

## 2 Learning About Machine Learning With Dihiggs 3 Production at the LHC

---

4 B. Tannenwald<sup>a</sup> C. Neu,<sup>a</sup> A. Li,<sup>a</sup> G. Buehlmann,<sup>a</sup> A. Cuddeback,<sup>a</sup> L. Hatfield,<sup>a</sup> R.  
5 Parvatam,<sup>a</sup> C. Thompson<sup>a</sup>

6 <sup>a</sup> *University of Virginia, 248 McCormick Road, Charlottesville, VA, USA*

7 *E-mail:* [benjamin.tannenwald@cern.sh](mailto:benjamin.tannenwald@cern.ch)

8 ABSTRACT: Many new domains of high energy physics analysis are starting to explore  
9 machine learning techniques. Powerful methods can be used to identify and measure rare  
10 processes from previously insurmountable backgrounds. One of the most profound Standard  
11 Model signatures still to be discovered at the LHC is the pair production of Higgs bosons  
12 through the Higgs self-coupling. The small cross section of this process makes detection very  
13 difficult even for the decay channel with the largest branching fraction ( $hh \rightarrow b\bar{b}b\bar{b}$ ). This  
14 paper benchmarks a variety of approaches (boosted decision trees, neural networks with  
15 straightforward and novel architectures, semi-supervised algorithms) against one another  
16 to catalogue the various techniques available to high energy physicists as the era of the  
17 HL-LHC approaches.

---

18 **Contents**

19	<b>1 Introduction</b>	<b>1</b>
20	<b>2 Dihiggs Physics</b>	<b>2</b>
21	2.1 Double Higgs Production	2
22	2.2 Event Reconstruction	3
23	<b>3 Supervised Learning</b>	<b>4</b>
24	3.1 Boosted Decision Tree	4
25	3.2 Random Forest	5
26	3.3 Feed Forward Neural Network	6
27	3.4 Convolutional Neural Network	7
28	3.5 Energy Flow Network	10
29	<b>4 Semi-Supervised Learning</b>	<b>12</b>
30	4.1 $k$ -Means Clustering	13
31	4.2 Autoencoder	13
32	<b>5 Results</b>	<b>15</b>
33	<b>6 Conclusions</b>	<b>15</b>

---

34 **1 Introduction**

35 The use of machine learning (ML) techniques in high energy particle physics has rapidly  
36 increased since its first use at the Tevatron [1]. The proliferation of techniques and applica-  
37 tions has touched nearly every segment of analysis and reconstruction [2] and will be vital  
38 to understand the full dataset of the LHC and data from future colliders.

39 Common approaches involve using linear techniques like decision trees as well as non-  
40 linear approaches like neural networks. These techniques are then used to reconstruct  
41 objects like leptons and jets, to tag objects like b-quarks or boosted decays, and classify  
42 different processes. Many models are built with the same kinematic input features physicists  
43 typically use while other architectures rely on emergent features produced in a more abstract  
44 phase-space. Regardless of the approach, the rise of so many approaches raises interesting  
45 questions about what types of information are best to feed to our networks– things that  
46 are best for physicists to learn from might not be optimal for sophisticated computing  
47 algorithms.

48 The goal of this paper is to explore a wide range of current ML techniques for appli-  
49 cability at identifying dihiggs production at the HL-LHC. Observing dihiggs production is  
50 necessary to measure the self-coupling of the Higgs boson and fully understand the nature

of electroweak symmetry breaking. The difficulty in measuring diHiggs production lies in the tiny cross-section of even the largest branching fraction ( $hh \rightarrow b\bar{b}b\bar{b}$ ) and the relative abundance of similarly reconstructed QCD events. Section 2 deals with the physics relevant for diHiggs production and the QCD background. Sections 3 and 4 summarize the various ML methods tested, and Section 5 compares the results from the various approaches.

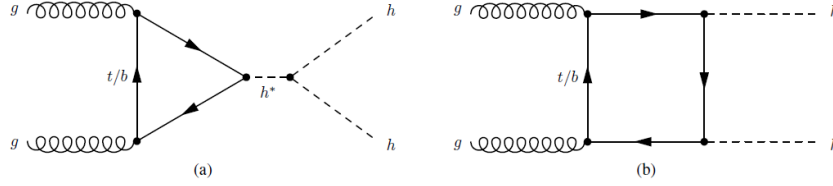
## 2 Dihiggs Physics

### 2.1 Double Higgs Production

The Higgs boson is an essential part of the Standard Model (SM) of particle physics and is a product of the mechanism responsible for electroweak symmetry breaking. Along with the interaction of the Higgs with the other particles of the Standard Model, the SM predicts the interaction of the Higgs boson with itself at tree-level (self-interaction). This mechanism contributes to non-resonant Higgs boson pair production together with quark-loop contributions via Yukawa-type interactions. Figure 1 shows a schematic diagram of non-resonant Higgs boson pair production. Since the production cross section for Higgs boson pair production is extremely small within the SM,

$$\sigma_{hh}(13 \text{ TeV}) = 33 \text{ fb},$$

any significant enhancement would indicate the presence of new physics.



