# A **GPU-based** framework for multi-variate analysis in particle physics

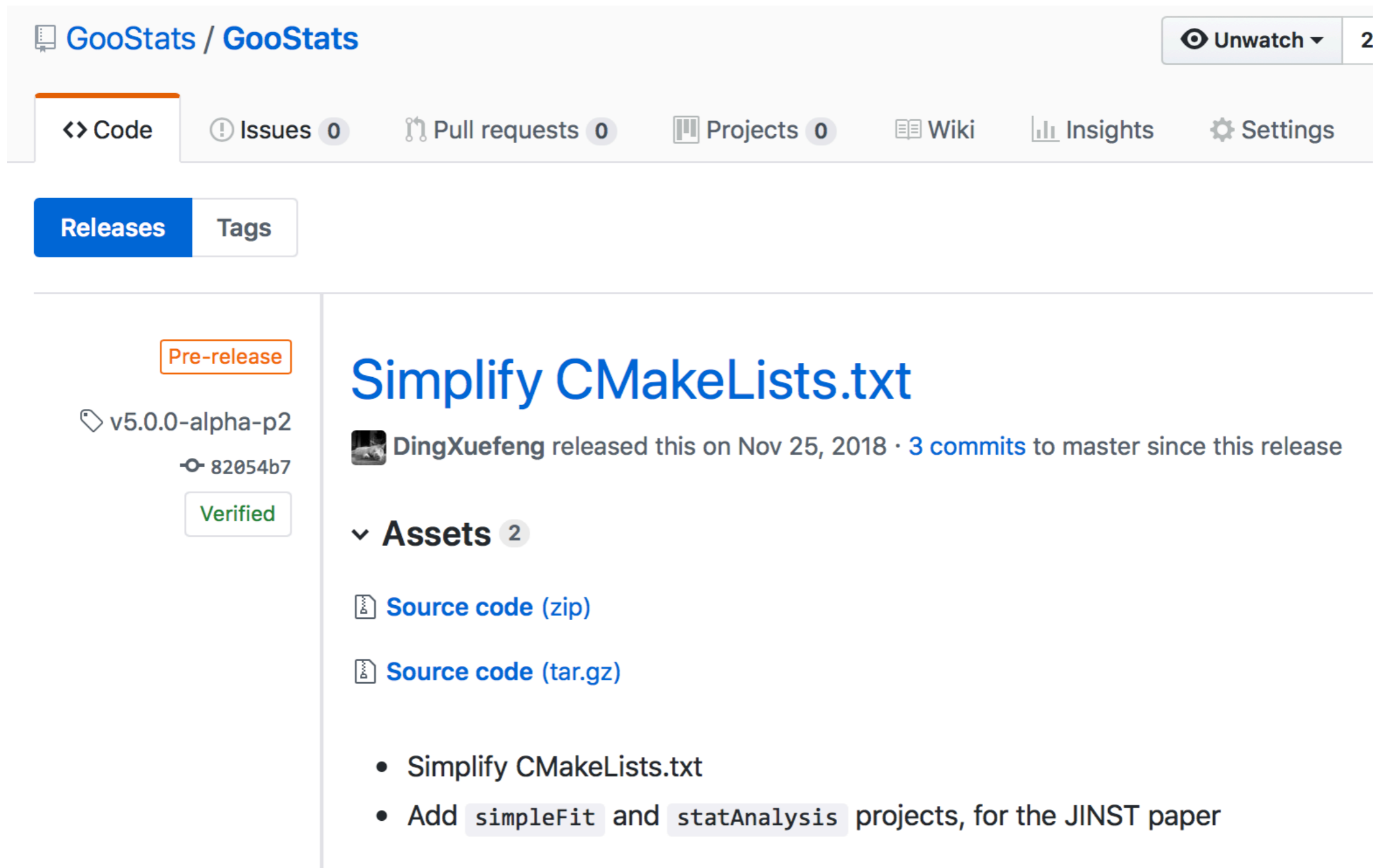25 Jan 2019, 12:10 (CET) @503-1-001, 12'+3'

