

Machine Learning Community White Paper

See page 19 for the full author list

August 23, 2017

Contents

5	1 Introduction	3
6	1.1 Motivation	3
7	1.2 Brief Overview of Machine Learning Algorithms in HEP	3
8	1.3 Structure of the Document	4
9	2 Machine Learning Applications and R&D	4
10	2.1 Detector Simulation	4
11	2.2 Real Time Analysis and Triggering	4
12	2.3 Object Reconstruction, Identification, and Calibration	5
13	2.4 Sustainable Matrix Element Method	6
14	2.5 Learning the Standard Model	8
15	2.6 Theory Applications	8
16	2.7 Monitoring Detectors, Hardware Anomalies and Preemptive Maintenance	9
17	2.8 Computing Resource Optimization and Control of Networks and Production Workflows	9
18	2.9 Collaborative Benchmark Datasets	9
19	3 Machine Learning Software and Tools	10
20	3.1 Software Methodology	10
21	3.2 I/O and Programming Languages	10
22	3.3 Software Interfaces to Acceleration Hardware	10
23	3.4 Parallelization and Interactivity	10
24	3.5 Internal and external ML tools	11
25	3.5.1 Machine Learning Data Formats	11
26	3.5.2 Desirable HEP-ML software and data format attributes	12
27	3.5.3 Interfaces and middleware	12
28	4 Computing and Hardware Resources	12
29	4.1 Resource Requirements	13
30	4.2 Resource Evaluation	13
31	4.2.1 High Performance Computing	13
32	4.2.2 Field Programmable Gate Arrays	13
33	4.2.3 Opportunistic Resources	14
34	4.2.4 Data Storage and Availability	14
35	4.2.5 Software Distribution and Deployment	14
36	4.2.6 Machine Learning As a Service	14
37	5 Training the community	14
38	6 Collaborating with other communities	15
39	6.1 Introduction	15
40	6.2 Machine Learning Challenges	15
41	6.3 Collaborative Benchmark Datasets	15
42	6.4 ML Academic outreach	15
43	6.5 Science outreach	16
44	6.6 Industry Engagement	16
45	6.6.1 CERN OpenLab and research-industry collaborative initiatives	16
46	6.7 ML community at large outreach	16

47	7 Roadmap	17
48	A Author list	19

1 Introduction

One of the main objectives of particle physics in the post-Higgs boson discovery era is to exploit the the full physics potential of the Large Hadron Collider (LHC) and its upgrade, the high luminosity LHC (HL-LHC) and of the current and future neutrino experiments. The HL-LHC will deliver data from 100 times the luminosity compared to the LHC, bringing quantitatively and qualitatively new challenges due to event size, 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.

1.1 Motivation

The experimental high-energy physics (HEP) program revolves around two main objectives: probing the Standard Model with increasing precision and searching for new physics. Both tasks require the identification of rare signals in immense backgrounds. Substantially increased levels of pile-up 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 estimate 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 data.

1.2 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*). Specific application areas of machine learning in HEP are detailed in Chapter 2.

Machine learning methods are designed to exploit large datasets in order to reduce complexity and identify new features in data. The 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 *model* is *trained* for *classification or regression* using signal and background events (or *instances*). As typical for all ML algorithms, training is the most human- and CPU- time consuming, while the application of the model to classification of new events, the so called *inference* stage, is relatively inexpensive. BDTs and NNs are typically used for particle identification, or in analysis to identify specific final states. They are also used for *regression*, where a continuous function is learned, for example to obtain the best estimate of a particle energy based on the measurements from several detectors.

Neural Networks have been used in HEP for some time; however, improvements in training algorithms and computing power have led in the last decade to the so-called Deep Learning revolution, which has made 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 deep learning neural networks used in HEP: fully-connected (FCN), convolutional (CNN) and recurrent (RNN). Additionally, neural networks are used in the context of Generative Models, when a Neural Network is trained to mimic multidimensional distributions to generate any number of new instances. Variational AutoEncoders and more recent Generative Adversarial Networks are two examples of such Generative Models used in HEP.

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

99 1.3 Structure of the Document

100 Applications of machine learning algorithms motivated by HEP drivers are detailed in Section 2. Section 3
101 focuses on the machine learning software in HEP and discusses the interplay between internally and externally
102 developed machine learning tools. Recent progress in machine learning was made possible in part by emergence
103 of suitable hardware for training complex models. In Section 4, the resource requirements of training and
104 applying machine learning algorithms in HEP are discussed. Section 6 focuses on outreach and collaboration
105 with the machine learning community. Section 5 discusses ways of training the HEP community in machine
106 learning. Section 7 presents the roadmap for the near future.

107 2 Machine Learning Applications and R&D

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

112 2.1 Detector Simulation

113 Particle discovery relies on the ability to accurately compare the observed detector response data with expecta-
114 tions based on the hypotheses of the Standard Model or models of new physics. While the processes of subatomic
115 particle interactions with matter are known, it is intractable to compute the detector response analytically. As
116 a result, Monte Carlo simulation tools, such as GEANT [1], have been developed to simulate the propagation
117 of particles in detectors to compare with the data.

118 For the HL-LHC, on the order of trillions of simulated collisions are needed in order to achieve the statistical
119 accuracy of the simulations to perform precision hypothesis testing. However, such simulations are highly
120 computationally expensive. For example, simulating the detector response of a single LHC proton-proton
121 collision event take on the order of several minutes [Citation]. Particularly time consuming is the simulation of
122 particles incident on the dense material of a calorimeter, the detector used to measure the energy deposited by
123 the particles. Radiative and nuclear interactions result in the production of a multitude of secondary particles,
124 collectively referred to as a shower. The high interaction probability and resulting high multiplicity of particles
125 passing through the dense material, make the simulation of such processes highly expensive. This problem is
126 further compounded when the particle showers overlap, as in the core of a jet of particles produced by high
127 energy quarks and gluons.

128 Fast simulation is the process of replacing the slowest components of the simulation chain with computational
129 efficient approximations. Often such approximations have been done by simplified parameterizations or particle
130 shower look-up tables. These are computationally fast but often suffer from insufficient accuracy for high
131 precision physics measurements and searches.

132 Recent progress in high fidelity fast generative models, such as Generative Adversarial Networks (GANs)
133 and Variational Auto-encoders, which are able to sample high dimensional feature distributions by learning
134 from existing data samples, offer a promising alternative for simulation.

135 A simplified first attempt at using such techniques saw orders of magnitude increase in simulation speed
136 over existing fast simulation techniques [2], but has not yet reached the required accuracy. Developing these
137 techniques for realistic detector models and understanding how to reach the required accuracy is still needed.
138 The fast advancement in the ML community of such techniques makes this a highly promising avenue to pursue.

139 Although fast simulation is a necessity, some data analyses will require the highest fidelity simulations from
140 GEANT. There is a large number of parameters that can be used to tune various aspects of the simulation
141 properties. Performing such tuning over high dimensional parameter space is highly non-trivial. Again, machine
142 learning may offer a solution. Modern optimization techniques, such as Bayesian Optimization, allow global
143 optimization of the simulator without the detailed knowledge of its internal details. Applying such techniques
144 to simulation tuning may further improve the output of the simulations.

145 2.2 Real Time Analysis and Triggering

146 The traditional approach to data analysis in particle physics assumes that the interesting events recorded by a
147 detector can be selected in real-time (a process known as “triggering”) with a reasonable efficiency, and that once
148 selected, these events can be affordably stored and distributed for further selection and analysis later. However,
149 the enormous production cross-section and luminosities of the LHC mean that these assumptions break down.¹

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

150 In particular there are whole classes of events, for example beauty and charm hadrons or low-mass dark matter
151 signatures, which are so abundant that it is not affordable to store all the events for later analysis. In order to
152 fully exploit the physics reach of the LHC, it will increasingly be necessary to perform more of the data analysis
153 in real-time.

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

162 Real-time analysis poses specific challenges to machine learning algorithm design, in particular how to
163 maintain insensitivity to detector performance which may vary over time. For example, the LHCb experiment
164 uses neural networks for fast fake-track and clone rejection. It will be important that these approaches maintain
165 performance for higher detector occupancy for the full range of tracks used in physics analyses. Another related
166 application is speeding up the reconstruction of beauty, charm, and other lower mass hadrons, where traditional
167 track combinatorics and vertexing techniques may become too computationally expensive.

