

Accelerating Full and Fast Simulation of the CMS Experiment

Natascha Krammer^{a,*} on behalf of the CMS Collaboration

^a*Institute of High Energy Physics, Austrian Academy of Science,
Vienna, Austria*

E-mail: natascha.krammer@oeaw.ac.at

Monte Carlo Simulation data for the CMS experiment can be produced using two software tools. The first, the Full Simulation (FullSim), is a more precise tool based on Geant4 detector simulation. The second, the Fast Simulation, provides a faster but still reliable tool and is based on parametric particle-material interactions. FullSim for the LHC Run-3 shows significant computing performance improvements compared to LHC Run-2. The major modifications of FullSim are the migration from Geant4 version 10.4.3 to 10.7.2, the new software package DD4hep for geometry description, the migration from CentOS7 to the new software platform AlmaLinux8 and using the LTO (Link time optimization) build flag. The challenging CMS detector upgrade plan for HL (High Luminosity)-LHC requires extra efforts due to the increased luminosity and the new and complex detector geometry. FullSim plans to meet the requirements for HL-LHC, which includes continues migration to newer versions of Geant4, the current one is 11.1.2, as well as physics improvements including machine learning (ML) techniques to reduce compute capacity needs. Major progresses of Fast Simulation are reached by a more efficient treatment of the generator particles as they propagate through the detectors. Recent developments include the implementation of an increasing more accurate shower generation, improved track finding and tuning of physics processes. This contribution reports the current Full and Fast Simulation performance innovations and further plans to fulfill the significant higher Monte Carlo Simulation demands in LHC Run-3 and for HL-LHC. Very promising software developments using ML for higher accuracy (Refinement Fast Simulation) and speed up (FlashSim) for the simulation run will be explained.

*The Eleventh Annual Conference on Large Hadron Collider Physics (LHCP2023)
22-26 May 2023
Belgrade, Serbia*

*Speaker

1. Introduction

Future LHC runs will reach new frontiers in energy and luminosity, leading to a greatly increased rate of data taking and pile up. To prepare for the upcoming experimental phase, the High-luminosity LHC, it is crucial to accelerate the simulation step. Due to the major upgrade of the CMS experiment [1], in particular the High Granularity Calorimeter (HGCal) [2], FullSim expects to be 2-3 times slower because of more complex geometry and more precise physics.

2. Accelerating Full and Fast Simulation

For Full and Fast Simulation, a wide range of improvements have been developed and introduced. New features in FullSim were implemented in Run-3 and Phase-2 [3] [4]: (i) the migration to DD4hep geometry description, (ii) the use of the LTO (Like Time Optimization) build method, (iii) the use of faster computations and less instructions: Geant4 Gamma General Process and VDT (VectorisD maTh) for fast and auto-vectorisable mathematical functions, (iv) the Geant4 Transportation with MSC (multiple scattering), (v) custom tracking managers to simplify the e-gamma transport in Geant4, (vi) G4HepEm external library, which focuses on the EM shower generation for GPU usage. In addition the operating systems was upgraded from CentOS7 to AlmaLinux8. The Geant4 version was migrated from 10.4.3 to 10.7.2 (CMSSW 11_3_X) and to 11.1.2 (CMSSW 13_1_X). Fast Simulation (FastSim) software and framework optimization achieves a more efficient handling of the generator particles through the detectors. In addition, there is an ongoing effort of R&D of GPU usage for simulation such as Accelerated demonstrator of electromagnetic Particle Transport (AdePT) [5] and the Celeritas [6] project targeting the computationally intensive HL-LHC runs. The success of the improvements are shown in Figure 1 for the average CPU time per event for the standard model (SM) process $t\bar{t}$ and the beyond SM process T1tttt ($pp \rightarrow \text{gluino} + \text{gluino}$, $\text{gluino} \rightarrow t\bar{t} + \text{lightest neutralino}$) for single threaded jobs [7].

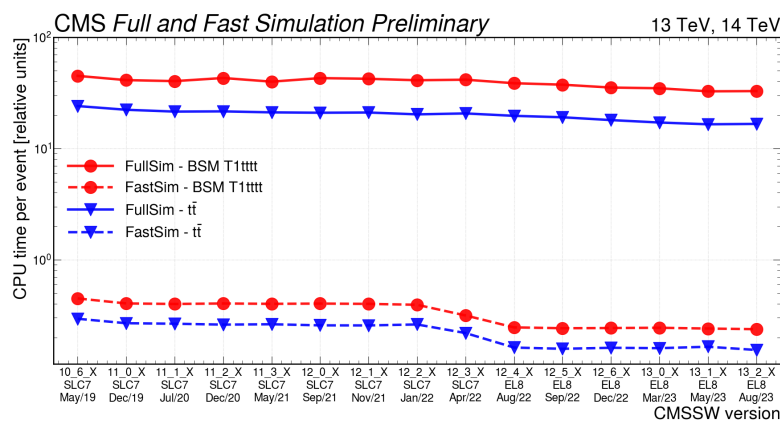


Figure 1: FullSim and FastSim CPU time evolution for the past 4 years. FastSim compared to FullSim is for BSM T1tttt 136 times and for $t\bar{t}$ 110 times faster using the latest CMSSW version.