**Figure 1.** Leading order Feynman diagrams for non-resonant production of Higgs boson pairs in the Standard Model through (a) the Higgs boson self-coupling and (b) the Higgs-fermion Yukawa interaction.

Many extensions of the SM predict the existence of additional scalar bosons which may have mass larger than twice the Higgs mass and can decay into a Higgs boson pair. Searching for resonances in the  $hh$  mass spectrum can help us discover or limit exotic models which predict the presence of such particles. More importantly, measuring the SM diHiggs cross-section (or placing limits on its magnitude) allow us to probe the self-coupling of the Higgs field and better understand the mechanism behind electroweak symmetry breaking.

The following work is focused on techniques for distinguishing non-resonant (SM-like) Higgs boson pair production where both Higgs bosons decay via  $h \rightarrow b\bar{b}$ . The choice of using the  $4b$  decay mode provides the largest possible amount of signal events but requires powerful background reduction techniques due to the large production cross-section of fully hadronic QCD processes. All results are quoted for simulated events produced by  $pp$  collisions with a center-of-mass energy of 14 TeV and scaled to the full design luminosity of

the HL-LHC (an integrated luminosity of  $3000 \text{ fb}^{-1}$ ). Simulated samples were produced using ROOT v6.12/04 [3] and Madgraph v2.7.0 [4]. Events were then showered using Pythia v8.2.44 [5] and reconstructed with Delphes v3.0 [6] using the v02 approximation of the upgraded Phase-II CMS detector.

Both the signal and background samples were generated with minimal pileup addition sampled from a Poisson distribution with an expectation value of zero additional vertices. An additional generator-level cut requiring total hadronic energy greater than 300 GeV was applied when generation background events. All code used to set up the generation environment and produce events is publically available [7]. A summary of the sample generation details is shown in Table 1.

Name	Process	$\sigma_{\text{eff}} [\text{fb}]$	$N_{\text{events}}$
Dihiggs	$pp > hh, h > b\bar{b}$	12.4	$1 \cdot 10^6$
QCD	$pp > b\bar{b}b\bar{b}$	441866.0	$4 \cdot 10^6$

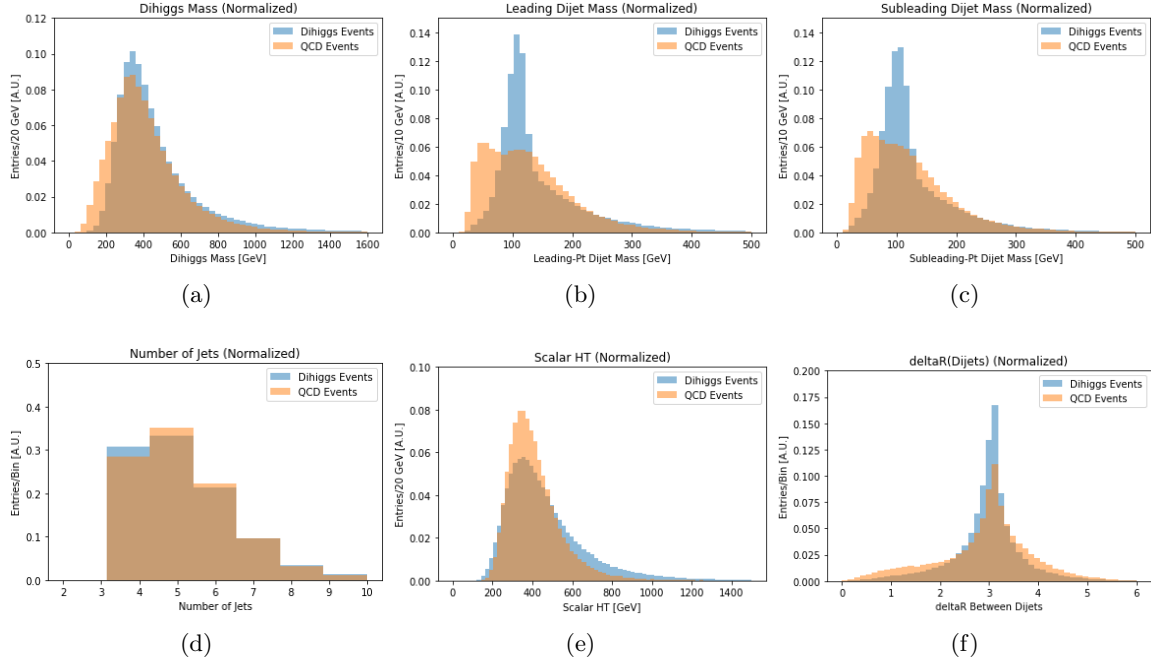
**Table 1.** Madgraph processes, effective cross-sections, and number of generated events for the signal and background samples used in this paper. The effective cross-sections differ slightly from the total theoretical cross-sections due to branching fractions and generation-level cuts on hadronic activity.

