

# Machine Learning in High Energy Physics Community White Paper

May 17, 2019

**Abstract:** Machine learning has been applied to several problems in particle physics research, beginning with applications to high-level physics analysis in the 1990s and 2000s, followed by an explosion of applications in particle and event identification and reconstruction in the 2010s. In this document we discuss promising future research and development areas for machine learning in particle physics. We detail a roadmap for their implementation, software and hardware resource requirements, collaborative initiatives with the data science community, academia and industry, and training the particle physics community in data science. The main objective of the document is to connect and motivate these areas of research and development with the physics drivers of the High-Luminosity Large Hadron Collider and future neutrino experiments and identify the resource needs for their implementation. Additionally we identify areas where collaboration with external communities will be of great benefit.

**Editors:** Sergei Gleyzer<sup>30</sup>, Paul Seyfert<sup>13</sup>, Steven Schramm<sup>32</sup>

**Contributors:** Kim Albertsson<sup>1</sup>, Piero Altoc<sup>2</sup>, Dustin Anderson<sup>3</sup>, John Anderson<sup>4</sup>, Michael Andrews<sup>5</sup>, Juan Pedro Araque Espinosa<sup>6</sup>, Adam Aurisano<sup>7</sup>, Laurent Basara<sup>8</sup>, Adrian Bevan<sup>9</sup>, Wahid Bhimji<sup>10</sup>, Daniele Bonacorsi<sup>11</sup>, Bjorn Burkle<sup>12</sup>, Paolo Calafiura<sup>10</sup>, Mario Campanelli<sup>9</sup>, Louis Capps<sup>2</sup>, Federico Carminati<sup>13</sup>, Stefano Carrazza<sup>13</sup>, Yi-Fan Chen<sup>4</sup>, Taylor Childers<sup>14</sup>, Yann Coadou<sup>15</sup>, Elias Coniavitis<sup>16</sup>, Kyle Cranmer<sup>17</sup>, Claire David<sup>18</sup>, Douglas Davis<sup>19</sup>, Andrea De Simone<sup>20</sup>, Javier Duarte<sup>21</sup>, Martin Erdmann<sup>22</sup>, Jonas Eschle<sup>23</sup>, Amir Farbin<sup>24</sup>, Matthew Feickert<sup>25</sup>, Nuno Filipe Castro<sup>6</sup>, Conor Fitzpatrick<sup>26</sup>, Michele Floris<sup>13</sup>, Alessandra Forti<sup>27</sup>, Jordi Garra-Tico<sup>28</sup>, Jochen Gemmler<sup>29</sup>, Maria Girone<sup>13</sup>, Paul Glayshe<sup>18</sup>, Sergei Gleyzer<sup>30</sup>, Vladimir Vava Gligorov<sup>31</sup>, Tobias Golling<sup>32</sup>, Jonas Graw<sup>2</sup>, Lindsey Gray<sup>21</sup>, Dick Greenwood<sup>33</sup>, Thomas Hacker<sup>34</sup>, John Harvey<sup>13</sup>, Benedikt Hegner<sup>13</sup>, Lukas Heinrich<sup>17</sup>, Ulrich Heintz<sup>12</sup>, Ben Hooberman<sup>35</sup>, Johannes Junggeburth<sup>36</sup>, Michael Kagan<sup>37</sup>, Meghan Kane<sup>38</sup>, Konstantin Kanishchev<sup>8</sup>, Przemysław Karpiński<sup>13</sup>, Zahari Kassabov<sup>39</sup>, Gaurav Kaul<sup>40</sup>, Dorian Kcira<sup>3</sup>, Thomas Keck<sup>29</sup>, Alexei Klimentov<sup>41</sup>, Jim Kowalkowski<sup>21</sup>, Luke Kreczko<sup>42</sup>, Alexander Kurepin<sup>43</sup>, Rob Kutschke<sup>21</sup>, Valentin Kuznetsov<sup>44</sup>, Nicolas Köhler<sup>36</sup>, Igor Lakomov<sup>13</sup>, Kevin Lannon<sup>45</sup>, Mario Lassnig<sup>13</sup>, Antonio Limosani<sup>46</sup>, Gilles Louppe<sup>17</sup>, Aashrita Mangu<sup>47</sup>, Pere Mato<sup>13</sup>, Helge Meinhard<sup>13</sup>, Dario Menasce<sup>48</sup>, Lorenzo Moneta<sup>13</sup>, Seth Moortgat<sup>49</sup>, Meenakshi Narain<sup>12</sup>, Mark Neubauer<sup>35</sup>, Harvey Newman<sup>3</sup>, Sydney Otten<sup>50</sup>, Hans Pabst<sup>40</sup>, Michela Paganini<sup>51</sup>, Manfred Paulini<sup>5</sup>, Gabriel Perdue<sup>21</sup>, Uzziel Perez<sup>52</sup>, Attilio Picazio<sup>53</sup>, Jim Pivarski<sup>54</sup>, Harrison Prosper<sup>55</sup>, Fernanda Psihas<sup>56</sup>, Alexander Radovic<sup>57</sup>, Ryan Reece<sup>58</sup>, Aurelius Rinkevicius<sup>44</sup>, Eduardo Rodrigues<sup>7</sup>, Jamal Rorie<sup>59</sup>, David Rousseau<sup>60</sup>, Aaron Sauers<sup>21</sup>, Steven Schramm<sup>32</sup>, Ariel Schwartzman<sup>37</sup>, Horst Severini<sup>61</sup>, Paul Seyfert<sup>13</sup>, Filip Siroky<sup>62</sup>, Konstantin Skazytkin<sup>43</sup>, Mike Sokoloff<sup>7</sup>, Graeme Stewart<sup>63</sup>, Bob Stienen<sup>64</sup>, Ian Stockdale<sup>65</sup>, Giles Strong<sup>6</sup>, Wei Sun<sup>4</sup>, Savannah Thais<sup>51</sup>, Karen Tomko<sup>66</sup>, Eli Upfal<sup>12</sup>, Emanuele Usai<sup>12</sup>, Andrey Ustyuzhanin<sup>67</sup>, Martin Vala<sup>68</sup>, Sofia Vallecorsa<sup>69</sup>, Justin Vassel<sup>56</sup>, Mauro Verzetti<sup>70</sup>, Xavier Vilasís-Cardona<sup>71</sup>, Jean-Roch Vlimant<sup>3</sup>, Ilija Vukotic<sup>72</sup>, Sean-Jium Wang<sup>30</sup>, Gordon Watts<sup>73</sup>, Michael Williams<sup>74</sup>, Wenjing Wu<sup>75</sup>, Stefan Wunsch<sup>29</sup>, Kun Yang<sup>4</sup>, Omar Zapata<sup>76</sup>

<sup>1</sup> Lulea University of Technology

<sup>2</sup> NVidia

<sup>3</sup> California Institute of Technology

<sup>4</sup> Google

<sup>5</sup> Carnegie Mellon University

<sup>6</sup> LIP Lisboa

<sup>7</sup> University of Cincinnati

<sup>8</sup> Università e INFN, Padova

<sup>9</sup> University of London

<sup>10</sup> Lawrence Berkeley National Laboratory

<sup>11</sup> Università e INFN, Bologna

<sup>12</sup> Brown University

<sup>13</sup> CERN

<sup>14</sup> Argonne National Laboratory  
<sup>15</sup> CPPM Aix Marseille Univ CNRS/IN2P3  
<sup>16</sup> Universitaet Freiburg  
<sup>17</sup> New York University  
<sup>18</sup> Deutsches Elektronen-Synchrotron  
<sup>19</sup> Duke University  
<sup>20</sup> SISSA Trieste Italy  
<sup>21</sup> Fermi National Accelerator Laboratory  
<sup>22</sup> RWTH Aachen University  
<sup>23</sup> Universität Zürich  
<sup>24</sup> University of Texas at Arlington  
<sup>25</sup> Southern Methodist University  
<sup>26</sup> Ecole Polytechnique Federale de Lausanne  
<sup>27</sup> University of Manchester  
<sup>28</sup> University of Cambridge  
<sup>29</sup> Karlsruher Institut für Technologie  
<sup>30</sup> University of Florida  
<sup>31</sup> LPNHE, Sorbonne Université et Université Paris Diderot, CNRS/IN2P3, Paris  
<sup>32</sup> Université de Genève  
<sup>33</sup> Louisiana Tech University  
<sup>34</sup> Purdue University  
<sup>35</sup> University of Illinois at Urbana-Champaign  
<sup>36</sup> Max Planck Institut für Physik  
<sup>37</sup> SLAC National Accelerator Laboratory  
<sup>38</sup> SoundCloud  
<sup>39</sup> University of Milan  
<sup>40</sup> Intel  
<sup>41</sup> Brookhaven National Laboratory  
<sup>42</sup> University of Bristol  
<sup>43</sup> Russian Academy of Sciences  
<sup>44</sup> Cornell University  
<sup>45</sup> University of Notre Dame  
<sup>46</sup> University of Melbourne  
<sup>47</sup> University of California Berkeley  
<sup>48</sup> Università & INFN, Milano Bicocca  
<sup>49</sup> Vrije Universiteit Brussel  
<sup>50</sup> University of Amsterdam and Radboud University Nijmegen  
<sup>51</sup> Yale University  
<sup>52</sup> University of Alabama  
<sup>53</sup> University of Massachusetts  
<sup>54</sup> Princeton University  
<sup>55</sup> Florida State University  
<sup>56</sup> Indiana University  
<sup>57</sup> College of William and Mary  
<sup>58</sup> University of California, Santa Cruz  
<sup>59</sup> Rice University  
<sup>60</sup> Université de Paris Sud 11  
<sup>61</sup> University of Oklahoma  
<sup>62</sup> Masaryk University  
<sup>63</sup> University of Glasgow  
<sup>64</sup> Radboud Universiteit Nijmegen  
<sup>65</sup> Altair Engineering  
<sup>66</sup> Ohio Supercomputer Center  
<sup>67</sup> Yandex School of Data Analysis  
<sup>68</sup> Technical University of Kosice  
<sup>69</sup> Gangneung-Wonju National University  
<sup>70</sup> University of Rochester  
<sup>71</sup> University of Barcelona  
<sup>72</sup> University of Chicago  
<sup>73</sup> University of Washington  
<sup>74</sup> Massachusetts Institute of Technology

<sup>75</sup> Chinese Academy of Sciences  
<sup>76</sup> OProject and University of Antioquia

