

# Algorithm Challenges in the LHC Data Analysis

Alfio Lazzaro, CERN openlab



Second International Workshop for Future Challenges in Tracking  
and Trigger Concepts

CERN

July 7<sup>h</sup>–8<sup>th</sup>, 2011

# From the experiment to the discovery

## □ 5 main phases:

- **Events Acquisition**, online (High Level Trigger)
- **Reconstruction**
- **MC Simulation**
- Data analysis: **event selection** and **results extraction**

## □ Quite distinct problems

- Efforts should be consider differently for each use cases
- Of course there are overlaps...

## □ Ideally we would like to have data analysis directly during the reconstruction (or inside the online part)

- Reduce the amount of data to collect and the time for analysis (a ringing bell when a good event is found)
- Dangerous: remove potential good data for discoveries

- The goal of physics experiments is to provide results (described in papers) that can be compared to theory predictions
  - All phases in an experiment are important
  - First phases can be centralized (common experiment, common simulation, common reconstruction...)
    - Several experts involved, saving a lot of resources
  - Last phases are very chaotic
    - Very easy to do “mistakes”
    - Every analysis wants the power to analyze data independently
    - “Competition” of several groups on a given measurement

□ We can distinguish two main categories:

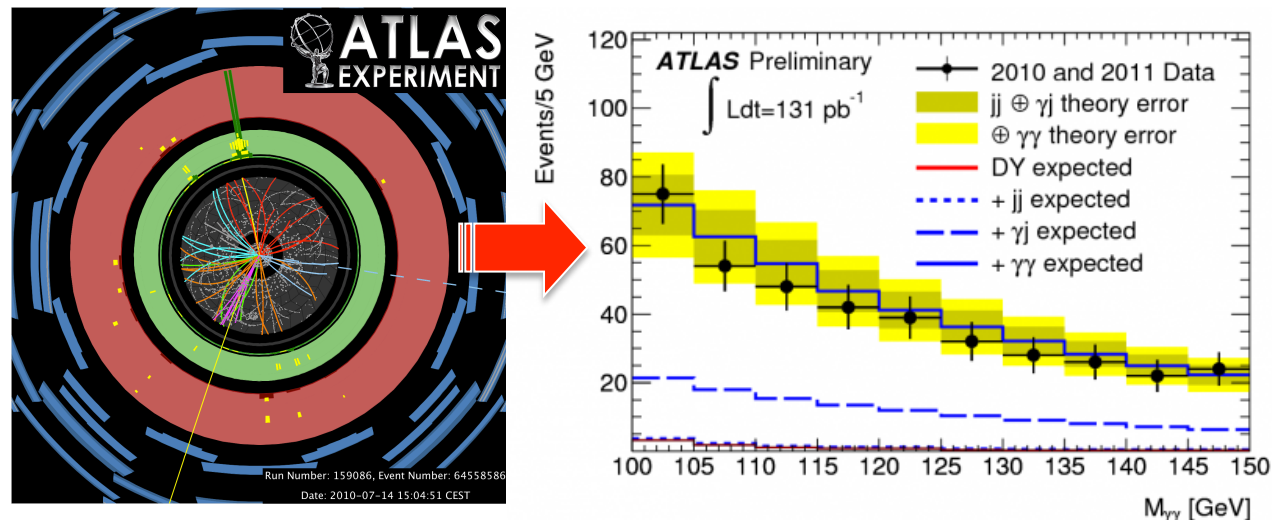
- **High precision measurements:** we know the physics phenomena and we want to compare the theory prediction with respect to results from data
  - Large sample of data to reduce the statistical errors
  - Improve experiments and data analysis techniques to reduce systematic errors
  - Push theory and experiment to the limits
- **Discovery measurement:** find new phenomena
  - Few events are enough to claim a discovery, but in most cases there is a tiny probability (cross section) to produce them and there is a huge contamination from other (well-know) events
  - We don't know where the new physics is, so in principle we want to collect and analyze everything!
  - Improve limits on the discoveries and hopefully move to precise measurements after the discovery

# Looking for signal events

- “If your experiment needs statistics, you ought to have done a better experiment.”

*Ernest Rutherford*

- Crucial to have a good discrimination between interesting (signal) events and the rest (background)
  - Data analysis techniques play a crucial role in this “war”



## □ Targets (2 main areas)

- Events selection: set of cuts applied on discriminant variables
- Signal/Background discrimination and parameters estimations

## □ Techniques

- Cuts Optimization (Bump Hunter), Fisher Discriminant, Neural Networks (NN), Boosted Decision Trees (BDT)...
- Maximum Likelihood (ML) fits

## □ Complexity

- Simple 1D fit / Cut&Count analyses
- Multivariate analyses for signal/background discrimination

# Common way to proceed (1)

- ❑ Efficient trigger (online selection) to reduce background events
- ❑ Reconstruction using framework (offline selection) for data and MC simulations
  - Runs on GRID
  - I/O bound (access to data)
  - Keep the useful information for a more aggressive selection, usually producing data for several analysis
- ❑ Further selection with more aggressive cuts
  - Event selection well performed in parallel using PROOF (data parallelism)
  - Local farm (Tier3)
  - Keep a small fraction of the initial sample

## Common way to proceed (2)

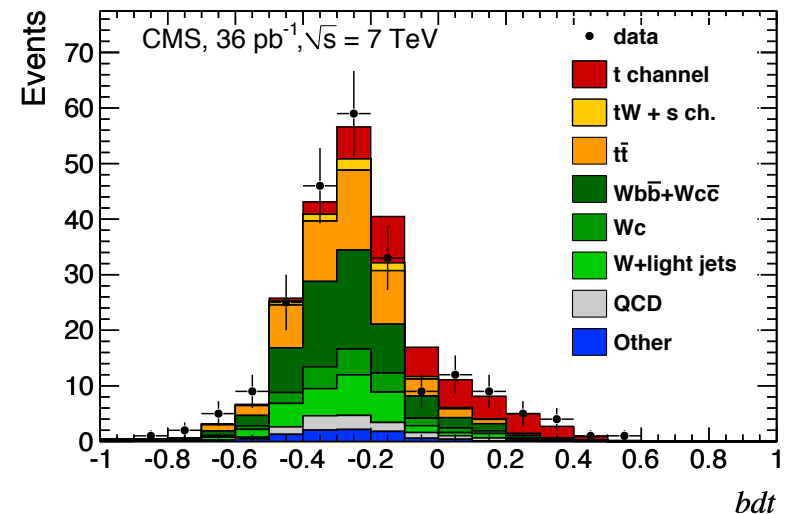
- The last steps are analysis- and user-dependent
  - In the **first phase** of the experiment it is reasonable to think that with small samples and (as usual for new experiments and in case of search for new phenomena) **simple analysis will be used** (events counting)
    - Reduce systematic errors estimations
  - Efforts will be concentrated to have results in a **reasonable time schedule**
  - Ideally you want to run on small systems at “home”: **laptops and desktops**
  - Long internal review inside the collaboration before publishing the results