3. Refining Fast Simulation using ML techniques

FastSim compared to FullSim has a major advantage in speed with the price of decreased accuracy in some of the final observables. The refining method using ML techniques [8] was developed to increase the accuracy and render FastSim suitable for a larger number of analyses throughout the collaboration. The refining method is to use the analysis observables simulated by the FastSim chains and compare them to the corresponding FullSim output. A fully-connected feed-forward neural network (NN) is trained to establish a refined version of the FastSim data sample, which is more accurate with respect to the FullSim sample. The current refining FastSim method focus on jet flavor tagging for four DeepJet discriminators (B(b+bb+lep), CvB(c/(c+b+bb+lep)), CvL(c/(c+uds+g)), QG(g/(g+uds))) in the CMS NanoAOD data analysis format. The DeepJet algorithm [9] is a multi-class NN trained to distinguish jets originating from b, c, light quarks and gluons. The network has six output nodes (b, bb, lep, c, uds, g) activated with a softmax function.

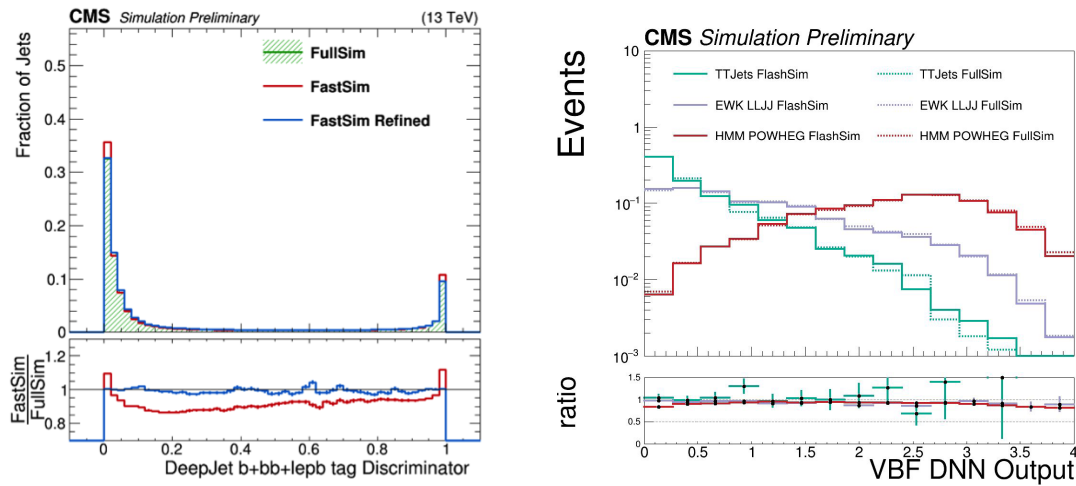
Two training samples are prepared by first generating a single sample of gluino pair-production, and then processing each with the FullSim and FastSim workflows. Jets are then identified which are matched in δR to associate the generator-level jet, the reconstructed jet processed with FullSim, and the reconstructed jet processed with FastSim, leading to a jet triplet. The application of the ResNet-like (Deep Residual Learning) architecture [10] results in a good approximation of FastSim to FullSim output and needs only to apply a residual correction. The pre-processing transforms the input variables/parameters and the post-processing transforms back. The refined DeepJet discriminators (B, CvB, CvL, QG) have to be constructed such that the sum of the original DeepJet output nodes (b, bb, lep, c, uds, g) is equal to unity. The network is implemented using the PyTorch package [11]. Two loss terms are used in the training: a primary loss MMD (distributed-based) [12] compares ensembles of jets not jet-jet pairs to cope with independent stochastic in both simulations chains and an additional loss term MSE/Huber (output-target pair-based) for the correction for deterministic FastSim biases. The two loss terms are combined via MDMM algorithm [14]. ResNet-like regression NN can be used as post-hoc refinement layer to FastSim output, which is now implemented in the official FastSim software and results in a considerably improved agreement with FullSim output, illustrated in Figure 2a.

4. FlashSim - a ML simulation framework

A new simulation framework, named *FlashSim* [15] uses the advantage of ML for a faster and accurate simulation. It is based on new Normalizing Flows (NF) [16] using generic ML generative techniques, which directly produce CMS NanoAOD format samples from generator level information. The advantage of the NanoAOD format for analysis is to reduce the number of variables to simulate from several thousands down to few hundreds. The difference to typical AI/ML sample generation, e.g. image generation, is that we need to condition the generation on some previous information. Not a generic CMS event should be simulated, but the CMS event corresponding to some given generator level input. For the simulation the natural factorization is to go one by one on the various objects and to use generator level representation of those objects as conditioning information. The simulation of each object is a so-called functional unit, which is a transformation implemented by the NF algorithm. As input only the relevant physical information

for the simulation of its target is taken into account and the various units are independent at first order. It may still be necessary to include additional correlation runs in a chain to also access not only generator level information but also reconstruction information of previous units. The advantage of this general, flexible simulator is to be not tailored to a specific analysis.