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

173 2.3 Object Reconstruction, Identification, and Calibration

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

180 Particles are observed in a detector through the energy they deposit when traversing material, which is
181 subsequently digitized. Reconstruction is the process of converting the raw digital signals in the detector into
182 the physical properties of particles. Particle Physics detectors are usually composed of several sub-detectors,
183 each taking advantage of specific interaction mechanism to detect passage of a specific type of particle and
184 measuring its properties. There is a variety of sub-detector technologies, but most belong to one of three
185 categories:

- 186 • Tracking Detectors: These detectors measure the trajectory of charge particles by spatially locating ioniza-
187 tion. Usually trackers are placed in a magnetic field, so that the particle momentum can be inferred from
188 the curvature of the trajectory. Very precise tracking detectors, such as those that employ silicon, provide
189 sufficient spatial resolution to enable locating the particle creation and/or decay point. The ionization
190 also allows identifying the particle type
- 191 • Calorimeters: These detectors measure the particle energy by causing them to interact and lose the energy
192 in material and counting secondary particles. Highly segmented calorimeters measure the profile of the
193 energy deposition and identify the particle type
- 194 • Particle Identification: These detectors are aimed at determining a specific particle type using a variety
195 of techniques

196 Algorithmic reconstruction typically involves several steps that turn the data from the detector electronics
197 (*raw* measurements) into higher level data objects, corresponding to the physical particles that were detected
198 (*features*):

- 199 • Feature Extraction: the signal from the passage of particles through a detector element, e.g. a calorimeter
200 cell, is observed above noise in the raw electronic output associated with the element. This signal is then
201 characterized
- 202 • Pattern Recognition: The pattern of signals in geometrically adjacent detector elements is associated
203 with the passage of a signal or group of particles. In calorimeters, this step is commonly referred to as
204 clustering.

- Object Characterization: Properties of the objects are measured. In tracking detectors, this step means fitting a pattern of “hits” to a helix. In calorimeters, this step extracts the energy, location, and other properties of the cluster that for example characterize the shape of the cluster.
- Combined reconstruction: Objects in different are associated together to create a particle candidate.

Machine learning can in principle be applied at any of these steps. For example, 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 been able to extract better performance with deep neural networks.

An active area of research is performing particle identification and extracting particle properties on the output of feature extraction with DNNs, in particular 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 a 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 photon 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 is 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 intensity. A current effort called HEP.TrkX investigates deep learning algorithms such as long-term short-term (LSTM) networks for track pattern recognition on many-core processors.

2.4 Sustainable Matrix Element Method

The Matrix Element (ME) Method [3–6] 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 SM measurements and Higgs boson searches [7–12] and at the LHC for measurements in the Higgs and top quark sectors of the SM [13–19]. The ME method is based on *ab initio* calculation 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 ξ with theory parameters α . One can compute $\mathcal{P}_\xi(\mathbf{x}|\alpha)$ 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

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

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, $\sigma_\xi^{\text{fiducial}}(\alpha)$ is the total cross section for the process ξ (with α) times the detector acceptance, $d\Phi(\mathbf{y})$ is the phase space density factor, $\mathcal{M}_\xi(\mathbf{y}|\alpha)$ 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 Eqn. 1 in a number of ways to search for new phenomena at particle colliders. For measurement of model parameters α , one would maximize the likelihood function for observed events $\mathcal{L}(\alpha)$ given by

$$\mathcal{L}(\alpha) = \prod_i \sum_k f_k \mathcal{P}_{\xi_k}(\mathbf{x}_i|\alpha) \quad (2)$$

where f_k are the fractions of (non-interfering) processes contributing to the data. For new particle searches, one can (using Bayes’ Theorem [20]) compute for a hypothesized signal S the probability $P(S|\mathbf{x})$ given by

$$P(S|\mathbf{x}) = \frac{\sum_i \beta_{S_i} \mathcal{P}_{S_i}(\mathbf{x}|\alpha_{S_i})}{\sum_i \beta_{S_i} \mathcal{P}(\mathbf{x}|\alpha_{S_i}) + \sum_j \beta_{B_j} \mathcal{P}(\mathbf{x}|\alpha_{B_j})} \quad (3)$$

where, S_i and B_j , denote all signal and background processes relevant to the considered phase space and β are the *a priori* expected process fractions. According to the Neyman-Pearson Lemma [21], Eqn. 3 is the optimal discriminant function for S in the presence of B and can be used to extract a signal fraction in the data.

The ME method brings in 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

239 of a hypothesized process, including all correlations, and (3) has a clear physical meaning in terms transition
240 probabilities within the framework of quantum field theory.

241 One drawback to the ME Method is that it has traditionally relied on LO matrix elements, although nothing
242 limits the ME method to LO calculations. Techniques that accomodate initial-state QCD radiation within the
243 LO ME framework using transverse boosting and dedicated transfer functions to integrate over the transverse
244 momentum of initial-state partons have been developed [22]. Another challenge is development of the transfer
245 functions which rely on tediously hand-crafted fits to full simulated Monte-Carlo events.

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

251 Accurate evaluation of Eqn. 1 is computationally challenging for two reasons: (1) it involves high-dimensional
252 integration over a large number of events, signal and background hypotheses, and systematic variations and (2)
253 it involves sharply-peaked integrands² over a large domain in phase space. In reference to point (1), the matrix
254 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
255 are not completely measured experimentally (e.g. neutrinos), one must integrate over the missing information
256 which increases the dimensionality of the integration. In reference to point (2), a clever technique to re-map the
257 phase space in order to reduce the sharpness of integrate in that space in an automated way (MADWEIGHT [23])
258 is often used in conjunction with a matrix element calculation package (MADGRAPH_aMCNLO [24]). In prac-
259 tice, evaluation of definite integrals by the ME approach invokes techniques such as importance sampling (see
260 VEGAS [25, 26] and FOAM [27]) or recursive stratified sampling (see MISER [28]) Monte Carlo integration.
261 Acceleration of some of these techniques on modern computing architectures has been achieved, for example
262 concurrent phase space sampling in VEGAS on GPUs.

263 Despite the attractive features of the ME method and promise of further optimization and parallelization
264 of the evaluation of Eqn. 1, the computational burden of the ME technique will continue to limit its range of
265 applicability for practical data analysis without new and innovative approaches. The primary idea put forward in
266 this Section is to utilize modern *machine learning techniques to dramatically speed up the numerical evaluation*
267 *of Eqn. 1* and therefore broaden the applicability of the ME method to the benefit of the HL-LHC physics
268 program.

269 Applying neural networks to numerical integration problems is plausible but not new (see [29–31], for
270 example). The technical challenge is to design a network which is sufficiently rich to encode the complexity of
271 the ME calculation for a given process over the phase space relevant to the signal process. Deep Neural Networks
272 (DNNs) are strong candidates for networks with sufficient complexity to achieve good approximation of Eqn. 1,
273 possibly in conjunction with smart phase-space mapping such as described in [23]. Promising demonstration
274 of the power of Boosted Decision Trees [32, 33] and Generative Adversarial Neural Networks [34] for improved
275 Monte Carlo integration can be found in [35]. Once a set of DNNs representing of definite integrals of the form
276 of Eqn. 1 to good approximation are generated, evaluation of the ME method calculations via the DNNs will
277 be very fast. These DNNs can be thought of as preserving the essence of ME calculations in a way that allows
278 for fast forward execution. The net result is that the DNNs can enable the ME method to be both *nimble* and
279 *sustainable*, neither of which is true today.

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

293 There are several activities which are proposed to further develop the idea of a Sustainable Matrix Element
294 Method. The first is to establish a cross-experiment group interested in developing the ideas presented in this
295 Section, along with a common software project for ME calculations in the spirit of [36]. This area is very
296 well-suited for impactful collaboration with computer scientists and those working in machine learning. Using a
297 few test cases (e.g. $t\bar{t}$ or $t\bar{t}h$ production), evaluation of DNN choices and configurations, developing methods for

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

298 DNN training from full ME calculations and direct comparisons of the integration accuracy between Monte Carlo
299 and DNN-based calculations should be undertaken. More effort should also be placed in developing compelling
300 applications of the ME method for HL-LHC physics. In the longer term, the possibility of Sustainable-Matrix-
301 Element-Method-as-a-Service (SMEMaaS) where shared software and infrastructure could be used through a
302 common API, is proposed.