# Contents

|           |                                                                                            |           |
|-----------|--------------------------------------------------------------------------------------------|-----------|
| <b>1</b>  | <b>Preface</b>                                                                             | <b>6</b>  |
| <b>2</b>  | <b>Introduction</b>                                                                        | <b>6</b>  |
| 2.1       | Motivation . . . . .                                                                       | 6         |
| 2.2       | Brief Overview of Machine Learning Algorithms in HEP . . . . .                             | 6         |
| 2.3       | Structure of the Document . . . . .                                                        | 7         |
| <b>3</b>  | <b>Machine Learning Applications and R&amp;D</b>                                           | <b>7</b>  |
| 3.1       | Simulation . . . . .                                                                       | 7         |
| 3.2       | Real Time Analysis and Triggering . . . . .                                                | 8         |
| 3.3       | Object Reconstruction, Identification, and Calibration . . . . .                           | 8         |
| 3.4       | End-To-End Deep Learning . . . . .                                                         | 9         |
| 3.5       | Sustainable Matrix Element Method . . . . .                                                | 9         |
| 3.6       | Matrix Element Machine Learning Method . . . . .                                           | 10        |
| 3.7       | Learning the Standard Model . . . . .                                                      | 11        |
| 3.8       | Theory Applications . . . . .                                                              | 12        |
| 3.9       | Uncertainty Assignment . . . . .                                                           | 12        |
| 3.10      | Monitoring of Detectors, Hardware Anomalies and Preemptive Maintenance . . . . .           | 13        |
| 3.11      | Computing Resource Optimization and Control of Networks and Production Workflows . . . . . | 13        |
| <b>4</b>  | <b>Collaborating with other communities</b>                                                | <b>13</b> |
| 4.1       | Introduction . . . . .                                                                     | 13        |
| 4.2       | Academic Outreach and Engagement . . . . .                                                 | 14        |
| 4.3       | Machine Learning Challenges . . . . .                                                      | 14        |
| 4.4       | Collaborative Benchmark Datasets . . . . .                                                 | 14        |
| 4.5       | Industry Engagement . . . . .                                                              | 15        |
| 4.6       | Machine Learning Community-at-large Outreach . . . . .                                     | 15        |
| <b>5</b>  | <b>Machine Learning Software and Tools</b>                                                 | <b>16</b> |
| 5.1       | Software Methodology . . . . .                                                             | 16        |
| 5.2       | I/O and Programming Languages . . . . .                                                    | 16        |
| 5.3       | Software Interfaces to Acceleration Hardware . . . . .                                     | 17        |
| 5.4       | Parallelization and Interactivity . . . . .                                                | 17        |
| 5.5       | Internal and External ML tools . . . . .                                                   | 17        |
| 5.5.1     | Machine Learning Data Formats . . . . .                                                    | 17        |
| 5.5.2     | Desirable HEP-ML Software and Data Format Attributes . . . . .                             | 19        |
| 5.5.3     | Interfaces and Middleware . . . . .                                                        | 19        |
| <b>6</b>  | <b>Computing and Hardware Resources</b>                                                    | <b>19</b> |
| 6.1       | Resource Requirements . . . . .                                                            | 20        |
| 6.2       | Graphical Processing Units . . . . .                                                       | 21        |
| 6.3       | Cloud TPUs . . . . .                                                                       | 21        |
| 6.4       | High Performance Computing . . . . .                                                       | 21        |
| 6.5       | Field Programmable Gate Arrays . . . . .                                                   | 21        |
| 6.6       | Opportunistic Resources . . . . .                                                          | 21        |
| 6.7       | Data Storage and Availability . . . . .                                                    | 22        |
| 6.8       | Software Distribution and Deployment . . . . .                                             | 22        |
| 6.9       | Machine Learning As a Service . . . . .                                                    | 22        |
| <b>7</b>  | <b>Training the community</b>                                                              | <b>22</b> |
| <b>8</b>  | <b>Roadmap</b>                                                                             | <b>22</b> |
| 8.1       | Timeline . . . . .                                                                         | 22        |
| 8.2       | Steps to Deployment . . . . .                                                              | 23        |
| <b>9</b>  | <b>Conclusions</b>                                                                         | <b>23</b> |
| <b>10</b> | <b>Acknowledgements</b>                                                                    | <b>23</b> |

**A Appendix** **24**

A.1 Matrix Element Methods . . . . . 24

# 1 Preface

To outline the challenges in computing that high-energy physics will face over the next years and strategies to approach them, the HEP software foundation has organised a Community White Paper (CWP) [1]. In addition to the main document, several more detailed documents were worked out by different working groups. The present document focusses on the topic of machine learning. The goals are to define the tasks at the energy and intensity frontier that can be addressed during the next decade by research and development of machine learning applications.

Machine learning in particle physics is evolving fast, while the contents of this community white paper were mainly compiled during community meetings in spring 2017 that took place at several workshops on machine learning in high-energy physics: S2I2 and [2–5]. The contents of this document thus reflect the state of the art at these events and does not attempt to take later developments into account.

## 2 Introduction

One of the main objectives of particle physics in the post-Higgs boson discovery era is to exploit the full physics potential of both the Large Hadron Collider (LHC) and its upgrade, the high luminosity LHC (HL-LHC), in addition to present and future neutrino experiments. The HL-LHC will deliver an integrated luminosity that is 20 times larger than the present LHC dataset, bringing quantitatively and qualitatively new challenges due to event size, data volume, and complexity. The physics reach of the experiments will be limited by the physics performance of algorithms and computational resources. Machine learning (ML) applied to particle physics promises to provide improvements in both of these areas.

Incorporating machine learning in particle physics workflows will require significant research and development over the next five years. Areas where significant improvements are needed include:

- **Physics performance** of reconstruction and analysis algorithms;
- **Execution time** of computationally expensive parts of event simulation, pattern recognition, and calibration;
- **Realtime implementation** of machine learning algorithms;
- **Reduction of the data footprint** with data compression, placement and access.

### 2.1 Motivation

The experimental high-energy physics (HEP) program revolves around two main objectives that go hand in hand: probing the Standard Model (SM) with increasing precision and searching for new particles associated with physics beyond the SM. Both tasks require the identification of rare signals in immense backgrounds. Substantially increased levels of pile-up collisions from additional protons in the bunch at the HL-LHC will make this a significant challenge.

Machine learning algorithms are already the state-of-the-art in event and particle identification, energy estimation and pile-up suppression applications in HEP. Despite their present advantage, machine-learning algorithms still have significant room for improvement in their exploitation of the full potential of the dataset.

### 2.2 Brief Overview of Machine Learning Algorithms in HEP

This section provides a brief introduction to the most important machine learning algorithms in HEP, introducing key vocabulary (in *italic*).

Machine learning methods are designed to exploit large datasets in order to reduce complexity and find new features in data. The current most frequently used machine learning algorithms in HEP are Boosted Decision Trees (BDTs) and Neural Networks (NN).

Typically, variables relevant to the physics problem are selected and a machine learning *model* is *trained* for *classification* or *regression* using signal and background events (or *instances*). Training the model is the most human- and CPU-time consuming step, while the application, the so called *inference* stage, is relatively inexpensive. BDTs and NNs are typically used to classify particles and events. They are also used for regression,

where a continuous function is learned, for example to obtain the best estimate of a particle’s energy based on the measurements from multiple detectors.

Neural Networks have been used in HEP for some time; however, improvements in training algorithms and computing power have in the last decade led to the so-called deep learning revolution, which has had a significant impact on HEP. Deep learning is particularly promising when there is a large amount of data and features, as well as symmetries and complex non-linear dependencies between inputs and outputs.

There are different types of DNN used in HEP: fully-connected (FCN), convolutional (CNN) and recurrent (RNN). Additionally, neural networks are used in the context of Generative Models, where a Neural Network is trained to reproduce the multidimensional distribution of the training instances set. Variational AutoEncoders (VAE) and more recent Generative Adversarial Networks (GAN) are two examples of such generative models used in HEP.

A large set of machine learning algorithms is devoted to time series analysis and prediction. They are in general not relevant for HEP data analysis where events are independent from each other. However, there is more and more interest in these algorithms for Data Quality, Computing and Accelerator Infrastructure monitoring, as well as those physics processes and event reconstruction tasks where time is an important dimension.

## 2.3 Structure of the Document

Applications of machine learning algorithms motivated by HEP drivers are detailed in Section 3, while Section 4 focuses on outreach and collaboration with the machine learning community. Section 5 focuses on the machine learning software in HEP and discusses the interplay between internally and externally developed machine learning tools. Recent progress in machine learning was made possible in part by emergence of suitable hardware for training complex models, thus in Section 6 the resource requirements of training and applying machine learning algorithms in HEP are discussed. Section 7 discusses strategies for training the HEP community in machine learning. Finally, Section 8 presents the roadmap for the near future.

# 3 Machine Learning Applications and R&D

This chapter describes the science drivers and high-energy physics challenges where machine learning can play a significant role in advancing the current state of the art. These challenges are selected because of their relevance and potential and also due to similarity with challenges faced outside the field. Despite similarities, major R&D work will go in adapting and evolving such methods to match the particular HEP requirements.

## 3.1 Simulation

Particle discovery relies on the ability to accurately compare the observed detector response data with expectations based on the hypotheses of the Standard Model or models of new physics. While the processes of subatomic particle interactions with matter are known, it is intractable to compute the detector response analytically. As a result, Monte Carlo simulation tools, such as GEANT [6], have been developed to simulate the propagation of particles in detectors to compare with the data. The dedicated CWP on detector simulation [7] discusses the challenges of simulations in great detail. This section focuses on the machine learning related aspects.

For the HL-LHC, on the order of trillions of simulated collisions are needed in order to achieve the required statistical accuracy of the simulations to perform precision hypothesis testing. However, such simulations are highly computationally expensive. For example, simulating the detector response of a single LHC proton-proton collision event takes on the order of several minutes. A particularly time consuming step is the simulation of particles incident on the dense material of a calorimeter. The high interaction probability and resulting high multiplicity in the so-called showers of particles passing through the detector material make the simulation of such processes very expensive. This problem is further compounded when particle showers overlap, as is frequently the case in the core of a jet of particles produced by high energy quarks and gluons.

Fast simulations replace the slowest components of the simulation chain with computationally efficient approximations. Often such approximations have been done by using simplified parametrizations or particle shower look-up tables. These are computationally fast but often suffer from insufficient accuracy for high precision physics measurements and searches.

Recent progress in high fidelity fast generative models, such as GANs and VAEs, which learn to sample from high dimensional feature distributions by minimizing an objective that measures the distance between the generated and actual distribution, offer a promising alternative for simulation. A simplified first attempt at using such techniques saw orders of magnitude increase in simulation speed over existing fast simulation techniques [8], but such generative models have not yet reached the required accuracy partly due to inherent shortcomings of the methods and the instability in training of the GANs. Developing these techniques for realistic detector models and understanding how to reach the required accuracy is still needed. The fast advancement in the ML community of such techniques makes this a highly promising avenue to pursue.

Orthogonal to the reduction of demand of computing resources with fast simulations, machine learning can also contribute to other aspects of the simulation. Event generators have a large number of parameters that can be used to tune various aspects of the simulated events. Performing such tuning over many-dimensional parameter space is highly non-trivial and may require generating many data samples in the process to test parameter space points. Modern machine learning optimization techniques, such as Bayesian Optimization, allow for global optimization of the generator without detailed knowledge of its internal details [9]. Applying such techniques to simulation tuning may further improve the output of the simulations.

### 3.2 Real Time Analysis and Triggering

The traditional approach to data analysis in particle physics assumes that the interesting events recorded by a detector can be selected in real-time (a process known as *triggering*) with a reasonable efficiency, and that once selected, these events can be affordably stored and distributed for further selection and analysis at a later point in time. However, the enormous production cross-section and luminosity of the LHC mean that these assumptions break down.<sup>1</sup> In particular there are whole classes of events, for example beauty and charm hadrons or low-mass dark matter signatures, which are so abundant that it is not affordable to store all of the events for later analysis. To exploit the full information the LHC delivers, it will increasingly be necessary to perform more of the data analysis in real-time [10].

This topic is discussed in some detail in the Reconstruction and Software Triggering chapter [11], but it is also an important driver of machine learning applications in HEP. Machine learning methods offer the possibility to offset some of the cost of applying reconstruction algorithms, and may be the only hope of performing the real-time reconstruction that enables real-time analysis in the first place. For example, the CMS experiment uses boosted decision trees in the Level 1 trigger to approximate muon momenta. One of the challenges is the trade-off in algorithm complexity and performance under strict inference time constraints. In another example, called the HEP.TrkX project, deep neural networks are trained on large resource platforms and subsequently perform fast inference in online systems.

Real-time analysis poses specific challenges to machine learning algorithm design, in particular how to maintain insensitivity to detector performance which may vary over time. For example, the LHCb experiment uses neural networks for fast fake-track and clone rejection and already employs a fast boosted decision tree for a large part of the event selection in the trigger [12]. It will be important that these approaches maintain performance for higher detector occupancy for the full range of tracks used in physics analyses. Another related application is speeding up the reconstruction of beauty, charm, and other lower mass hadrons, where traditional track combinatorics and vertexing techniques may become too computationally expensive.

In addition, the increasing event complexity particularly in the HL-LHC era will mean that machine learning techniques may also become more important to maintaining or improving the efficiency of traditional triggers. Examples of where ML approaches can be useful are the triggering of electroweak events with low-energy objects; improving jet calibration at a very early stage of reconstruction allowing jet triggers thresholds to be lowered; or supernovae and proton decay triggering at neutrino experiments.

### 3.3 Object Reconstruction, Identification, and Calibration

The physical processes of interest in high energy physics experiments occur on time scales too short to be observed directly by particle detectors. For instance, a Higgs boson produced at the LHC will decay within approximately  $10^{-22}$  seconds and thus decays essentially at the point of production. However, the decay products of the initial particle, which are observed in the detector, can be used to infer its properties. Better knowledge of the properties (e.g. type, energy, direction) of the decay products permits more accurate reconstruction of the initial physical process. Event reconstruction at large is discussed in [11] and also disusses the following

<sup>1</sup>They may well also break down in other areas of high-energy physics in due course.



applications of machine learning.

