

Maximum Likelihood Fits on GPUs

S. Jarp, A. Lazzaro, J. Leduc,
A. Nowak, F. Pantaleo
CERN openlab

Openlab Minor review meeting
November 2nd, 2010



- We have a sample composed by N events, belonging to s different specie (signals, backgrounds), and we want to extract the number of events for each species and other parameters
- We use the *Maximum Likelihood fit technique* to estimate the values of the free parameters, **minimizing** the Negative Log-Likelihood (NLL) function

$$NLL = \sum_{j=1}^s n_j - \sum_{i=1}^N \left(\ln \sum_{j=1}^s n_j \mathcal{P}_j(x_i; \theta_j) \right)$$

j species (signals, backgrounds)

n_j number of events

\mathcal{P}_j probability density function (PDF)

θ_j Free parameters in the PDFs

- 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

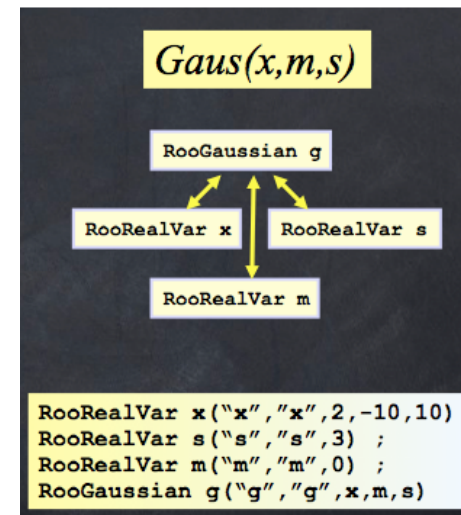
$$\left. \frac{\partial NLL}{\partial \hat{\theta}} \right|_{\hat{\theta}_0} \approx \frac{NLL(\hat{\theta}_0 + \hat{d}) - NLL(\hat{\theta}_0 - \hat{d})}{2\hat{d}}$$

2 function calls
per each
parameter

- The calculation of the covariance matrix of the free parameters (which means 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*

- RooFit is a Maximum Likelihood fitting package (W. Verkerke and D. Kirkby) for the *NLL* calculation
 - Inside ROOT (details at <http://root.cern.ch/drupal/content/roofit>)
 - Allows to build complex models and declare the likelihood function
 - 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