**Xuefeng Ding[1,2]**

1. Gran Sasso Science Institute, L'Aquila, Italy
2. INFN Sezione di Milano, Milan, Italy

PHYSTAT-nu 2019
@ CERN, Geneva, Switzerland 25 January 2019

# Outline

- The GooStats framework

- Motivation: Borexino spectrum fit

- Package design

- Systematic uncertainty estimation

- Conclusion

# GooStats hosted on GitHub

## https://github.com/GooStats/GooStats.git

- A convenient GPU multivariate analysis framework
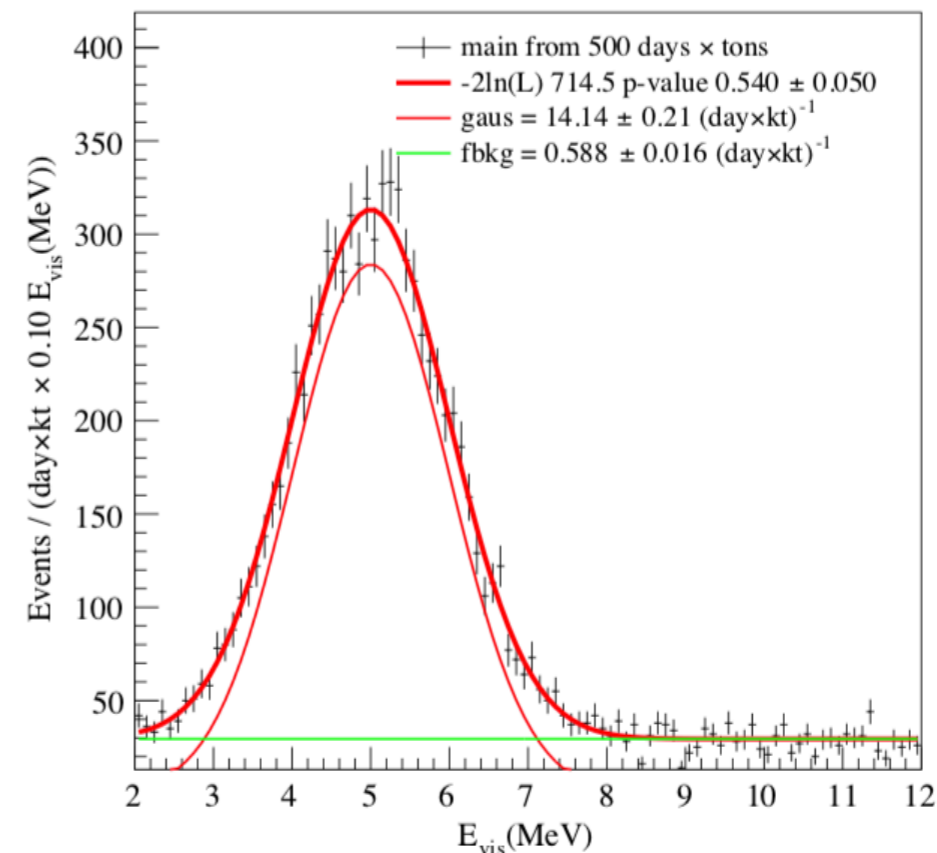
- Easy spectrum fitting



**Figure 7.** Left: screen shot of fit result summary. Right: produced figure in the file format of pdf.

- Speed up: around 200 times



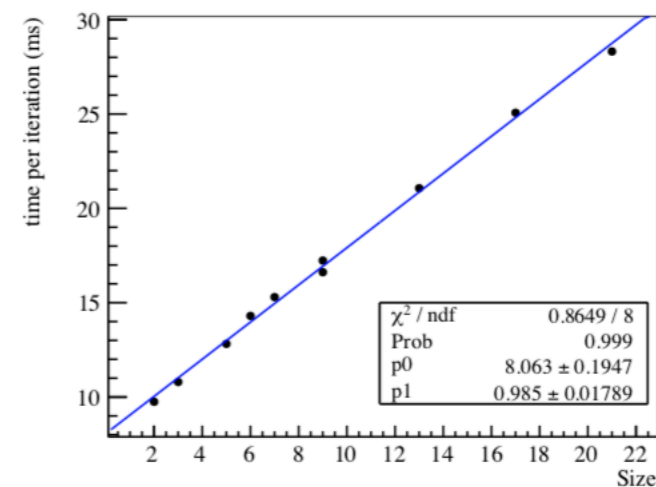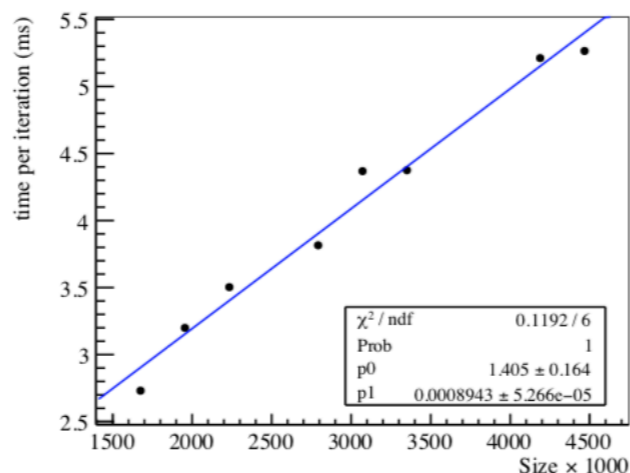**Table 1.** Comparison of fitting time between GooStats and original software used by the Borexino collaboration. $T_{tot}$: total execution time. $N$: the number of iterations taken to converge in MINUIT. $T_{it}$: average execution time per iteration. Speed up: $T_{it}(CPU)/T_{it}(GPU)$.