Experiments have trained ML algorithms on the features from combined reconstruction algorithms to perform particle identification for decades. In the past decade BDTs have been one of the most popular techniques in this domain. More recently, experiments have focused on extracting better performance with deep neural networks.

An active area of research is the application of DNNs to the output of feature extraction in order to perform particle identification and extracting particle properties [13]. This is particularly true for calorimeters or time projection chambers (TPCs), where the data can be represented as a 2D or 3D image and the problems can be cast as computer vision tasks, in which neural networks are used to reconstruct images from pixel intensities. These neural networks are adapted for particle physics applications by optimizing network architectures for complex, 3-dimensional detector geometries and training them on suitable signal and background samples derived from data control regions. Applications include identification and measurements of electrons and photons from electromagnetic showers, jet properties including substructure and b-tagging, taus and missing energy. Promising deep learning architectures for these tasks include convolutional, recurrent and adversarial neural networks. A particularly important application is to Liquid Argon TPCs (LArTPCs), which are the chosen detection technology for the flagship neutrino program.

For tracking detectors, pattern recognition is the most computationally challenging step. In particular, it becomes computationally intractable for the HL-LHC. The hope is that machine learning will provide a solution that scales linearly with LHC collision density. A current effort called HEP.TrkX investigates deep learning algorithms such as long short-term memory (LSTM) networks for track pattern recognition on many-core processors.

### 3.4 End-To-End Deep Learning

The vast majority of analyses at the LHC use high-level features constructed from particle four-momenta, even when the analyses make use of machine learning. A high-profile example of such variables are the seven, so-called MELA variables, used in the analysis of the final states  $H \rightarrow ZZ \rightarrow 4\ell$ . While a few analyses, first at the Tevatron, and later at the LHC, have used the four-momenta directly, the latter are still high-level relative to the raw data. Approaches based on the four-momenta are closely related to the Matrix Element Method, which is described in the next section.

Given recent spectacular advances in image recognition based on the use of raw information, we are led to consider whether there is something to be gained by moving closer to using raw data in LHC analyses. This so-called *end-to-end deep learning* approach uses low level data from a detector together with deep learning algorithms [14, 15]. One obvious challenge is that low level data, for example, detector hits, tend to be both high-dimensional and sparse. Therefore, there is interest in also exploring automatic ways to compress raw data in a controlled way that does not necessarily rely on domain knowledge.

### 3.5 Sustainable Matrix Element Method

The Matrix Element (ME) Method [16–19] is a powerful technique which can be utilized for measurements of physical model parameters and direct searches for new phenomena. It has been used extensively by collider experiments at the Tevatron for standard model (SM) measurements and Higgs boson searches [20–25] and at the LHC for measurements in the Higgs and top quark sectors of the SM [26–32]. A few more details on the ME method are given in Appendix A.1.

The ME method has several unique and desirable features, most notably it (1) does not require training data being an *ab initio* calculation of event probabilities, (2) incorporates all available kinematic information of a hypothesized process, including all correlations, and (3) has a clear physical meaning in terms of the transition probabilities within the framework of quantum field theory.

One drawback to the ME Method is that it has traditionally relied on leading order (LO) matrix elements, although nothing limits the ME method to LO calculations. Techniques that accommodate initial-state QCD radiation within the LO ME framework using transverse boosting and dedicated transfer functions to integrate over the transverse momentum of initial-state partons have been developed [33]. Another challenge is development of the transfer functions which rely on tediously hand-crafted fits to full simulated Monte-Carlo events.

The most serious difficulty in the ME method that has limited its applicability to searches for beyond-the-SM physics and precision measurements is that it is very *computationally intensive*. If this limitation is overcome, it would enable more widespread use of ME methods for analysis of LHC data. This could be particularly important for extending the new physics reach of the HL-LHC which will be dominated by increases in integrated luminosity rather than center-of-mass collision energy.

The application of the ME method is computationally challenging for two reasons: (1) it involves high-dimensional integration over a large number of events, signal and background hypotheses, and systematic variations and (2) it involves sharply-peaked integrands<sup>2</sup> over a large domain in phase space. Therefore, despite the attractive features of the ME method and promise of further optimization and parallelization, the computational burden of the ME technique will continue to limit its range of applicability for practical data analysis without new and innovative approaches. The primary idea put forward in this section is to utilize modern *machine learning techniques to dramatically speed up the numerical evaluations in the ME method* and therefore broaden the applicability of the ME method to the benefit of HL-LHC physics.

Applying neural networks to numerical integration problems is plausible but not new (see [34–36], for example). The technical challenge is to design a network which is sufficiently rich to encode the complexity of the ME calculation for a given process over the phase space relevant to the signal process. Deep Neural Networks (DNNs) are strong candidates for networks with sufficient complexity to achieve good approximations, possibly in conjunction with smart phase-space mapping such as described in [37]. Promising demonstration of the power of Boosted Decision Trees [38, 39] and Generative Adversarial Networks [40] for improved Monte Carlo integration can be found in [41]. Once a set of DNNs representing definite integrals is generated to good approximation, evaluation of the ME method calculations via the DNNs will be very fast. These DNNs can be thought of as preserving the essence of ME calculations in a way that allows for fast forward execution. They can enable the ME method to be both *nimble* and *sustainable*, neither of which is true today.

The overall strategy is to do the expensive full ME calculations as infrequently as possible, ideally once for DNN training and once more for a final pass before publication, with the DNNs utilized as a good approximation in between. A future analysis flow using the ME method with DNNs might look something like the following: One performs a large number of ME calculations using a traditional numerical integration technique like VEGAS [42, 43] or FOAM [44] on a large CPU resource, ideally exploiting acceleration on many-core devices. The DNN training data is generated from the phase space sampling in performing the full integration in this initial pass, and DNNs are trained either *in situ* or *a posteriori*. The accuracy of the DNN-based ME calculation can be assessed through this procedure. As the analysis develops and progresses through selection and/or sample changes, systematic treatment, etc., the DNN-based ME calculations are used in place of the time-consuming, full ME calculations to make the analysis nimble and to preserve the ME calculations. Before a result using the ME method is published, a final pass using full ME calculation would likely be performed both to maximize the numerical precision or sensitivity of the results and to validate the analysis evolution via the DNN-based approximations.

There are several activities which are proposed to further develop the idea of a Sustainable Matrix Element Method. The first is to establish a cross-experiment group interested in developing the ideas presented in this section, along with a common software project for ME calculations in the spirit of [45]. This area is very well-suited for impactful collaboration with computer scientists and those working in machine learning. Using a few test cases (e.g.  $t\bar{t}$  or  $t\bar{t}h$  production), evaluation of DNN choices and configurations, developing methods for DNN training from full ME calculations and direct comparisons of the integration accuracy between Monte Carlo and DNN-based calculations should be undertaken. More effort should also be placed in developing compelling applications of the ME method for HL-LHC physics. In the longer term, the possibility of Sustainable-Matrix-Element-Method-as-a-Service (SMEMaaS), where shared software and infrastructure could be used through a common API, is proposed.

### 3.6 Matrix Element Machine Learning Method

The matrix element method is based in the fact that the physics of particle collisions is encoded in the distribution of the particles' four-momenta and with their flavors. As noted in the previous section, the fundamental task is to approximate the left-hand side of Eq. (5) for all (exclusive) final states of interest. In the matrix element method, one proceeds by approximating the right-hand side of Eq. (5). But, since the goal is to compute  $\mathcal{P}_\xi(\mathbf{x}|\alpha)$ , and given that billions of fully simulated events will be available, and that the simulations use exactly the same inputs as in the matrix element method, namely, the matrix elements, parton distribution

<sup>2</sup>a consequence of imposing energy/momentum conservation in the processes

functions, and transfer (or response) functions, one can ask whether a more direct machine learning approach can be developed to approximate  $\mathcal{P}_\xi(\mathbf{x}|\boldsymbol{\alpha})$  without the need to execute the calculation of the right-hand side of Eq. (5) explicitly. We believe the answer is yes, provided that a key advantage of the matrix element method can be replicated, namely, the fact that the method provides a function that depends explicitly on the model parameters  $\boldsymbol{\alpha}$ . Simulated events are typically simulated at fixed values of the model parameters. In order to replicate the advantage of the matrix element method, it would seem necessary to simulate events over an ensemble of model parameter points. But, that is a huge computational burden. The question is: is there a way to sidestep? Perhaps.

The matrix element method is technically feasible because there are now many codes that provide access to the square of the matrix elements as a function of  $\boldsymbol{\alpha}$ . Consequently, it would be possible to build a parametrized set of simulated events by reweighting each simulated event using the weighting function

$$w(\boldsymbol{\alpha}, \boldsymbol{\alpha}_0) = \frac{|\mathcal{M}_\xi(\mathbf{y}|\boldsymbol{\alpha})|^2}{|\mathcal{M}_\xi(\mathbf{y}|\boldsymbol{\alpha}_0)|^2}, \quad (1)$$

where  $\boldsymbol{\alpha}_0$  denotes the values of parameters at which the events were simulated. In practice, in order to keep the dynamic range of the weights within reasonable bounds, one presumably would simulate sets of events at a few reasonable choices of the parameters and use the weights to interpolate between these fixed parameter points.

What could one do with these parametrized simulated events? One could take advantage of the mathematical fact (shown more than quarter century ago) that, given a sufficiently expressive parametrized function  $f(x, \omega)$ , with parameters  $\omega$ , fitted by minimizing either the quadratic loss or cross entropy using data comprising two classes of objects of equal size, for example, signals with density  $s(x)$  assigned a target of unity and backgrounds with density  $b(x)$  assigned a target of zero, then asymptotically — that is, for arbitrarily large training samples,

$$f(x, \omega) = \frac{s(x)}{s(x) + b(x)}. \quad (2)$$

This result was derived in the context of neural networks. However, it is in fact entirely independent of the nature of the function  $f(x, \omega)$ .

We can exploit this result to approximate  $\mathcal{P}_\xi(\mathbf{x}|\boldsymbol{\alpha})$  directly using, for example, DNNs. For each simulated event, in the training set, one could sample  $\boldsymbol{\alpha}$  from a known distribution  $q(\boldsymbol{\alpha})$ , thereby yielding an ensemble of triplets  $\{(\mathbf{x}, \boldsymbol{\alpha}, w(\boldsymbol{\alpha}, \boldsymbol{\alpha}_0))\}$ . Call this ensemble the “signal”. Sample  $\mathbf{x}$  from another known distribution  $p(\mathbf{x})$  and for each  $\mathbf{x}$  sample  $\boldsymbol{\alpha}$  from  $q(\boldsymbol{\alpha})$ . Call the ensemble of pairs  $\{(\mathbf{x}, \boldsymbol{\alpha})\}$  the “background”. From Eq. (2), we have

$$f(x, \omega) = \frac{s(\mathbf{x}, \boldsymbol{\alpha})}{s(\mathbf{x}, \boldsymbol{\alpha}) + p(\mathbf{x})q(\boldsymbol{\alpha})}, \quad (3)$$

from which we find

$$\mathcal{P}_\xi(\mathbf{x}|\boldsymbol{\alpha}) = \frac{s(\mathbf{x}, \boldsymbol{\alpha})}{q(\boldsymbol{\alpha})} = p(\mathbf{x}) \left( \frac{f}{1-f} \right). \quad (4)$$

If this could be made to work, it would be a direct machine learning alternative to the matrix element method, which incorporates the advantages of the former with the added advantage that the transfer function is automatically incorporated without the need to model it explicitly and the calculation of  $\mathcal{P}_\xi(\mathbf{x}|\boldsymbol{\alpha})$  would be fast because it would be given directly in terms of the DNN,  $f(x, \omega)$ .

### 3.7 Learning the Standard Model

New physics may manifest itself as unusual or rare events. One approach is to accurately identify the Standard Model processes and search for anomalies. Classifying the Standard Model events is a challenging task, as it consists of many complicated physics processes. Multi-class machine learning algorithms are well-suited for this classification problem. Once an event is classified as likely to be from a known physics process, it can be filtered out, and remaining events can be further analyzed for hints of new physics. Additionally, unsupervised machine learning techniques can be applied to remaining events to cluster them together. This approach would also be useful in identifying detector problems. Another possibility is to use unsupervised learning to estimate the distributions of Standard Model events in the high-dimensional feature space. The comparison of such distributions with the measured datasets allows one to assess the compatibility of data with the background or to spot differences. The latter case would represent a hint of the presence of new physics events in the data. Density estimation techniques may be particularly useful in this task.