303 2.5 Learning the Standard Model

304 New physics may manifest itself as unusual or rare events. One approach is to accurately identify the Standard
305 Model processes and search for anomalies. Classifying the Standard Model events is a challenging task, as it
306 consists of many complicated physics processes. Multi-class machine learning algorithms are well-suited for this
307 classification problem. Once an event is classified as likely a known physics process it can be filtered out and
308 remaining events can be further analyzed for hints of new physics. Additionally, unsupervised machine learning
309 techniques can be applied to remaining events to cluster them together. This approach would also be useful in
310 identifying detector problems.

311 2.6 Theory Applications

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

316 Making progress towards the objectives of the HL-LHC program (see section 1) requires not only obtaining
317 the experimental measurements of the physical processes but also reliable theory inputs to compare to. This
318 becomes increasingly challenging as the experimental data gets more precise. There are numerous examples of
319 phenomenologically relevant processes where the experimental uncertainty is comparable to the estimate of the
320 theoretical uncertainty of the corresponding calculation.

321 Furthermore, the theory does not predict the value of all the inputs required for the computations (for
322 example the value of the strong coupling constant α_S at the Z mass), and there are situations where the
323 equations resulting from theory cannot be solved to describe the Physics adequately, and the corresponding
324 theory inputs must be obtained from data instead. A more complex example is the determination of Parton
325 Distribution Functions (PDFs): Quantum Chromodynamics (QCD) describes the proton collisions at high
326 energy in terms of *partons* (e.g. quarks and gluons), but it is not possible to calculate directly from QCD the
327 momentum carried by each quark or gluon within a proton since QCD is not solvable in its confined regime. Our
328 lack of theoretical knowledge about the characterization of partons within a proton is embedded into a suitable
329 definition of the Parton Distribution Functions (approximately the momentum densities of each of the partons)
330 The PDFs then need to be determined from experimental data. The NNPDF collaboration uses Machine
331 Learning techniques to obtain a PDF determination that is accurate enough to be suitable for high-precision
332 collider data comparison. The NNPDF fitting procedure is described in full details in [37].

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

339 While the dataset is small, each experimental point has an indirect relation to the PDFs, as it is the result of
340 the convolution of one or two PDFs with the corresponding partonic cross section. Code has been developed to
341 compute these convolutions APFELgrid [38]. Future research directions include the possibility of using standard
342 ML frameworks to express efficiently the PDF fitting problem. The uncertainties of the theory calculations need
343 to be taken into account as well in the fits. A fully systematic treatment of theory errors in PDFs is a topic
344 of research where Machine Learning could play an important role. The dominant uncertainties in the data are
345 no longer statistical and instead arise from correlated systematics. Determining those systematics accurately
346 is non-trivial on the side of the experimental analyses and can have a major impact on the resulting PDFs.
347 The problem grows more complex when ML techniques for which there is no simple recipe to estimate the
348 uncertainty are used extensively in the experimental analysis. Taking full advantage of these advanced methods
349 requires interdisciplinary research and communication on topics such as developing regularization schemes for
350 experimental covariance matrices.

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

353 2.7 Monitoring Detectors, Hardware Anomalies and Preemptive Maintenance

354 Data-taking of current complex HEP detectors is continuously monitored by physicists taking shifts to monitor
355 the quality of the incoming data. Typically, hundreds of histograms have been defined by experts and shifters
356 are alerted when an unexpected deviation with respect to a reference occurs. It regularly happens that a new
357 type of problem is unseen in a timely manner because it has not been foreseen by the expert.

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

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

366 2.8 Computing Resource Optimization and Control of Networks and Production 367 Workflows

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

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

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

383 Networks are going to play a crucial role in data exchange and data delivery to scientific applications in HL-
384 LHC era. The network-aware application layer and configurations may significantly affect experiment's daily
385 operations. ML applications can be in network security in identifying anomalies in network traffic; predicting
386 network congestion; bug detection via analysis of self-learning networks, and WAN path optimization based on
387 user access patterns.

388 2.9 Collaborative Benchmark Datasets

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

393 These benchmark datasets could be built based on public simulation engine, or released by experiments
394 within the bounds of their data access policy. Even a small subset of an experiment simulated data can be
395 the base of a very valuable benchmark dataset. For example, the CMS experiment has released a significant
396 amount of its simulated and collected data via the CERN Open Data Portal [**Opendata:2017**].

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

- 398 • Simplify the dataset as much as possible
- 399 • Document the dataset to make it understandable by a non-HEP expert
- 400 • Create methodology and metrics for evaluation proposed solutions, and document them.
- 401 • Prepare an integration plan for incoming ideas and solutions
- 402 • Feedback results of successful applications

403 **3 Machine Learning Software and Tools**

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

408 **3.1 Software Methodology**

409 Presently, there are two machine learning software methodologies in high-energy physics. The first approach
410 focuses on HEP-developed ML toolkits, such as the Toolkit for Multivariate Analysis (TMVA) in ROOT, while
411 the second approach relies on externally developed software, of which there are many examples. Historically,
412 a variety of approaches and competition among them has led to important breakthroughs in the field. On the
413 other hand, having too many choices increases repetition and leads community segmentation and possible issues
414 with reproducibility.

415 **3.2 I/O and Programming Languages**

416 The sheer amounts of data accumulated by HEP experiments require a close look at data access optimization.
417 To train and applying ML techniques on these data, efficient I/O becomes critical, especially for training. I/O
418 performance is very dependent on data formats. Moreover, support for reading data in different formats is
419 required for certain use-cases.

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

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

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

427 **3.3 Software Interfaces to Acceleration Hardware**

428 Modern machine learning software significantly benefits from using hardware accelerators such as GPUs. At
429 the same time, ML users should not be forced to write platform-dependent code. Various interfaces to different
430 hardware architectures are needed in order to make efficient use of provided computing resources. Emergence of
431 the Open Computing Language (OpenCL) allows programming of high-level interfaces that can run on various
432 hardware platforms.

433 Machine learning tools often provide different sets of APIs to develop and train the models in one language,
434 and various bindings to use trained models in another programming languages. This is a convenient model for
435 many HEP applications, such as the trigger system, where application latency puts stringent requirements on
436 the software and hardware used.

437 **3.4 Parallelization and Interactivity**

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

443 Often one needs to produce many different machine learning models, for example while tuning hyper-
444 parameters or performing k-fold cross-validation, and distribution of these algorithms is key to the reduction of
445 the overall training time.

446 ML algorithm inference significantly benefits from parallelization as well. For example, in particle physics
447 trigger systems, the stringent latency requirements impose constraints on the type of algorithms that can be
448 easily parallelized in the hardware.

449 Availability of interactive frameworks, for example Jupyter notebooks, allows for rapid prototype develop-
450 ment and testing of ML tools. Such frameworks also ease the connection between the description of models and
451 the data, providing straightforward means of visualizing models and data. HEP has started to exploring inter-
452 active frameworks, such as the Service for Web Based Analysis (SWAN). One of the challenges is availability of
453 adequate hardware resources for these systems.

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 conversion is from one data-format to another is straightforward. The following notations 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 Learn	R	Spark ML	VW	libFM	RGF	Torch
ROOT [T , C]	\checkmark									
CSV [F]		\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\times	\times	\times	\checkmark
libSVM [M]							\times	\checkmark	\times	
VW [M]							\checkmark			
RGF [M]									\checkmark	
NumPy [R]		\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\times	\times	\times	\checkmark
Avro [S , R]					\checkmark	\checkmark				
Parquet [S , C]					\checkmark	\checkmark				
HDF5 [S]										\checkmark
R df [R]					\checkmark					

3.5 Internal and external ML tools

Internally developed tools, such as the Toolkit for Multivariate Analysis (TMVA), have been developed to apply a variety of machine learning algorithms to HEP challenges. Currently, most published HEP analyses with machine learning have made use of TMVA. There are also tools developed in HEP, such as NeuroBayes 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 is 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.

