

TOWARDS HARDWARE ACCELERATION FOR PARTON DENSITIES ESTIMATION

Stefano Carrazza
Juan Cruz-Martinez
Jesús Urtasun-Elizari
Emilio Villa

*TIF Lab, Dipartimento di Fisica, Università degli Studi di Milano and INFN Sezione di Milano,
Via Celoria 16, 20133, Milano, Italy*

Abstract

In this proceedings we describe the computational challenges associated to the determination of parton distribution functions (PDFs). We compare the performance of the convolution of the parton distributions with matrix elements using different hardware instructions. We quantify and identify the most promising data-model configurations to increase PDF fitting performance in adapting the current code frameworks to hardware accelerators such as graphics processing units.

1 Introduction

The determination of parton distribution functions (PDFs) is a particular topic which strongly relies on three dynamic and time dependent factors: new experimental data, higher order theoretical predictions and fitting methodology. In this environment, there are two main tasks for PDF fitters such as the NNPDF collaboration [1–3], MMHT [4] or CTEQ [5]. The first task consists in maintaining and organizing a workflow which incrementally implements new features proposed by the respective experimental and theoretical communities. The second task of a PDF fitter corresponds to investigate new numerical and efficient approaches to PDF fitting methodology. While the former relies almost exclusively on external groups and communities, the later is under full control of the PDF fitting collaborations, and in most cases it reflects the differences between them.

As a real example of the previous description we show in figure 1 a non-exhaustive timeline for PDF determinations with QED corrections. Since 2004 we observe at least three different fitting approaches to the determination of the photon PDF, starting from model based approach where the photon PDF is modelled by an ad-hoc distribution [6, 7], then to a data driven approach where the photon PDF is extracted directly from data [8–10], and finally to a more precise procedure involving theory calculations [11–15]. Improvements in terms of PDF quality and physics content are however accompanied by

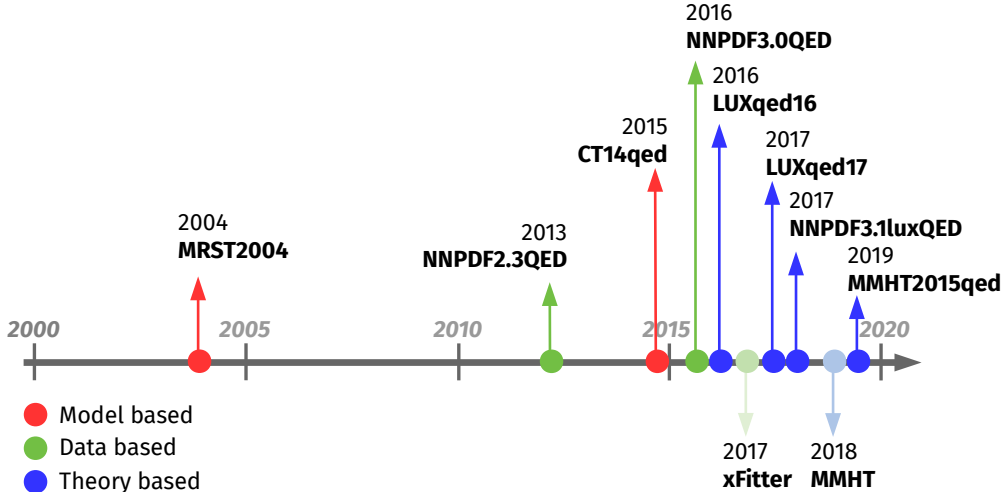


Figure 1: Illustrative timeline of PDF releases with QED corrections since 2004. Arrows pointing up refer to publications providing PDF sets as deliverables.

a negative performance trend associated to longer fitting times and increased computational cost due to larger datasets and increasingly complex procedures. This in turn makes more difficult the task of developing and testing novel fitting procedures.

In this proceedings we describe a new approach to deal with this growing trend in computational resources for PDF determination. We focus our discussion on the new methodology recently presented by the N3PDF team within the NNPDF collaboration and summarized in the next paragraphs.