### 3.8 Theory Applications

The theoretical physics community has a number of challenges where machine learning can have an impact. These include areas of theoretical model optimization with hundreds of parameters, searches for new models, understanding and estimation of the parton distribution functions and possibly quantum machine learning. The following details one such application: learning of the parton distribution functions with machine learning.

Making progress towards the objectives of the HL-LHC program (see Section 2) requires not only obtaining experimental measurements of the physical processes but also reliable theory inputs to compare to. This becomes increasingly challenging as the experimental data gets more precise. There are numerous examples of phenomenologically relevant processes where the experimental uncertainty is comparable to the estimate of the theoretical uncertainty of the corresponding calculation.

Furthermore, the theory does not predict the value of all the inputs required for the computations (for example the value of the strong coupling constant  $\alpha_S$  evaluated at the  $Z$  mass), and there are situations where the equations resulting from theory cannot be solved to describe the physics adequately, and the corresponding theory inputs must be obtained from data instead. A more complex example is the determination of Parton Distribution Functions (PDFs): Quantum Chromodynamics (QCD) describes the proton collisions at high energy in terms of *partons* (e.g. quarks and gluons), but it is not possible to calculate directly from QCD the momentum carried by each quark or gluon within a proton since QCD is not solvable in its confined regime. Our lack of theoretical knowledge about the characterization of partons within a proton is embedded into a suitable definition of the Parton Distribution Functions (approximately the momentum densities of each of the partons). The PDFs then need to be determined from experimental data. The NNPDF collaboration uses machine learning techniques to obtain a PDF determination that is accurate enough to be suitable for high-precision collider data comparison. The NNPDF fitting procedure is described in full details in [46].

The idea is to combine data from all relevant physical processes and fit a neural network representing each PDF. The difficulty of the procedure stems from the fact that multiple experimental inputs need to be combined to obtain a PDF fit. Each of these inputs adjoins only indirect constraints on the PDFs, leaving some regions of the PDF completely unconstrained by data. NNPDF fit includes around 50 datasets from different physical processes, and results that are not always consistent among themselves. Therefore, it is crucial to propagate the uncertainty of the experimental inputs into uncertainty on the PDFs via Monte Carlo Dropout [47].

While the dataset is small, each experimental point has an indirect relation to the PDFs, as it is the result of the convolution of one or two PDFs with the corresponding partonic cross section. Code has been developed to compute these convolutions called APFELgrid [48]. Future research directions include the possibility of using standard ML frameworks to efficiently express the PDF fitting problem. The uncertainties of the theory calculations need to be taken into account as well in the fits. A fully systematic treatment of theory errors in PDFs is a topic of research where machine learning could play an important role. The dominant uncertainties in the data are no longer statistical and instead arise from correlated systematics. Determining those systematics accurately is non trivial on the side of the experimental analyses and can have a major impact on the resulting PDFs. The problem grows more complex when ML techniques for which there is no simple recipe to estimate the uncertainty are used extensively in the experimental analysis. Taking full advantage of these advanced methods requires interdisciplinary research and communication on topics such as developing regularization schemes for experimental covariance matrices.

In conclusion, it is not only important to obtain the best fit PDF, but also a reliable estimation of the uncertainty, which in turn requires controlling the uncertainty of the experimental and theoretical inputs.

### 3.9 Uncertainty Assignment

Until fairly recently, little attention has been paid to the problem of assigning a measure of uncertainty to the output of machine learning algorithms. However, there is a growing recognition that the lack of such measures is a serious deficiency that needs to be addressed. This is particularly problematic in particle physics where accurate uncertainty assignments are crucial for assessing the quality of quantitative results. Uncertainty assignment is especially urgent if machine learning methods are to be used for regression.

Standard methods exist in the statistics literature to quantify the uncertainty when an explicit likelihood function is available. Unfortunately, however, the overlap between the machine learning and statistics communities has been minimal at best. The problem of uncertainty offers many opportunities for fruitful collaborations between physicists, statisticians, computer scientists, and machine learning practitioners, and is something that

ought to be vigorously encouraged. One area ripe for collaboration is in the use of Bayesian methods to quantify uncertainty using, for example, Hamiltonian Monte Carlo (HMC) sampling of posterior densities. One serious issue is that HMC is computationally prohibitive for DNNs. Therefore, new and more effective ways to perform large scale Bayesian calculations that go beyond HMC need to be developed, and it would be particularly useful to develop methods that can take advantage of multicore machines as discussed in Section 6. One possibility may be to explore deep probabilistic programming [49]

### 3.10 Monitoring of Detectors, Hardware Anomalies and Preemptive Maintenance

Data-taking of current complex HEP detectors is continuously monitored by physicists taking shifts to check the quality of the incoming data. Typically, hundreds of histograms have been defined by experts and shifters are alerted when an unexpected deviation with respect to a reference occurs. It is a common occurrence for a new type of problem to remain unseen for a non-negligible period of time because such a situation had not been foreseen by the expert.

A whole class of ML algorithms called anomaly detection can be useful in such situations. They are able to learn from data and produce an alert when a deviation is seen. By monitoring many variables at the same time such algorithms are sensitive to subtle signs forewarning of imminent failure, so that preemptive maintenance can be scheduled. Such techniques are already used in industry applications.

One challenge is that normal drifts in environmental conditions can induce drifts in the data. Beyond just reporting a problem, the natural next step is to connect an anomaly detection algorithm to the appropriate action: restart an online computer, contact an on-call expert, or similar. In the long term, the hardware and data structures of future detectors should be designed to facilitate the operation of anomaly detection algorithms.

### 3.11 Computing Resource Optimization and Control of Networks and Production Workflows

Data operations is one of the significant challenges for the upcoming HL-LHC. In the current infrastructure, LHC experiments rely on in-house solutions for managing the data. While these approaches work reasonably well today, machine learning can help automate and improve the overall system throughput and reduce operational costs.

Machine learning can be applied in many areas of computing infrastructure, workflow and data management. For example, dataset placement optimization and reduction of transfer latency can lead to a better usage of site resources and an increased throughput of analysis jobs. One of the current examples is predicting the “popularity” of a dataset from dataset usage, which helps reduce disk resource utilization and improve physics analysis time turn-over.

Data volume in data transfers is one of the challenges facing the current computing systems as thousand of users need to access thousands of datasets across the Grid. There is an enormous amount of metadata collected by application components, such as information about failures, file accesses, etc. Resource utilization optimization based on this data, including Grid components and software stack layers, can improve overall operations. Understanding the data transfer latencies and network congestion may improve operational costs of hardware resources.

Networks are going to play a crucial role in data exchange and data delivery to scientific applications in HL-LHC era. The network-aware application layer and configurations may significantly affect experiment’s daily operations. ML can be applied to network security in order to identify anomalies in network traffic; predict network congestion; detect bugs via analysis of self-learning networks, and optimize WAN paths based on user access patterns.

## 4 Collaborating with other communities

### 4.1 Introduction

Discovery science provides a challenge that attracts brilliant minds eager to push the boundaries of scientific understanding of nature. Particle physics has a rich problem domain that offers avenues for intellectual reward. The goal is to achieve a vibrant collaboration between the data science and high-energy physics communities



by finding a common language and working together to further science. Both communities can benefit from such collaboration. The HEP community can explore new research directions and applications of machine learning, novel algorithms, and direct collaboration on HEP challenges. The ML community can benefit from a diverse set of particle physics problems with unique challenges in scale and complexity, and a large community of researchers that can expand the machine learning horizon by contributing to solving problems relevant to both communities. **For example, the treatment of systematic uncertainties is an important topic for the HEP and ML communities.** By working together on common challenges, the two fields can further progress in solving such problems.

**There are a number of existing examples of collaboration between HEP and ML that have produced fruitful results through mostly local connections (e.g. [8, 50, 51]).** The HEP community should continue such collaborations and look for additional collaborations with the ML community.

The HEP community needs to define its challenges in a language that the ML community can understand. This may involve only retaining necessary information with clear and concise explanations as to its relevance. Machine learning likewise has a significant amount of domain knowledge. Ideas and solutions provided by both communities should be presented in an understandable way for scientists without in-depth knowledge.

## 4.2 Academic Outreach and Engagement

Direct collaboration between HEP researchers and computer scientists working in the area of machine learning is an important possible driver of innovation in the area of machine learning applications described in Section 3. It is important for high-energy physics to engage and collaborate with the academic community focused on new machine learning algorithms and applications as they will naturally be interested in applying new ideas to interesting and complex data provided by HEP.

Conferences and workshops are a core aspect of the academic ML community, and organizing or contributing to key conferences is a means of gaining interest. Organizing sessions or mini-workshops within major ML conferences, such as NIPS, would increase the familiarity of HEP within the ML community and jump-start future collaborations. This has been explored in single cases [52] but is not an established, regular workshop series. At the same time inviting ML experts to HEP workshops, as done at [53] and the DS@HEP series [2, 54, 55], can foster greater long-term collaboration. There should be coordinated efforts to:

- Organize workshops and conferences open to external collaborators to discuss applications, algorithms and tools
- Organize thematic workshops around topics relevant to HEP

HEP should reach out to other scientific communities with similar challenges, for example astrophysics/cosmology, medium energy nuclear physics and computational biology. This can lead to more active partnerships to better collaborate on ideas, techniques, and algorithms.

## 4.3 Machine Learning Challenges

To engage the wider ML community, **challenges such as the Higgs Boson Challenge (2014) or the Flavor Physics Challenge [52, 53] have been organized on Kaggle.** These types of challenges draw considerable attention from the machine learning community and additional similar challenges should be organized in the future.

Organization of a challenge requires a well documented dataset, a starting-kit and an evaluation metric to rank the solutions. This forces the organizers to simplify the problem as much as possible, while retaining its intrinsic complexity. The drawback of challenges is that once they are launched, participants' priority is winning the challenge and not eventual collaboration with HEP. It is important to foresee upfront a way to integrate incoming solutions, for example via forums and post-challenge workshops (like [53]) where a diversity of competitive algorithms can be presented. The challenge dataset and evaluation metric should be released publicly so that further developments can continue.

## 4.4 Collaborative Benchmark Datasets

There is a strong incentive for HEP to develop public benchmark datasets, beyond just challenges. Access to a dataset makes the discussion much more concrete and productive. Within the HEP community, common datasets enable comparisons of algorithms with much better accuracy, which is very useful for research and development. The same benchmark datasets can also be used for teaching, tutorials and training.

These benchmark datasets could be built based on public simulation engines or released by experiments within the bounds of their data access policy. Even a small subset of an experiment’s simulated data can be the base of a very valuable benchmark dataset. For example, the CMS experiment has released a significant amount of its simulated and collected data via the CERN Open Data Portal [56].

To be maximally useful, the subsequent guidelines should be followed when designing a benchmark dataset:

- Simplify the dataset as much as possible while conserving all methodical difficulties relevant to solving the problem
- Document the dataset to make it understandable by a non-HEP expert
- Create methodologies and metrics for the evaluation of proposed solutions, and document them
- Prepare an integration plan for incoming ideas and solutions
- Feedback results of successful applications

The HEP community should organize and curate a variety of such benchmark datasets covering its current physics drivers and make them publicly available. To improve the reproducibility of results and algorithm comparisons, some of the data used for evaluation of the solutions should be kept private. Additionally, after investing heavily into producing highly-detailed and realistic simulations, the HEP community can provide the machine learning community with labeled datasets with high statistical power to test algorithms and develop novel ideas.

## 4.5 Industry Engagement

Industry has been focused on the development and adoption of machine learning techniques. In addition to algorithm and software development, one of the promising areas is the adoption of dedicated specialized hardware and high performance co-processors. GPUs, FPGAs, and high core count co-processors all have the potential to dramatically increase performance of machine learning applications relevant to the HEP community. One of the challenges is gaining the human expertise for development and implementation. Industry interactions bring specific technology opportunities and access to specialized expertise that can be difficult to hire and support internally.

