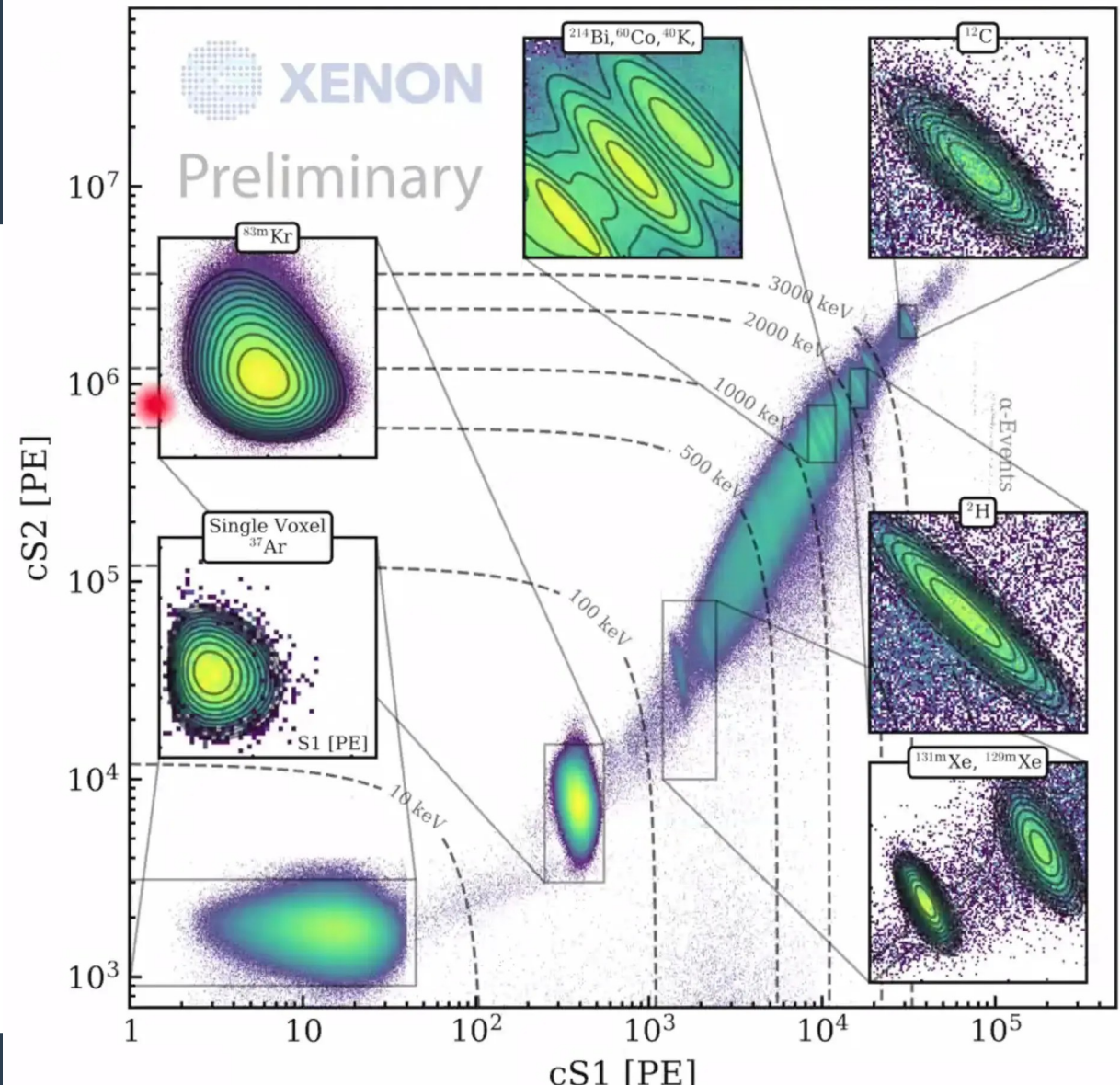  - Position in the detector
  - Status of the vetoes

# High Level Analysis

# Multi-Dimensional Data



- **High level data are nativelly highly-dimensional**
- **An *event* correspond to an interaction in the detector**
  - Several information are available
    - Primary data (energy, position, topology)
    - Nuisance data (veto, metadata)