2 A deep learning approach to PDFs

In Ref. [16] we presented a new approach to PDF fits based on deep learning techniques in the context of the NNPDF methodology. We implement a new efficient computing framework based on graph generated models for PDF parametrization and gradient descent optimization called `n3fit`. The best model configuration is derived from a robust cross-validation mechanism through a hyperparametrization tune procedure. From a practical point of view the `n3fit` code uses Keras [17] and TensorFlow (TF) [18] as backends.

From a technical perspective, one of the most relevant achievements of `n3fit` is the reduction of computational time required to obtain PDF sets based on the NNPDF3.1 dataset [19]. In figure 2 we show the running time (in hours) required to fit 100 PDF replicas using the new `n3fit` fitting code and we compare it to the latest NNPDF3.1 algorithm. On the left plot, we show a fit to DIS-only data, while in the right plot we have a global fit. In both cases we observe an improvement of between one to two orders of magnitude, *e.g.* the new `n3fit` code takes in average one hour to complete a global fit whereas the old code could take more than 40 hours. These improvements are partially due to the new minimizer (based on gradient descent instead of on genetic algorithms) in combination with multi-threading CPU calculations when executing the TensorFlow graph model.

The great performance improvement observed with `n3fit` suggests that we may find new code strategies which take advantage of hardware acceleration. At this point, one may ask if it is possible

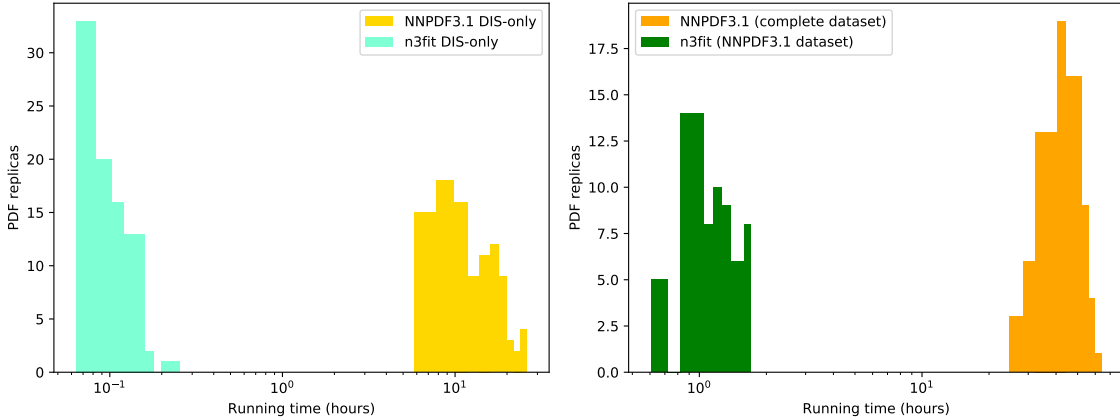


Figure 2: Fitting time distribution per replica PDF for NNPDF and `n3fit` codes for DIS-only (left plot) and global datasets (right plot).

improve the performance from the default TensorFlow graph optimization for CPU and eventually use hardware accelerators such as graphic processing units (GPUs), field-programmable gate arrays (FPGAs), or tensor processing units (TPUs).

3 Hardware accelerating PDF determination

The first step towards faster fits consists in profiling the code and isolate the most time consuming operations during PDF fits. Fortunately, the answer to this question is simple and involves the computation of physical observables through PDF and matrix elements convolutions,

$$\sigma^N = \sum_{i,j,\alpha,\beta} f_\alpha(x_i) f_\beta(x_j) \text{FK}_{ij\alpha\beta}^N, \quad (1)$$

where $f_\alpha(x_i)$ stands for the PDF of a particular flavor species α evaluated in the point x_i of the grid in x . $\text{FK}_{ij\alpha\beta}^N$ is a Fast Kernel (FK) table which contains the information about the partonic cross section, following the description presented in [16]. In the case of hadronic observables the evaluation of predictions produces a vector of N observables, σ^N , by building a neural network that generates the PDFs sampling from a grid in x , representing the fractions of momenta that a particular parton could carry, and then convoluting the result with an FK table containing the partonic information.