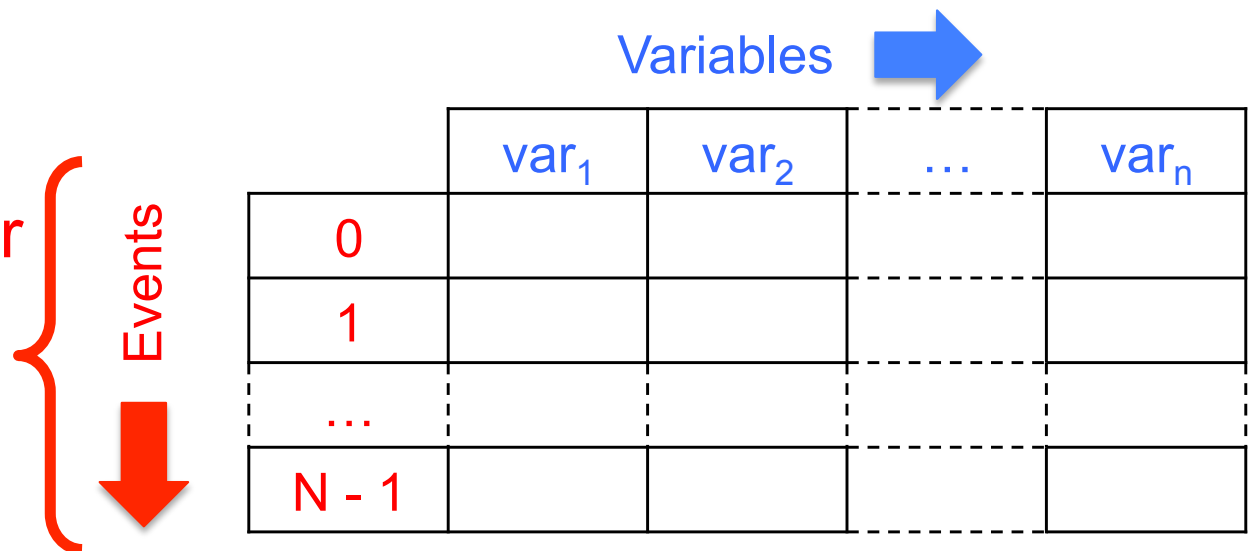


- On top of RooFit developed another package for advanced data analysis techniques, RooStats
 - Limits and intervals on Higgs mass and New Physics effects

Likelihood Function calculation in RooFit

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** `evaluate()` which define the function
 - Each PDF implements the method `evaluate()`
 - 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*

Parallel execution over the events (as it is already implemented)



- Two algorithms implemented:
 1. RooFit Event-based (CPU Implementation), described before
 - Parallelization at event level, using fork
 - Not shared resources
 2. PDF-Event-based Algorithm
 - GPU Implementation (CUDA)
 - CPU Implementation (OpenMP)



NEW

Note: everything done in double precision

New approach to the *NLL* calculation:

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, i.e. **not need a virtual function** (drawback: require more memory to store temporary results: 1 double per each event and PDF)
 - Apply normalization
3. Combine the arrays of results (composite PDFs)
4. Calculation of the *NLL*

Parallelization splitting calculation of each PDF over the events

- Particularly suitable for thread parallelism on GPU, requiring one thread for each PDF/event
- Possible benefit from vectorization on the CPU

- PCs
 - CPU: Nehalem @ 3.2GHz: 4 cores – 8 hw-threads
 - OS: SLC5 64bit - GCC 4.3.4
 - ROOT trunk (October 11th, 2010)
- GPU: ASUS nVidia GTX470 PCI-e 2.0
 - Commodity card (for gamers)
 - Architecture: GF100 (Fermi)
 - Memory: 1280MB DDR5
 - Core/Memory Clock: 607MHz/837MHz
 - Maximum # of Threads per Block: 1024
 - Number of SMs: 14
 - CUDA Toolkit 3.1 06/2010
 - Developer Driver 256.40
 - Power Consumption 200W
 - Price ~\$340



- 1D PDFs commonly used in HEP:
 - Symmetric and Asymmetric Gaussian
 - Breit-Wigner
 - Crystal Ball Function
 - Argus
 - Generic Polynomial
 - Chi Square
- Composition of PDFs:
 - Sum of two or more PDFs
 - Product of two or more PDFs
 - Multivariate PDFs
- **Very easy to build complex models (via composition) and add new PDFs**

CPU (existing code from RooFit)

```
Double_t Pdf::evaluate() const // virtual method
{
    return <function to be evaluated>(<data>,<pars>); // data and pars are data members of the Pdf Class
}
```

GPU

```
Double_t Pdf::evaluate() const // virtual method
{
    return evaluateLocal(<data>,<pars>); // data and pars are data members of the Pdf Class
}

//inside class definition
#ifdef __USECUDA__
    friend void KernelEvaluatePdf(const Pdf* pdf, const Double_t* data,
        <par>, Double_t *results, const UInt_t N);

    __host__ __device__
#endif
inline Double_t Pdf::evaluateLocal(<data>,<pars>) const // inline method
{
    return <function to be evaluated>(<data>,<pars>); // data and pars are local to the method
}
```

GPU code (Kernel implementation)

```
#ifdef __USECUDA__

Bool_t Pdf::evaluate(const Data& data) const // data links events on GPU and CPU
{
    // set pointer to GPU for the results of the evaluation
    ...
    // pars are local to the method
    KernelEvaluatePdf<<<NUM_BLOCKS,NUM_THREADS>>>(this,<pointer to GPU data>,pars,<pointer to GPU results>,data.size());

    return kTRUE;
}

__global__ void KernelEvaluatePdf(const Pdf* pdf, const Double_t* data,
    <pars>, Double_t *results, const UInt_t N)
{
    UInt_t idx = blockIdx.x * blockDim.x + threadIdx.x;
    if (idx<N) {
        results[idx] = pdf->evaluateLocal(data[idx],<pars>);
    }
}

#endif
```