In a real-world scenario analysis test the feasibility of the model is performed with generated datasets of $t\bar{t}$, Drell-Yan, EWK LLJJ and signal ($H \rightarrow \mu^+\mu^-$) processes. Higgs to di-muon VBF (vector boson fusion) analysis [17] uses events with muons and jets simulated by FlashSim to verify the usability of the approach at the analysis level. The flash-simulated quantities are calculated using Deep Neural Network (DNN) to separate signal from background and are compared with FullSim results. Multiple derived input variables are used by DNN, including some that are correlating the di-jets part of the VBF event with the di-muon part of the Higgs decay. The consistency of FlashSim and FullSim results are shown in Figure 2b.



(a) The distribution of the DeepJet discriminator $B(b+bb+lepb)$ for FullSim, FastSim and the refined version of FastSim.

(b) Flash and Full Simulation comparison of analysis DNN output using muons and jets information of $t\bar{t}$ and EWK LLJJ samples and signal events.

Figure 2: Test sample results of Refining FastSim (a) using ResNet-like (Deep Residual Learning) architecture and FlashSim (b) using Deep Neural Network techniques.

Conclusions As result of these efforts CPU time performance for $t\bar{t}$ /BSM T1tttt processes of Run-3 versus Run-2 was improved by 32/27% for FullSim and 48/47% for FastSim. The ML-based refinement for FastSim has shown a clear improvement in agreement with the FullSim output. FlashSim is on the way of building a complete ML simulation framework and considerably reduce the time taken by simulation tasks.

References

- [1] CMS Collaboration, The CMS experiment at the CERN LHC, JINST 3 S08004 (2008)
- [2] CMS Collaboration, The Phase-2 Upgrade of the CMS Endcap Calorimeter, CERN-LHCC-2017-023, CMS-TDR-019 (2017)

- [3] V. Ivanchenko, S. Banerjee, G. Hugo, S. L. Meo, I. Osborne, K. Pedro, D. Piparo, D. Sorokin, N. Srimanobhas, C. Vuosalo, CMS Full Simulation for Run 3, EPJ Web Conf. 251 03016 (2021)
- [4] N. Srimanobhas, S. Banerjee, J. Hahnfeld, V. Ivantchenko, N. Krammer, S. Muzaffar, K. Pedro, D. Piparo, Full Simulation of CMS for Run-3 and Phase-2, CERN-CMS-CR-2023/152 (2023)
- [5] G. Amadio, J. Apostolakis, P. Buncic et al., Offloading electromagnetic shower transport to GPUs, J. Phys.: Conf. Ser. 2438 012055 (2023)
- [6] S. C. Tognini, P. Canal, T. M. Evans et al., Celeritas: GPU-accelerated particle transport for detector simulation in High Energy Physics experiments, FERMILAB-FN-1159-SCD, arXiv 2203.09467 (2022)
- [7] CMS Collaboration, CPU performance evolution of Full and Fast simulations from Run-2 to Run-3/CMSSW_13_2_0, CERN-CMS-DP-2023-063 (2023)
- [8] S. Bein, P. Connor, K. Pedro, P. Schleper, M. Wolf, Refining fast simulation using machine learning, CERN-CMS-CR-2023/128 (2023)
- [9] E. Bols, J. Kieseler, M. Verzetti, M. Stoye, A. Stakia, Jet flavour classification using DeepJet, JINST 15 P12012, arXiv:2008.10519 (2020)
- [10] K. He, X. Zhang, S. Ren, J. Sun, Deep Residual Learning for Image Recognition, 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR) pp. 770-778 (2016)
- [11] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga et al., PyTorch: An Imperative Style, High-Performance Deep Learning Library, Part of Advances in Neural Information Processing Systems (NeurIPS) 32 pp. 8024–8035 (2019)
- [12] A. Gretton, K.M. Borgwardt, M.J. Rasch, B. Schölkopf, A. Smola, A Kernel Two-Sample Test, Journal of Machine Learning Research 13 723-773 (2012)
- [13] J. Platt, A. Barr, Constrained Differential Optimization, Neural Information Processing Systems Vol.0 612-621 (1987)
- [14] K. Crowson, mdmm, Software version 0.1.3, <https://github.com/crowsonkb/mdmm> (2021)
- [15] F. Vaselli, FlashSim: accelerating HEP simulation with an end-to-end Machine Learning framework, CERN-CMS-CR-2023/090 (2023)
- [16] G. Papamakarios, E. Nalisnick, D.J. Rezende, S. Mohamed, B. Lakshminarayanan, Normalizing flows for probabilistic modeling and inference, Journal of Machine Learning Research 22(57) 1-64, arXiv:1912.02762 (2021)
- [17] CMS Collaboration, Evidence for Higgs boson decay to a pair of muons, JHEP 01 148, arXiv:2009.04363 (2021)