Given that TensorFlow relies on symbolic computation and graph generation to represent a model, as a first step we investigate if the memory usage it requires is higher than the one needed by a custom code specialized in convolutions. We wrote a custom operator in C++ for TensorFlow that performs the convolution and its corresponding gradient, directly without graph evaluations. In table 1 we cross-check the implementation by looking at examples of convolution and gradient computation for DIS and hadronic observables. The ratio between both implementations confirms excellent numerical agreement. The memory usage for the default TensorFlow implementation and the custom convolution code are shown in table 2. We observe a reduction of 3.2 GB and 5.9 GB of resident memory when using our custom operator, when loading all NNPDF3.1 hadronic and global data respectively. The reduction of memory usage is a great benefit because it gives the possibility to run PDF fits in consumer level hardware and, most importantly, load all data in the limited memory space available in hardware accelerators. Once

		TensorFlow [pb]	Custom [pb]	Ratio
DIS	Convolution	1.9207904	1.9207904	1.0000000
		2.4611666	2.4611664	0.9999999
		1.3516952	1.3516952	1.0000000
	Gradient	1.8794115	1.8794115	1.0000000
		1.505316	1.505316	1.0000000
		2.866085	2.866085	1.0000000
Hadronic	Convolution	8.142365	8.142366	1.0000001
		8.947762	8.947762	1.0000000
		7.4513326	7.4513316	0.9999999
	Gradient	18.525095	18.525095	1.0000000
		19.182995	19.182993	0.9999999
		19.551006	19.551004	0.9999999

Table 1: Example of convolution and gradient computations for both DIS and hadronic observables. The ratio between TensorFlow and the custom computation confirms excellent numerical agreement.

		TensorFlow	Custom Convolution	Difference (TensorFlow - custom)
Hadronic	Virtual	17.7 GB	13.8 GB	3.9 GB
	RES	12.1 GB	8.39 GB	3.2 GB
Global	Virtual	23.5 GB	19.7 GB	3.8 GB
	RES	18.4 GB	12.5 GB	5.9 GB

Table 2: Memory usage after model generation and fit for both hadronic and global (DIS + hadronic).

the memory saving is obtained, the performance can also be improved by multi-threading our custom operator on the CPU.

After the memory usage analysis, we carried out a time performance comparison running the convolution both on CPU and GPU. Shifting the computation from CPU to GPU one can take advantage of the parallelization over the increased number of cores. In table 3 we show the overall running time for several examples of toy PDF and FK tables based on hadronic observables. The numbers include the computation time as well as the time required for the memory transfer to the GPU.

The performance of Advanced Vector Extensions (AVX) on CPU, OpenCL [20] on GPU and TF both on CPU and GPU are compared. As it is shown in table 3, AVX and TF running on CPU are faster up to a certain number of columns (for the FK table) or rows (for the PDFs matrices). Once the size of the operation is big enough, AVX is over one order of magnitude slower than OpenCL and even two orders of magnitude slower than TF on GPU. TF on CPU is more resilient than our AVX implementation and it is competitive with the GPU convolutions for a larger range of dimensions.

According to these results, for given dimensions of the FK table and the PDF matrices, it is convenient to carry out the PDF fits on GPU devices, parallelizing, for instance, over the different PDF replicas, allowing the fit of all of them to run simultaneously. In other words, hardware accelerators become competitive tools for PDF fitting once the penalty introduced by the memory transfer between devices is overcome.

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
$35721 \times 5 \cdot 10^4$	10×35721	5.33	8.93	7.83	5.80
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

Table 3: Time performances achieved with AVX, TensorFlow (both on CPU and GPU) and OpenCL for the given sizes of the FK table and the PDF matrix. In green is highlighted the lowest value within each row. Time is given in seconds. The base FK table used in this comparison consists on 49 flavour combinations and an grid in x of size 27 (35721 elements in total).

4 Outlook and future developments

The results presented in this proceedings strongly suggest that PDF fits can benefit from hardware accelerators such as GPUs, and in future, FPGAs or TPUs thanks to the possibility of offloading the most time-consuming tasks to the accelerator. However, we should notice that in order to achieve performance improvements some precautions are required by defining the sizes of FK tables and the number of PDFs we would like to convolute simultaneously. In future work we are planning to extend the `n3fit` framework to support GPU hardware.