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.

3.5.1 Machine Learning Data Formats

Unfortunately, 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: ROOT file format is very flexible, though requires a significant investment to properly use. Table 1 summarizes the current machine learning toolkits and file formats they use.

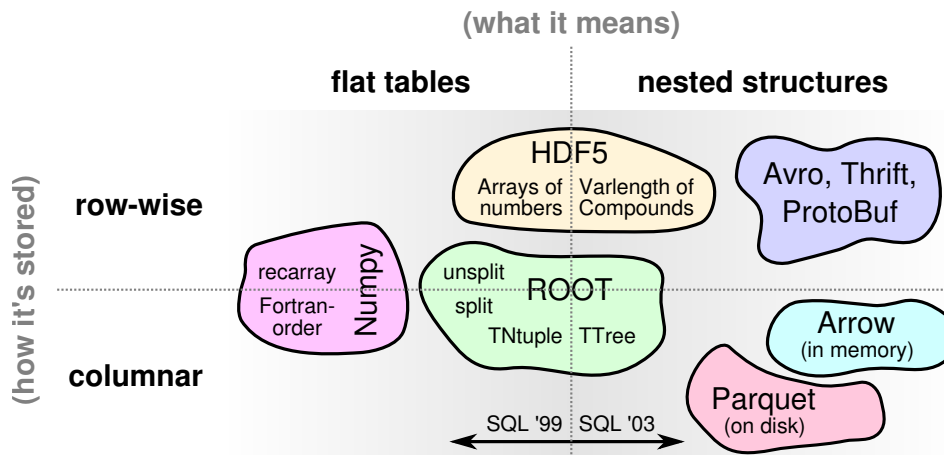


Figure 1: Existing data-formats used by ML communities

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

PyROOT	Python extension module that allows the user to interact with ROOT data/classes [39].
root_numpy	The interface between ROOT and NumPy supported by the Scikit-HEP community [40].
c2numpy	Pure C-based code to convert ROOT data into Numpy arrays which can be used in C/C++ frameworks [41].
root4j	The hep.io.root package contains a simple Java interface for reading ROOT files. This tool has been developed based on freehep-rootio [42].
root2npy	The go-hep contains a reading ROOT files. This tool has been developed based on freehep-rootio [42].
root2hdf5	Converts ROOT files containing TTrees into HDF5 files containing HDF5 tables [43].

472 3.5.2 Desirable HEP-ML software and data format attributes

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

476 HEP machine learning applications require high performance and flexible algorithms to address the variety
477 of use cases. Some applications, such as triggering, also have to work under tight latency constraints of the
478 order of a few microseconds and below. The data sets are extremely large, which comes with I/O challenges
479 described in section 3.2. This is expected to become even more challenging, as the LHC continues to ramp-up
480 and deliver increasingly large amounts of data.

481 As discussed in section 3.3, Machine Learning tools use a number of languages. To use them it will be
482 important to offer adequate support. C++ converters or similar tools are also needed to make sure the training
483 result can be efficiently evaluated.

484 Advantages of using the **external tools** are the size of the community that uses and supports them, being
485 able to easily keep up with progress in the industry and profit from the forefront of the ML research.

486 It should also be noted that some of the recent industrial efforts to develop and maintain ML-tools rely on
487 resources far beyond that of basic research. The deep learning tools of the previous and current generations
488 constitute a demonstration of corresponding quality.

489 A disadvantage of using external tools: too many choices that are not guaranteed to be supported over the
490 lifetime of particle physics experiments, difficulty of adaptation to HEP specific requirements not among the
491 priorities of the ML community.

492 Advantages of using **internal tools** are decisions about long-term support remain in the community, can be
493 adapted to specific needs of HEP. Disadvantages include challenges in incorporating new algorithms and ideas
494 on a timely basis and possible lack of resources for long-term maintenance.

495 3.5.3 Interfaces and middleware

496 One approach to bridge the gap between internal and external tools is by building interfaces. A number of
497 interfaces have already been built between TMVA and external machine learning tools, allowing for their use
498 and direct comparison between their performance. Currently, interfaces to R, scikit-learn, keras and tensorflow
499 have been developed.

500 Another approach is building middleware solutions that export HEP-specific formats like ROOT to formats
501 used by external machine learning tools. Existing middleware solutions are shown in Table 2.

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

505 4 Computing and Hardware Resources

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

513 At present, training of machine learning algorithms is done using dedicated or private resources. These vary
514 in configuration and processing power, depending on the size of the data and complexity of the algorithm. For
515 a given event, the evaluation of algorithms is performed on a single core producing a single discriminator or
516 regressor output. In order to progress to evaluation of complex machine learning, more computing power is
517 needed in both the training and evaluation stages, as larger amounts of data are needed to feed models with tens
518 or hundreds of thousands of parameters. This implies the expansion of the current computing model to include
519 architectures that are well suited to machine learning tasks, such as many integrated core (MIC), graphical
520 processing units (GPUs) and tensor processing units (TPUs). This is a fundamental departure from the single-
521 core or few-core jobs. These architectures provide a significant computational speed improvement for both
522 training and evaluation of ML algorithms, but require dedicated hardware, drivers, and software configuration.
523 Similarly, the locality and bandwidth of large data stores will need to be optimized in order to avoid
524 bottlenecks in training and evaluation for analysis. Data placement and the need to use dedicated hardware
525 indicate that a transition to HPC, or HPC-like, architectures may be needed to achieve the desired performance.
526 Due to significant synergy with the direction of industry in this respect, use of commercially available resources
527 should be considered for future high-energy physics computing models.

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

531 4.1 Resource Requirements

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

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

540 4.2 Resource Evaluation

541 It can be considered that a typical HEP application would require up to 1 GPU-week to train a single model.
542 To obtain the best results and understand the performance of the model, an average of 100 hyper-parameter
543 points optimization is typically performed. A single project could therefore easily require up to a full GPU-year
544 for training.

545 The speed up of the training process can be obtained by means of faster and more capable hardware,
546 parallelization of single training and over multiple nodes. Resources, such as GPUs, TPUs and MIC need to be
547 evaluated in the context of realistic benchmark particle physics applications.

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

551 4.2.1 High Performance Computing

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

559 4.2.2 Field Programmable Gate Arrays

560 Field Programmable Gate Arrays (FPGAs) allow an efficient and low-latent application of machine learning
561 algorithms directly at the level of hardware, as desirable for the high-energy physics trigger systems. The following
562 ML algorithms are more suitable for FPGAs due to their simpler parallelization: boosted decision trees, random
563 forests and decision rule ensembles. For example the CMS experiment currently uses boosted decision trees in
564 FPGAs in the trigger system to estimate muon momenta. Further research and development is needed in this
565 area to apply more advanced machine learning techniques like deep learning directly in the hardware. One of the
566 challenges is the limited availability of floating point operations gates and the precision needed to maintain the

567 best performance. The possibility of coupling the FGPA's with a CPU with significant random-access memory
568 (RAM) allows the shift of some of these operations to RAM.

569 **4.2.3 Opportunistic Resources**

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

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

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

581 **4.2.4 Data Storage and Availability**

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

587 Availability of data at PByte/EByte scale represents another challenge for ML community. A good solution
588 must provide access to a large data volume for hundred or thousand of users simultaneously. As discussed in
589 the Data Storage chapter [reference], data movement might need to be automatized to make the training data
590 available transparently at high speed local storage with use of automatic caching mechanism. The success of
591 Apache Spark and Google BigQuery platforms may serve as a model. Data streaming, transformation and
592 readout in mini-batches may be required to train models over large data sets. This is in addition to the regular
593 HEP workflows described in the Data Storage and Networking chapter [reference]

594 **4.2.5 Software Distribution and Deployment**

595 To efficiently use the resources described in previous subsections, machine learning software needs to be available
596 on the computing resources. Platforms, such as CERNVM File System (cvmfs), are very useful for software
597 distribution that does not require local installation. Additional tools like docker containers for application
598 shipping can be useful in providing homogeneous software environments across the different systems. Another
599 challenge is the software layer needs to be agnostic to the hardware back-end.