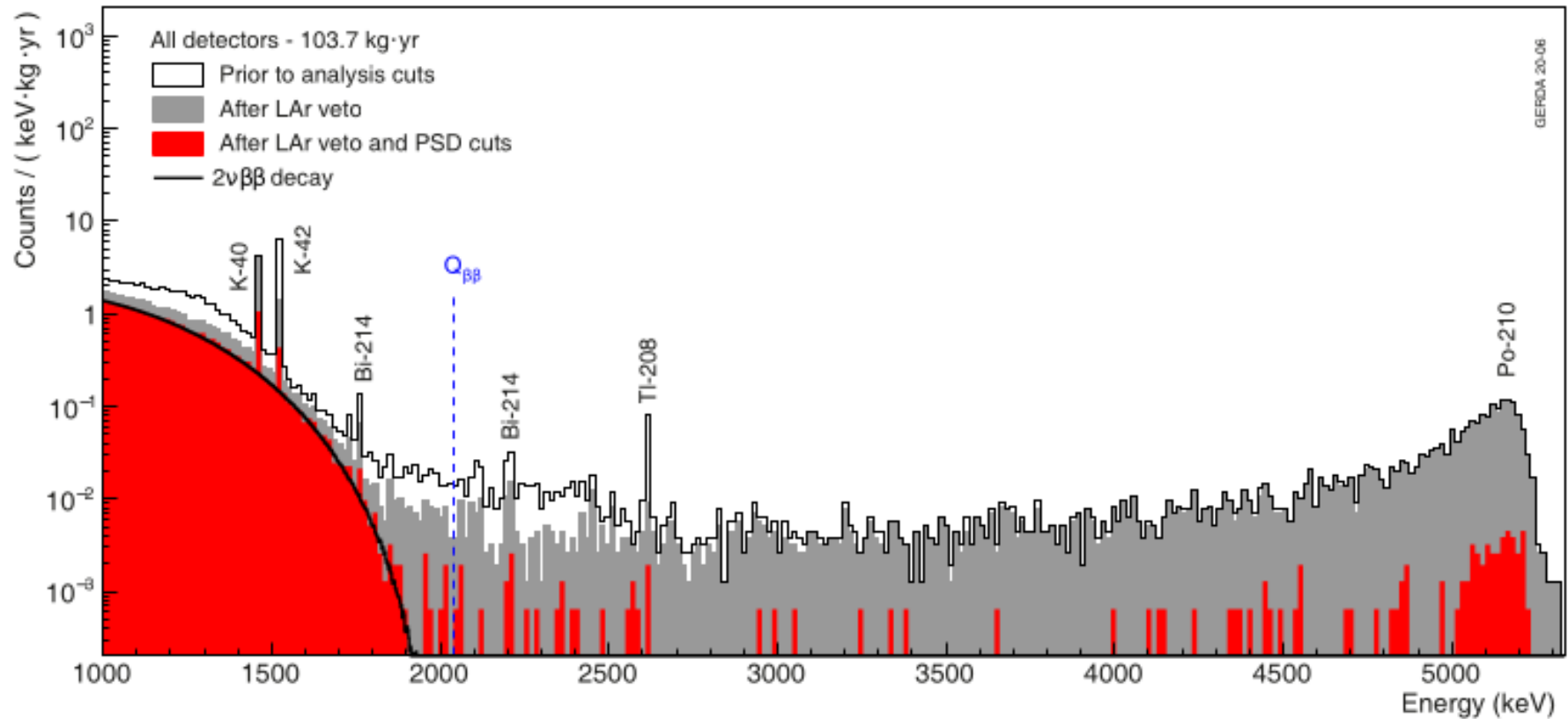- **However we typically do not employ multi-dimensional regression**

- **Typically we populate histograms with our primary data**
- **Search peaks (or other features) in these histograms**