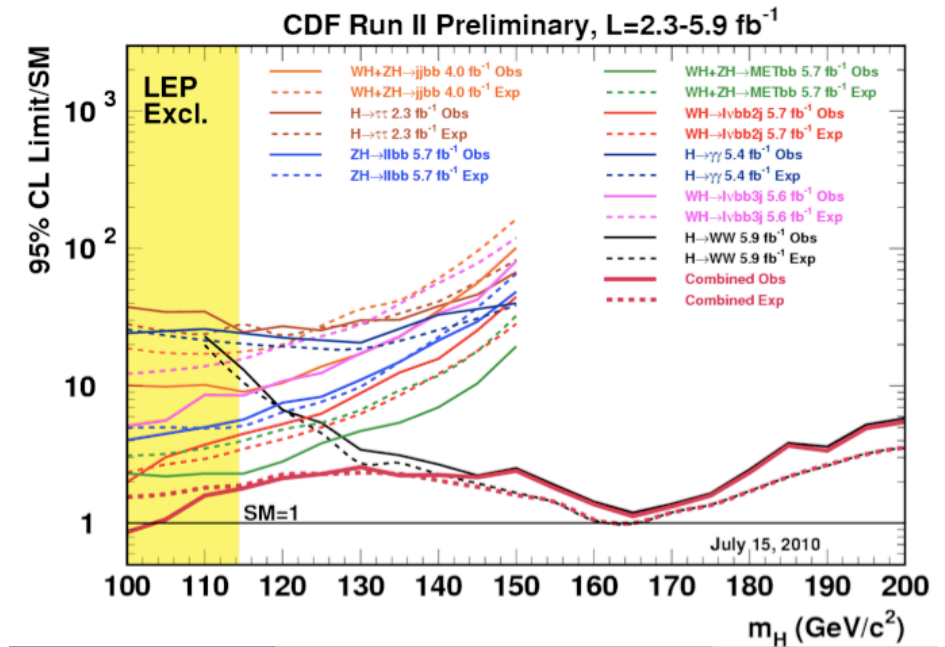
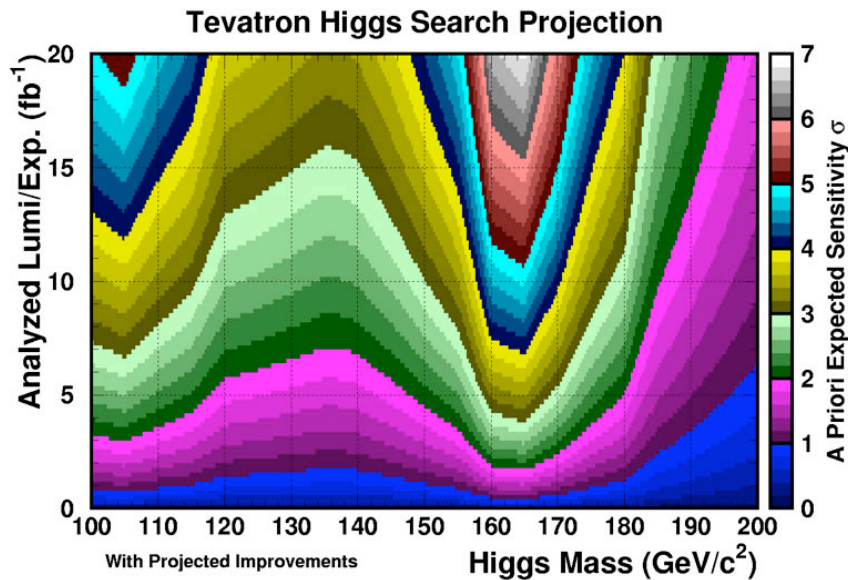
# Common way to proceed (3)

- Use MC or data-driven techniques to understand the background events
  - Tuning of the MC to data using standard “candles” events (well known physics events used as control sample)
  - Understand the detector and the systematics
  - After initial measurements, start the exploration of new “territory”
    - Use “golden mode” channel
    - Keep low systematics (statistical dominated)
    - 5 Gaussian standard deviations (sigmas) level to claim the discovery

Single top analysis at CMS  
(arXiv:1106.3052)  
3.5 sigmas



- Reduce limits up to 5 sigmas to rule out a model



- Several approaches
  - frequentist or bayesian
  - Hopefully the measurement is model-independent

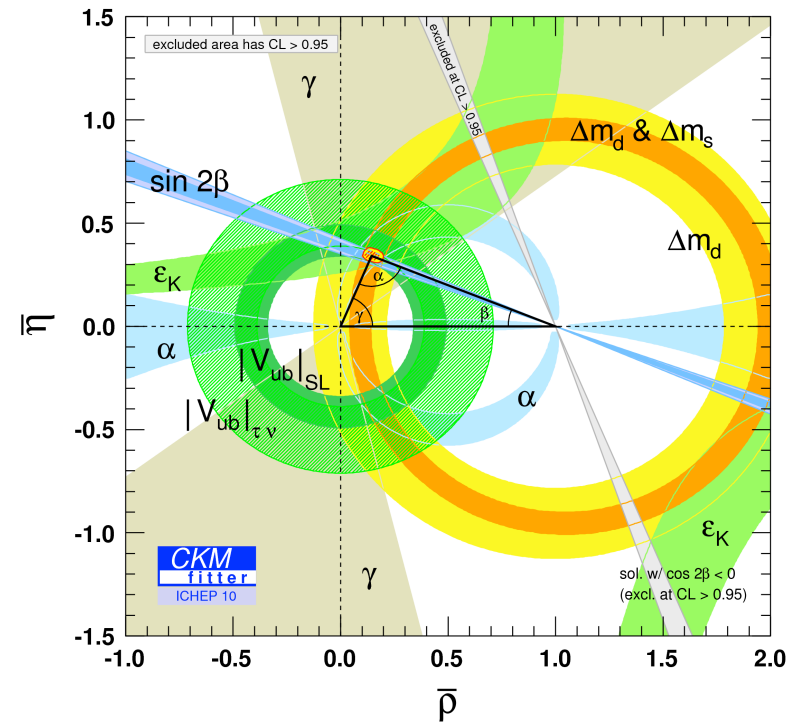
- Hope to find the Higgs...