600 **4.2.6 Machine Learning As a Service**

601 Current cloud providers rely on machine learning as a service model allowing for efficient use of common
602 resources and use interactive machine learning tools. Machine Learning As a Service is not yet widely used
603 in HEP, but examples of successful publications which used Machine Learning As a Service exist, e.g. [44].
604 Specialized HEP services for interactive analysis, such as CERN's Service for Web-based Analysis (SWAN) may
605 play an important role in adoption of machine learning tools in HEP workflows. In order to use these tools
606 more efficiently, sufficient and appropriately tailored hardware resources described in this chapter are needed.

607 **5 Training the community**

608 In order to address the communication barrier and to speak the same language, the HEP community should
609 be trained in ML concepts and terminology as part of a standard curriculum. The training should focus on
610 well-maintained and well-documented software packages. It should provide lectures on general ML concepts
611 and hands-on tutorials on specific tools based on concrete examples.

612 Being able to apply machine learning to practical HEP problems requires the understanding of basic ML
613 concepts and algorithms. For this, regular data science lecture series and seminars are very useful. At the
614 University level, courses dedicated to machine learning applications in physics research is an excellent way to
615 train undergraduate and graduate students. For example, "Deep Learning in Physics Research" course with
616 60 participants consisting of 12 lectures and exercises which are performed on 20 GPU's of the VISPA internet
617 platform [45].

618 Experiments currently have training activities for newcomers that focus on analysis software and introduction
619 to domain knowledge. Machine learning should next be incorporated into the incoming collaborators training
620 efforts of the experiments.

621 As discussed in the training chapter, ensuring the development and availability of resources for knowledge
622 transfer is likewise essential to ML.

623 **6 Collaborating with other communities**

624 **6.1 Introduction**

625 Discovery science provides a challenge that attracts brilliant minds eager to push the boundaries of scientific
626 understanding of nature. Particle physics has a rich problem domain that offers avenues for intellectual reward.
627 The goal is to achieve vibrant collaboration between data science and high-energy physics communities by
628 finding a common language and working together to further science.

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

636 There are a number of existing examples of collaboration between HEP and ML that have produced fruitful
637 results through mostly local connections (e.g. [2, 46]). The HEP community should continue such collaborations
638 and look for additional collaborations with ML.

639 Domain knowledge can present a barrier to collaboration. The HEP community needs to define its chal-
640 lenges in a language that the ML community can understand. This may involve stripping the domain knowledge
641 entirely, or retaining necessary information with clear and concise explanations as to its relevance. Machine
642 learning likewise has a significant amount of domain knowledge. Ideas and solutions provided by both commu-
643 nities should be presented in an understandable way for scientists without in-depth knowledge.

644 **6.2 Machine Learning Challenges**

645 To engage the wider ML community, challenges such as the Higgs Boson Challenge (2014) or the Flavor Physics
646 Challenge [47, 48] have been organized on Kaggle. These types of challenges draw considerable attention from
647 the Machine Learning community and more of such challenges should be organized in the future.

648 Organization of a challenge requires a well documented dataset, a starting-kit and an evaluation metric to
649 rank the solutions. This forces the organizers to simplify the problem as much as possible, while retaining its
650 intrinsic complexity.

651 The drawback of challenges is that once they are launched, participants priority is winning the challenge and
652 not eventual collaboration with HEP. It is important to foresee upfront a way to integrate incoming solutions, for
653 example via forums and post-challenge workshops where a diversity of competitive algorithms can be presented.
654 The challenge dataset and evaluation metric should be released publicly so that further developments can
655 continue.

656 **6.3 Collaborative Benchmark Datasets**

657 As discussed in section 2.9, collaborative benchmark datasets can be useful for developing ML challenges and
658 for collaboration with the ML community. The HEP community should organize and curate a variety of such
659 benchmark datasets covering its current physics drivers and make it publicly available. To improve reproducibil-
660 ity of results and algorithm comparisons, some of the data used for evaluation of the solutions should be kept
661 private.

662 Additionally, after investing heavily into producing highly-detailed and realistic simulations, the HEP com-
663 munity can provide the machine learning community with labeled datasets with high statistical power to test
664 algorithms and develop novel ideas.

665 **6.4 ML Academic outreach**

666 Conferences and workshops are a core aspect of the academic ML community, and organizing or contributing
667 to key conferences is a means of gaining interest. Organizing sessions or mini-workshops within major ML
668 conferences, such as NIPS, would increase the familiarity of HEP within the ML community and jump-start

669 future collaborations. This has been explored in single cases [47] but is not an established, regular workshop
670 series. At the same time inviting ML experts to HEP workshops as done at [48] and the DS@HEP series [49–51],
671 can foster greater long-term collaboration.

- 672 • Organize workshops and conferences open to external collaborators to discuss the applications, algorithms
673 and tools
- 674 • Organize thematic workshops around topics relevant to HEP

675 6.5 Science outreach

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

679 6.6 Industry Engagement

680 Industry has been focused on development and adoption of machine learning techniques. In addition to
681 algorithm and software development, one of the promising areas is the adoption of dedicated specialized hardware
682 and high performance co-processors. GPUs, FPGAs, and high core count co-processors all have the potential
683 to dramatically increase performance of machine learning applications relevant to HEP applications.

684 One of the challenges is gaining the human expertise for development and implementation. Industry brings
685 specific technology opportunities and access to specialized expertise that can be difficult to hire and support
686 internally.

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

694 6.6.1 CERN OpenLab and research-industry collaborative initiatives

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

701 6.7 ML community at large outreach

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

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

- 709 • [Linear Digressions](#) (co-hosted by former ATLAS Ph.D. Katie Malone)
- 710 • [Partially Derivative](#)
- 711 • [Talking Machines](#)
- 712 • [Data Skeptic](#)
- 713 • [Becoming a Data Scientist Podcast](#)
- 714 • [Not So Standard Deviations](#)
- 715 • [This Week in ML & AI](#)

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

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

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

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

727 7 Roadmap

728 The incorporation of ML into particle physics experiments must respect two time lines: the schedule of LHC
729 and funding agencies, and the experiments' need for extensive validation of the algorithms.

730 The current LHC schedule has Run 3 starting in 2021 and the HL-LHC, if approved, starting in 2026. As
731 software processes and algorithms are re-imagined, their implementation must fit into these time frames if they
732 are to maximize their benefit to the physics. To fit this schedule, a newly proposed implementation would need
733 to show a demonstration in 2018 to prove viability. Two years later, in 2020, the idea needs to attain a level of
734 maturity to be included in the HL-LHC Technical Design Report. The project should then be further refined
735 towards a large scale test around the middle of Run 3, about 2022. Run 3 is scheduled to end in late 2023. The
736 project must then be adapted to the HL-LHC software and physics analysis environment as it will be relied on
737 by the experiment.

738 The path of taking a ML idea from conception to community-wide acceptance and deployment will entail
739 several stages, as appropriate. There are ample opportunities to make the process more efficient. For example,
740 in many steps having common datasets, as discussed in Section 2.9 will likely accelerate the progress.

- 741 1. Problem formulation and dataset preparation: Problem formulation is the first step in building an ML
742 algorithm. The inputs and desired output need to be established. The training and validation datasets
743 must be identified and simulated. In many cases, these datasets are large, and resources must be identified
744 to possibly create and store the data. In most cases, the data needs to be processed into a form suitable for
745 input into the algorithm. Since these steps are often lengthy, common datasets with well-defined problems
746 will be very helpful
- 747 2. Feasibility/Demonstration: Given a dataset, appropriate ML algorithms need to be investigated and
748 evaluated for ability to solve the problem. In some cases, such studies can be preformed on simplified
749 datasets
- 750 3. First application: An application of the solution to one or few specific physics analysis where the ML tech-
751 nique significantly improves the physics result. Here the incorporation of the technique into the computing
752 workflow will likely be very specific to the application and require significant manual intervention
- 753 4. Scaling/Optimization: Evolving from a demonstration to a general solution requires use of realistic
754 datasets with full detector simulation, noise, etc. Furthermore, the solution will also require optimization
755 to achieve nominal physics and computing performance. A good practice would be to apply the solution to
756 a specific physics analysis. This stage will likely require significant computing resources to scale solutions
757 to the full detector and datasets
- 758 5. Integration/Validation: The solution needs to be incorporated into the experimental software and workflow
759 and validated.