- **Unbinned likelihood are not common**
- **Histogramming is often abused**
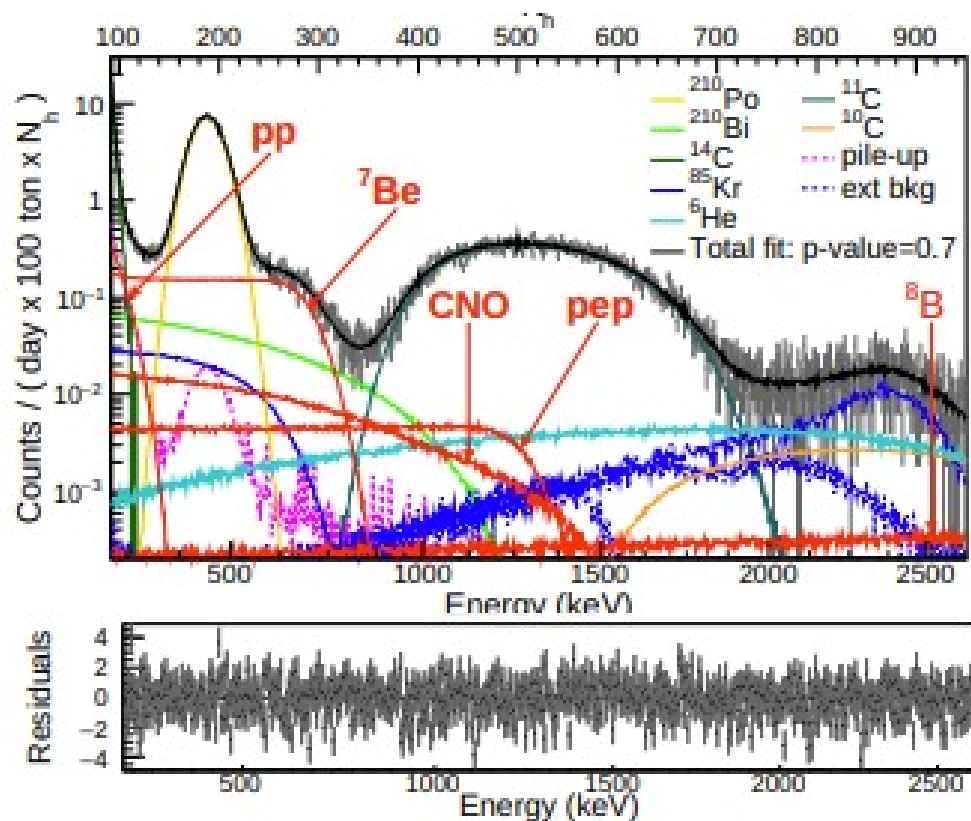  - Even to calculate the with of a peak we typically use histograms