- ❑ There is **not a common framework** for data analysis (like in reconstruction) for different analysts'
  - In general everybody wants the “power” to obtain the final results, i.e. **his own version of data analysis code**
- ❑ This means a “plethora” of programs
  - Not always based on the same base-code (different languages, Matlab, different algorithms...)
- ❑ Advantage: possible to make comparisons to spot bugs out
- ❑ **Disadvantage: “sometimes” all the versions are not well optimized**
- ❑ Inside the collaborations there is a general agreement to use common software (based on ROOT) from the beginning

# Combinations and Global fits

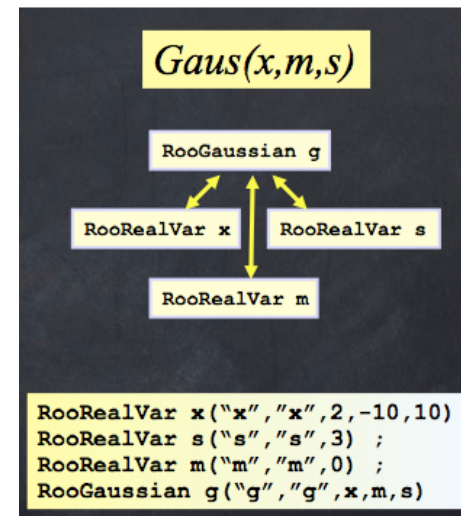
- Combine measurements from different experiments
  - Higgs combination for ATLAS + CMS & CDF + D0
  
- Combine several measurements to test an entire theory
  - EW model
  - SUSY models
  - CKM model

CKM fits



- RooFit is commonly used in High Energy Physics experiments to define the likelihood functions (W. Verkerke and D. Kirkby)
  - Inside ROOT. Details at <http://root.cern.ch/drupal/content/roofit>
  - Mathematical concepts are represented as C++ objects

| Mathematical concept |                                      |   | RooFit class    |
|----------------------|--------------------------------------|---|-----------------|
| variable             | $x$                                  | → | RooRealVar      |
| function             | $f(x)$                               | → | RooAbsReal      |
| PDF                  | $f(x)$                               | → | RooAbsPdf       |
| space point          | $\vec{x}$                            | → | RooArgSet       |
| integral             | $\int_{x_{\min}}^{x_{\max}} f(x) dx$ | → | RooRealIntegral |
| list of space points |                                      | → | RooAbsData      |



- On top of RooFit developed another package for advanced data analysis techniques, RooStats
  - Limits and intervals on Higgs mass and New Physics effects

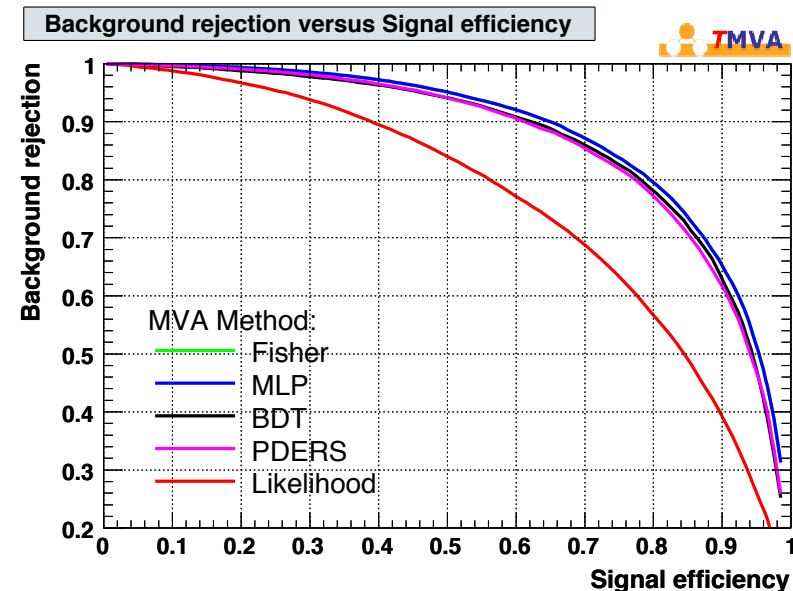
## ❑ Toolkit for Multivariate Data Analysis with ROOT

- Details at <http://tmva.sourceforge.net/>

## ❑ Several techniques, e.g.

- Function discriminant analysis (FDA)
- Artificial neural networks (three different MLP implementations)
- Boosted/Bagged decision trees
- Predictive learning via rule ensembles (RuleFit)
- Support Vector Machine (SVM)

## ❑ Working on parallelization of the package



# RooFit optimization and parallelization

- In the following slides I will focus on an R&D project we are doing for improvement the likelihood calculation in RooFit
  - Biased from my experience in the Babar and Atlas experiments. However, **data analysis is not our goal, so we don't focus on any specific analysis**
    - Strong collaboration with physics collaborators to have wide coverage of different analyses
- Our way to proceed:
  - Understanding the current version of the algorithm
  - Rewriting the algorithm so that we can improve it
    - Optimizations, vectorization, numerical accuracy
  - Apply parallelization
  - Porting the algorithm on accelerators (see tomorrow Yngve's talk)

# Likelihood-based techniques

- Data are a collection of independent events
  - an event consists of the measurement of a set of variables (energies, masses, spatial and angular variables...) recorded in a brief span of time by the physics detectors
- Introducing the concept of probability  $\mathcal{P}$  (= Probability Density Function, PDF) for a given event to be signal or background, we can combine this information for all events in the *likelihood function*

$$\mathcal{L} = \prod_{i=1}^N \mathcal{P}(\hat{x}_i | \hat{\theta})$$

$N$  number of events

$\hat{x}_i$  set of variables for the event  $i$

$\hat{\theta}$  set of parameters

- Several data analysis techniques requires the evaluation of  $\mathcal{L}$  to discriminate signal versus background events



- It allows to estimate free parameters over a data sample, by minimizing the corresponding Negative Log-Likelihood ( $NLL$ ) function (extended likelihood)

$$NLL = \sum_{j=1}^s n_j - \sum_{i=1}^N \left( \ln \sum_{j=1}^s n_j \mathcal{P}_j(\hat{x}_i | \hat{\theta}_j) \right)$$