| | CPU | | | GPU | | | |
|---|---|---|---|---|---|---|---|
| Type | AMD Opteron(TM) Processor 6238 | | | nVidia Tesla K20m | | | |
| Size | $T_{tot}$ (s) | $N$ | $T_{it}$ (ms) | $T_{tot}$ (s) | $N$ | $T_{it}$ (ms) | speed up |
| 400 | 27.6 | 1128 | 24.4 | 0.181 | 1346 | 0.135 | 181 |
| 350 | 29.4 | 1331 | 22.1 | 0.156 | 1294 | 0.121 | 183 |
| 300 | 22.5 | 1239 | 18.2 | 0.243 | 1995 | 0.122 | 149 |

- Motivation: spectral-fitting tool for Borexino analysis

# Analytical Response Function

$$M : f(E) \mapsto g(\text{charge}) = \int_0^{E_{\text{end}}} \mathrm{d}E \cdot f(E) \cdot \text{RPF}\left[\text{charge}; \mu(E), \text{var}(\mu)\right]$$

- Analytical shape of spectrum of mono-energetic events

  - **Momentum based approximation**

  - Match the average ( energy scale + non-linearity model )

  - Match the variance ( energy resolution model )

  - … (—> simplified)

  - More: "Mask", "pile-up" etc…

- We can simplify because

  - Borexino response is simple: small FV in center, low energies => no irregular tail

  - We are not sensitive.. => small systematics

  - **Fit full MC to get the bias introduced in simplification**

$$\mathcal{L}_{\mathrm{MV}}\left(\vec{\theta}\right) = \mathcal{L}_{\mathrm{TFC-sub}}\left(\vec{\theta}\right) \cdot \mathcal{L}_{\mathrm{TFC-tagged}}\left(\vec{\theta}\right) \cdot \mathcal{L}_{\mathrm{RD}}\left(\vec{\theta}\right) \cdot \mathcal{L}_{\mathrm{PS}}\left(\vec{\theta}\right)$$

- Scaling factor introduced to remove bias.