## 2.2 Event Reconstruction

The first step in reconstructing the  $4b$  system is to reconstruct and select b-jet candidates. Jets are clustered using an anti- $k_T$  algorithm with a radius of  $R=0.4$ . To be selected for use in event reconstruction, a jet must have  $p_T > 20 \text{ GeV}$  and an absolute value of  $|\eta| < 2.5$ . Delphes uses an internal  $b$ -tagging efficiency parameterization to predict whether jets are tagged, and an event is reconstructed only if at least 4 jets are  $b$ -tagged unless otherwise specified. The properties and kinematics of selected jets are shown in Fig 2. This strict requirement of having at least 4  $b$ -tags in an event helps to reduce contributions from QCD and reduce the combinatoric ambiguity in event reconstruction.

Once events with at least 4  $b$ -tags are selected, there is a choice about how to reconstruct the diHiggs system. Several reconstruction methods were tested for pairing b-jets to find an optimal algorithm for correctly pairing Higgs boson constituents. Two algorithms were selected for use in the following sections: the first iterates through all selected jets in an event and returns the two pairs with closest dijet masses to one another and the second returns the two jet pairs that minimize the difference between the individual candidate pairs and a Higgs boson mass of 125 GeV. Unless otherwise specified, the method that selects dijets with masses closest to each other is used when training techniques that require reconstructed events. Fig 2 shows a selection of distributions describing the diHiggs system using this reconstruction algorithm.

Reconstructed variables include the masses and momentum of the four- and two-body Higgs candidates as well as the angular separations between the two Higgs candidates and their constituent jets. Additional event-level variables like the number of selected jets, the number of  $b$ -tagged jets, and the missing transverse energy in the event were also con-



**Figure 2.** Sample of event kinematics and reconstructed diHiggs system from QCD and diHiggs simulation. Distributions are normalized to the same area to compare shapes.

112 sidered as inputs to various algorithms. All possible variables were evaluated using the  
 113 Kolmogorov-Smirnov (KS) test for individual separation power between signal and back-  
 114 ground. Variables were sorted in descending order of KS separability. Each algorithm is  
 115 trained on a subset that balances minimizing the number of variables without sacrificing  
 116 performance.

### 117 3 Supervised Learning

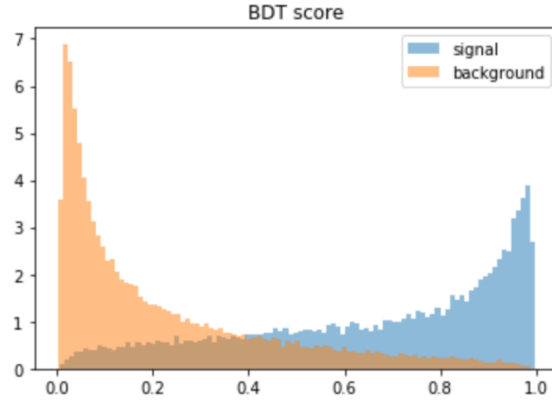
118 Searches for specific signatures or interactions in collider data can in general be thought of  
 119 as a classification problem - some known signal process must be identified and separated  
 120 from some known and well-modeled set of background processes. Any iterative algorithm  
 121 can then improve its ability to properly identify signal from background by comparing its  
 122 predictions to the true known classifications and adjusting its internal parameters. This  
 123 type of approach is known as supervised machine learning, and it is particularly relevant for  
 124 measuring diHiggs decays.

#### 125 3.1 Boosted Decision Tree

126 Boosted Decision Trees (BDTs) have a long history in high energy physics from enabling the  
 127 first observation of single top production at the Tevatron [1] to helping in the discovery of  
 128 the Higgs boson at the LHC [8, 9]. A decision tree functions by making a series of cuts (or  
 129 decisions) that maximize the separation between signal and background events in a single

dimension. Each cut produces a branch in the tree containing independent populations. The depth of the tree sets the number of decisions a tree will make, and a well-designed tree will have end-nodes that efficiently separate and properly identify the constituent classes. Any series of cuts for identifying events will inevitably misclassify some events, and there are many strategies for improving the results. A boosted decision tree attempts to improve the classification by creating a new set of data from the improperly classified events and training a new decision tree on these inputs. Each step of re-training with misclassified events is called a *boost*, and the total prediction for an event is the weighted sum of predictions from the original tree and the boosted trees where each sequential boost receives a smaller weight in the sum.

The BDT trained for diHiggs detection was built using the xgboost package [10]. The top seventeen reconstructed and event-level variables ranked by KS separability (discussed in Section 2.2) were used in training. The hyperparameters describing the boosted decision tree were optimized for maximum  $S/\sqrt{B}$ . The optimal hyperparameters were found to be as follows: multiplicative boost factor of 0.1, maximum tree depth of 9, gamma (minimum loss reduction needed for further partition) of 1.1, and an L2 regularization term of 8.28.



**Figure 3.** Signal predictions of the trained BDT for signal and background samples independent from the training sets.