There are specific areas of development where industry has expressed interest in collaborating with HEP. Automated resource provisioning, data placement, and scheduling are similar to industrial applications to improve efficiency. Applications such as data quality monitoring, detector health monitoring and preventative maintenance can be automated using techniques developed for other industrial quality control applications. There are two more forward looking areas that coincide with HEP physics drivers, namely computer vision techniques for object identification and real-time event classification. These present a challenge to industry due to its complexity and benefit outside of HEP.

CERN OpenLab is a public-private partnership that accelerates the development of cutting-edge solutions for the LHC community and wider scientific research. CERN OpenLab has established the infrastructure to maintain non-disclosure agreements, to arrange ownership of intellectual property, and to provide an interface between CERN and industry. As part of its upcoming phases, OpenLab plans to explore machine learning applications for the benefit of LHC experiments computing and the HL-LHC. Such initiatives and industry partnerships should be supported in the future.

## 4.6 Machine Learning Community-at-large Outreach

Another form of engagement is using the communications mediums to broadcast our challenges and attract interested collaborators. There are a variety of channels which can be leveraged to increase the visibility of our problems and research opportunities in the ML community. These can be popular forums such as reddit, personal or official blogs, social media, or direct contact with influential personalities.

Podcasts have shown to be a great vehicle for reaching a large audience. Listeners are keen to consume material that is outside of their immediate problem domain in a way that is easy to digest. There is an abundance of machine learning podcasts with a large base of listeners that can be targeted for outreach:

- [Linear Digressions](#) (co-hosted by former ATLAS Ph.D. Katie Malone)
- [Partially Derivative](#)

- [Talking Machines](#)
- [Data Skeptic](#)
- [Becoming a Data Scientist Podcast](#)
- [Not So Standard Deviations](#)
- [This Week in ML & AI](#)

Another form of engagement is through outreach-style blog posts to explain HEP challenges in a way that is easy to understand by the public.

Another outreach opportunity is to make HEP related presentations at machine learning Meetups across the world to generate awareness, engage community, foster cross pollination of ideas between HEP and industry. Some popular ML meetups are:

- NYC: <https://www.meetup.com/NYC-Machine-Learning/>
- Berlin: <https://www.meetup.com/Advanced-Machine-Learning-Study-Group/>
- SF: <https://www.meetup.com/SF-Bayarea-Machine-Learning/>

In conclusion, existing outreach efforts should be expanded to attract greater collaboration between the HEP and ML communities. By understanding and speaking the same language, the two communities can better collaborate and find solutions to present and future challenges.

## 5 Machine Learning Software and Tools

Machine learning does not exist without software. There are a large variety of algorithms written in different programming languages and general software frameworks that combine many classes of methods into one package. The following sections focus on specific topics and challenges related to machine learning software design in HEP.

### 5.1 Software Methodology

Presently, there are two machine learning software methodologies in high-energy physics. The first approach is to implement abstract ML algorithms in HEP-developed toolkits, such as the Toolkit for Multivariate Analysis (TMVA) in ROOT. This provides on site support to physicists and dedicated development for HEP data formats and applications. The second approach is to rely on externally developed software, of which there are many examples. This often requires tedious and repetitive work to adapt HEP data formats to external software, breaks analysis workflows and introduces difficulties in the analysis software development. Historically, a variety of approaches and competition among them has led to important breakthroughs in the field. On the other hand, having too many choices increases repetition and leads community segmentation and possible issues with reproducibility. There is some convergence of the two approaches. Converters have been written for the reintegration of some externally trained models into HEP tools, like [57]. Complementarily, interfaces between HEP and external tools have also been developed.

### 5.2 I/O and Programming Languages

The sheer amount of data accumulated by HEP experiments requires data access optimization because. Such I/O performance is very dependent on data formats. Moreover, support for reading data in different formats is required for certain use-cases.

Exploration of new file systems and methods to improve I/O limitations are important and the following R&D studies should take place:

- Explore new file systems to assess I/O limitations;
- Use alternative industry approaches such as Google BigQuery to explore various data access patterns;
- Explore parallel data processing platforms such as Apache Spark for ML training.

Although particle physics has been reliant on C++ over the past decade, the machine learning community has explored other programming languages, in particular the python-based ecosystem.



### 5.3 Software Interfaces to Acceleration Hardware

Modern machine learning software significantly benefits from the use of hardware accelerators such as GPUs. At the same time, ML users should not be hindered in their development of new applications by writing platform-dependent code. Various interfaces to different hardware architectures are needed in order to make efficient use of the available computing resources. The emergence of the Open Computing Language (OpenCL) allows programming of high-level interfaces that can run on various hardware platforms.

Machine learning tools often provide different sets of APIs to develop and train the models in one language, and various bindings to use trained models in other programming languages. This is a convenient model for many HEP applications, such as the trigger system, where (1) application latency puts stringent requirements on the software and hardware used and (2) the general purpose training platform might be different from the highly specific deployment platform.

### 5.4 Parallelization and Interactivity

Training ML algorithms takes a significant amount of time and parallelization at various levels is desired, such as the parallelization of the computations within a single model. Another type of parallelism is data parallelism that targets the processing phase of the training with data partitioning and model training using distributed workers. Frameworks like Apache Spark and ideas such as batch training offer promise in this area.

One often needs to produce many different machine learning models, for example while tuning hyper-parameters or performing k-fold cross-validation, and distribution of these algorithms is key to the reduction of the overall training time. ML algorithm inference significantly benefits from parallelization as well. Although event selection in a particle physics trigger system is “embarrassingly parallel [58]”, parallel processing frameworks are developed to reduce the memory footprint. Complex models might well be too big to reside once per CPU core in memory and need to be hosted in shared memory. At the same time the stringent latency requirements impose constraints on the type of algorithms that can be easily deployed.

The availability of interactive frameworks, for example Jupyter notebooks, allows for rapid prototype development and testing of ML tools. Such frameworks also ease the connection between the description of models and the data, providing straightforward means of visualizing models and data. HEP has started to exploring interactive frameworks, such as the Service for Web Based Analysis (SWAN) [59]. One of the challenges is availability of adequate hardware resources for these systems.

### 5.5 Internal and External ML tools

Internally developed ML software, such as the Toolkit for Multivariate Analysis (TMVA, [60]), has been developed to apply a variety of machine learning algorithms to HEP challenges. Currently, many published HEP analyses with machine learning have made use of TMVA. There is also software developed in HEP, such as NeuroBayes [61, 62] and RuleFit, that have gained popularity outside of HEP.

At the same time, the ML landscape has evolved and many different ML tools have emerged and gained popularity. There are a growing number of published results based on externally developed tools. The latter, often developed directly by industry for specific applications, are constantly undergoing development, incorporating the latest algorithms from academia. Currently, both internal and external tools are used by the HEP community. TMVA has also undergone significant development in recent years. In addition, there are smaller tools developed in the HEP community, extending either internal or external ML tools for specific use cases and applications within HEP experiments, such as `hep_ml` [63] for training ML classifiers in HEP specific tasks, or `tmva-branch-adder` [64] and `lwttn` [57] to facilitate classifier inference tasks.

This begs the question: what aspects of ML development and use should the HEP community focus on in the next 5-10 years? There are several aspects to consider including data formats, community size, and interfaces.

#### 5.5.1 Machine Learning Data Formats

HEP and ML communities currently make use of different data formats. HEP heavily relies on the ROOT software framework for data storage, data processing, and data analysis. The machine learning community uses a large variety of formats, as shown in Figure 1. This figure also shows the relationship of machine learning data formats with ROOT: the ROOT file format is very flexible, though it requires a significant time investment to learn the proper use. Table 1 summarizes the current machine learning toolkits and file formats they support.

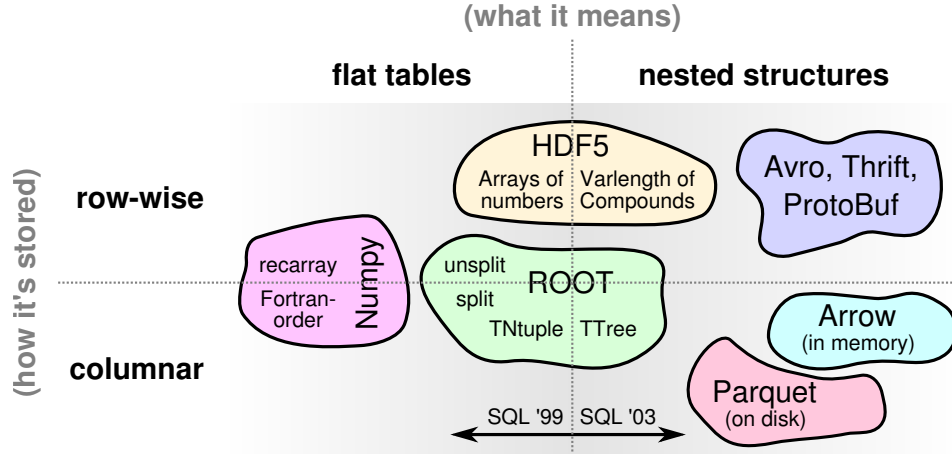


Figure 1: Existing data-formats used by ML communities.

Table 1: This table lists various data formats (rows) and ML tools (columns). The  $\checkmark$  indicates that there is a native solution, while  $\times$  means that the conversion from one data-format to another is straightforward (the conversion mechanisms from Table 2 are omitted here). The following notation has been used to denote the data-formats: **T** Trees, **F** flat tables, **M** sparse matrices, **R** row-wise arrays, **C** column-wise arrays **S** static data structures

|                | TMVA         | TensorFlow                                         | Theano       | Scikit<br>Learn | R            | Spark<br>ML  | VW           | libFM        | RGF          | Torch        |
|----------------|--------------|----------------------------------------------------|--------------|-----------------|--------------|--------------|--------------|--------------|--------------|--------------|
| ROOT [T, C]    | $\checkmark$ | through conversion into other formats, see Table 2 |              |                 |              |              |              |              |              |              |
| CSV [F]        |              | $\checkmark$                                       | $\checkmark$ | $\checkmark$    | $\checkmark$ | $\checkmark$ | $\times$     | $\times$     | $\times$     | $\checkmark$ |
| libSVM [M]     |              |                                                    |              |                 |              |              | $\times$     | $\checkmark$ | $\times$     |              |
| VW [M]         |              |                                                    |              |                 |              |              | $\checkmark$ |              |              |              |
| RGF [M]        |              |                                                    |              |                 |              |              |              |              | $\checkmark$ |              |
| NumPy [R]      | [65]         | $\checkmark$                                       | $\checkmark$ | $\checkmark$    | $\checkmark$ | $\checkmark$ | $\times$     | $\times$     | $\times$     | $\checkmark$ |
| Avro [S, R]    |              |                                                    |              |                 | $\checkmark$ | $\checkmark$ |              |              |              |              |
| Parquet [S, C] |              |                                                    |              |                 | $\checkmark$ | $\checkmark$ |              |              |              |              |
| HDF5 [S]       |              | $\times$                                           | $\times$     | $\times$        |              |              |              |              |              | $\checkmark$ |
| R df [R]       |              |                                                    |              |                 | $\checkmark$ |              |              |              |              |              |

### 5.5.2 Desirable HEP-ML Software and Data Format Attributes

A desirable data format should have the following attributes: high read speed for efficient training, sparse readability without loading the entire dataset into RAM, compression, and common use by the machine learning community.

HEP machine learning applications require highly performant and flexible algorithms to address the variety of use cases. Some applications, such as triggering, also have to work under tight latency constraints of the order of a few microseconds and below. The data sets are extremely large, which comes with I/O challenges as described in Section 5.2. This is expected to become even more challenging, as the LHC continues to ramp-up and deliver increasingly large amounts of data.

As discussed in Section 5.3, machine learning tools use a number of languages. To use these tools it will therefore be important to offer adequate support. In production software environments, ML classifiers will be part of larger C++ or python applications. Classifiers from external tools will probably fit perfectly into python applications, but for C++ applications, restrictions to toolkits with C++ API or the use of C++ converters may be needed to make sure the training result can be evaluated.