- **Middle Layer** between **GooFit** and **User module**

# Parameter Synchronization

- Synchronize fit parameters in joint analysis

- Supported by tree type internal data structure

# TTree Output

- All fit parameters saved in TTree output automatically

- More quantities can be added as lambda functions in the output builder

```
root [1] .ls
TFile**          test_tree.root
 TFile*          test_tree.root
  KEY: TTree     fit_result;1     Fit result of GooStats
root [2] fit_result->Show(0)
======> EVENT:0
 default.NReactor = 3.02987
 default.NReactor_err = 0.0159146
 default.Nbkg     = 1.02405
 default.Nbkg_err = 0.0162459
 default.U235     = 0.5
 default.U235_err = 0
 default.U238     = 0.2
 default.U238_err = 0
 default.Pu239    = 0.1
 default.Pu239_err = 0
 default.U241     = 0.2
 default.U241_err = 0
 default.LY       = 1300
 default.LY_err   = 0
 default.qc1      = 2.78788
 default.qc1_err  = 0
 default.qc2      = -0.528003
 default.qc2_err  = 0
 default.v1       = 0.3
 default.v1_err   = 0
 default.vT       = 5
 default.vT_err   = 0
 default.Reactor_dEvis = 1078.28
 default.Reactor_dEvis_err = 13.2291
 chi2             = 390.448
 NDF              = 397
 likelihood       = 1883.96
```
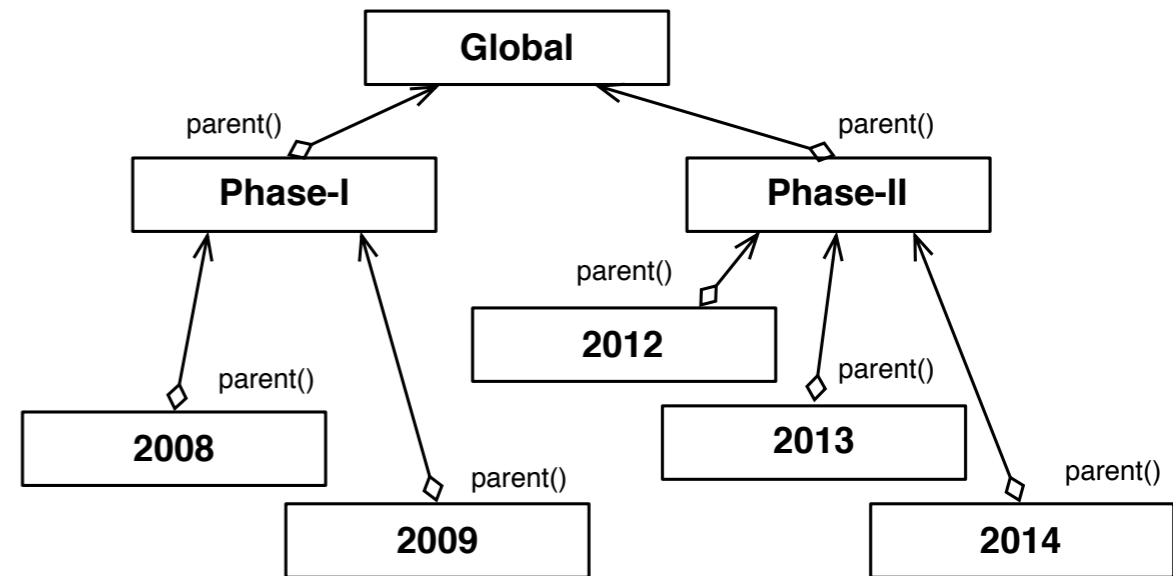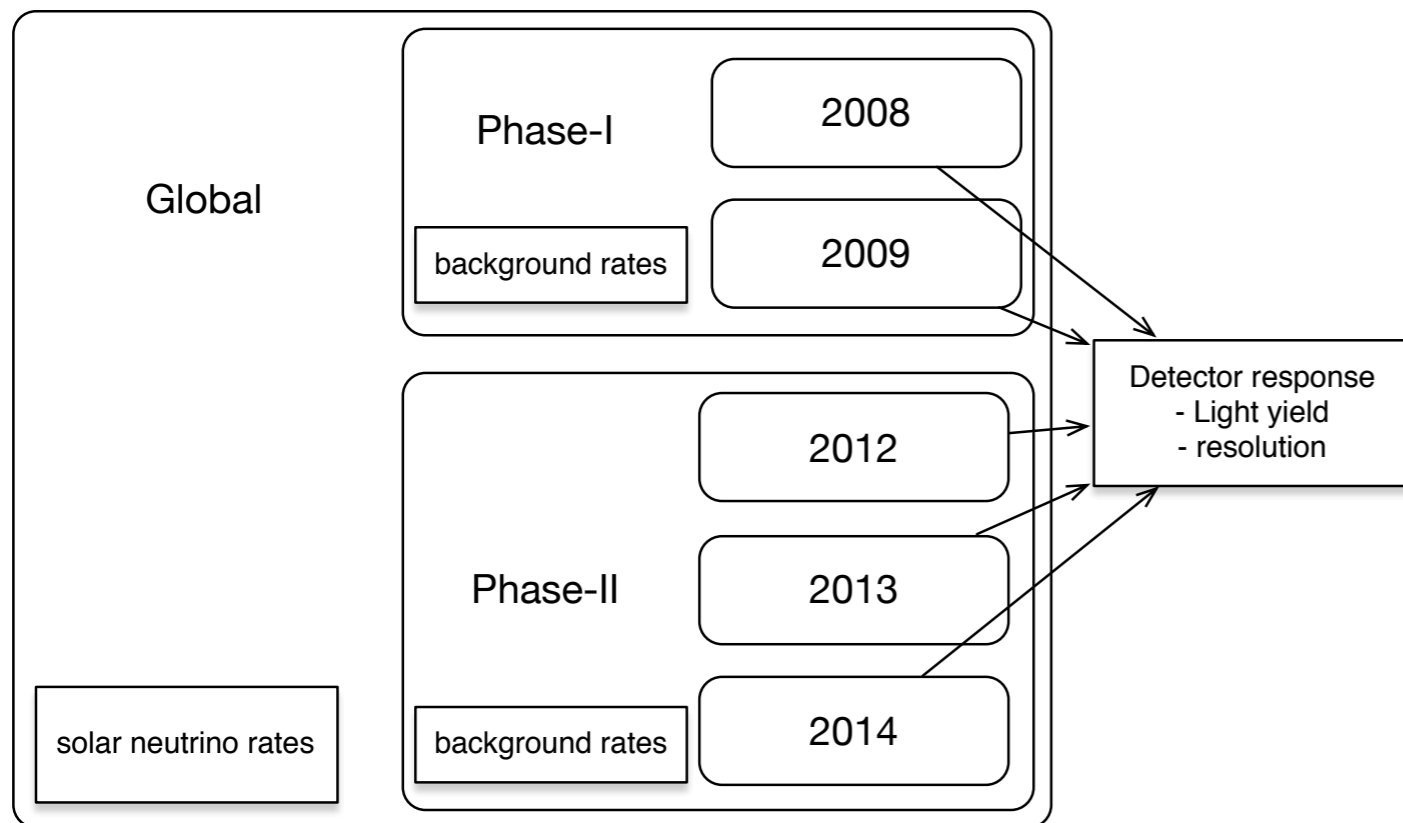
# Statistical analysis modules

