

Performing PDF fits on GPU: a preliminary search for the best solution

Emilio Villa

Supervisor: Dr. Stefano Carrazza

Co-supervisor: Dr. Juan Cruz-Martinez

University of Milan

3 December 2019



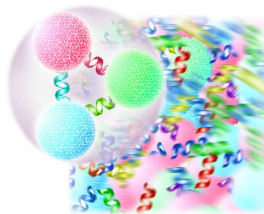
Structure of the presentation

- Parton Distribution Functions (PDFs)
- Neural-Network PDF (NNPDF) collaboration
- Speeding up the fits: APPLgrid, APFEL and APFELgrid
- Going beyond the Fast Kernel method: convolution on GPU
- Performance analysis
- Conclusions and outlook



Parton Distribution Functions

- **QCD** is the theory of strong interactions between *quarks* and *gluons* (*partons*), the elementary constituents of *hadrons*.
- The **factorization theorem** allows to determine the hadronic cross sections in terms of a convolution between the partonic cross sections and the PDFs.
- PDFs are the probability densities of finding a parton with a certain momentum fraction x inside the parent hadron at the energy scale Q^2 of the scattering process.



Parton Distribution Functions

The expression of a typical cross section $pp \rightarrow X$ reads:

$$\sigma_{pp \rightarrow X} = \sum_s \sum_k \int dx_1 dx_2 \hat{\sigma}^{(k)(s)} \alpha_s^{k+k_{LO}}(Q^2) F^{(s)}(x_1, x_2, Q^2)$$

$F^{(s)}$ represents the partonic density of the s -th subprocess:

$$F^{(s)}(x_1, x_2, Q^2) = \sum_{i,j} C_{ij}^{(s)} f_i(x_1, Q^2) f_j(x_2, Q^2)$$

f_i and f_j are the **PDFs**

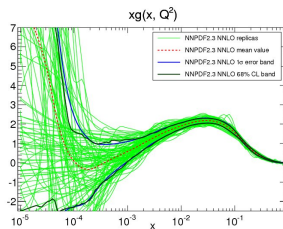
$C_{ij}^{(s)}$ counts the PDFs combinations that contribute to the s -th subprocess.



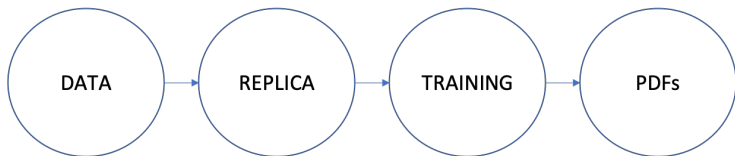
Parton Distribution Functions

PDFs are **non perturbative** objects and, at the moment, they cannot be determined from first principles. However, it is possible to determine them by **fit** to data.

- data \rightarrow hadronic cross sections;
- theory \rightarrow partonic cross sections;
- fit \rightarrow PDFs.



Over the years, the NNPDF collaboration has developed a fitting methodology based on **Monte Carlo methods** and the usage of **neural networks** as interpolating functions.



- huge number of convolutions \rightarrow **40 [h/PDF replica]**
- need to speed up the convolutions:
 - **reduction** of the number of operations required during the fit
 - increase the performances at **code level**.



APPLgrid and APFEL

Packages such as APPLgrid and APFEL, by means of interpolating techniques, allow to **lower** significantly the number of operations needed at the time of fitting.