# χ² test
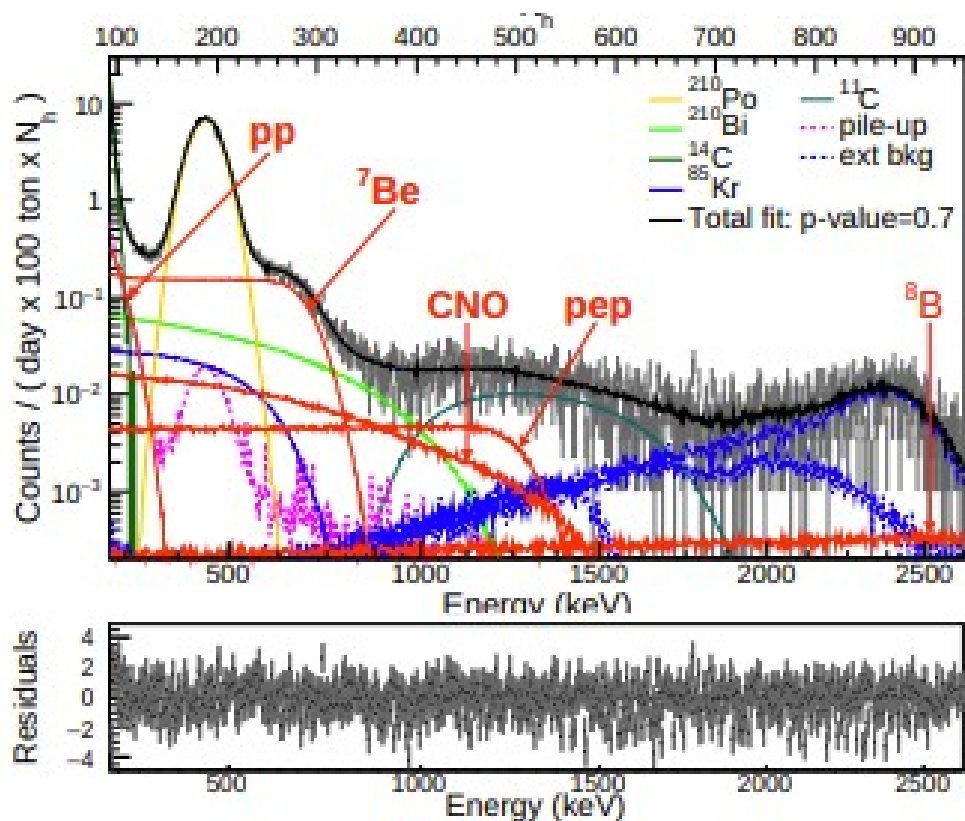

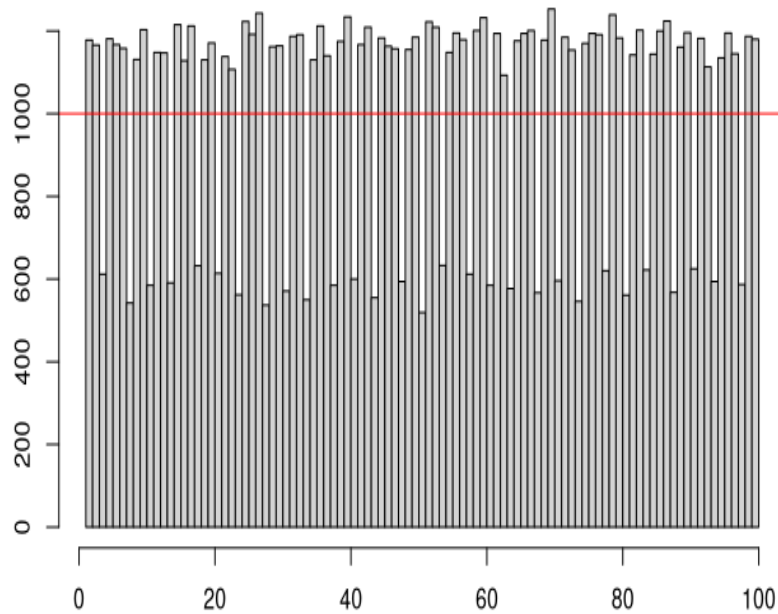
- **The research of features on histograms is based on the χ²**
  - We assume each bin is an independent random variable
  - We use the Neyman χ²

- **The minimization of the χ² provides the model parameters and the goodness of fit**
- **It is a very handy process**
  - Not immune to several quirks

# Multivariate fit of Borexino data
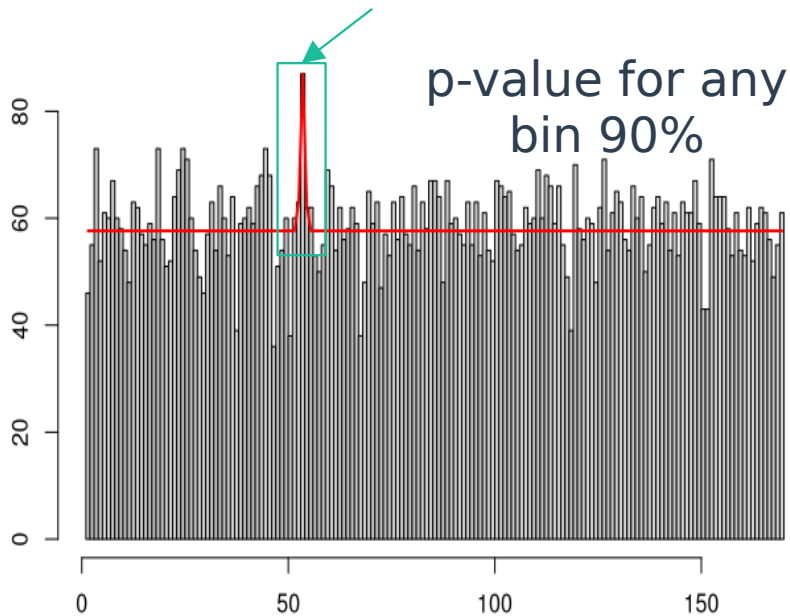
# χ² not always is benevolent

Rounding problems

- **Not always the bin contents are independent**
  - Or the binning is wrong
  - # of degree of freedom is overestimated
    - P-value wrong
- **Look elsewhere effect**
  - Random peaks can appear
- **...**
- **CERN requires 5 σ evidence**