- Adding tasks of statistical analysis by registering new modules inheriting abstract classes

- A few pre-installedt modules: ScanPar, DiscoveryTest..

```cpp
PrepareData *data = new PrepareData();
SimpleFit *fit = new SimpleFit();
DiscoveryTest *discovery = new DiscoveryTest();

ana->registerModule(inputManager);
ana->registerModule(data);
ana->registerModule(fit);
ana->registerModule(discovery);
ana->registerModule(outManager);
```

## Solar Neutrinos

## JUNO style

# Two ways of estimating $\sigma_{sys}$

- create an ensemble of models according to experiment precision

  - Fit the data with varying models through their coordinates

  - Generate pseudo-experiments with varying models and take the width of distribution of best-fit

# Method 1: fit with varying models

- **Full MC** during each iteration of the fit, vary **kb / absorption length spectrum etc.** **re-simulate and produce new pdf on the fly** —> when one day computer is fast enough

  - ~200, 000 CPU x years per fit

- **Semi-analytical:** analytical non-linearity model + response Matrix

  - ~30 minutes per fit

- **Full analytical**

  - ~2 hours per fit

- Systematic uncertainties of LY can be included by scaling the response matrix

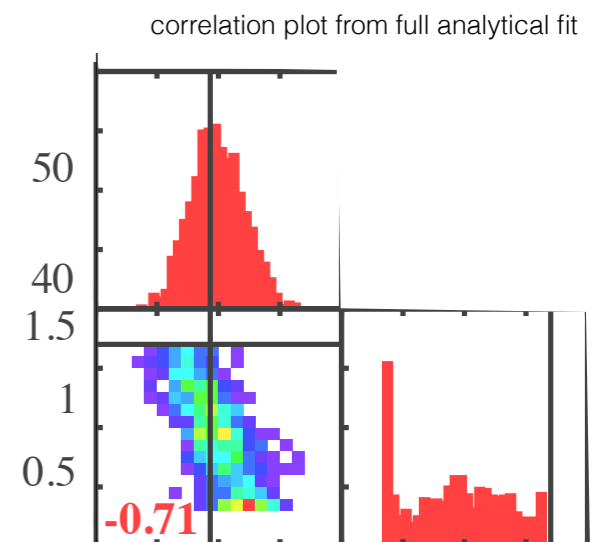- This could be dangerous if the interested parameters is correlated with the resolution parameters

- When LY is scaled, the resolution should also be changed



default.Major from 582.51 days × tons
-2ln(L) 8270.6 p-value 1.000 ± 0.000
$^{210}$Bi = 16.7 ± 2.2 (day×100t)$^{-1}$
$^{11}$C = 2.09 ± 0.30 (day×100t)$^{-1}$
$^{14}$C = 40.00 ± 0.99 (sec×100t)$^{-1}$ [p] 40.0 ± 1.0
Ext $^{214}$Bi = 1.86 ± 0.33 (day×100t)$^{-1}$
Ext $^{40}$K = 0.5 ± 1.5 (day×100t)$^{-1}$
Ext $^{208}$Tl = 3.35 ± 0.16 (day×100t)$^{-1}$
$^{85}$Kr = 6.4 ± 2.1 (day×100t)$^{-1}$
$^{210}$Po = 259.12 ± 0.74 (day×100t)$^{-1}$
$\nu(^{7}Be)_{862}$ = 46.3 ± 1.3 (day×100t)$^{-1}$
$\nu(^{7}Be)_{384}$ = 1.979 ± 0.057 (day×100t)$^{-1}$
$\nu(CNO)$ = 4.92 ± 0.56 (day×100t)$^{-1}$ [p] 4.92 ± 0.56
$\nu(pep)$ = 2.80 ± 0.48 (day×100t)$^{-1}$
$\nu(pp)$ = 133.8 ± 10.8 (day×100t)$^{-1}$