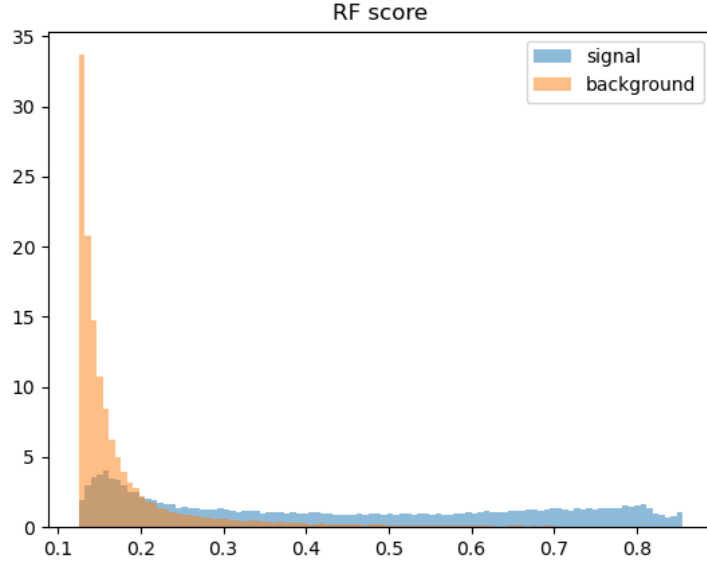
The predictions from the optimized BDT are shown in Figure 3. A maximum significance of  $1.84 \pm 0.09$  was obtained, yielding 986 signal events and  $2.8 \cdot 10^5$  background events.

### 3.2 Random Forest

Random Forest algorithms share a similar tree structure with BDTs, but they leverage ensembles of independent decision trees as opposed to iteratively improving the predictions of a single tree using misclassified events. Each tree in a random forest is ‘grown’ using a random sampling of input variables and training events. The randomness of the sampling ensures each tree yields a unique but correlated prediction compared to the other trees in the forest. The class prediction of the forest is the majority vote of the constituent trees.

156 Tuning the hyperparameters of a random forest requires optimizing the number of trees  
 157 in the forest, the variable sub-sampling used to produce each tree, and the depth of the  
 158 constituent trees.

159 The random forest trained for diHiggs classification uses the reconstruction algorithm  
 160 that selects dijet pairs consistent with a Higgs mass hypothesis, and the top seventeen  
 161 reconstructed and event-level variables were used as input. The random forest was imple-  
 162 mented using the 'XGBRFClassifier' functionality from xgboost [10]. An optimal forest was  
 163 obtained by individually varying each hyperparameter over a reasonable range and selecting  
 164 the best performing model. The optimal set of hyperparameters consisted of training with  
 165 300 constituent trees, a maximum tree depth of 20, column sub-sampling rate of 0.8, and  
 166 an L1 regularization term of 1.175. The best significance for the random forest approach  
 167 was found to be  $S/\sqrt{B} = 2.04 \pm 0.07$  when requiring a prediction score  $> 0.70$ . Prediction  
 168 results are shown in Figure 4.



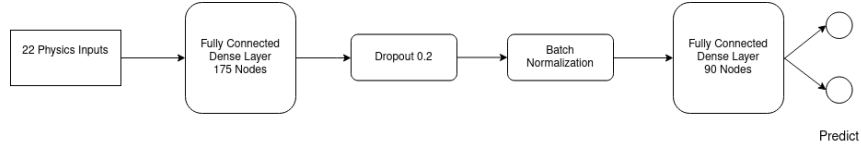
**Figure 4.** Output score on the testing dataset with the fully trained random forest classifier.

### 169 3.3 Feed Forward Neural Network

170 Fully connected or feed-forward neural networks (NN) have a long history in high energy  
 171 physics. One of the earliest applications of this type of approach was in a search for top  
 172 quark production using the CDF experiment at the Tevatron. The fundamental element  
 173 of a feed-forward neural network is called a 'layer'. Multiple layers are stacked together  
 174 connecting input variables to a predicted outcome which is evaluated against a known target  
 175 value. A fully-connected network can have multiple internal (or hidden) layers between  
 176 the input and output layers, and each hidden layer is composed of a series of trainable

177 activation functions and weights that allow the network to identify and iteratively combine  
178 important features of the input space. A vector of relevant physics-level information (e.g.  
179 mass of two highest  $p_T$  jets, angle between measured objects, etc) is constructed for each  
180 event, and these vectors are then ‘fed forward’ through multiple layers in order to predict  
181 whether the event comes from a signal diHiggs process or a background QCD process. A  
182 function (called the loss function) is chosen to quantify the difference between the model  
183 prediction and target values. The loss calculated after a single training iteration is used to  
184 adjust the internal network weights in the next training iteration through a process called  
185 backpropagation. The model is fully trained once the improvement in the loss between  
186 iterations falls beneath some user-defined threshold.