$s$  species, i.e. signals and backgrounds

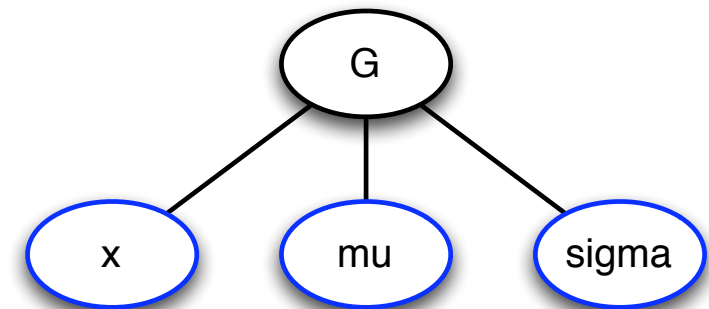
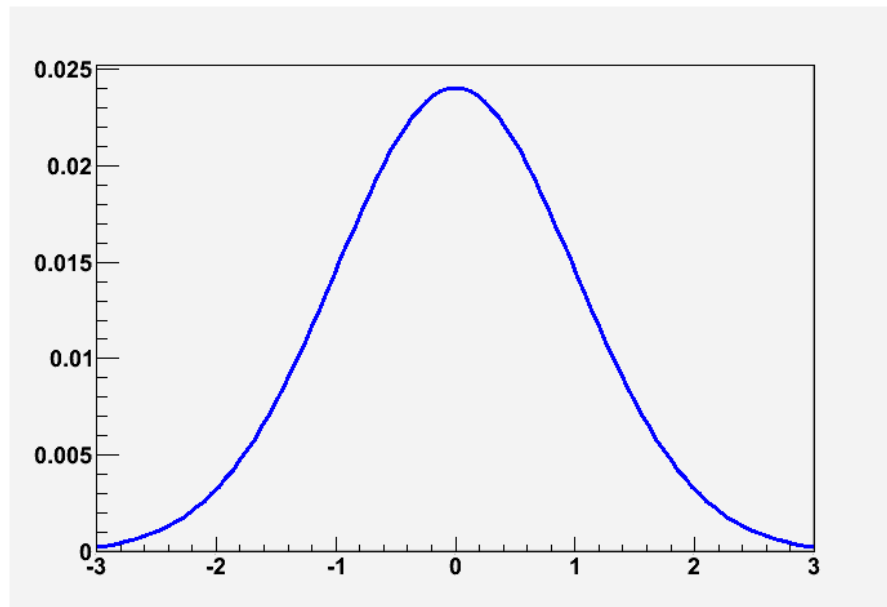
$n_j$  number of events belonging to the species  $j$

- The procedure of minimization can require several evaluation of the  $NLL$ 
  - Depending on the complexity of the function, the number of observables, the number of free parameters, and the number of events, the entire procedure can require long execution time
  - Mandatory to speed-up the execution

- In most cases PDFs can be factorized as product of the  $n$  PDFs of each variable (i.e. case of uncorrelated variables)

$$\mathcal{P}_j(\hat{x}_i | \hat{\theta}_j) = \prod_{v=1}^n \mathcal{P}_j^v(x_i^v | \hat{\theta}_j)$$

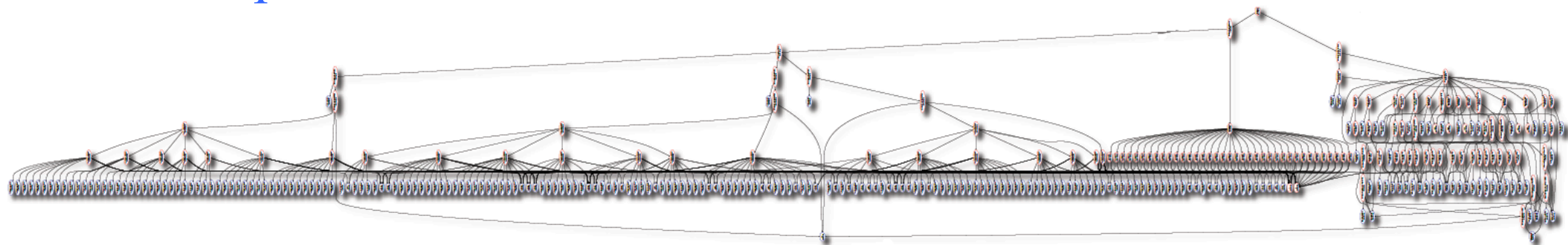
Gaussian  
 $G(x | \mu, \sigma)$



- In most cases PDFs can be factorized as product of the  $n$  PDFs of each variable (i.e. case of uncorrelated variables)

$$\mathcal{P}_j(\hat{x}_i | \hat{\theta}_j) = \prod_{v=1}^n \mathcal{P}_j^v(x_i^v | \hat{\theta}_j)$$

Combined Atlas & CMS Higgs analysis:  
12 variables  
50 free parameters



- Numerical minimization of the *NLL* using MINUIT (F. James, Minuit, Function Minimization and Error Analysis, CERN long write-up D506, 1970)
- MINUIT uses the gradient of the function to find local minimum (MIGRAD), requiring
  - The calculation of the gradient of the function for each free parameter, naively

$$\frac{\partial NLL}{\partial \hat{\theta}} \Big|_{\hat{\theta}_0} \approx \frac{NLL(\hat{\theta}_0 + \hat{d}) - NLL(\hat{\theta}_0 - \hat{d})}{2\hat{d}}$$

2 function calls  
per each  
parameter

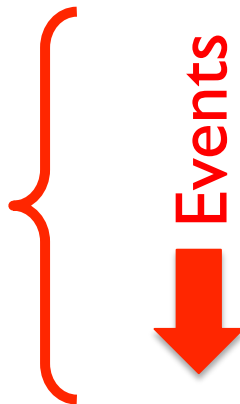
- The calculation of the covariance matrix of the free parameters, i.e. evaluation of the second order derivatives
- The minimization is done in several steps moving in the Newton direction: each step requires the calculation of the gradient
  - ⇒ Several calls to the *NLL*


- ❑ We developed a **new algorithm** for the likelihood function evaluation to be added in RooFit
  - We don't replace the current RooFit algorithm, which is used for results checking
  - Very chaotic situation: users can implement any kind of model
  - No need to change the user code to use the new implementation, i.e. same interface (use a simple flag to switch to the new algorithm)
- ❑ The new algorithm is optimized to run on the CPU
  - Used as reference for the GPU implementation: “fair” comparison
- ❑ All data in the calculation are in double precision floating point numbers
- ❑ Our target is to use commodity systems (e.g. laptops or desktops), easily accessible to data analysts
  - Of course we tests also on server systems

# Likelihood Function evaluation in RooFit (1)

1. Read the values of the variables for each event
2. Make the calculation of PDFs for each event
  - Each PDF has a common interface declared inside the class RooAbsPdf with a **virtual method** which defines the function
  - Automatic calculation of the normalization integrals for each PDF
  - Calculation of composite PDFs: sums, products, extended PDFs
3. Loop on all events and make the calculation of the *NLL*
  - *A single loop for all events*

Parallel execution over the events (*by fork*), with final reduction of the contributions



Variables 

|     | var <sub>1</sub> | var <sub>2</sub> | ... | var <sub>n</sub> |
|-----|------------------|------------------|-----|------------------|
| 1   |                  |                  |     |                  |
| 2   |                  |                  |     |                  |
| ... |                  |                  |     |                  |
| N   |                  |                  |     |                  |