- Data are copied on the GPU once
- Results for each PDF are resident only on the GPU
 - Arrays of results are allocated on the **global memory once** and they are deallocated at the end of the fitting procedure
 - **Minimize CPU ↔ GPU communication**
 - Only the final results are copied on the CPU for the final sum to compute *NLL*
- Device algorithm performance with a linear polynomial PDF and 1,000,000 events
 - **45 GFLOPS and 3.5 GB/s CPU ↔ GPU data transfer**

1,000,000 events and 1000 iterations

PDF Name	Formula	CPU vs GPU time ratio	kernels execution time portion
Gaussian	$\exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$	11.3	21.3%
Asymmetric Gaussian	$\begin{cases} \exp\left\{-\frac{(x-\mu)^2}{2\sigma_l^2}\right\} & x \leq \mu \\ \exp\left\{-\frac{(x-\mu)^2}{2\sigma_r^2}\right\} & x > \mu \end{cases}$	11.9	21.9%
Breit-Wigner	$\frac{1}{(x-x_0)^2 + \Gamma^2/4}$	2.4	13.4%
χ^2	$\frac{1}{2^{k/2}\Gamma(\frac{k}{2})} x^{\frac{k}{2}-1} \exp\{-x/2\}, k = 5$	27.6	70.0%
Argus	$\sqrt{\left(1 - \frac{x^2}{c^2}\right)} \exp\left\{-\frac{1}{2}\eta\left(1 - \frac{x^2}{c^2}\right)\right\}$	25.3	50.1%

- CPU algorithm is the event-based (RooFit) in sequential
- GPU time includes data transfer time (data and results)
 - A significant portion of time, limiting the scalability
 - More complex PDF => Bigger portion of time spent in evaluation VS time for data transfers

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

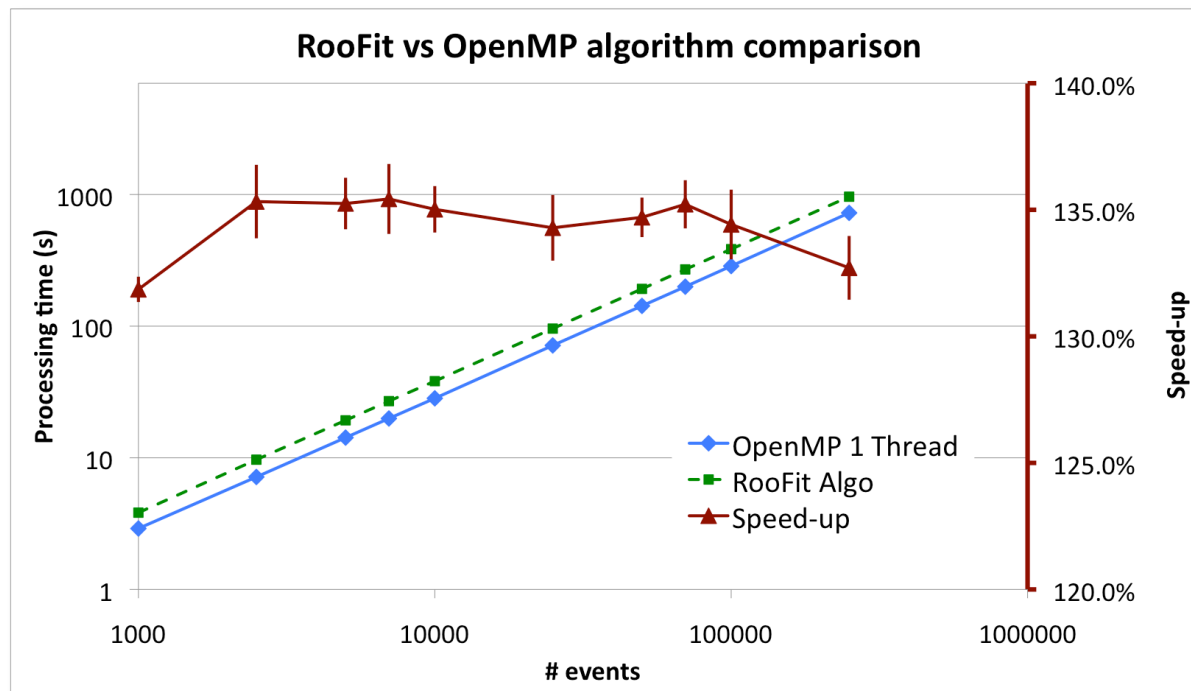
17 PDFs in total, 3 variables, 4 components, 35 parameters