# χ² not always is benevolent

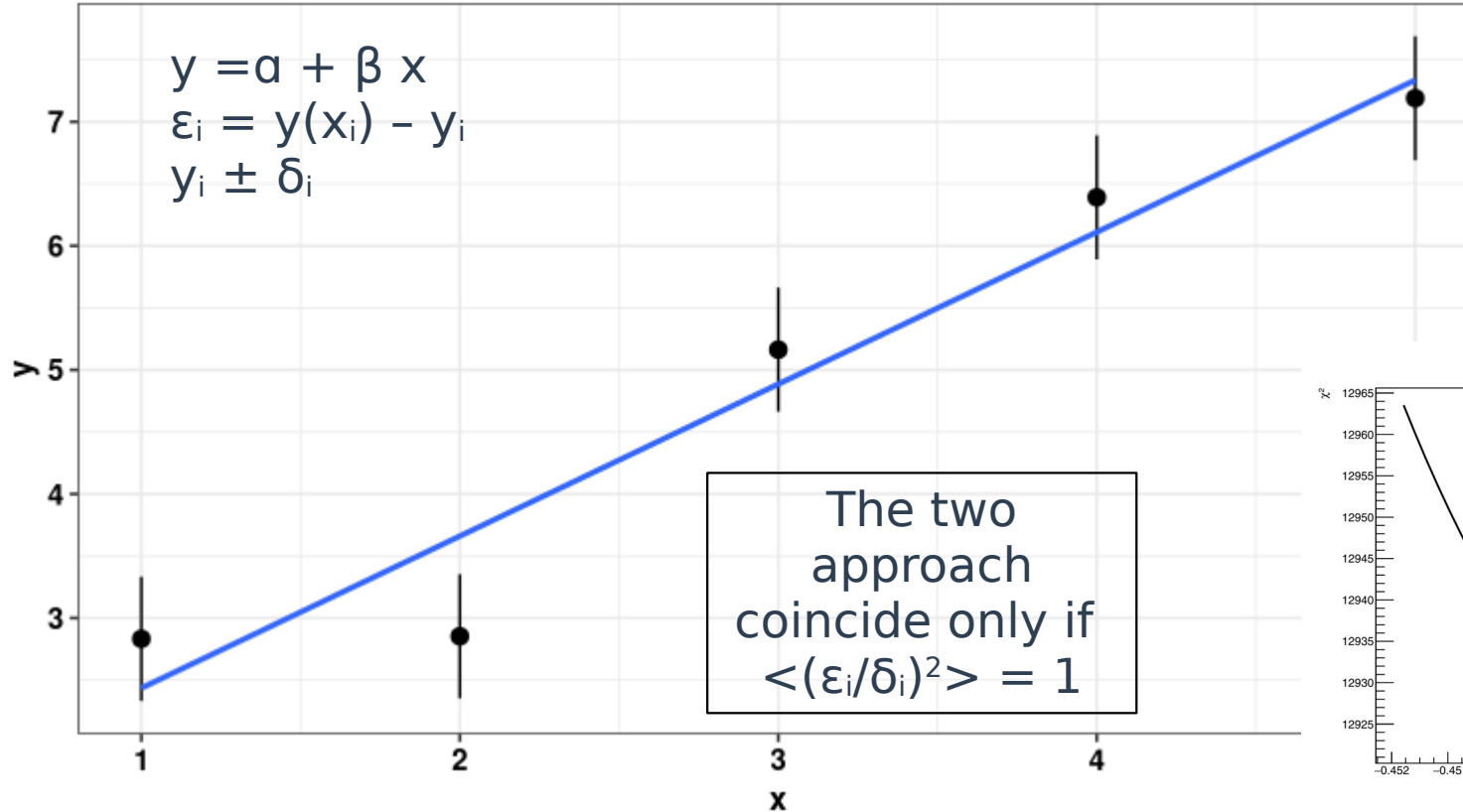p-value of fluctuation .5%

p-value for any bin 90%



Look elsewhere

- **Not always the bin contents are independent**
  - Or the binning is wrong
  - # of degree of freedom is overestimated
    - P-value wrong
- **Look elsewhere effect**
  - Random peaks can appear
- **...**
- **CERN requires 5 σ evidence**