760 As an example, consider the simulation physics driver. An effort has recently started to build generative
761 models that can significantly accelerate simulation of particle showers in calorimeters. These early efforts are
762 based on simplified datasets specifically created for this problem, without the complications of realistic data
763 and limited to a small section of calorimeters. The first papers[2] use GANs to generate calorimetric data
764 which are reasonably faithful, but still require tuning. The next step involves exploration of DNN architecture
765 and systematic hyper-parameter scans on HPCs to achieve the required performance. The technique can be
766 applied to searches at LHC that involve boosted objects, where the required simulation samples require CPU

767 intensive full GEANT-based simulation and are therefore limited in statistics due to resource limitations. The
768 process of employing the new technique in a publication will illicit scrutiny by the full experiment, effectively
769 validating the technique. Once the technique is accepted, it can be generalized beyond this first application
770 and then incorporated into the experiment's software for use by others. Finally, as the technique is applied to
771 an increasing number of physics analyses, the technique will be incorporated into the experiment's production
772 workflows.

773 **A Author list**

- 774 • Aaron Sauers
- 775 • Aashrita Mangu (CS)
- 776 • Adam Aurisano (NOvA)
777 University of Cincinnati (US)
- 778 • Adrian Bevan (ATLAS)
779 University of London (GB)
- 780 • Alessandra Forti (ATLAS)
781 University of Manchester (GB)
- 782 • Alexander Kurepin (ALICE)
783 Russian Academy of Sciences (RU)
- 784 • Alexander Radovic (NOvA)
785 College of William and Mary (GB)
- 786 • Alexei Klimentov (ATLAS)
787 Brookhaven National Laboratory (US)
- 788 • Amir Farbin (ATLAS)
789 University of Texas at Arlington (US)
- 790 • Andrey Ustyuzhanin (Yandex, LHCb)
791 Yandex School of Data Analysis (RU)
- 792 • Antonio Limosani (ATLAS)
793 University of Melbourne (AU)
- 794 • Ariel Schwartzman (ATLAS)
795 SLAC National Accelerator Laboratory (US)
- 796 • Attilio Picazio (ATLAS)
797 University of Massachusetts (US)
- 798 • Aurelius Rinkevicius (CMS)
799 Cornell University (US)
- 800 • Ben Hooberman (ATLAS)
801 University of Illinois at Urbana-Champaign (US)
- 802 • Benedikt Hegner (SFT)
803 CERN
- 804 • Bob Stienen
805 Radboud Universiteit Nijmegen (NL)
- 806 • Claire David (ATLAS)
807 Deutsches Elektronen-Synchrotron (DESY) (DE)
- 808 • Conor Fitzpatrick (LHCb)
809 Ecole Polytechnique Federale de Lausanne (CH)
- 810 • Daniele Bonacorsi (CMS)
811 Universita e INFN, Bologna (IT)
- 812 • Dario Menasce (CMS and INFN)
813 Universita & INFN, Milano
814 Bicocca (IT)
- 815 • David Rousseau (ATLAS)
816 Universite de Paris
817 Sud 11 (FR)

- 818 • Dick Greenwood (ATLAS)
819 Louisiana Tech University (US)
- 820 • Dorian Kcira (CMS)
821 California Institute of Technology (US)
- 822 • Douglas Davis
- 823 • Dustin Anderson (CMS)
824 California Institute of Technology (US)
- 825 • Eduardo Rodrigues (LHCb)
826 University of Cincinnati (US)
- 827 • Elias Coniavitis (ATLAS)
828 Universitaet Freiburg (DE)
- 829 • Federico Carminati (SFT)
830 CERN
- 831 • Fernanda Psihas (NOvA)
832 Indiana University (US)
- 833 • Filip Siroky (CMS)
834 Masaryk University (Czech Republic)
- 835 • Gabriel Perdue (MINERvA)
836 Fermilab (US)
- 837 • Gaurav Kaul (Intel)
- 838 • Giles Strong (CMS)
839 LIP-Lisbon (Portugal)
- 840 • Gilles Louppe (ATLAS)
841 New York University (US)
- 842 • Gordon Watts (ATLAS)
843 University of Washington (US)
- 844 • Graeme Stewart (ATLAS)
845 University of Glasgow (Scotland)
- 846 • Hans Pabst (Intel)
- 847 • Harvey Newman (CMS)
848 California Institute of Technology (US)
- 849 • Helge Meinhard (CERN)
- 850 • Horst Severini (ATLAS)
851 University of Oklahoma (US)
- 852 • Ian Stockdale
- 853 • Igor Lakomov (ALICE, CERN)
- 854 • Ilija Vukotic (ATLAS)
855 University of Chicago (US)
- 856 • Jamal Rorie (CMS)
857 Rice University (US)
- 858 • Javier Duarte (CMS)
859 California Institute of Technology (US)
- 860 • Jean-Roch Vlimant (CMS)
861 California Institute of Technology (US)
- 862 • Jim Kowalkowski (Fermilab)

- 863 • Jim Pivarski (CMS)
864 Princeton University (US)
- 865 • Jochen Gemmler (Belle2)
866 KIT/IEKP (DE)
- 867 • Johannes Junggeburth
868 Max Planck Institut für Physik (DE)
- 869 • John Harvey (CERN)
- 870 • Jonas Eschle (LHCb)
871 Universität Zürich (CH)
- 872 • Jonas Graw
- 873 • Jordi Garra-Tico (LHCb)
874 University of Cambridge (GB)
- 875 • Juan Pedro Araque Espinosa (ATLAS)
876 LIP Lisboa (Portugal)
- 877 • Karen Tomko (Ohio Supercomputer Center, US)
- 878 • Kevin Lannon (CMS)
879 University of Notre Dame (US)
- 880 • Kim Albertsson (ATLAS)
881 Lulea University of Technology (Sweden)
- 882 • Konstantin Kanishchev (AMS-02)
883 Università INFN, Padova (Italy)
- 884 • Konstantin Skazytkin (ALICE)
885 Russian Academy of Sciences (RU)
- 886 • Kyle Cranmer (ATLAS)
887 New York University
- 888 • Laurent Basara (RE1)
889 Università INFN, Padova (Italy)
- 890 • Lindsey Gray (CMS)
891 Fermilab (US)
- 892 • Lorenzo Moneta (ROOT, CMS)
893 CERN
- 894 • Louis Capps
- 895 • Lukas Heinrich (ATLAS)
896 New York University
- 897 • Luke Kreczko (CMS, LZ)
898 University of Bristol (UK)
- 899 • Maria Girone (CERN openlab, CERN)
- 900 • Mario Campanelli (ATLAS)
901 University of London (UK)
- 902 • Mario Lassnig (ATLAS)
903 CERN
- 904 • Mark Neubauer (ATLAS)
905 University of Illinois at Urbana-Champaign (US)
- 906 • Martin Erdmann (CMS)

- 907 • Martin Vala (ALICE)
908 Technical University of Kosice (Slovakia)
- 909 • Matthew Feickert (ATLAS)
910 Southern Methodist University (US)
- 911 • Mauro Verzetti (CMS)
912 University of Rochester (US)
- 913 • Meghan Kane (SoundCloud, formerly @MIT)
- 914 • Michael Andrews (CMS)
915 Carnegie-Mellon University (US)
- 916 • Michael Kagan (ATLAS)
917 SLAC (US)
- 918 • Michael Williams (LHCb)
919 MIT (US)
- 920 • Michela Paganini (ATLAS)
921 Yale University (US)
- 922 • Michele Floris (ALICE)
923 CERN
- 924 • Mike Sokoloff (LHCb)
925 University of Cincinnati (US)
- 926 • Nicolas Köhler
927 Max-Planck-Institut für Physik (DE)
- 928 • Nuno Filipe Castro (ATLAS)
929 LIP-Lisbon (Portugal)
- 930 • Paolo Calafiura (ATLAS)
931 Lawrence Berkeley National Lab (US)
- 932 • Paul Glaysner (ATLAS)
933 Deutsches Elektronen-Synchrotron (DESY) (DE)
- 934 • Paul Seyfert (LHCb)
935 CERN
- 936 • Pere Mato (SFT, LHCb)
937 CERN
- 938 • Piero Altoe (Nvidia)
- 939 • Przemysław Karpiński (CERN openlab)
940 CERN
- 941 • Rob Kutschke (Mu2e, Intensity Frontier)
942 Fermilab (US)
- 943 • Ryan Reece (ATLAS)
944 University of California, Santa Cruz (US)
- 945 • Savannah Thais (ATLAS)
946 Yale University (US)
- 947 • Sean-Jiun Wang (CMS)
948 University of Florida (US)
- 949 • Sergei Gleyzer (CMS)
950 University of Florida (US)
- 951 • Seth Moortgat (CMS)
952 Vrije Universiteit Brussel (Belgium)