187 The NN trained for diHiggs detection was built using the Keras [11] and Tensorflow [12]  
188 packages. The top twenty-two most separable reconstructed and event-level variables are  
189 used as the input variables for the NN. The complete network structure consists of the input  
190 layer, two hidden layers, and a single-node output layer. The first hidden layer contains 175  
191 nodes with an L2 kernel regularizer ( $\lambda = 10^{-4}$ ). The second hidden layer contains 90  
192 nodes with no kernel regularizer. A batch normalization layer and a dropout (0.2) function  
193 are placed in between the two hidden layers to prevent over-fitting. Both hidden layers use a  
194 rectified linear (ReLU) activation function, while the output layer uses a sigmoid activation  
195 function. Several models were trained by individually tuning each hyperparameter over a  
196 reasonable range in order to produce a final optimized model. A schematic flowchart of the  
197 network structure is shown in Figure 5.



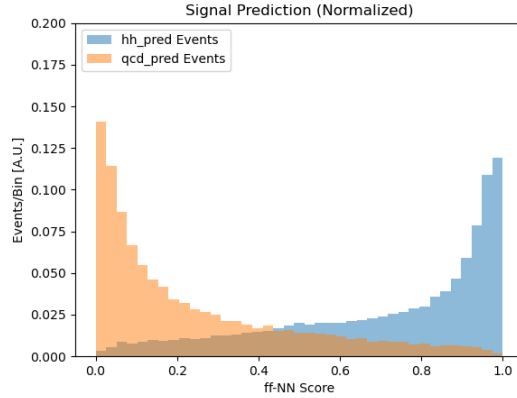
**Figure 5.** Structure of the feed-forward neural network. The input variables are fed through two fully connected dense layers to classify events. One dropout layer and one batch normalization layer help mitigate over-fitting during training.

198 The NN was trained for 25 epochs before the minimal loss-improvement threshold was  
199 met, and the results are shown in Figure 6. The trained model obtained a maximum  $S/\sqrt{B}$   
200  $= 2.2 \pm 0.07$  when considering events with a signal prediction score  $> 0.94$ . This phase-space  
201 has a signal yield of 1520.1 events and a background yield of 478872.3 events.

### 202 3.4 Convolutional Neural Network

203 Convolutional Neural Networks (CNNs) are neural networks that use assumptions about  
204 the local relationships between neighboring pixels in order to predict the content of an  
205 input image. For this analysis, content prediction is simplified to a general classification of  
206 whether the image comes from a diHiggs or QCD event. The fundamental elements of any  
207 convolutional network are convolutional layers and pooling layers. Convolutional layers  
208 use filters that perform linear combinations of neighboring pixels within the filter size,



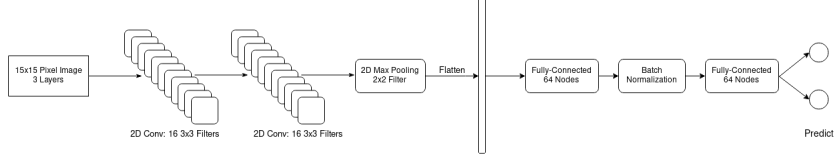


**Figure 6.** Final predictions of the feed-forward network for signal and background samples.

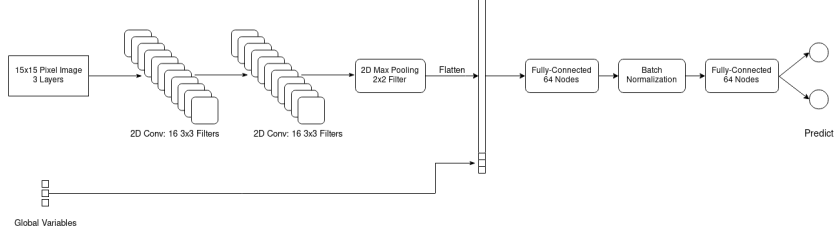
209 and pooling layers aggregate information by grouping neighboring pixels by either their  
 210 maximum or average values. After some number of these layers, the output is flattened  
 211 into a one-dimensional vector, and this flattened vector is pushed through a set of feed-  
 212 forward layers in order to make a final output prediction.

213 Many previous papers have explored the training convolutional networks using low-  
 214 level measured quantities (e.g. tracks and calorimeter deposits) in high energy physics  
 215 in the context of object identification [13]. This paper extends the application to event-  
 216 level identification. Using low-level quantities removes the need to reconstruct objects  
 217 like jets since no higher-level quantities (e.g. dijet masses) are needed; only the detector-  
 218 level measurements are required for classification. The performance of four convolutional  
 219 networks are studied in the context of diHiggs identification. The first network uses a 3-  
 220 layer image composed of energy/momentum weighted tracks, electromagnetic calorimeter  
 221 deposits, and hadronic calorimeter deposits. The second network uses the same three layers  
 222 but appends additional global information to the flattened vector after image processing  
 223 and before the fully connected layers. Figures 7 and 8 depict both network structures. The  
 224 third and fourth networks follow the same pattern as the previous two but with the addition  
 225 of two image layers corresponding to impact parameter-weighted track information (both  
 226 longitudinal and transverse).