- G: Gaussian
- AG: Asymmetric Gaussian
- BW: Breit-Wigner
- P: Polynomial

Note: all PDFs have analytical normalization integral

Event-based VS PDF-event-base performance

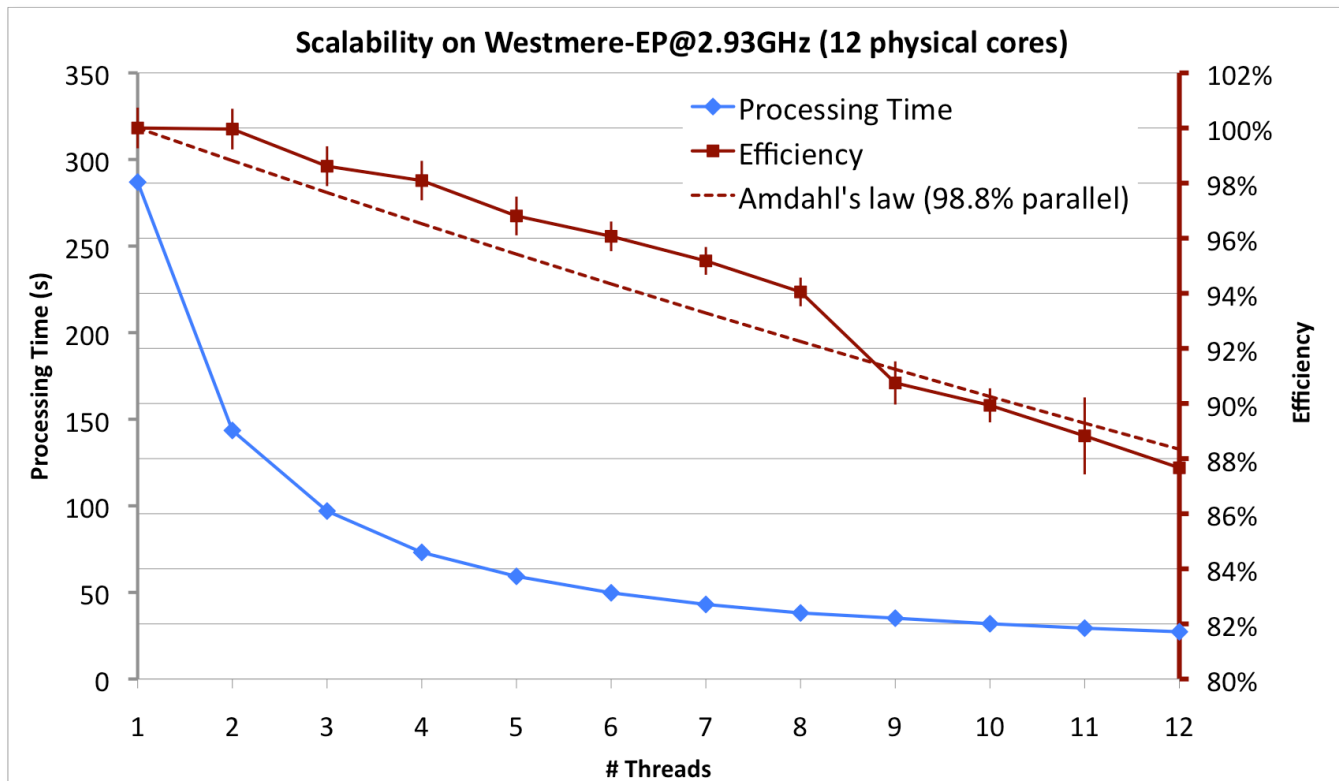
- ❑ Driven by the GPU implementation, we implemented a corresponding CPU implementation
 - take benefit from the code optimizations (due to migration from C++ to C)
 - ❑ No virtual functions
 - ❑ Inlining of the `evaluate` function
 - ❑ Data organized in C arrays, perfect for vectorization
 - it can be easily parallelized using OpenMP