Advantages of using the **external tools** are the size of the community that uses and supports them, being able to easily keep up with progress in the industry and profit from the forefront of the ML research. It should also be noted that some of the recent industrial efforts to develop and maintain ML-tools rely on resources far beyond that of basic research. The deep learning tools of the previous and current generations constitute a demonstration of corresponding quality.

Disadvantages of using external tools are that, they is often no guaranteed support over the lifetime of particle physics experiments, and it can be difficult to adapt them to HEP specific requirements which may not be among the priorities of the ML community.

Advantages of using **internal tools** are that decisions about long-term support remain in the community, and the tools can be adapted to the specific needs of HEP. Disadvantages include that any new algorithm and idea first needs to be ported before their use can be evaluated.

### 5.5.3 Interfaces and Middleware

One approach to bridge the gap between internal and external tools is by building interfaces. Some researchers prefer to convert their data to the file formats used by external tools and to work exclusively with external tools. This has the advantage of working as close as possible with the ML community's tools and its documentation, and is as close to what is used by ML researchers as possible. This may come with the difficulty of converting back to a HEP format, when analyzing the data further with a fitting framework [66–68]. At the same time, interfaces have been built between TMVA and external machine learning tools, allowing for their use and direct comparison between their performance. Currently, interfaces to R, scikit-learn, keras and tensorflow have been developed. Those have the advantage of providing a homogeneous interface and require little training overhead for those already knowledgeable in using TMVA.

A more general approach to file format conversion is to build middleware solutions that export HEP-specific formats like ROOT to formats used by external machine learning tools. Existing middleware solutions are shown in Table 2.

Approaches to bridge the different languages and data formats inside and outside of HEP include providing interfaces or building middleware solutions that translate HEP-specific data formats to external ML tools. It is a topic of current research to determine the most efficient solution.

## 6 Computing and Hardware Resources

A typical high-energy physics data model consists of a hierarchy of increasingly refined data stores. Each store provides a refined view of a list of “events”, the self contained records that capture the state of the detector at the time when a particle interaction occurs. At the bottom of the hierarchy is the raw data, a byte-stream of the readout from detector electronics. At the top of the hierarchy are the “high-level” physics objects, such as electrons or jets, providing descriptive information about the quality and topology of physics events. The data stores are typically processed by independent copies of identical code processed in batch computing queues.

Table 2: Middleware solutions translating the ROOT data format to other formats.

|                    |                                                                                                                                             |
|--------------------|---------------------------------------------------------------------------------------------------------------------------------------------|
| <b>PyROOT</b>      | Python extension module that allows the user to interact with ROOT data/classes. [69]                                                       |
| <b>root_numpy</b>  | The interface between ROOT and NumPy supported by the Scikit-HEP community. [65]                                                            |
| <b>root_pandas</b> | The interface between ROOT and Pandas dataframes supported by the DIANA/HEP project. [70]                                                   |
| <b>uprooot</b>     | A high throughput I/O interface between ROOT and NumPy. [71]                                                                                |
| <b>c2numpy</b>     | Pure C-based code to convert ROOT data into Numpy arrays which can be used in C/C++ frameworks. [72]                                        |
| <b>root4j</b>      | The hep.io.root package contains a simple Java interface for reading ROOT files. This tool has been developed based on freehep-rootio. [73] |
| <b>root2npy</b>    | The go-hep package contains a reading ROOT files. This tool has been developed based on freehep-rootio. [73]                                |
| <b>root2hdf5</b>   | Converts ROOT files containing TTrees into HDF5 files containing HDF5 tables. [74]                                                          |

The result of this processing is filtered data and extracted physics parameters.

At present, training of machine learning algorithms is done using dedicated or private resources. These vary in configuration and processing power, depending on the size of the data and complexity of the algorithm. For a given event, the evaluation of algorithms is performed on a single core producing a single discriminator or regressor output. In order to progress to evaluation of complex machine learning algorithms, more computing power is needed in both the training and evaluation stages, as larger amounts of data are needed to feed models with tens or hundreds of thousands of parameters. This implies the expansion of the current computing model to include architectures that are well suited to machine learning tasks, such as many integrated core (MIC), graphical processing units (GPUs) and Tensor Processing Units (TPUs). This is a fundamental departure from the single-core or few-core jobs. These architectures provide a significant computational speed improvement for both training and evaluation of ML algorithms, but require dedicated hardware, drivers, and software configuration.

Similarly, the locality and bandwidth of large data stores will need to be optimized in order to avoid bottlenecks in training and evaluation for analysis. Data placement and the need to use dedicated hardware indicate that a transition to HPC, or HPC-like, architectures may be needed to achieve the desired performance. Due to significant synergy with the direction of industry in this respect, use of commercially available resources should be considered for future high-energy physics computing models.

In the following subsections we will discuss the resource needs for the physics drivers mentioned earlier: fast simulation, real-time analysis, object and event reconstruction and particle identification. The limitations of the current computing model are discussed as well as how those physics driver needs can be met in the future.

## 6.1 Resource Requirements

The popularity of deep learning methods is to a large extent due to the possibility of training these models in a reasonable amount of time with large scale parallelism. In particular, the training stage requires repeated simultaneous access to many data elements and specialized hardware has been developed for training deep learning models.

In contrast, inference can be an operation applied to a single data element at a time and needs to be performed only once. Inference has less demands on I/O and is limited only by the computing power and model complexity. Because inference has real-time applications in high-energy physics, latency and throughput constraints are the main challenges.

A typical HEP application can require up to 1 GPU-week to train a single model. To obtain the best results and understand the performance of the model, an average of 100 hyper-parameter points optimization is typically performed. A single project could therefore easily require up to a full GPU-year for training. The speed up of the training process can be obtained by means of faster and more capable hardware, parallelization of single training and by splitting training over multiple nodes. Resources, such as GPUs, TPUs and MIC need to be evaluated in the context of realistic benchmark particle physics applications.

If different ML techniques can achieve equivalent physics performance but require different processing power, it is important to quantify what is driving the performance gains and what level of performance-processing power trade-off is maintainable to achieve the required physics goals.

## 6.2 Graphical Processing Units

GPUs have been extremely useful in speeding up the training of complex deep learning models. Many researchers in HEP are currently relying on private or university GPU clusters to perform machine learning training. Unfortunately, there is no centralized GPU resource available for general HEP usage. Availability of greater and more modern GPU resources would significantly reduce the training time and assist many on-going *R&D* efforts in the community.

## 6.3 Cloud TPUs

Cloud TPUs systems are built around Google's TPUs, which are custom silicon chips for machine learning. Cloud TPU Pods are multi-rack supercomputers that can include more than 1,000 TPU chips as well as many host machines. Cloud TPUs and Cloud TPU Pods are optimized to accelerate the dense linear algebra that is commonly found in cutting-edge deep neural network models, but they can potentially be used for other purposes as well. Current efforts in the HEP community are studying the use of Cloud TPUs for possible acceleration of the training stage. Cloud TPUs and Cloud TPU Pods are available to the public via Google Cloud [75].

## 6.4 High Performance Computing

Resource-rich many-core processors such as MIC, GPUs, and TPUs are vital to the optimization of the training time of most modern machine learning algorithms, including deep neural networks, generative adversarial networks, autoencoders, etc. Availability of High Performance Computing (HPC) resources equipped with many-core processors and high-performance network storage are essential to distributed running of large-scale machine learning algorithms. Current efforts to bring and expand the availability of HPC resources in high-energy physics computing will be vital to the successful progress of the application of machine learning techniques to current and future experiments.

## 6.5 Field Programmable Gate Arrays

Field Programmable Gate Arrays (FPGAs) provide an efficient and low-latency application of machine learning algorithms directly at the level of hardware, as desirable for HEP trigger systems. The following ML algorithms are more suitable for FPGAs due to their simpler parallelization: boosted decision trees, random forests and decision rule ensembles. For example the CMS experiment currently uses boosted decision trees in FPGAs in the trigger system to estimate muon momenta. Further research and development is needed in this area to apply more advanced machine learning techniques like deep learning directly in the hardware. One of the challenges is the limited availability of floating point operations gates and the precision needed to maintain the best performance. The possibility of coupling the FPGAs with a CPU with significant random-access memory (RAM) allows the shift of some of these operations to RAM.

## 6.6 Opportunistic Resources

The current HEP computing model is based on a tiered structure where computing resources are mostly large data centers providing CPU resources for collaboration. Although existing resources are gradually moving towards supporting GPUs, it is unlikely to reach all HEP computing centers in the near future. Therefore opportunistic resources are a possible option for training machine learning applications.

Currently, cloud solutions provided by the industry run ML workflows on dedicated hardware and offer interfaces for training machine learning models. The scientific community should work closely with cloud providers to harmonize our analysis computing needs and data access patterns with their business models. Costs of the cloud resources should be compared with the costs of procuring these resources independently.

In order to make the best use of resources available to the community, all resources should ideally be made available through a unique work queue. That implies some uniformization of the software stack, and several specific requirements in the resource management system, especially in terms of data movement.

## 6.7 Data Storage and Availability

Data storage limitations will have a major impact on machine learning applications. Presently, to train machine learning algorithms, it has been possible to take advantage of increases in statistics of Monte-Carlo simulated events needed for other use cases. Further machine learning progress may require more simulated data than what is available today. How to produce and store these additional large amounts of data is a challenge that needs to be overcome.

Availability of data at PByte/EByte scale represents another challenge for the ML community. A good solution must provide access to a large data volume for hundred or thousand of users simultaneously. Additionally, data movement might need to be automatized to make the training data available transparently at high speed local storage with use of an automatic caching mechanism. The success of Apache Spark and Google BigQuery platforms may serve as a model. In addition to the regular HEP workflows, data streaming, transformation and readout in mini-batches may be required to train models over large data sets.

## 6.8 Software Distribution and Deployment

To efficiently use the resources described in previous subsections, machine learning software needs to be available on the computing resources. Platforms, such as the CERNVM File System (cvmfs), are very useful for distributing software instead of requiring local installation. Additional tools like docker containers for application shipping can be useful in providing homogeneous software environments across the different systems. Another challenge is that the software layer needs to be agnostic to the hardware back-end.

## 6.9 Machine Learning As a Service

Current cloud providers rely on the machine learning as a service model, allowing for the efficient use of common resources, and make use of interactive machine learning tools. Machine learning as a service is not yet widely used in HEP, but examples of successful publications which used machine learning as a service exist, e.g. [76]. Specialized HEP services for interactive analysis, such as CERN's Service for Web-based Analysis (SWAN) may play an important role in adoption of machine learning tools in HEP workflows. In order to use these tools more efficiently, sufficient and appropriately tailored hardware resources described in this chapter are needed. In addition, the evolution of current HEP reconstruction software model to support off-loading ML tasks to heterogeneous computing resources, as well as ML task pooling, should be studied.

# 7 Training the community

The CWP on training [77] discusses the training needs of the HEP community in detail. In order to address the communication barrier and to speak the same language, training in machine learning concepts and terminology have to be part of a standard curriculum. In addition to lectures also hands-on tutorials on specific tools are necessary. Being able to apply machine learning to practical HEP problems requires the understanding of basic machine learning concepts and algorithms. For this, regular data science lecture series and seminars, like [78], are very useful. Also university level courses dedicated to machine learning applications in physics research are an excellent way to train undergraduate and graduate students. Experimental collaborations currently have training activities for newcomers that focus on analysis software and introduction to domain knowledge [79]. Machine learning should next be incorporated into the incoming collaborators' training efforts of the experiments.

# 8 Roadmap

In this chapter we discuss the roadmap towards implementation of the research and development areas described in Section 3.

## 8.1 Timeline

The incorporation of machine learning into particle physics experiments must respect two primary time lines: the schedules of HL-LHC and funding agencies, and the experiments' need for extensive validation of the algorithms.

The current LHC schedule has Run 3 starting in 2021 and the HL-LHC starting in 2026. As software processes and algorithms are re-imagined, their implementation must fit into these time frames if they are to maximize their benefit to the particle physics community. To fit this schedule, a newly proposed implementation would

need to have a demonstration in 2018 to prove viability. Two years later, in 2020, the proposed idea needs to attain a level of maturity to be included in the HL-LHC Technical Design Report. The project should then be further refined towards a large scale test around the middle of Run 3, about 2022. Run 3 is scheduled to end in late 2023, after which point the project must then be adapted to the HL-LHC software and physics analysis environment as it will be required by the experiment.