npmts_dt1 (LY=529.1±2.0 p.e./MeV, [2.00±0.74,3.00±1.84,7±0

correlation plot from full analytical fit

```
            ν(⁷Be)
48.5 (-2.1 +2.3)
inj 48.00 ± 0.07

         β-RPF(0)
 0.9 (-0.4 +0.3)
inj 1.49 ± 0.01
```
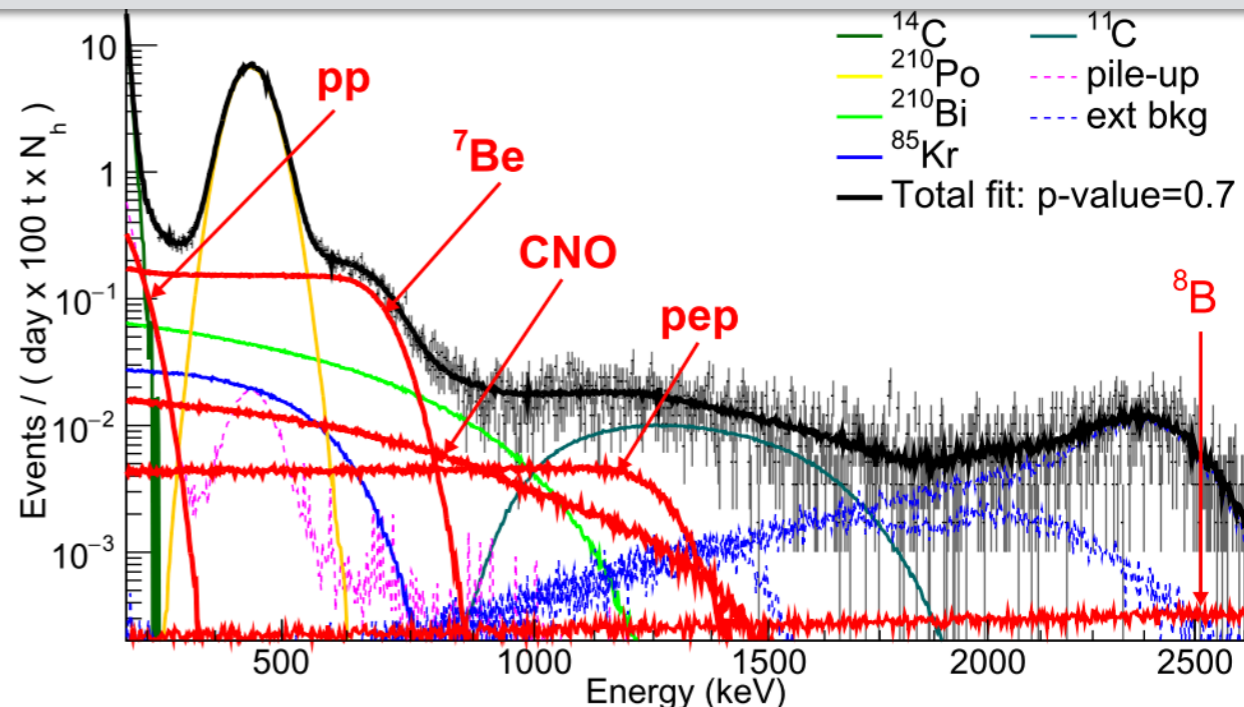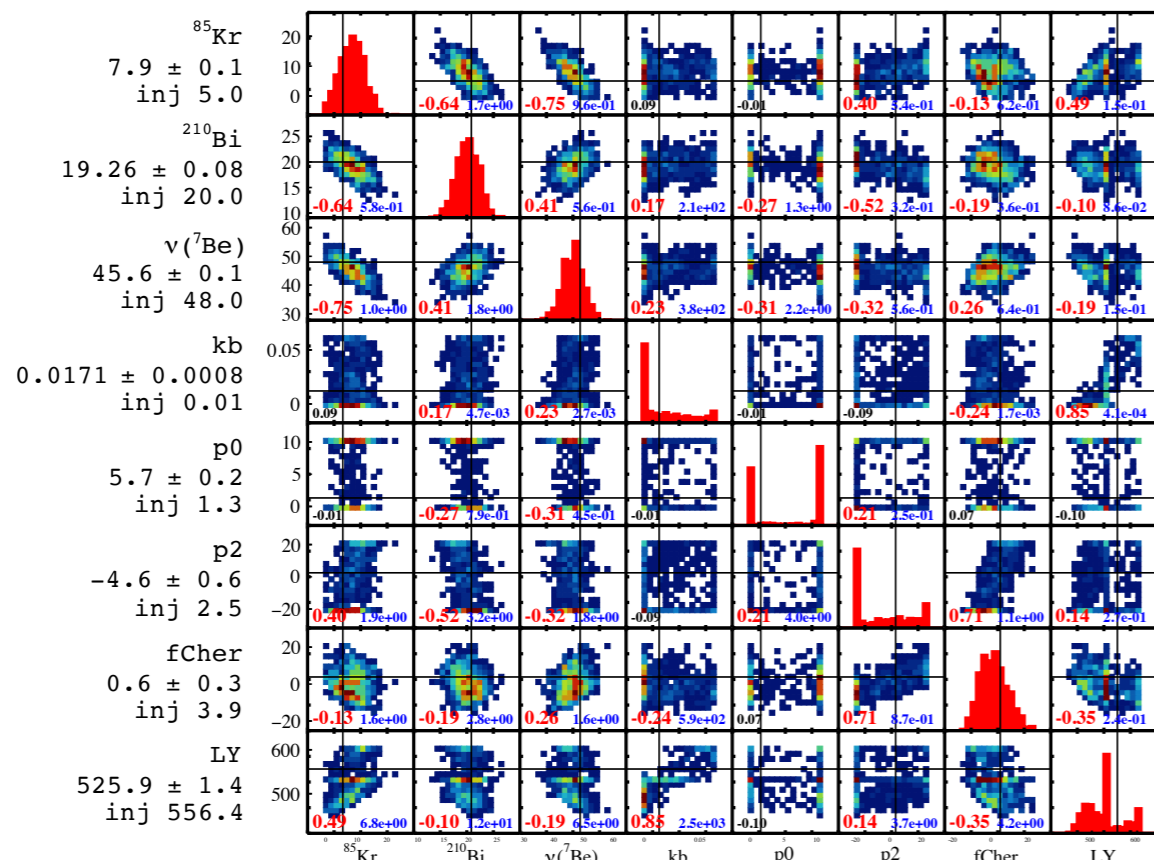