227 In order to produce coherent images, the center of mass and the center of momentum  
 228 for each event are calculated. All constituents are then boosted longitudinally into the center  
 229 of mass of the event and rotated in phi to the center of momentum. At the end, each image  
 230 layer corresponds to a 31x31 pixel grid centered on the total activity in the event. Figure 9  
 231 shows an average QCD image and Figure 10 shows an average signal image. While the  
 232 layers in each figure closely resemble one another, they do contain different information,  
 233 and variations are visible. More importantly, the average QCD image differs significantly  
 234 from the average diHiggs image. The scalar structure of the Higgs boson leads to the more  
 235 circular forward/backward symmetry seen in Figure 10 compared to Figure 9. The variance  
 236 in pixel intensity is also smaller for the average diHiggs events reflecting the more balanced



**Figure 7.** Structure of the nominal convolutional neural network. The input images are fed through two convolutional layers and a single max-pooling layer before being flattened into a one-dimensional vector. The flattened vector is then fed through one fully connected layer, a batch normalization layer, and a final fully connected layer before a final prediction is made.



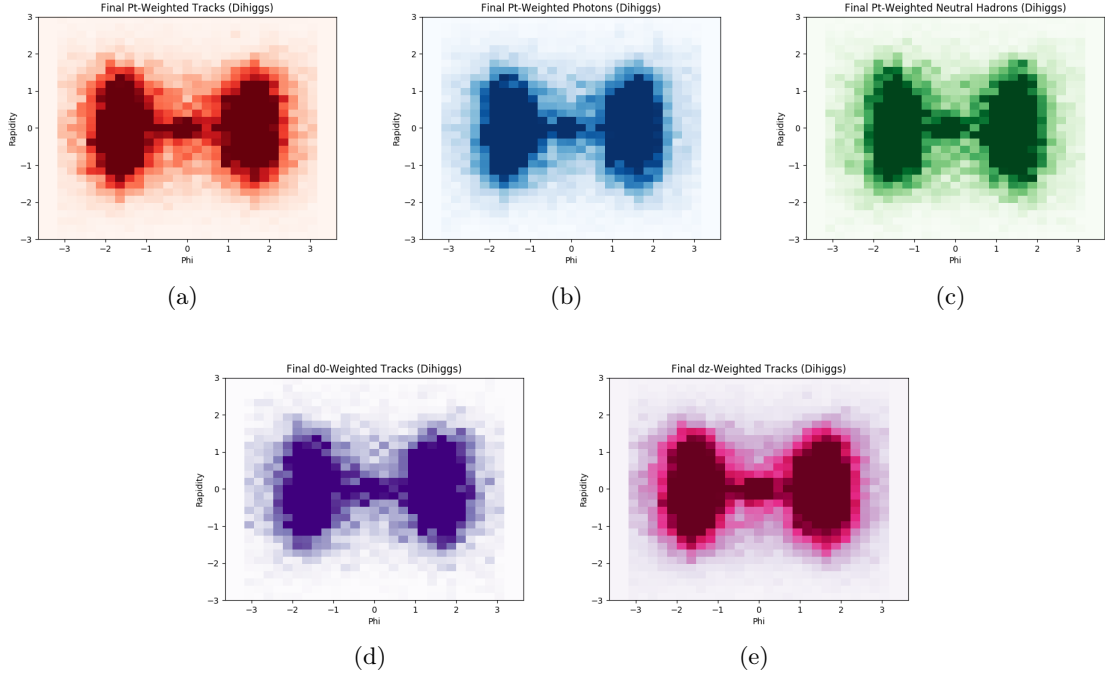
**Figure 8.** Structure of the hybrid convolutional neural network. The input images are fed through two convolutional layers and a single max-pooling layer before being flattened into a one-dimensional vector. Scaled user-specified variables (e.g.  $H_T$ ) are then concatenated with the flattened image vector. The concatenated vector is then fed through one fully connected layer, a batch normalization layer, and a final fully connected layer before a final prediction is made.

kinematics of double Higgs production.

As shown in Figures 7 and 8, the CNN network structure uses two sequential 2D convolutional layers each with 16 3x3 filters, one max-pooling layer with a 2x2 window, a flattening of the outputs, two 64-node fully connected hidden layers, and one output layer for making the final prediction. As previously described, two of the networks append additional high level variables (scalar sum of transverse hadronic energy, number of jets, and number of  $b$ -tags) after the flattening and before the image information is fed through the fully connected layers. Optimal significances for each network are shown in Table 2 and the output signal predictions and ROC for the 5-color network with high-level inputs are shown in Figure 11.