# Model Driven vs Data Driven

$$y = \alpha + \beta x$$
$$\varepsilon_i = y(x_i) - y_i$$
$$y_i \pm \delta_i$$

$$s_{\widehat{\beta}} = \sqrt{\frac{\frac{1}{n-2}\sum_{i=1}^{n}\widehat{\varepsilon}_i^2}{\sum_{i=1}^{n}(x_i - \bar{x})^2}}$$

$$s_{\widehat{\alpha}} = s_{\widehat{\beta}}\sqrt{\frac{1}{n}\sum_{i=1}^{n}x_i^2} =$$

The two approach coincide only if $\langle(\varepsilon_i/\delta_i)^2\rangle = 1$
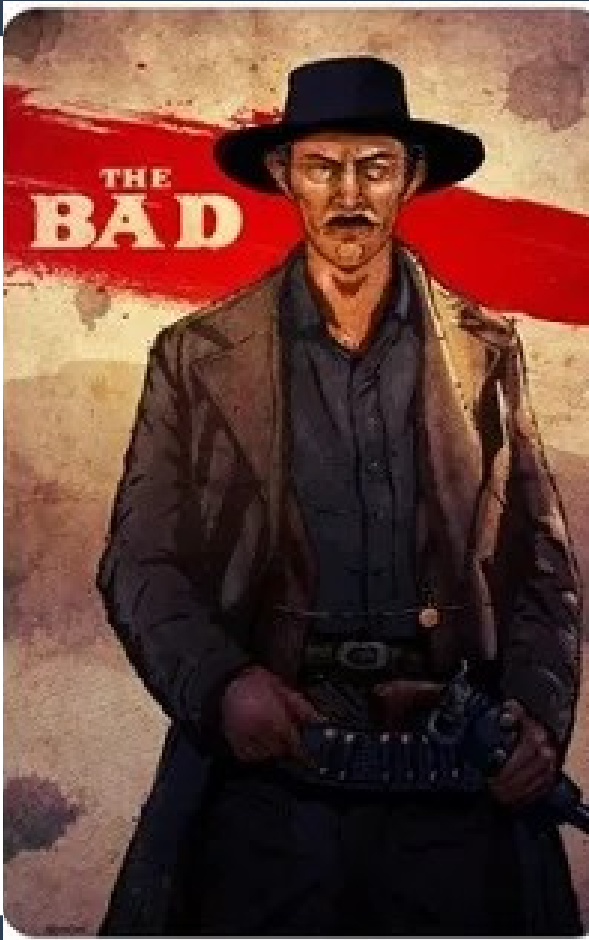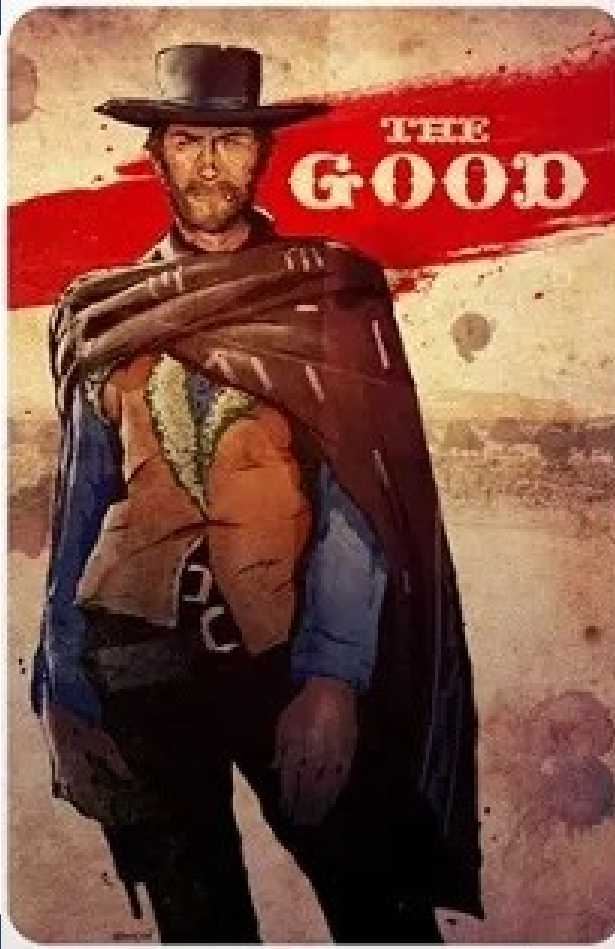
Chi-squared scan