# Likelihood Function evaluation in RooFit (2)

$$\text{Ex: } \mathcal{P} = \mathcal{P}_A(a_i) \mathcal{P}_B(b_i)$$

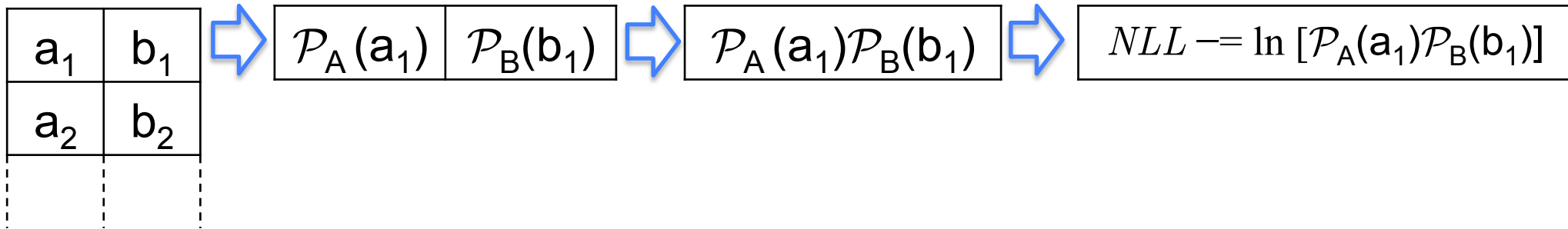
|           |
|-----------|
| $NLL = 0$ |
|-----------|

|       |       |
|-------|-------|
| $a_1$ | $b_1$ |
| $a_2$ | $b_2$ |

Vertical dashed lines extend downwards from the bottom of the table.

# Likelihood Function evaluation in RooFit (2)

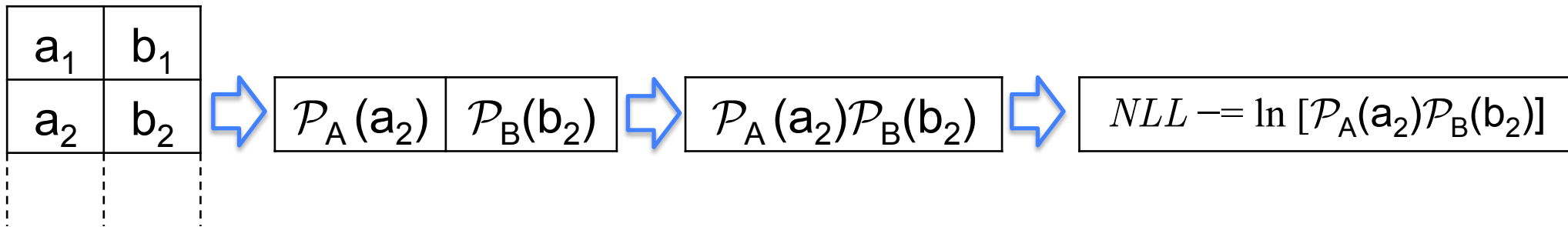
$$\text{Ex: } \mathcal{P} = \mathcal{P}_A(a_i) \mathcal{P}_B(b_i)$$





# Likelihood Function evaluation in RooFit (2)

Ex:  $\mathcal{P} = \mathcal{P}_A(a_i) \mathcal{P}_B(b_i)$



# Likelihood Function evaluation in RooFit (3)

Looping over all events and do the accumulation on  $NLL$

- Data are stored in something like ROOT TTree (RooTreeDataStore)
  - Very inefficient. At then our variables are simple float/double/int values
  - It breaks any possible vectorization
  - No thread safe, parallelization done with a fork, i.e. no shared memory
- In the C++ OO spirit, there is a common interface (RooAbsReal) and then virtual methods in all derivate classes
  - Each PDF calls virtual methods to access parameters, the observables, the integral value for the normalization, calculation of the ln's, ...
  - In case of composite PDFs (e.g. sums, products) it requires the call to virtual method of corresponding PDFs
  - **A lot of virtual function calls!**
- If the PDF doesn't change in the minimization, they are precalculated for all events and stored as a standard variable in the dataset
  - Not efficient way for caching the values of the PDFs
  - It doesn't take in account caching of constant values of the PDF inside a single minimization iteration

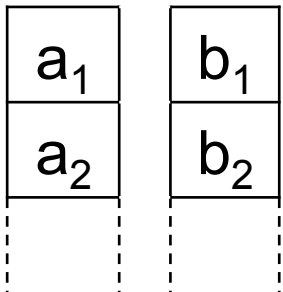
# Likelihood Function evaluation in RooFit (4)

- PDFs are considered as independent entities, i.e. a PDFs doesn't know if it is called inside a minimization process, from a mother composite PDF, or with a direct call
  - A PDF is not responsible to read the corresponding data
  - The PDF provides a single result for a given values of the data and parameters
  - In case of calculation which gives errors (e.g. negative probability), we get a warning message for the given values of the data and parameters
- Parallelization with a fork increases the memory footprint with the number of threads, but data are read-only!
  - Still it is easy to implement and it gives good scalability
- **At the end, we are doing the evaluation of functions (PDFs) over a vector of read-only data!**
  - **Suitable for loop parallelism (note functions can be very complex!)**

# New algorithm and parallelization (I)

1. Read all events and store in arrays in memory
2. For each PDF make the calculation on all events
  - ❑ Corresponding array of results is produced for each PDF
  - ❑ Evaluation of the function inside the local PDF
3. Combine the arrays of results (composite PDFs)
4. Loop over the final array of results to calculate  $NLL$  (final reduction)

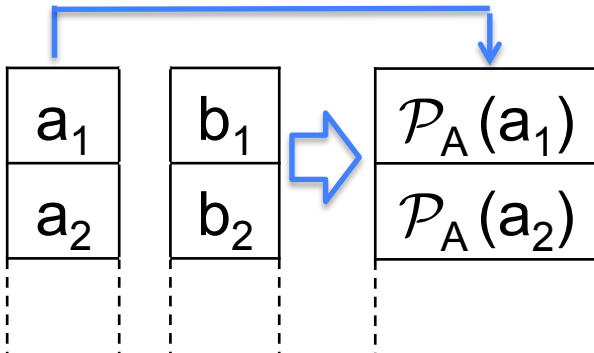
$$\text{Ex: } \mathcal{P} = \mathcal{P}_A(a_i) \mathcal{P}_B(b_i)$$