Acknowledgements

S.C. acknowledges the NVIDIA Corporation for the donation of a Titan V GPU used for this research. S.C., J.C.M. and J.U. are supported by the European Research Council under the European Unions Horizon 2020 research and innovation Programme (grant agreement number 740006). S.C. is supported by the UNIMI Linea 2A grant “New hardware for HEP”.

References

1. R. Abdul Khalek *et al.* [NNPDF Collaboration], arXiv:1906.10698 [hep-ph].
2. R. Abdul Khalek *et al.* [NNPDF Collaboration], arXiv:1905.04311 [hep-ph].
3. R. D. Ball *et al.* [NNPDF Collaboration], Eur. Phys. J. C **78** (2018) no.5, 408 doi:10.1140/epjc/s10052-018-5897-7 [arXiv:1802.03398 [hep-ph]].

4. L. A. Harland-Lang, A. D. Martin, P. Motylinski and R. S. Thorne, *Eur. Phys. J. C* **75** (2015) no.5, 204 doi:10.1140/epjc/s10052-015-3397-6 [arXiv:1412.3989 [hep-ph]].
5. T. J. Hou *et al.*, arXiv:1908.11238 [hep-ph].
6. A. D. Martin, R. G. Roberts, W. J. Stirling and R. S. Thorne, *Eur. Phys. J. C* **39** (2005) 155 doi:10.1140/epjc/s2004-02088-7 [hep-ph/0411040].
7. C. Schmidt, J. Pumplin, D. Stump and C. P. Yuan, *Phys. Rev. D* **93** (2016) no.11, 114015 doi:10.1103/PhysRevD.93.114015 [arXiv:1509.02905 [hep-ph]].
8. R. D. Ball *et al.* [NNPDF Collaboration], *Nucl. Phys. B* **877** (2013) 290 doi:10.1016/j.nuclphysb.2013.10.010 [arXiv:1308.0598 [hep-ph]].
9. V. Bertone and S. Carrazza, *PoS DIS 2016* (2016) 031 doi:10.22323/1.265.0031 [arXiv:1606.07130 [hep-ph]].
10. F. Giuli *et al.* [xFitter Developers' Team], *Eur. Phys. J. C* **77** (2017) no.6, 400 doi:10.1140/epjc/s10052-017-4931-5 [arXiv:1701.08553 [hep-ph]].
11. A. Manohar, P. Nason, G. P. Salam and G. Zanderighi, *Phys. Rev. Lett.* **117** (2016) no.24, 242002 doi:10.1103/PhysRevLett.117.242002 [arXiv:1607.04266 [hep-ph]].
12. A. V. Manohar, P. Nason, G. P. Salam and G. Zanderighi, *JHEP* **1712** (2017) 046 doi:10.1007/JHEP12(2017)046 [arXiv:1708.01256 [hep-ph]].
13. V. Bertone *et al.* [NNPDF Collaboration], *SciPost Phys.* **5** (2018) no.1, 008 doi:10.21468/SciPostPhys.5.1.008 [arXiv:1712.07053 [hep-ph]].
14. R. Nathvani, R. Thorne, L. Harland-Lang and A. Martin, *PoS DIS 2018* (2018) 029 doi:10.22323/1.316.0029 [arXiv:1807.07846 [hep-ph]].
15. L. A. Harland-Lang, A. D. Martin, R. Nathvani and R. S. Thorne, arXiv:1907.02750 [hep-ph].
16. S. Carrazza and J. Cruz-Martinez, *Eur. Phys. J. C* **79** (2019) no.8, 676 doi:10.1140/epjc/s10052-019-7197-2 [arXiv:1907.05075 [hep-ph]].
17. F. Chollet *et al.*, *Keras* (2015) <https://keras.io>
18. M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G.S. Corrado, A. Davis, J. Dean, M. Devin *et al.*, *TensorFlow: Large-scale machine learning on heterogeneous systems* (2015), software available from tensorflow.org, <http://tensorflow.org/>
19. R. D. Ball *et al.* [NNPDF Collaboration], *Eur. Phys. J. C* **77** (2017) no.10, 663 doi:10.1140/epjc/s10052-017-5199-5 [arXiv:1706.00428 [hep-ph]].
20. J.E. Stone, D. Gohara and G. Shi, *Computing in Science Engineering* (2010) 12 no.3, 66 doi:10.1109/MCSE.2010.69