- 953 • Sofia Vallecorsa (SFT)
954 Gangneung-Wonju National University (South Korea)
- 955 • Stefan Wunsch (CMS)
956 KIT - Karlsruhe Institute of Technology (DE)
- 957 • Stefano Carrazza (CERN)
- 958 • Steven Schramm (ATLAS)
959 Université de Genève (CH)
- 960 • Taylor Childers (ATLAS)
961 Argonne National Laboratory (US)
- 962 • Thomas Keck (Belle2)
963 KIT - Karlsruhe Institute of Technology (DE)
- 964 • Tom Hacker
- 965 • Uzziel Perez (CMS)
- 966 • Valentin Kuznetsov (CMS)
967 Cornell University (US)
- 968 • Vladimir Vava Gligorov (LHCb)
969 Centre National de la Recherche Scientifique (FR)
- 970 • Wahid Bhijmi (Daya-Bay)
- 971 • Wenjing Wu (ATLAS)
972 Chinese Academy of Sciences (China)
- 973 • Xavier Vilasís-Cardona (LHCb)
974 University of Barcelona (Spain)
- 975 • Omar Zapata (<http://oproject.org>)
- 976 • Zahari Kassabov
977 University of Turin
978 University of Milan

References

- [1] S. Agostinelli et al. “GEANT4: A Simulation toolkit.” *Nucl. Instrum. Meth.* A506 (2003), pp. 250–303. DOI: [10.1016/S0168-9002\(03\)01368-8](https://doi.org/10.1016/S0168-9002(03)01368-8).
- [2] 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](https://arxiv.org/abs/1705.02355) [[hep-ex](#)].
- [3] 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](https://doi.org/10.1143/JPSJ.57.4126).
- [4] 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](https://doi.org/10.1016/j.nima.2010.09.024). arXiv: [1003.1316](https://arxiv.org/abs/1003.1316) [[hep-ex](#)].
- [5] I. Volobouev. “Matrix Element Method in HEP: Transfer Functions, Efficiencies, and Likelihood Normalization.” *ArXiv e-prints* (Jan. 2011). arXiv: [1101.2259](https://arxiv.org/abs/1101.2259) [[physics.data-an](#)].
- [6] 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](https://doi.org/10.1103/PhysRevD.96.015021). arXiv: [1705.02563](https://arxiv.org/abs/1705.02563) [[hep-ph](#)].
- [7] 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](https://doi.org/10.1038/nature02589). arXiv: [hep-ex/0406031](https://arxiv.org/abs/hep-ex/0406031) [[hep-ex](#)].
- [8] 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](https://doi.org/10.1103/PhysRevD.75.031105). arXiv: [hep-ex/0612060](https://arxiv.org/abs/hep-ex/0612060) [[hep-ex](#)].
- [9] 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](https://doi.org/10.1103/PhysRevLett.100.201801). arXiv: [0801.4806](https://arxiv.org/abs/0801.4806) [[hep-ex](#)].
- [10] 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](https://doi.org/10.1103/PhysRevLett.104.061803). arXiv: [1001.4468](https://arxiv.org/abs/1001.4468) [[hep-ex](#)].
- [11] V. M. Abazov et al. “Observation of Single Top Quark Production.” *Phys. Rev. Lett.* 103 (2009), p. 092001. DOI: [10.1103/PhysRevLett.103.092001](https://doi.org/10.1103/PhysRevLett.103.092001). arXiv: [0903.0850](https://arxiv.org/abs/0903.0850) [[hep-ex](#)].
- [12] 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](https://doi.org/10.1103/PhysRevLett.103.092002). arXiv: [0903.0885](https://arxiv.org/abs/0903.0885) [[hep-ex](#)].
- [13] S. Chatrchyan et al. “Observation of a new boson at a mass of 125 GeV with the CMS experiment at the LHC.” *Phys. Lett.* B716 (2012), pp. 30–61. DOI: [10.1016/j.physletb.2012.08.021](https://doi.org/10.1016/j.physletb.2012.08.021). arXiv: [1207.7235](https://arxiv.org/abs/1207.7235) [[hep-ex](#)].
- [14] 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](https://doi.org/10.1103/PhysRevD.89.092007). arXiv: [1312.5353](https://arxiv.org/abs/1312.5353) [[hep-ex](#)].
- [15] 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](https://doi.org/10.1103/PhysRevD.91.012006). arXiv: [1408.5191](https://arxiv.org/abs/1408.5191) [[hep-ex](#)].
- [16] 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.* B758 (2016), pp. 321–346. DOI: [10.1016/j.physletb.2016.05.005](https://doi.org/10.1016/j.physletb.2016.05.005). arXiv: [1511.06170](https://arxiv.org/abs/1511.06170) [[hep-ex](#)].
- [17] 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](https://doi.org/10.1140/epjc/s10052-015-3454-1). arXiv: [1502.02485](https://arxiv.org/abs/1502.02485) [[hep-ex](#)].
- [18] 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](https://doi.org/10.1140/epjc/s10052-015-3543-1). arXiv: [1503.05066](https://arxiv.org/abs/1503.05066) [[hep-ex](#)].
- [19] 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.* B756 (2016), pp. 228–246. DOI: [10.1016/j.physletb.2016.03.017](https://doi.org/10.1016/j.physletb.2016.03.017). arXiv: [1511.05980](https://arxiv.org/abs/1511.05980) [[hep-ex](#)].
- [20] M. Bayes and M. Price. “An Essay towards Solving a Problem in the Doctrine of Chances. By the Late Rev. Mr. Bayes, F. R. S. Communicated by Mr. Price, in a Letter to John Canton, A. M. F. R. S.” *Philosophical Transactions* 53 (1763), pp. 370–418. DOI: [10.1098/rstl.1763.0053](https://doi.org/10.1098/rstl.1763.0053). eprint: <http://rstl.royalsocietypublishing.org/content/53/370.full.pdf+html>. URL: <http://rstl.royalsocietypublishing.org/content/53/370.short>.