- Linear increase with the number of events (as expected)
- Speed-up of 34% (almost flat over the number of events), just optimizing the algorithm! (not parallelization)

PDF-event-base scalability with OpenMP

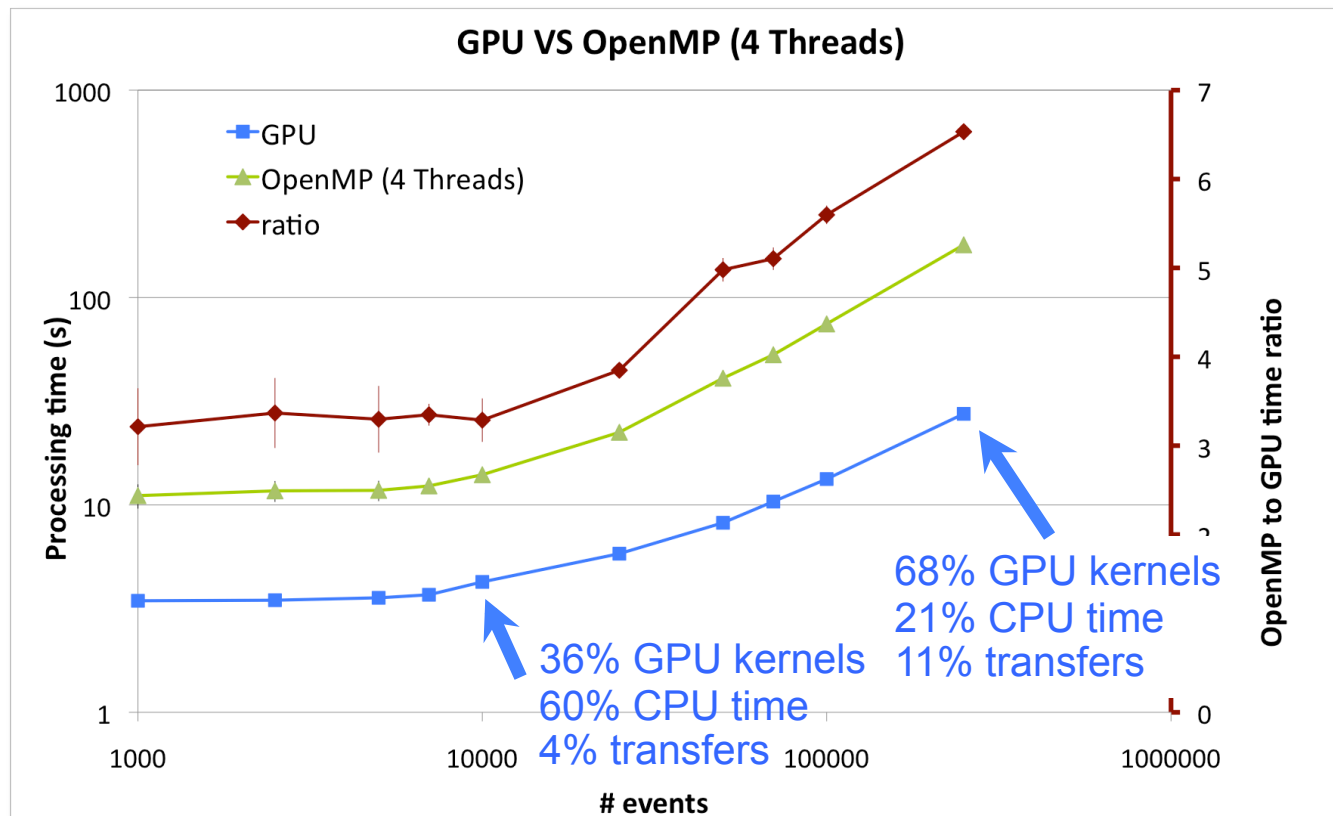
- ❑ Test done on the Westmere-EP @ 2.93 GHz
 - ❑ 12 cores / 24 threads
- ❑ 100,000 events
- ❑ 98.8% of the sequential execution can be parallelized (1.2% required for initialization of the arrays for data and results and normalization integrals calculation)