# New algorithm and parallelization (I)

1. Read all events and store in arrays in memory
2. For each PDF make the calculation on all events
  - ❑ Corresponding array of results is produced for each PDF
  - ❑ Evaluation of the function inside the local PDF
3. Combine the arrays of results (composite PDFs)
4. Loop over the final array of results to calculate  $NLL$  (final reduction)

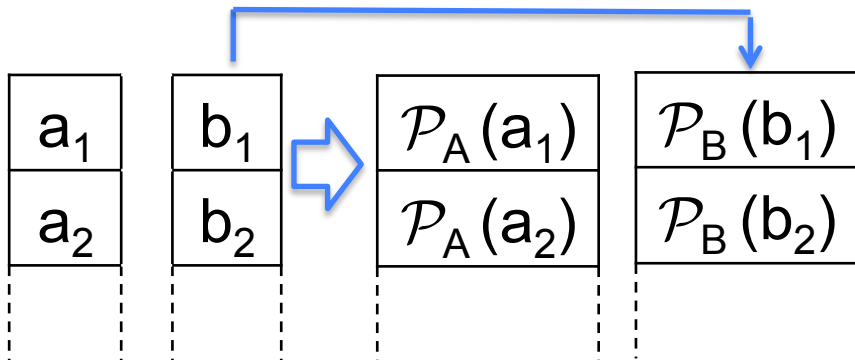
$$\text{Ex: } \mathcal{P} = \mathcal{P}_A(a_i) \mathcal{P}_B(b_i)$$



# New algorithm and parallelization (I)

1. Read all events and store in arrays in memory
2. For each PDF make the calculation on all events
  - ❑ Corresponding array of results is produced for each PDF
  - ❑ Evaluation of the function inside the local PDF
3. Combine the arrays of results (composite PDFs)
4. Loop over the final array of results to calculate  $NLL$  (final reduction)

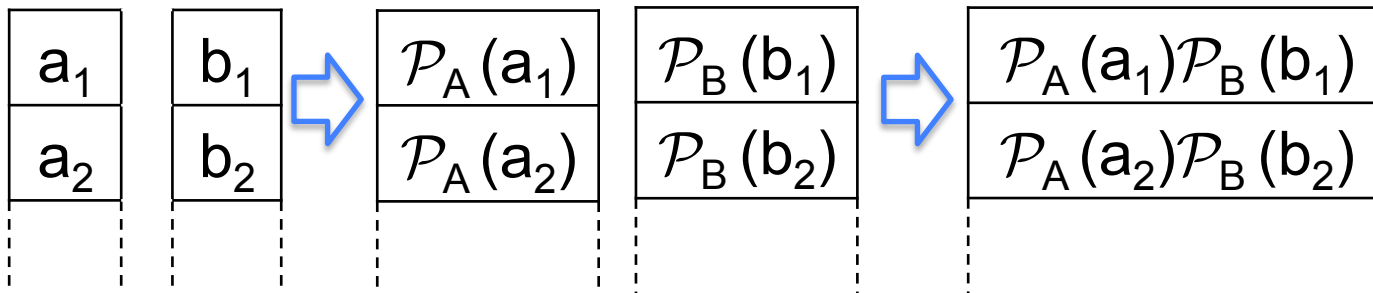
Ex:  $\mathcal{P} = \mathcal{P}_A(a_i) \mathcal{P}_B(b_i)$



# New algorithm and parallelization (I)

1. Read all events and store in arrays in memory
2. For each PDF make the calculation on all events
  - ❑ Corresponding array of results is produced for each PDF
  - ❑ Evaluation of the function inside the local PDF
3. Combine the arrays of results (composite PDFs)
4. Loop over the final array of results to calculate  $NLL$  (final reduction)

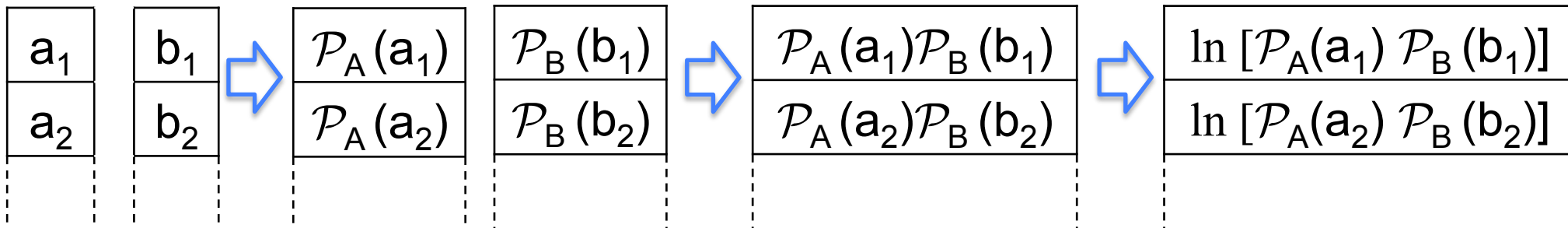
Ex:  $\mathcal{P} = \mathcal{P}_A(a_i) \mathcal{P}_B(b_i)$



# New algorithm and parallelization (I)

1. Read all events and store in arrays in memory
2. For each PDF make the calculation on all events
  - ❑ Corresponding array of results is produced for each PDF
  - ❑ Evaluation of the function inside the local PDF
3. Combine the arrays of results (composite PDFs)
4. Loop over the final array of results to calculate  $NLL$  (final reduction)

Ex:  $\mathcal{P} = \mathcal{P}_A(a_i) \mathcal{P}_B(b_i)$

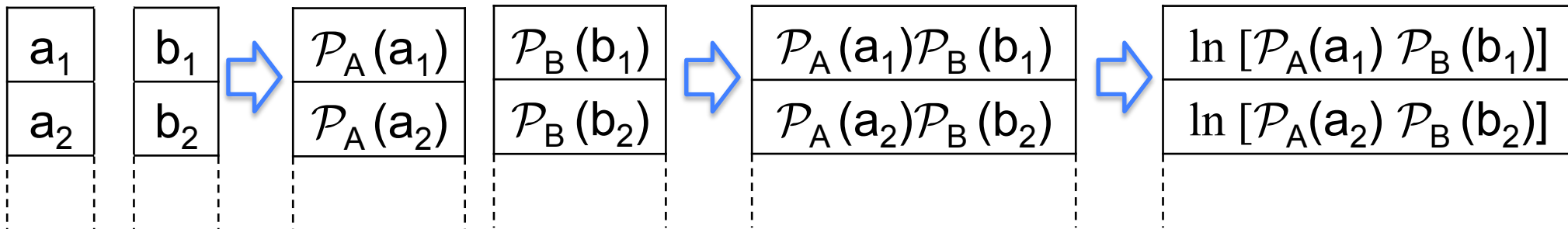




# New algorithm and parallelization (I)

1. Read all events and store in arrays in memory
2. For each PDF make the calculation on all events
  - ❑ Corresponding array of results is produced for each PDF
  - ❑ Evaluation of the function inside the local PDF
3. Combine the arrays of results (composite PDFs)
4. Loop over the final array of results to calculate  $NLL$  (final reduction)