## 8.2 Steps to Deployment

The path of taking a machine learning idea from conception to community-wide acceptance and deployment will entail several stages, as appropriate. There are ample opportunities to make the process more efficient.

1. **Problem formulation and data set preparation:** Problem formulation is the first step in building a machine learning algorithm. The inputs and desired output need to be established. The training and validation data sets must be identified and simulated. In many cases, these data sets are large, and resources must be identified to possibly create and store the data. In most cases, the data needs to be processed into a form suitable for input into the algorithm. As discussed in Section 4.4, common benchmark samples will later on facilitate the comparison of different approaches.
2. **Feasibility and demonstration:** Given a dataset, appropriate machine learning algorithms need to be investigated and evaluated for their suitability to solve the problem.
3. **First application:** An application of the solution to one or a few specific physics analysis examples where the ML technique significantly improves the result. The incorporation of the technique into the computing work-flow will likely be very specific to the application and require significant manual intervention.
4. **Scaling and optimization:** Evolving from a demonstration to a general solution requires the use of realistic data sets with full detector simulation. Furthermore, the solution will also require optimization to achieve nominal physics and computing performance. This stage will likely require significant computing resources to scale solutions to the full detector and data sets.
5. **Integration and Validation:** The solution needs to be incorporated into the experimental software and work-flow and must be validated.

As an example, an effort has recently started to build generative models that can significantly accelerate the simulation of particle showers in calorimeters. These early efforts are based on simplified data sets specifically created for this problem, without the complications of realistic data and limited to a small calorimeter sections. The first papers [8] use generative adversarial networks to generate calorimetric data which are reasonably faithful but still require tuning. The next step involves exploration of neural network architectures and systematic hyper-parameter scans to achieve the required performance.

The process of employing the new technique in a publication will elicit scrutiny by the full experiment, effectively validating the technique. Once the technique is accepted, it can be generalized beyond this first application and then incorporated into the experiment’s software for use by others. Finally, as the technique is applied to an increasing number of physics analyses, the technique will be incorporated into the experiment’s production work-flows.

## 9 Conclusions

Machine learning algorithms are already state of the art in many areas of particle physics and will likely be called on to take on a greater role in solving upcoming data analysis and event reconstruction challenges. In this document we have outlined the promising areas of research and development applications of machine learning in particle physics and focused on addressing the most important science drivers. We identified the need for greater collaboration with external communities in machine learning and a need to train the particle physics community in machine learning. We also identified how these prerequisites for successful incorporation of machine learning applications into high-energy physics can be met and provided an example roadmap for the implementation of machine learning applications into the workflows of particle physics experiments.

## 10 Acknowledgements

Vladimir Vava Gligorov acknowledges funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme under grant agreement No 724777 “RECEPT”. Omar Zapata is supported by Sostenibilidad-UdeA and COLCIENCIAS through the grant No 111577657253.

# A Appendix

## A.1 Matrix Element Methods

The ME method is based on *ab initio* calculations of the probability density function  $\mathcal{P}$  of an event with observed final-state particle momenta  $\mathbf{x}$  to be due to a physics process  $\xi$  with theory parameters  $\boldsymbol{\alpha}$ . One can compute  $\mathcal{P}_\xi(\mathbf{x}|\boldsymbol{\alpha})$  by means of the factorization theorem from the corresponding partonic cross-sections of the hard-scattering process involving parton momenta  $\mathbf{y}$  and is given by

$$\mathcal{P}_\xi(\mathbf{x}|\boldsymbol{\alpha}) = \frac{1}{\sigma_\xi^{\text{fiducial}}(\boldsymbol{\alpha})} \int d\Phi(\mathbf{y}_{\text{final}}) dx_1 dx_2 \frac{f(x_1)f(x_2)}{2sx_1x_2} |\mathcal{M}_\xi(\mathbf{y}|\boldsymbol{\alpha})|^2 \delta^4(\mathbf{y}_{\text{initial}} - \mathbf{y}_{\text{final}}) W(\mathbf{x}, \mathbf{y}) \quad (5)$$

where and  $x_i$  and  $\mathbf{y}_{\text{initial}}$  are related by  $y_{\text{initial},i} \equiv \frac{\sqrt{s}}{2}(x_i, 0, 0, \pm x_i)$ ,  $f(x_i)$  are the parton distribution functions,  $\sqrt{s}$  is the collider center-of-mass energy,  $\sigma_\xi^{\text{fiducial}}(\boldsymbol{\alpha})$  is the total cross section for the process  $\xi$  (with  $\boldsymbol{\alpha}$ ) times the detector acceptance,  $d\Phi(\mathbf{y})$  is the phase space density factor,  $\mathcal{M}_\xi(\mathbf{y}|\boldsymbol{\alpha})$  is the matrix element (typically at leading-order (LO)), and  $W(\mathbf{x}, \mathbf{y})$  is the probability density (aka “transfer function”) that a selected event  $\mathbf{y}$  ends up as a measured event  $\mathbf{x}$ . One can use calculations of Eq. 5 in a number of ways (e.g. likelihood functions) to search for new phenomena at particle colliders.

As stated in Sect. 3.5, the ME method has three notable features: it (1) does not require training data being an *ab initio* calculation of event probabilities, (2) incorporates all available kinematic information of a hypothesized process, including all correlations, and (3) has a clear physical meaning in terms of the transition probabilities within the framework of quantum field theory.

In reference to point (1), the matrix element  $\mathcal{M}_\xi(\mathbf{y}|\boldsymbol{\alpha})$  in the method involves all partons in the  $n \rightarrow m$  process, so when the 4-momentum of particles are not completely measured experimentally (e.g. neutrinos), one must integrate over the missing information which increases the dimensionality of the integration. In reference to point (2), a clever technique to re-map the phase space in order to reduce the sharpness of integrate in that space in an automated way (MADWEIGHT [37]) is often used in conjunction with a matrix element calculation package (MADGRAPH\_aMCNLO [80]). In practice, evaluation of definite integrals by the ME approach invokes techniques such as importance sampling (see VEGAS [42, 43] and FOAM [44]) or recursive stratified sampling (see MISER [81]) Monte Carlo integration. Acceleration of some of these techniques on modern computing architectures has been achieved, for example concurrent phase space sampling in VEGAS on GPUs.



## References

- [1] A. A. Alves Jr et al. “A Roadmap for HEP Software and Computing R&D for the 2020s” (2017). arXiv: 1712.06982 [physics.comp-ph].
- [2] *DS@HEP*. 2017. URL: <https://indico.fnal.gov/conferenceDisplay.py?confId=13497>.
- [3] *IML Machine Learning Workshop*. 2017. URL: <https://indico.cern.ch/event/595059/>.
- [4] *18th International Workshop on Advanced Computing and Analysis Techniques in Physics Research*. 2017. URL: <https://indico.cern.ch/event/567550/>.
- [5] *HEP Software Foundation Workshop*. 2017. URL: <https://indico.cern.ch/event/613093>.
- [6] S. Agostinelli et al. “GEANT4: A Simulation toolkit.” *Nucl. Instrum. Meth.* A506 (2003), pp. 250–303. DOI: 10.1016/S0168-9002(03)01368-8.
- [7] J. Apostolakis et al. “HEP Software Foundation Community White Paper Working Group - Detector Simulation” (2018). Ed. by V. Elvira and J. Harvey. arXiv: 1803.04165 [physics.comp-ph].
- [8] M. Paganini, L. de Oliveira, and B. Nachman. “CaloGAN: Simulating 3D High Energy Particle Showers in Multi-Layer Electromagnetic Calorimeters with Generative Adversarial Networks” (2017). arXiv: 1705.02355 [hep-ex].
- [9] P. Ilten, M. Williams, and Y. Yang. “Event generator tuning using Bayesian optimization.” *JINST* 12.04 (2017), P04028. DOI: 10.1088/1748-0221/12/04/P04028. arXiv: 1610.08328 [physics.data-an].
- [10] S. Benson et al. “The LHCb Turbo Stream.” *Journal of Physics: Conference Series* 664.8 (2015), p. 082004. URL: <http://stacks.iop.org/1742-6596/664/i=8/a=082004>.
- [11] J. Albrecht et al. “HEP Community White Paper on Software trigger and event reconstruction” (2018). Ed. by V. V. Gligorov and D. Lange. arXiv: 1802.08638 [physics.comp-ph].
- [12] V. V. Gligorov and M. Williams. “Efficient, reliable and fast high-level triggering using a bonsai boosted decision tree.” *JINST* 8 (2013), P02013. DOI: 10.1088/1748-0221/8/02/P02013. arXiv: 1210.6861 [physics.ins-det].
- [13] D. Derkach et al. “Machine-Learning-based global particle-identification algorithms at the LHCb experiment.” *Proceedings, 18th International Workshop on Advanced Computing and Analysis Techniques in Physics Research (ACAT 2017): Seattle, United States of America, August 21-25, 2017*. URL: <https://indico.cern.ch/event/567550/contributions/2629676/>.
- [14] M. Andrews et al. “End-to-End Physics Event Classification with the CMS Open Data: Applying Image-based Deep Learning on Detector Data to Directly Classify Collision Events at the LHC” (2018). arXiv: 1807.11916 [hep-ex].
- [15] M. Andrews et al. “End-to-End Jet Classification of Quarks and Gluons with the CMS Open Data” (2019). arXiv: 1902.08276 [hep-ex].
- [16] K. Kondo. “Dynamical Likelihood Method for Reconstruction of Events With Missing Momentum. 1: Method and Toy Models.” *J. Phys. Soc. Jap.* 57 (1988), pp. 4126–4140. DOI: 10.1143/JPSJ.57.4126.
- [17] F. Fiedler et al. “The Matrix Element Method and its Application in Measurements of the Top Quark Mass.” *Nucl. Instrum. Meth.* A624 (2010), pp. 203–218. DOI: 10.1016/j.nima.2010.09.024. arXiv: 1003.1316 [hep-ex].
- [18] I. Volobouev. “Matrix Element Method in HEP: Transfer Functions, Efficiencies, and Likelihood Normalization.” *ArXiv e-prints* (Jan. 2011). arXiv: 1101.2259 [physics.data-an].
- [19] F. Elahi and A. Martin. “Using the modified matrix element method to constrain  $L_\mu - L_\tau$  interactions.” *Phys. Rev. D* 96.1 (2017), p. 015021. DOI: 10.1103/PhysRevD.96.015021. arXiv: 1705.02563 [hep-ph].
- [20] V. M. Abazov et al. “A precision measurement of the mass of the top quark.” *Nature* 429 (2004), pp. 638–642. DOI: 10.1038/nature02589. arXiv: hep-ex/0406031 [hep-ex].
- [21] A. Abulencia et al. “Precision measurement of the top quark mass from dilepton events at CDF II.” *Phys. Rev. D* 75 (2007), p. 031105. DOI: 10.1103/PhysRevD.75.031105. arXiv: hep-ex/0612060 [hep-ex].
- [22] T. Aaltonen et al. “First Measurement of ZZ Production in panti-p Collisions at  $\sqrt{s} = 1.96$ -TeV.” *Phys. Rev. Lett.* 100 (2008), p. 201801. DOI: 10.1103/PhysRevLett.100.201801. arXiv: 0801.4806 [hep-ex].
- [23] T. Aaltonen et al. “Inclusive Search for Standard Model Higgs Boson Production in the WW Decay Channel using the CDF II Detector.” *Phys. Rev. Lett.* 104 (2010), p. 061803. DOI: 10.1103/PhysRevLett.104.061803. arXiv: 1001.4468 [hep-ex].
- [24] V. M. Abazov et al. “Observation of Single Top Quark Production.” *Phys. Rev. Lett.* 103 (2009), p. 092001. DOI: 10.1103/PhysRevLett.103.092001. arXiv: 0903.0850 [hep-ex].