APPLgrid

$$W_{\alpha\beta,\tau}^{(p)(s)} = \mathcal{I}_\tau (Q^2) \int dx_1 dx_2 \hat{\sigma}^{(p)(s)} \mathcal{I}_\alpha (x_1) \mathcal{I}_\beta (x_2) \quad (1)$$

$$\sigma_{pp \rightarrow X} = \sum_s \sum_p \sum_{\alpha\beta\tau} \alpha_s^{p+P_{LO}} (Q^2) W_{\alpha\beta,\tau}^{(p)(s)} F_{\alpha\beta,\tau}^{(s)} \quad (2)$$

APFEL

$$f_i (x_\alpha, Q_\tau^2) = \sum_k \sum_\beta A_{\alpha\beta,ik}^\tau f_k (x_\beta, Q_0^2) \quad (3)$$

$$F_{\alpha\beta,\tau}^{(s)} = \sum_{i,j} \sum_{k,l} \sum_{\delta,\gamma} C_{ij}^{(s)} [A_{\alpha\delta,ik}^\tau f_k (x_\delta, Q_0^2) A_{\beta\gamma,jl}^\tau f_l (x_\gamma, Q_0^2)] \quad (4)$$



APFELgrid simplifies further the convolution:

$$\sigma_{pp \rightarrow X} = \sum_{k,l} \sum_{\delta,\gamma} \text{FK}_{kl,\delta\gamma} f_k(x_\delta, Q_0^2) f_l(x_\gamma, Q_0^2)$$

$\text{FK}_{kl,\delta\gamma}$ is referred to as **Fast Kernel table**

Pros:

- number of operations lowered;
- PDFs at the initial scale;
- dot product \rightarrow multi threading.

Cons:

- theory embedded in FK.



At the end, the convolution is reduced to a simple **matrix** \times **matrix** product:

$$\sigma_i = \sum_{j=1}^N \text{FK}_{ij} \text{PDF}_{jn}$$

$$i \rightarrow \text{data}, \quad j \rightarrow k l \delta \gamma, \quad n \rightarrow \text{PDFs}$$

- Ndata \rightarrow data (σ at different Q^2);
- Npdf \rightarrow PDFs convolved with the FK;
- N $\rightarrow N_k N_l N_\delta N_\gamma$, i.e. product of active partons and x -grid points.

Next step is to **speed up** the evaluation of such product.



- **SSE3** (Streaming SIMD Extensions 3);
- **AVX** (Advanced Vector Extension);
- **Eigen**;
- **OpenBLAS**;
- **GSL** (GNU Scientific Library);
- **MKL** (Math Kernel Library);
- **OpenCL** (Open Computing Language);
- **TensorFlow**.

Scalar Operation

$$A_1 \times B_1 = C_1$$

$$A_2 \times B_2 = C_2$$

$$A_3 \times B_3 = C_3$$

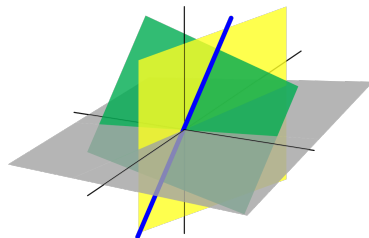
$$A_4 \times B_4 = C_4$$

SIMD Operation

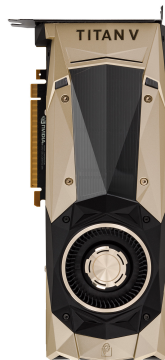
$$\begin{bmatrix} A_1 \\ A_2 \\ A_3 \\ A_4 \end{bmatrix} \times \begin{bmatrix} B_1 \\ B_2 \\ B_3 \\ B_4 \end{bmatrix} = \begin{bmatrix} C_1 \\ C_2 \\ C_3 \\ C_4 \end{bmatrix}$$



- **SSE3** (Streaming SIMD Extensions 3);
- **AVX** (Advanced Vector Extension);
- **Eigen**;
- **OpenBLAS**
- **GSL** (GNU Scientific Library);
- **MKL** (Math Kernel Library);
- **OpenCL** (Open Computing Language);
- **TensorFlow**.



- **SSE3** (Streaming SIMD Extensions 3);
- **AVX** (Advanced Vector Extension);
- **Eigen**;
- **OpenBLAS** (Open Basic Linear Algebra Subprograms);
- **GSL** (GNU Scientific Library);
- **MKL** (Math Kernel Library);
- **OpenCL** (Open Computing Language);
- **TensorFlow**.