Method	Best $S/\sqrt{B}$	AUC
Tracks+HCAL+ECAL	$1.77 \pm 0.01$	0.818
Tracks+HCAL+ECAL + high-level	$2.12 \pm 0.01$	0.846
Tracks+HCAL+ECAL+D0+DZ	$2.45 \pm 0.02$	0.863
Tracks+HCAL+ECAL+D0+DZ + high-level	$2.86 \pm 0.03$	0.882

**Table 2.** Normalized to full HL-LHC dataset of  $3000 \text{ fb}^{-1}$

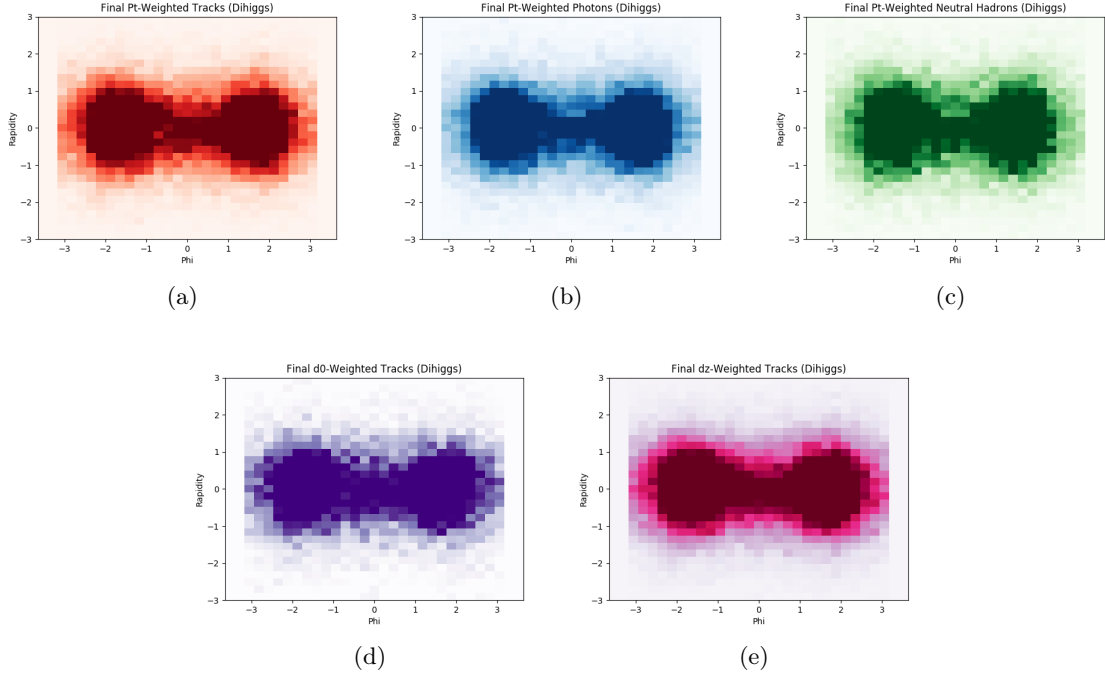


**Figure 9.** Average QCD image showing (a)  $p_T$ -weighted tracks, (b)  $E_T$ -weighted ECAL deposits, (c)  $E_T$ -weighted HCAL deposits, (d) transverse impact parameter-weighted tracks, (e) longitudinal impact parameter-weighted tracks.

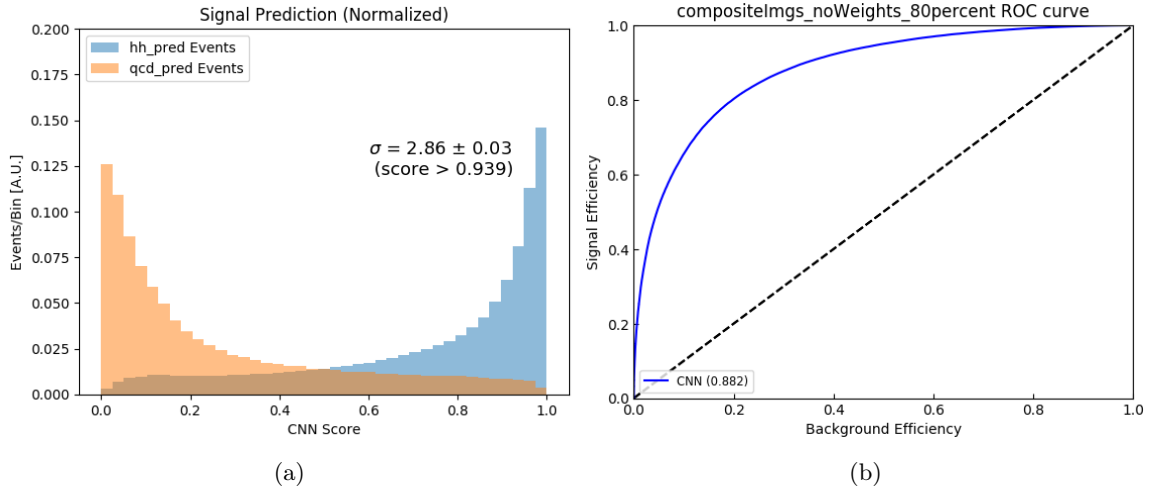
### 247 3.5 Energy Flow Network

248 Energy Flow Networks (EFN) and Particle Flow Networks (PFN) are algorithms that take  
 249 basic jet constituents information as input rather than reconstructed jets and multi-jet  
 250 composites, e.g. Higgs candidates. The EFN structure takes only the rapidity  $y$  and  
 251 azimuthal angle  $\phi$  of jet constituents as input, while the PFN takes the rapidity  $y$ , azimuthal  
 252 angle  $\phi$ , and transverse momentum  $p_T$  of jet constituents as input. Using the constituents  
 253 as input means no high level reconstruction is necessary when identifying events. Both  
 254 the EFN and PFN are two-component networks, and their internal structures are shown in  
 255 Figure 12. The implementations of the EFN and PFN used for diHiggs classification use 200  
 256 nodes for each hidden layer in network (a), 256 for latent space dimension and 300 nodes  
 257 for each hidden layer in network (b).