-0.71

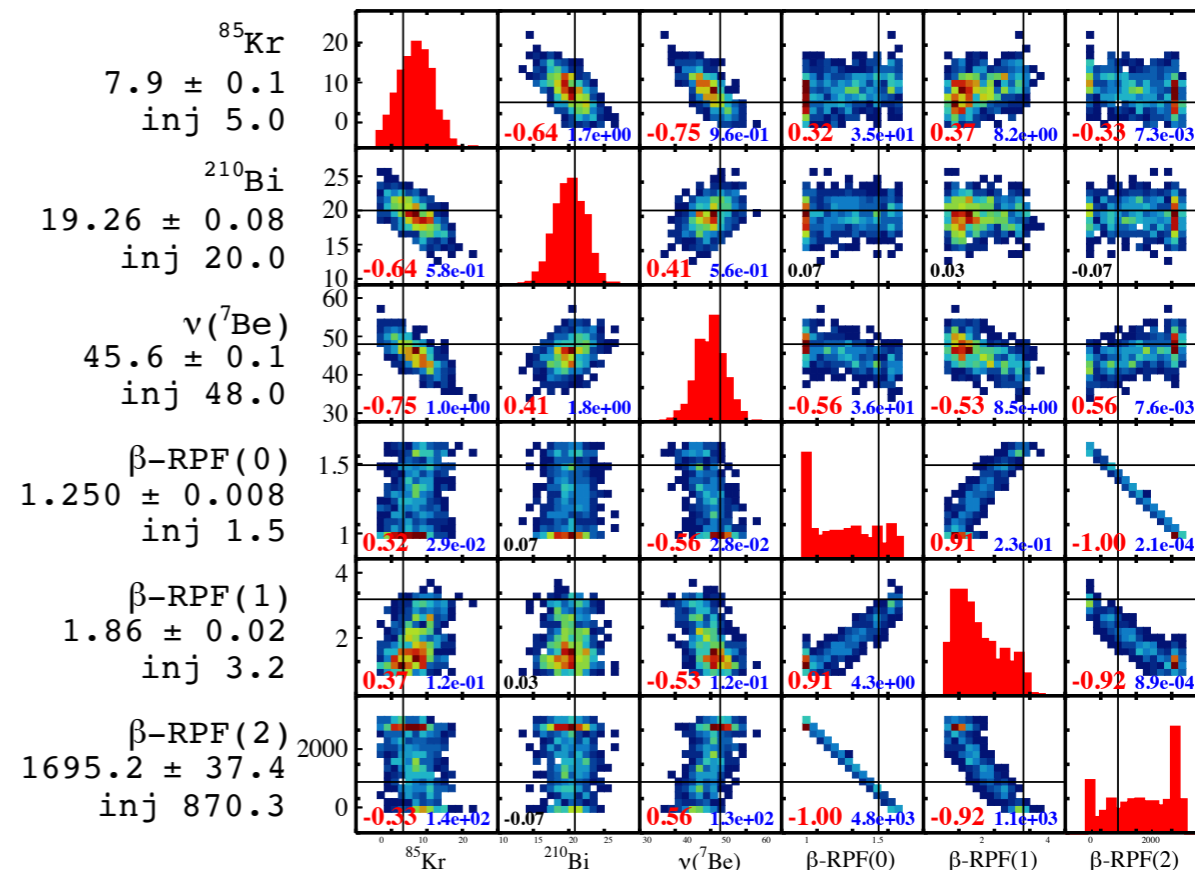- **Using full analytical response function, we can see the correlation with detector responses**