- Negligible increase in memory (arrays are shared)
- Scalability as expected
- Using SMT (hw-threading) with 24 threads we reach 110% in efficiency w.r.t 12 threads (+32% in case of ideal speed-up)

PDF-event-base: GPU VS OpenMP

- ❑ Fair comparison
 - ❑ Same algorithm
 - ❑ Algorithm on CPU optimized and parallelized (4 threads)
 - ❑ CPU does the final sum of the *NLL* and normalization integral calculations
- ❑ Check that the results are compatible: asymmetry less than 10^{-12}



- Speed-up increases with the dimension of the sample, taking benefit from the data streaming on GPU and the integral calculation only on the CPU
- ~3x for small samples, up to ~7x for large samples

- ✓ Implementation of the algorithm in CUDA to calculate the NLL on GPU, as part of the RooFit package
 - ❑ Require not so drastic changes in the existing RooFit code
 - ❑ New design of the algorithm for PDF-event parallelism
- ✓ The CUDA implementation “forces” us to develop an OpenMP implementation on the CPU of the same PDF-event algorithm
 - ❑ With 1 thread +34% better performance with respect to RooFit implementation
- ✓ In our test GPU implementation gives >3x speed-up (~7x for large samples) with respect to OpenMP with 4 threads
 - ❑ Note that our target is running fits at the user-level on the GPU of small systems (laptops), i.e. with small number of CPU cores
- ✓ This is a preliminary work (mainly by the summer student, Felice: 2.5 months work). Still a lot to do. Some examples:
 - ❑ Simultaneous fits with index variables
 - ❑ More complex tests
 - ❑ Parallelization of PDFs with numerical integrals
 - ❑ Further optimization on the GPU (better treatment of the memory)
- ✓ Last but not least: insert the code in the official RooFit/ROOT release

- ❑ Try to use OpenCL
 - ❑ The great benefit is the possibility to have hardware-independent code, i.e. GPUs (NVIDIA, AMD, Intel) and CPUs
 - ❑ There are issues related to the implementation that are to be investigated.
 - ❑ In contact with a guru, Tim Mattson (Intel)
- ❑ **It turns out that the new implementation of the algorithm (which is required to run on the GPU) gives better performance on the CPU and it is easy to parallelize (using OpenMP)**
 - ❑ **We will continue to improve this version. This is our first priority**
- ❑ We are working on the evaluation of the Knights Ferry (32 cores) and soon of the Single-Chip Cloud Computer (48 cores, no cache coherency), as part of the collaboration with Intel
 - ❑ Very promising architectures for massive parallelization with intensive calculations
 - ❑ It can be put in the general context of accelerators

Backup Slides

Physical structure: discrete behavior

- GTX 470 Fermi Card is a discrete device, made up of 14 stream multiprocessors
- As the device is being filled, the processing time does not follow a $O(N)$ growth
- As soon as the device is completely filled and the # of events is increased, the performance drops and we begin to watch a $O(N)$ behavior