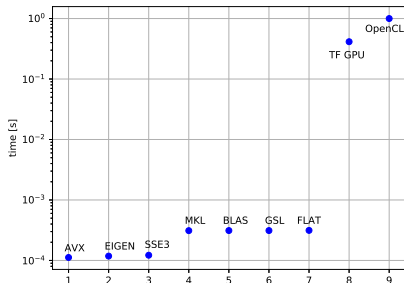
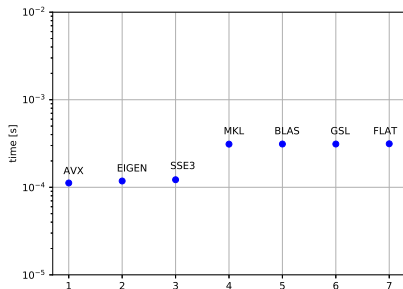
CUDA and OpenCL are very similar to each other, however:

- OpenCL can be executed over a **variety** of platforms, including CPUs, GPUs and other type of processors.
- CUDA executes **only** on NVIDIA GPUs.



First time benchmark

Convolution between an FK 8×35721 (atlas-Z0-rapidity.root) and a single PDF
(Ndata = 8 , Npdf = 1 , N = 35721)



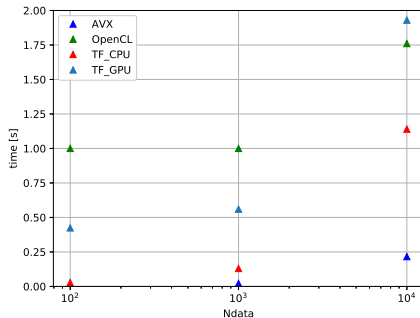
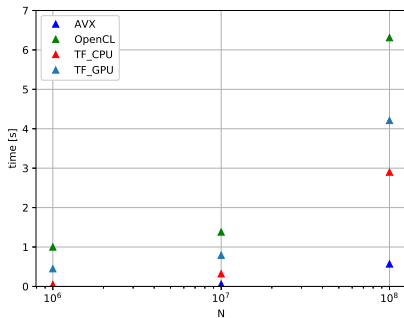
First time benchmark

- OpenCL, TF GPU → **slowest** methods;
- need to develop a **different approach**;
- convolution evaluated for **arbitrary** dimensions of FK and PDF matrices;
- test AVX and TensorFlow on **CPU**, OpenCL and TensorFlow on **GPU**.

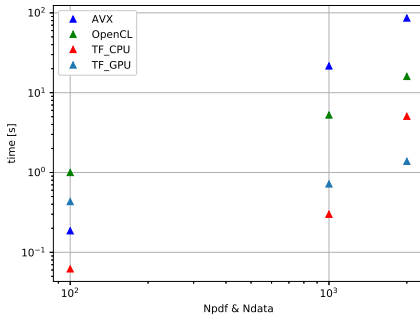
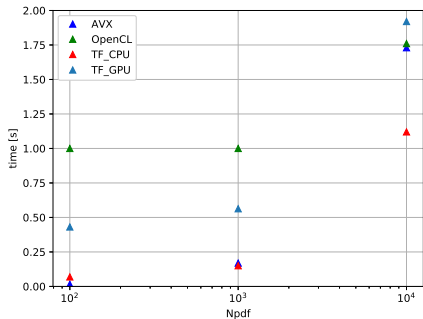
$$\begin{bmatrix} \text{FK}_{11} & \cdots & \text{FK}_{1,N} \\ \text{FK}_{21} & \cdots & \text{FK}_{2,N} \\ \vdots & \ddots & \vdots \\ \text{FK}_{Ndata,1} & \cdots & \text{FK}_{Ndata,N} \end{bmatrix} \times \begin{bmatrix} \text{PDF}_{11} & \cdots & \text{PDF}_{1,Npdf} \\ \text{PDF}_{21} & \cdots & \text{PDF}_{2,Npdf} \\ \vdots & \ddots & \vdots \\ \text{PDF}_{N,1} & \cdots & \text{PDF}_{N,Npdf} \end{bmatrix}$$