- 1032 [21] J. Neyman and E. S. Pearson. “On the Problem of the Most Efficient Tests of Statistical Hypotheses.”
1033 *Philosophical Transactions of the Royal Society of London A: Mathematical, Physical and Engineering*
1034 *Sciences* 231.694-706 (1933), pp. 289–337. ISSN: 0264-3952. DOI: [10.1098/rsta.1933.0009](https://doi.org/10.1098/rsta.1933.0009). eprint:
1035 <http://rsta.royalsocietypublishing.org/content/231/694-706/289.full.pdf>. URL: [http://](http://rsta.royalsocietypublishing.org/content/231/694-706/289)
1036 rsta.royalsocietypublishing.org/content/231/694-706/289.
- 1037 [22] J. Alwall, A. Freitas, and O. Mattelaer. “The Matrix Element Method and QCD Radiation.” *Phys. Rev.*
1038 *D*83 (2011), p. 074010. DOI: [10.1103/PhysRevD.83.074010](https://doi.org/10.1103/PhysRevD.83.074010). arXiv: [1010.2263](https://arxiv.org/abs/1010.2263) [hep-ph].
- 1039 [23] P. Artoisenet et al. “Automation of the matrix element reweighting method.” *JHEP* 12 (2010), p. 068.
1040 DOI: [10.1007/JHEP12\(2010\)068](https://doi.org/10.1007/JHEP12(2010)068). arXiv: [1007.3300](https://arxiv.org/abs/1007.3300) [hep-ph].
- 1041 [24] J. Alwall et al. “The automated computation of tree-level and next-to-leading order differential cross
1042 sections, and their matching to parton shower simulations.” *JHEP* 07 (2014), p. 079. DOI: [10.1007/](https://doi.org/10.1007/JHEP07(2014)079)
1043 [JHEP07\(2014\)079](https://doi.org/10.1007/JHEP07(2014)079). arXiv: [1405.0301](https://arxiv.org/abs/1405.0301) [hep-ph].
- 1044 [25] G. P. Lepage. “A new algorithm for adaptive multidimensional integration.” *Journal of Computational*
1045 *Physics* 27.2 (1978), pp. 192–203. ISSN: 0021-9991. DOI: [http://dx.doi.org/10.1016/0021-9991\(78\)](https://dx.doi.org/10.1016/0021-9991(78)90004-9)
1046 [90004-9](https://dx.doi.org/10.1016/0021-9991(78)90004-9). URL: <http://www.sciencedirect.com/science/article/pii/0021999178900049>.
- 1047 [26] T. Ohl. “Vegas revisited: Adaptive Monte Carlo integration beyond factorization.” *Comput. Phys. Com-*
1048 *munic.* 120 (1999), pp. 13–19. DOI: [10.1016/S0010-4655\(99\)00209-X](https://doi.org/10.1016/S0010-4655(99)00209-X). arXiv: [hep-ph/9806432](https://arxiv.org/abs/hep-ph/9806432) [hep-ph].
- 1049 [27] S. Jadach. “Foam: A general-purpose cellular Monte Carlo event generator.” *Computer Physics Com-*
1050 *munications* 152.1 (2003), pp. 55–100. ISSN: 0010-4655. DOI: [http://dx.doi.org/10.1016/S0010-](https://dx.doi.org/10.1016/S0010-4655(02)00755-5)
1051 [4655\(02\)00755-5](https://dx.doi.org/10.1016/S0010-4655(02)00755-5). URL: <http://www.sciencedirect.com/science/article/pii/S0010465502007555>.
- 1052 [28] W. H. Press and G. R. Farrar. “RECURSIVE STRATIFIED SAMPLING FOR MULTIDIMENSIONAL
1053 MONTE CARLO INTEGRATION.” *Submitted to: Comp.in Phys.* (1989).
- 1054 [29] Z. Zhe-Zhao, W. Yao-Nan, and W. Hui. “Numerical integration based on a neural network algorithm.” 8
1055 (Aug. 2006), pp. 42–48.
- 1056 [30] L. y. Xu and L. j. Li. “The New Numerical Integration Algorithm Based on Neural Network.” *Third*
1057 *International Conference on Natural Computation (ICNC 2007)*. Vol. 1. Aug. 2007, pp. 325–328. DOI:
1058 [10.1109/ICNC.2007.730](https://doi.org/10.1109/ICNC.2007.730).
- 1059 [31] L. Yan, J. Di, and K. Wang. “Spline Basis Neural Network Algorithm for Numerical Integration.” *Interna-*
1060 *tional Journal of Mathematical, Computational, Physical, Electrical and Computer Engineering* 7.3 (2013),
1061 pp. 458–461. ISSN: PISSN:2010-376X, EISSN:2010-3778. URL: <http://waset.org/Publications?p=75>.
- 1062 [32] J. Friedman, T. Hastie, and R. Tibshirani. “Additive logistic regression: a statistical view of boosting
1063 (With discussion and a rejoinder by the authors).” *Ann. Statist.* 28.2 (Apr. 2000), pp. 337–407. DOI:
1064 [10.1214/aos/1016218223](https://doi.org/10.1214/aos/1016218223). URL: <http://dx.doi.org/10.1214/aos/1016218223>.
- 1065 [33] J. H. Friedman. “Greedy function approximation: A gradient boosting machine.” *Ann. Statist.* 29.5 (Oct.
1066 2001), pp. 1189–1232. DOI: [10.1214/aos/1013203451](https://doi.org/10.1214/aos/1013203451). URL: [http://dx.doi.org/10.1214/aos/](http://dx.doi.org/10.1214/aos/1013203451)
1067 [1013203451](http://dx.doi.org/10.1214/aos/1013203451).
- 1068 [34] I. J. Goodfellow et al. “Generative Adversarial Networks.” *ArXiv e-prints* (June 2014). arXiv: [1406.2661](https://arxiv.org/abs/1406.2661)
1069 [stat.ML].
- 1070 [35] J. Bendavid. “Efficient Monte Carlo Integration Using Boosted Decision Trees and Generative Deep Neural
1071 Networks” (2017). arXiv: [1707.00028](https://arxiv.org/abs/1707.00028) [hep-ph].
- 1072 [36] *MeMEMta: Modular Matrix Element Implementation*. URL: <https://github.com/MoMEMta>.
- 1073 [37] R. D. Ball et al. “Parton distributions for the LHC Run II.” *JHEP* 04 (2015), p. 040. DOI: [10.1007/](https://doi.org/10.1007/JHEP04(2015)040)
1074 [JHEP04\(2015\)040](https://doi.org/10.1007/JHEP04(2015)040). arXiv: [1410.8849](https://arxiv.org/abs/1410.8849) [hep-ph].
- 1075 [38] V. Bertone, S. Carrazza, and N. P. Hartland. “APFELgrid: a high performance tool for parton density
1076 determinations.” *Comput. Phys. Commun.* 212 (2017), pp. 205–209. DOI: [10.1016/j.cpc.2016.10.006](https://doi.org/10.1016/j.cpc.2016.10.006).
1077 arXiv: [1605.02070](https://arxiv.org/abs/1605.02070) [hep-ph].
- 1078 [39] *PyROOT*. URL: <https://root.cern.ch/pyroot>.
- 1079 [40] *root numpy converter*. URL: https://github.com/scikit-hep/root%5C_numpy.
- 1080 [41] *c2numpy*. URL: <https://github.com/diana-hep/c2numpy>.
- 1081 [42] *root4j*. URL: <https://github.com/diana-hep/root4j>.
- 1082 [43] *root2hdf5*. URL: <http://www.rootpy.org/commands/root2hdf5.html>.
- 1083 [44] R. Aaij et al. “Search for the lepton flavour violating decay $\tau^- \rightarrow \mu^- \mu^+ \mu^-$.” *JHEP* 02 (2015), p. 121.
1084 DOI: [10.1007/JHEP02\(2015\)121](https://doi.org/10.1007/JHEP02(2015)121). arXiv: [1409.8548](https://arxiv.org/abs/1409.8548) [hep-ex].

- 1085 [45] URL: <https://vispa.physik.rwth-aachen.de>.
- 1086 [46] T. Likhomanenko, D. Derkach, and A. Rogozhnikov. “Inclusive Flavour Tagging Algorithm.” *J. Phys.*
1087 *Conf. Ser.* 762.1 (2016), p. 012045. DOI: [10.1088/1742-6596/762/1/012045](https://doi.org/10.1088/1742-6596/762/1/012045). arXiv: [1705.08707](https://arxiv.org/abs/1705.08707)
1088 [[hep-ex](https://arxiv.org/abs/1705.08707)].
- 1089 [47] *ALEPH Workshop @ NIPS 2015 – Applying (machine) Learning to Experimental Physics (ALEPH) and*
1090 *“Flavours of Physics» challenge”*. Neural Information Processing Systems (NIPS). 2015. URL: [http:](http://yandexdataschool.github.io/aleph2015/)
1091 [//yandexdataschool.github.io/aleph2015/](http://yandexdataschool.github.io/aleph2015/).
- 1092 [48] *Heavy Flavour Data Mining Workshop*. Zurich, 2016. URL: <https://indico.cern.ch/event/433556/>.
- 1093 [49] *Data Science @ LHC 2015 Workshop*. 2015. URL: <https://indico.cern.ch/event/395374/>.
- 1094 [50] *DS@HEP at the Simons Foundation*. 2016. URL: [https://indico.hep.caltech.edu/indico/conferenceDisplay.](https://indico.hep.caltech.edu/indico/conferenceDisplay.py?confId=102)
1095 [py?confId=102](https://indico.hep.caltech.edu/indico/conferenceDisplay.py?confId=102).
- 1096 [51] *DS@HEP*. 2017. URL: <https://indico.fnal.gov/conferenceDisplay.py?confId=13497>.