Ex:  $\mathcal{P} = \mathcal{P}_A(a_i) \mathcal{P}_B(b_i)$



Final reduction in  $NLL$

## New algorithm and parallelization (2)

- Parallelization splitting calculation of each PDF over the events (**data parallelism**) and over the independent PDFs (**task parallelism**)
- Data are organized in vector, which are shared in memory
  - **Perfect for vectorization**
- Call the PDFs once for all events
  - **Reduce dramatically the number of virtual function calls!**
  - Perfect for caching values over the iterations during the minimization
- Drawbacks
  - Require to handle arrays of temporary results: 1 value per each event and PDF
  - Memory footprint increases with the number of events and number of PDFs, but not with the number of threads!
  - Due to the vectorization, we cannot have warning messages for a given event, but only at the end of the loop for the calculation over all events

- ❑ First of all we added a new class to manage the data as vectors (based on map of `std::vector`'s, where the key is the name of the observable)
- ❑ We added a class to take in account the array of results (based on `std::vector`)
- ❑ The loop parallelism is implemented using OpenMP
  - An OpenMP pragma loop for each loop used in the evaluation of the function
- ❑ **Added new** methods to the PDF interface
  - Still the old interface is working
- ❑ Using **Intel compiler** for the auto-vectorization of the loops (using `svml` library by Intel)
  - GNU compiler cannot auto-vectorize complex functions (like `exp`'s), unless you use intrinsics...

```
// Inline method for the Gaussian PDF calculation,
// defined inside the class RooGaussian
inline double evaluateLocal(const double x,
                           const double mu,
                           const double sigma) const
{
    return std::exp(-0.5*std::pow((x-mu)/sigma,2));
}

// Virtual method for the calculation of the
// Gaussian PDF on a single event
// (this is the original RooFit algorithm)
virtual double evaluate() const
{
    return evaluateLocal(x,mu,sigma);
}

// Virtual method for the calculation of the
// Gaussian PDF on all events
// (new implemented algorithm)
virtual bool evaluate(const RooAbsData& data)
{
    // retrieve the data array of values for the variable
    const double *dataArray = data.GetDataArray(x.arg());
    // check if there is an array for the variable
    if (dataArray==0)
        return false;

    // retrieve the number of events
    int nEvents = data.GetEntries();
    // retrieve the array for the partial results
    double *resultsArray = GetResultsArray();
    double m_mu = mu;
    double m_sigma = sigma;

    // loop over the events to calculate the Gaussian
    #pragma omp parallel for
    for (int idx = 0; idx<nEvents; ++idx) {
        resultsArray[idx] = evaluateLocal(dataArray[idx],
                                          m_mu,m_sigma);
    }

    return true;
}
```

# OpenMP parallelization

- ❑ Very easy parallelization with OpenMP
- ❑ Take benefit from the code optimizations
  - ❑ Inlining of the functions, no virtual functions
  - ❑ Data organized in C arrays, perfect for vectorization
- ❑ Easily avoid race conditions, keep the parallel region limited inside each PDF

- ❑ The final reduction for the *NLL* evaluation done in parallel using block-wise algorithm
  - Numerical approximation w.r.t. sequential reduction, which are number of threads dependent
  - Minuit is very sensitive to these approximation
    - Of course differences are negligible, but still they can worry people (and they can be non deterministic)
- ❑ We implemented a parallel reduction based on double-double algorithm which reduces the approximations (Y. He and C. H. Q. Ding, *The Journal of Supercomputing*, 18, 259–277, 2001; P. Kornerup *at al.*, *IEEE Transactions on Computers*, 01 Feb. 2011)
  - We need to switch off any compiler optimization inside the reduction, using pragmas
- ❑ **Now the results are identical up to  $10^{-6}$ , no matter how many threads you are running**

$$\begin{aligned} n_a [f_{1,a} G_{1,a}(x) + (1 - f_{1,a}) G_{2,a}(x)] A G_{1,a}(y) A G_{2,a}(z) + \\ n_b G_{1,b}(x) B W_{1,b}(y) G_{2,b}(z) + \\ n_c A R_{1,c}(x) P_{1,c}(y) P_{2,c}(z) + \\ n_d P_{1,d}(x) G_{1,d}(y) A G_{1,d}(z) \end{aligned}$$

Model from B. Aubert *et. al.*,  
Phys. Rev. Lett. 98, 031801, 2007

17 PDFs in total, 3 variables, 4 components, 35 parameters

- G: Gaussian
- AG: Asymmetric Gaussian
- BW: Breit-Wigner
- AR: Argus function
- P: Polynomial

40% of the  
execution time  
is spent in exp's  
calculation

Note: all PDFs have analytical normalization integral, i.e. >98%  
of the sequential portion can be parallelized

# Test on CPU in sequential

- Dual socket Intel Westmere-based system: CPU (L5640) @ 2.27GHz (12 physical cores, 24 hardware threads in total), 10x4096MB DDR3 memory @ 1333MHz
- Linux 64bit, Intel C++ compiler version 12.0.2

| # Events                            | 10,000 | 25,000 | 50,000 | 100,000 |
|-------------------------------------|--------|--------|--------|---------|
| <b>RooFit</b>                       |        |        |        |         |
| # <i>NLL</i> evaluations            | 15810  | 14540  | 19041  | 12834   |
| Time (s)                            | 826.0  | 1889.0 | 5192.9 | 6778.9  |
| Time per <i>NLL</i> evaluation (ms) | 52.25  | 129.92 | 272.72 | 528.19  |
| <b>OpenMP (w/o vectorization)</b>   |        |        |        |         |
| # <i>NLL</i> evaluations            | 15237  | 17671  | 15761  | 11396   |
| Time (s)                            | 315.1  | 916.0  | 1642.6 | 2397.3  |
| Time per <i>NLL</i> evaluation (ms) | 20.68  | 51.84  | 104.22 | 210.36  |
| w.r.t. RooFit                       | 2.5x   | 2.5x   | 2.6x   | 2.5x    |
| <b>OpenMP (w/ vectorization)</b>    |        |        |        |         |
| # <i>NLL</i> evaluations            | 15304  | 17163  | 15331  | 12665   |
| Time (s)                            | 178.8  | 492.1  | 924.2  | 1536.9  |
| Time per <i>NLL</i> evaluation (ms) | 11.68  | 28.67  | 60.28  | 121.35  |
| w.r.t. RooFit                       | 4.5x   | 4.5x   | 4.4x   | 4.4x    |