Results obtained varying respectively N and Ndata



Results obtained varying respectively Npdf and Npdf & Ndata



Summary

Size of the PDF	Size of the FK Table	TensorFlow CPU [s]	AVX [s]	TensorFlow GPU [s]	OpenCL [s]
35721×1	8×35721	$1.10 \cdot 10^{-2}$	$1.57 \cdot 10^{-4}$	$4.14 \cdot 10^{-1}$	~ 1
$10^6 \times 1$	8×10^6	$4.70 \cdot 10^{-2}$	$5.00 \cdot 10^{-3}$	$4.49 \cdot 10^{-1}$	~ 1
$10^7 \times 1$	8×10^7	$3.20 \cdot 10^{-1}$	$5.70 \cdot 10^{-2}$	$7.90 \cdot 10^{-1}$	1.38
$10^8 \times 1$	8×10^8	2.90	$5.70 \cdot 10^{-1}$	4.21	6.31
35721×10^2	8×35721	$6.90 \cdot 10^{-2}$	$1.60 \cdot 10^{-2}$	$4.31 \cdot 10^{-1}$	~ 1
35721×10^3	8×35721	$1.50 \cdot 10^{-1}$	$1.69 \cdot 10^{-1}$	$5.63 \cdot 10^{-1}$	~ 1
35721×10^4	8×35721	1.12	1.73	1.92	1.76
35721×1	$10^2 \times 35721$	$2.80 \cdot 10^{-2}$	$2.43 \cdot 10^{-3}$	$4.24 \cdot 10^{-1}$	~ 1
35721×1	$10^3 \times 35721$	$1.30 \cdot 10^{-1}$	$2.14 \cdot 10^{-2}$	$5.60 \cdot 10^{-1}$	~ 1
35721×1	$10^4 \times 35721$	1.14	$2.16 \cdot 10^{-1}$	1.93	1.76
35721×10^2	$10^2 \times 35721$	$6.20 \cdot 10^{-2}$	$1.86 \cdot 10^{-1}$	$4.32 \cdot 10^{-1}$	~ 1
35721×10^3	$10^3 \times 35721$	$3.00 \cdot 10^{-1}$	21.61	$7.19 \cdot 10^{-1}$	5.25
$35721 \times 2 \cdot 10^3$	$2 \cdot 10^3 \times 35721$	5.06	86.13	1.38	15.97



Conclusions and outlook

Conclusions¹:

- PDF fits can benefit from **hardware accelerators** (i.e. GPUs and OpenCL compatible devices) thanks to the possibility of offloading the most time-consuming tasks to the accelerator;
- However, in order to achieve performance improvements some precautions are required by defining the sizes of FK tables and the number of PDFs that should be convoluted simultaneously.

Outlook:

- Solutions proposed in this work should be tested in a real PDF fit. This will be possible thanks to the future extension of the `n3fit` framework to support GPU hardware.

¹S. Carrazza, J. Cruz-Martinez, J. Elizari, E. Villa, *Towards hardware acceleration for parton densities estimation*, Proceeding of PHOTON 2019, preprint arXiv:1909.10547



Intel(R) Core(TM) i9-9980XE CPU 3.00GHz:

- # of cores: 18
- processor base frequency: 3.00 GHz

GPU Nvidia Titan V:

- # of streaming multiprocessors: 80
- # of CUDA cores: 5120
- base clock: 1.2 GHz
- total memory: 12288 MB
- memory clock: 850 MHz
- total memory bandwidth: 652.8 GB/s