- [25] T. Aaltonen et al. “First Observation of Electroweak Single Top Quark Production.” *Phys. Rev. Lett.* 103 (2009), p. 092002. DOI: 10.1103/PhysRevLett.103.092002. arXiv: 0903.0885 [hep-ex].
- [26] S. Chatrchyan et al. “Observation of a new boson at a mass of 125 GeV with the CMS experiment at the LHC.” *Phys. Lett. B* 716 (2012), pp. 30–61. DOI: 10.1016/j.physletb.2012.08.021. arXiv: 1207.7235 [hep-ex].
- [27] S. Chatrchyan et al. “Measurement of the properties of a Higgs boson in the four-lepton final state.” *Phys. Rev. D* 89.9 (2014), p. 092007. DOI: 10.1103/PhysRevD.89.092007. arXiv: 1312.5353 [hep-ex].
- [28] G. Aad et al. “Measurements of Higgs boson production and couplings in the four-lepton channel in pp collisions at center-of-mass energies of 7 and 8 TeV with the ATLAS detector.” *Phys. Rev. D* 91.1 (2015), p. 012006. DOI: 10.1103/PhysRevD.91.012006. arXiv: 1408.5191 [hep-ex].
- [29] V. Khachatryan et al. “Measurement of spin correlations in  $t\bar{t}$  production using the matrix element method in the muon+jets final state in pp collisions at  $\sqrt{s} = 8$  TeV.” *Phys. Lett. B* 758 (2016), pp. 321–346. DOI: 10.1016/j.physletb.2016.05.005. arXiv: 1511.06170 [hep-ex].
- [30] V. Khachatryan et al. “Search for a Standard Model Higgs Boson Produced in Association with a Top-Quark Pair and Decaying to Bottom Quarks Using a Matrix Element Method.” *Eur. Phys. J. C* 75.6 (2015), p. 251. DOI: 10.1140/epjc/s10052-015-3454-1. arXiv: 1502.02485 [hep-ex].
- [31] G. Aad et al. “Search for the Standard Model Higgs boson produced in association with top quarks and decaying into  $b\bar{b}$  in pp collisions at  $\sqrt{s} = 8$  TeV with the ATLAS detector.” *Eur. Phys. J. C* 75.7 (2015), p. 349. DOI: 10.1140/epjc/s10052-015-3543-1. arXiv: 1503.05066 [hep-ex].
- [32] G. Aad et al. “Evidence for single top-quark production in the  $s$ -channel in proton-proton collisions at  $\sqrt{s} = 8$  TeV with the ATLAS detector using the Matrix Element Method.” *Phys. Lett. B* 756 (2016), pp. 228–246. DOI: 10.1016/j.physletb.2016.03.017. arXiv: 1511.05980 [hep-ex].
- [33] J. Alwall, A. Freitas, and O. Mattelaer. “The Matrix Element Method and QCD Radiation.” *Phys. Rev. D* 83 (2011), p. 074010. DOI: 10.1103/PhysRevD.83.074010. arXiv: 1010.2263 [hep-ph].
- [34] Z. Zhe-Zhao, W. Yao-Nan, and W. Hui. “Numerical integration based on a neural network algorithm.” 8 (Aug. 2006), pp. 42–48.
- [35] L. y. Xu and L. j. Li. “The New Numerical Integration Algorithm Based on Neural Network.” *Third International Conference on Natural Computation (ICNC 2007)*. Vol. 1. Aug. 2007, pp. 325–328. DOI: 10.1109/ICNC.2007.730.
- [36] L. Yan, J. Di, and K. Wang. “Spline Basis Neural Network Algorithm for Numerical Integration.” *International Journal of Mathematical, Computational, Physical, Electrical and Computer Engineering* 7.3 (2013), pp. 458–461. ISSN: PISSN:2010-376X, EISSN:2010-3778. URL: <http://waset.org/Publications?p=75>.
- [37] P. Artoisenet et al. “Automation of the matrix element reweighting method.” *JHEP* 12 (2010), p. 068. DOI: 10.1007/JHEP12(2010)068. arXiv: 1007.3300 [hep-ph].
- [38] J. Friedman, T. Hastie, and R. Tibshirani. “Additive logistic regression: a statistical view of boosting (With discussion and a rejoinder by the authors).” *Ann. Statist.* 28.2 (Apr. 2000), pp. 337–407. DOI: 10.1214/aos/1016218223.
- [39] J. H. Friedman. “Greedy function approximation: A gradient boosting machine.” *Ann. Statist.* 29.5 (Oct. 2001), pp. 1189–1232. DOI: 10.1214/aos/1013203451.
- [40] I. J. Goodfellow et al. “Generative Adversarial Networks.” *ArXiv e-prints* (June 2014). arXiv: 1406.2661 [stat.ML].
- [41] J. Bendavid. “Efficient Monte Carlo Integration Using Boosted Decision Trees and Generative Deep Neural Networks” (2017). arXiv: 1707.00028 [hep-ph].
- [42] G. P. Lepage. “A new algorithm for adaptive multidimensional integration.” *Journal of Computational Physics* 27.2 (1978), pp. 192–203. ISSN: 0021-9991. DOI: [http://dx.doi.org/10.1016/0021-9991\(78\)90004-9](http://dx.doi.org/10.1016/0021-9991(78)90004-9). URL: <http://www.sciencedirect.com/science/article/pii/0021999178900049>.
- [43] T. Ohl. “Vegas revisited: Adaptive Monte Carlo integration beyond factorization.” *Comput. Phys. Commun.* 120 (1999), pp. 13–19. DOI: 10.1016/S0010-4655(99)00209-X. arXiv: hep-ph/9806432 [hep-ph].
- [44] S. Jadach. “Foam: A general-purpose cellular Monte Carlo event generator.” *Computer Physics Communications* 152.1 (2003), pp. 55–100. ISSN: 0010-4655. DOI: [http://dx.doi.org/10.1016/S0010-4655\(02\)00755-5](http://dx.doi.org/10.1016/S0010-4655(02)00755-5). URL: <http://www.sciencedirect.com/science/article/pii/S0010465502007555>.
- [45] S. Brochet et al. “MoMEMta, a modular toolkit for the Matrix Element Method at the LHC.” *Eur. Phys. J. C* 79.2 (2019), p. 126. DOI: 10.1140/epjc/s10052-019-6635-5. arXiv: 1805.08555 [hep-ph].

- [46] R. D. Ball et al. “Parton distributions for the LHC Run II.” *JHEP* 04 (2015), p. 040. DOI: 10.1007/JHEP04(2015)040. arXiv: 1410.8849 [hep-ph].
- [47] Y. Gal and Z. Ghahramani. “Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning.” *arXiv e-prints*, arXiv:1506.02142 (June 2015), arXiv:1506.02142. arXiv: 1506.02142 [stat.ML].
- [48] V. Bertone, S. Carrazza, and N. P. Hartland. “APFELgrid: a high performance tool for parton density determinations.” *Comput. Phys. Commun.* 212 (2017), pp. 205–209. DOI: 10.1016/j.cpc.2016.10.006. arXiv: 1605.02070 [hep-ph].
- [49] D. Tran et al. “Deep Probabilistic Programming.” *arXiv e-prints*, arXiv:1701.03757 (Jan. 2017), arXiv:1701.03757. arXiv: 1701.03757 [stat.ML].
- [50] T. Likhomanenko, D. Derkach, and A. Rogozhnikov. “Inclusive Flavour Tagging Algorithm.” *J. Phys. Conf. Ser.* 762.1 (2016), p. 012045. DOI: 10.1088/1742-6596/762/1/012045. arXiv: 1705.08707 [hep-ex].
- [51] *darkmachines*. URL: <http://www.darkmachines.org>.
- [52] *ALEPH Workshop @ NIPS 2015 – Applying (machine) Learning to Experimental Physics (ALEPH) and «Flavours of Physics» challenge*. Neural Information Processing Systems (NIPS). 2015. URL: <http://yandexdataschool.github.io/aleph2015/>.
- [53] *Heavy Flavour Data Mining Workshop*. Zurich, 2016. URL: <https://indico.cern.ch/event/433556/>.
- [54] *Data Science @ LHC 2015 Workshop*. 2015. URL: <https://indico.cern.ch/event/395374/>.
- [55] *DS@HEP at the Simons Foundation*. 2016. URL: <https://indico.hep.caltech.edu/indico/conferenceDisplay.py?confId=102>.
- [56] *CERN Open Data Portal*. URL: <https://opendata.cern.ch>.
- [57] D. H. Guest et al. *lwtnn/lwtnn: Version 2.7.1*. DOI: 10.5281/zenodo.1403888. URL: <https://doi.org/10.5281/zenodo.1403888>.
- [58] I. Foster. *Designing and Building Parallel Programs*. Section 1.4.4. Addison-Wesley, 1995. ISBN: 9780201575941.
- [59] E. Tejedor Saavedra. “Facilitating collaborative analysis in SWAN.” *23rd International Conference on Computing in High Energy and Nuclear Physics*. 2018. URL: <https://indico.cern.ch/event/587955/contributions/2935943/>.
- [60] A. Hoecker et al. “TMVA: Toolkit for Multivariate Data Analysis.” *PoS ACAT* (2007), p. 040. arXiv: physics/0703039.
- [61] M. Feindt. “A Neural Bayesian Estimator for Conditional Probability Densities” (2004). arXiv: physics/0402093.
- [62] M. Feindt and U. Kerzel. “The NeuroBayes neural network package.” *Nucl.Instrum.Meth.* A559 (2006), pp. 190–194. DOI: 10.1016/j.nima.2005.11.166.
- [63] A. Rogozhnikov. *hep\_ml python package*. URL: [https://github.com/arogozhnikov/hep\\_ml](https://github.com/arogozhnikov/hep_ml).
- [64] P. Seyfert. *pseyfert/tmva-branch-adder: Version 0.5.0*. DOI: 10.5281/zenodo.1149690. URL: <https://doi.org/10.5281/zenodo.1149690>.
- [65] N. Dawe et al. *scikit-hep/root-numpy: 4.7.3*. DOI: 10.5281/zenodo.842249. URL: <https://doi.org/10.5281/zenodo.842249>.
- [66] W. Verkerke and D. Kirkby. “The RooFit toolkit for data modeling” (2003). arXiv: physics/0306116.
- [67] L. Moneta et al. “The RooStats Project.” *PoS ACAT2010* (2010), p. 057. DOI: 10.22323/1.093.0057. arXiv: 1009.1003 [physics.data-an].
- [68] K. Cranmer et al. “HistFactory: A tool for creating statistical models for use with RooFit and RooStats” (2012).
- [69] *PyROOT*. URL: <https://root.cern.ch/pyroot>.
- [70] I. Babuschkin et al. *scikit-hep/root-pandas: 0.3.2*. DOI: 10.5281/zenodo.1188928. URL: <https://doi.org/10.5281/zenodo.1188928>.
- [71] J. Pivarski et al. *scikit-hep/uproot: 2.8.17*. DOI: 10.5281/zenodo.1219218. URL: <https://doi.org/10.5281/zenodo.1219218>.
- [72] *c2numpy*. URL: <https://github.com/diana-hep/c2numpy>.
- [73] *root4j*. URL: <https://github.com/diana-hep/root4j>.
- [74] *root2hdf5*. URL: <http://www.rootpy.org/commands/root2hdf5.html>.

- [75] *cloudgoogle*. URL: <https://cloud.google.com/tpu/>.
- [76] R. Aaij et al. “Search for the lepton flavour violating decay  $\tau^- \rightarrow \mu^- \mu^+ \mu^-$ .” *JHEP* 02 (2015), p. 121. DOI: 10.1007/JHEP02(2015)121. arXiv: 1409.8548 [hep-ex].
- [77] D. Berzano et al. “HEP Software Foundation Community White Paper Working Group - Training, Staffing and Careers.” 2018.
- [78] *Third Machine Learning in High Energy Physics Summer School*. 2017. URL: <https://indico.cern.ch/event/613571/>.
- [79] L. Bel. “The LHCb Starterkit.” *Proceedings of the 38th International Conference on High Energy Physics (ICHEP2016). 3-10 August 2016. Chicago, USA*. 2016, p. 334. URL: <http://pos.sissa.it/cgi-bin/reader/conf.cgi?confid=282>.
- [80] J. Alwall et al. “The automated computation of tree-level and next-to-leading order differential cross sections, and their matching to parton shower simulations.” *JHEP* 07 (2014), p. 079. DOI: 10.1007/JHEP07(2014)079. arXiv: 1405.0301 [hep-ph].
- [81] W. H. Press and G. R. Farrar. “Recursive stratified sampling for multidimensional Monte Carlo integration.” *Submitted to: Comp.in Phys.* (1989).