# Software Tools

- **ROOT is a software package developed at CERN**
  - With the main focus on particle physics
- **It has a nice data storage (columnar db-file)**
  - Allows to analyze large data set
    - Larger than memory size
- **It was written in C++**
  - Well before C++11
  - It uses a C++ interpreter (clang)
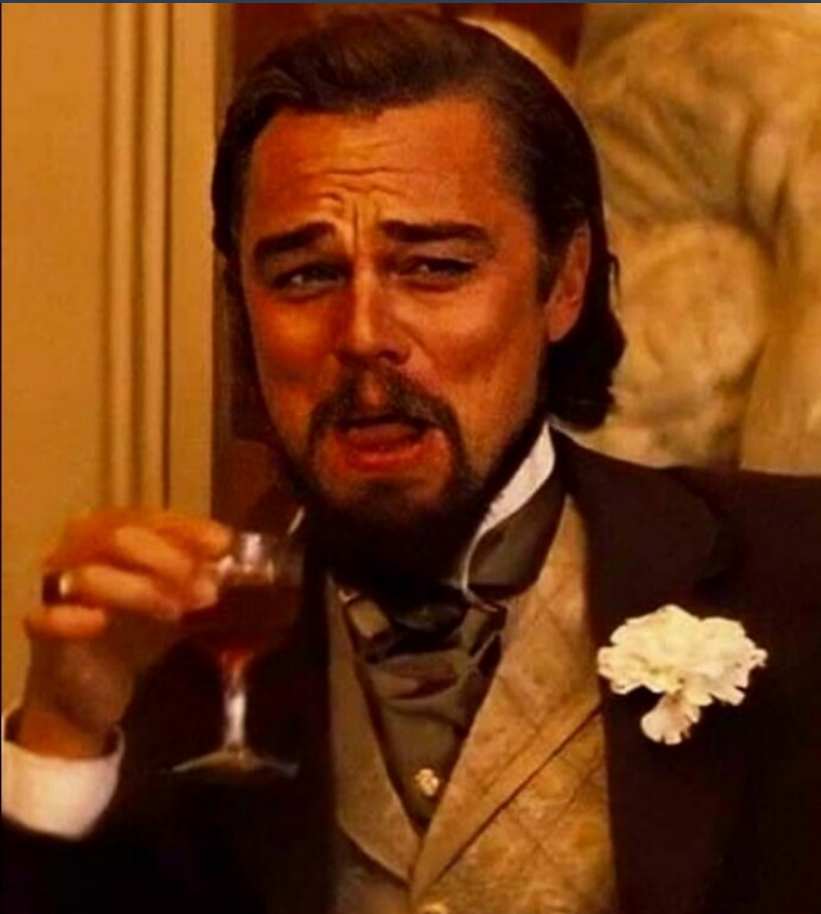  
  **Its syntax is ugly**

# C++

- **After Fortran77 CERN pushed the scientific community for the adoption of C++**
- **C++ is fine for the DAQ**
  - And for the intensive low-level analysis
- **C++ is not a scripting language**
  - Complex (bloated) syntax
  - Not vectorized + external iterator
  - Handling string is painful
  - Missing syntactic sugar
  - **Prone to memory leaks**

"There are only two kinds of languages: the ones people complain about and the ones nobody uses"
Bjarne Stroustrup

- **Python is acquiring popularity in particle physics**
  - It is easy to learn
  - Has advanced scientific libraries
    - ROOT library as well

- **It is much better than C++ but**
  - Missing large data-set operation
  - Many dislike the syntax
  - It is much slower than C++
    - Some compilation/optimization options available

- **Personally I use R**
- **It is not tailored for particle physics**
  - Histogram fitting is somehow missing
- **It forces me to study and to think differently**
- **No built-in support for large data set**
  - I use SPARK clusters with parallel map/reduce

# Conclusion

- **A rich environment that includes**
  - Theory
  - Experiments
  - Models
  - Data
  - Plus our idiomatic use of statistics
- **Large experiments are fueled by large international collaboration of physicists**

# Thank you