**4.5x faster!**

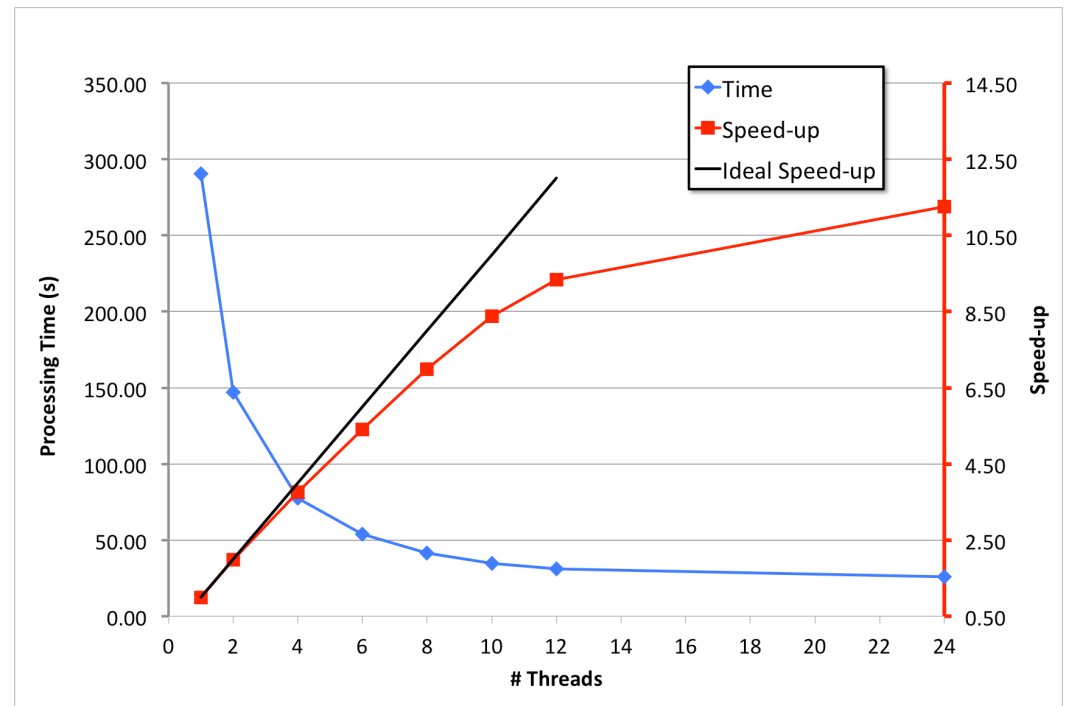


**Vectorization gives a 1.8x speed-up (SSE).**

**Additional 12% using AVX on Intel Sandy Bridge**

# Test on CPU in parallel

- ❑ Dual socket Intel Westmere-based system: CPU @ 2.67GHz (12 physical cores, 24 hardware threads in total), Turbo Mode ON, 10x4096MB DDR3 memory @ 1333MHz
- ❑ Linux 64bit, Intel C++ compiler version 12.0.2
- ❑ 100,000 events
- ❑ Data is shared, i.e. no significant increase in the memory footprint
  - Possibility to use Hyper-threading (about 20% improvement)
- ❑ Limited by the sequential part, OpenMP overhead, and memory access to data





- Scalability is limited by accessing the array of results
  - In particular the effect becomes important for PDFs with simple function, like polynomials and composite PDFs (add and prod)
  - We do pinning of the threads to the physical cores, taking in account the NUMA effect
  - However the performance depends on the cache memory available on the systems
    - Testing on a 4 core i7 desktop system (8 MB L3 cache) we reach a factor  $\sim 2x$  with 8 threads (using SMT)
- We solve this problem with different techniques
  - Merge the number of OpenMP parallel region and reuse the data (in particular for composite PDFs)
  - Do block-splitting, i.e. do full evaluation for small sub-groups of events
- Doing this optimization we are able to reach 4.6x on the 4 core i7 desktop system (8 threads with SMT)

- ❑ Generating events from PDFs and repeat on them the analysis (frequentist approach)
- ❑ Implemented in RooFit using PROOF
  - Not a parallel generator implemented (using pseudo-generators TRandom3 with different seeds)

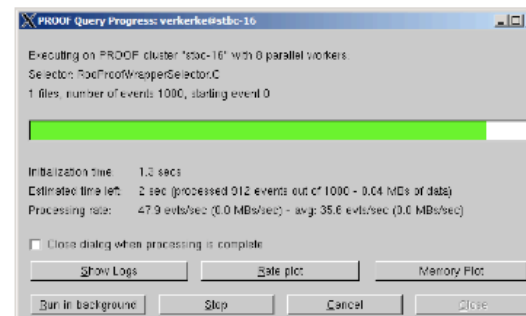
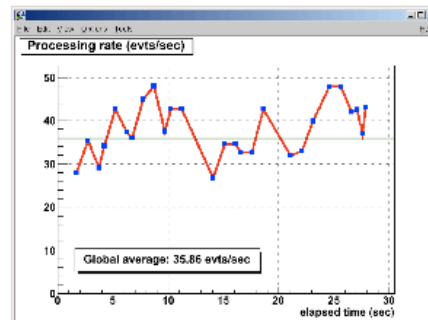
## Demo of parallelization with PROOF-lite

- Example – Factor 8 speed up on a dual-quad core box.
  - Works with out-of-the box ROOT distribution
  - Also: Graceful early termination when users presses 'Stop'

```

RooStudyManager mcs(*w,gfs) ;
mcs.run(1000) ; // inline running
mcs.runProof(1000,"") ; // empty string is PROOF-lite
mcs.prepareBatchInput("default",1000,kTRUE) ;

```



- Data analysis will be the major challenge in the next year at LHC
  - Squeeze all possible physics results from data
  - Analyses will increase their complexity in the next future
    - Intensive analysis period during the 2013-2014 shutdown
- Other experiments will have the same complexity of LHC
  - SuperB, Panda experiments...
- It is useful to converge on common efforts
  - Common softwares
  - Extending and improve existing software
    - Parallelization is mandatory in a lot of analysis
  - Better selection during the online can help a lot the final analysis

- In our work the OpenMP implementation required not so drastic changes in the existing RooFit code
  - In any case we added our implementation, so that users can use the original implementation for reference
- Optimization gives a great speed-up:  $\sim 5x$
- Note that our target is running at the user-level of small systems (laptops, desktops), i.e. with small number of CPU cores
- Very important to take under control numerical accuracy
  - We would like to try single precision in case of PDF evaluation, moving to double precision for the final reduction
    - Reduce memory footprint (half space for results)
    - Gain a factor possible  $2x$  from vectorization