[1] M. Agostini *et al.*, "First Simultaneous Precision Spectroscopy of $pp$, $^7Be$, and $pep$ Solar Neutrinos with Borexino Phase-II," pp. 1–8, Jul. 2017.
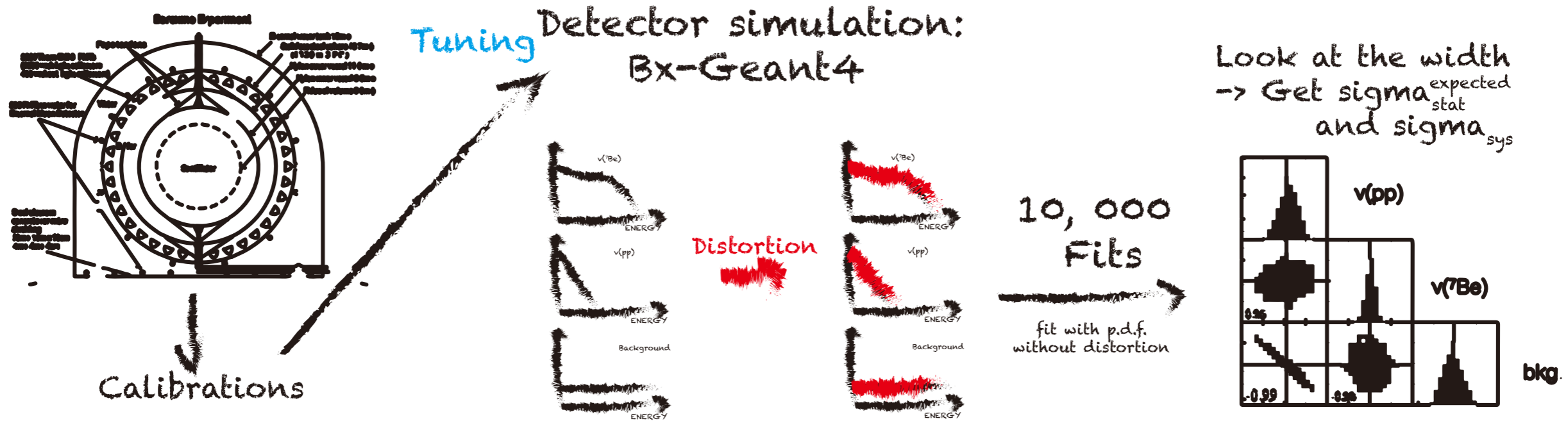
## correlation with NL



## correlation with res

- pseudo-experiment spectra without distortion —> **statistical sensitivity**

- pseudo-experiment spectra with distortion —> **statistical** **systematic uncertainty**

# Conclusion

- GooStats, a GPU based multivariate analysis framework is introduced.

- Statistical analysis module is easy to be implemented as needed. TTree/figure output provided.

- Full analytical response function has advantage that it treated the NL and resolution in a coherent way when evaluating the systematic uncertainties

- The systematic uncertainties can also be evaluated using Monte Carlo method by looking at change of width of best fit with/without distortion