258 The EFN/PFN networks were trained using four separate categories split by number  
 259 of jets and number of  $b$ -tags in order to test the network's dependence on higher-level jet  
 260 information. Independent networks were trained using: all events, only events with  $\geq 4$  jets,  
 261 only events with  $\geq 4$  jets and  $=2$   $b$ -tags, and only events with  $\geq 4$  jets and  $\geq 4$   $b$ -tags. In each  
 262 configuration, the number of signal and background events were adjusted to maintain an  
 263 equal proportion of each population in the training sample. L2 regularization and dropout  
 264 layers were added to minimize overfitting. The results obtained with the EFN are shown  
 265 in Table 3. The results of the PFN are shown in Table 4.

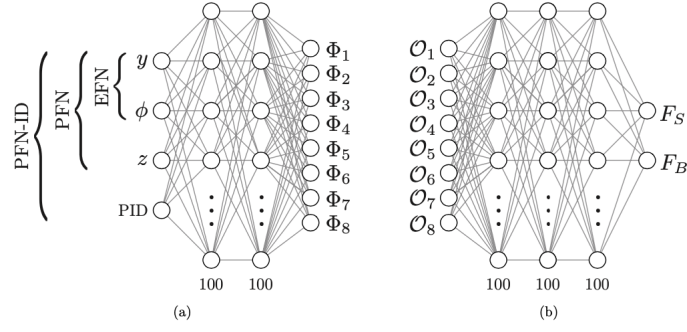


**Figure 10.** Average diHiggs image showing (a) pt-weighted tracks, (b) Et-weighted ECAL deposits, (c) Et-weighted HCAL deposits, (d) transverse impact parameter-weighted tracks, (e) longitudinal impact parameter-weighted tracks.



**Figure 11.** (Right) signal prediction for the 5-color convolutional network with additional high-level inputs, (Left) ROC for the 5-color convolutional network with additional high-level inputs.

266 The PFN provided a best significance of 1.618 when trained over all events without  
 267 any cuts on the number of jets or  $b$ -tags.



**Figure 12.** Network (a) takes jet constituents information as input and outputs latent space  $\Phi$  for each jet constituents. Network (b) takes  $\mathcal{O}$ , which is the linear combination of  $\Phi$ , as input and outputs final result.

Category	0PU		
	Best $S/\sqrt{B}$	$N_{\text{Signal}}$	$N_{\text{Background}}$
All Events	$1.407 \pm 0.006$	$1.89 \cdot 10^4$	$1.80 \cdot 10^8$
4Jets	$1.363 \pm 0.006$	$1.63 \cdot 10^4$	$1.43 \cdot 10^8$
4Jets 2BTags	$1.343 \pm 0.006$	$1.33 \cdot 10^4$	$9.95 \cdot 10^7$
4Jets 4BTags	$0.867 \pm 0.008$	3468.65	$1.60 \cdot 10^7$

**Table 3.** EFN results. Normalized to full HL-LHC dataset of  $3000 \text{ fb}^{-1}$

Category	0PU		
	Best $S/\sqrt{B}$	$N_{\text{Signal}}$	$N_{\text{Background}}$
All Events	$1.618 \pm 0.008$	$1.79 \cdot 10^4$	$1.21 \cdot 10^8$
4Jets	$1.580 \pm 0.008$	$1.32 \cdot 10^4$	$7.00 \cdot 10^7$
4Jets 2BTags	$1.574 \pm 0.009$	$1.32 \cdot 10^4$	$4.85 \cdot 10^7$
4Jets 4BTags	$0.903 \pm 0.009$	3297.34	$1.33 \cdot 10^7$

**Table 4.** PFN results. Normalized to full HL-LHC dataset of  $3000 \text{ fb}^{-1}$

## 268 4 Semi-Supervised Learning

269 While it makes sense to treat searches for new physics or rare signatures as a supervised  
270 classification problem, an alternative approach is to let an algorithm learn intrinsic features  
271 from an unlabeled dataset and then evaluate whether this self-learned information can be  
272 used to separate signal from background processes. The process of learning features from  
273 unlabelled datasets is called unsupervised learning. Unsupervised machine learning can be  
274 transformed into semi-supervised learning when the results of unsupervised learning are  
275 evaluated using known classification information.

## 276 4.1 $k$ -Means Clustering

277 K-means clustering is an unsupervised learning algorithm that finds unlabelled groupings in  
 278 the phase-space defined by the input variables. The number of clusters to fit is a user-defined  
 279 hyperparameter, and three different clusterings (15, 20, 40) were tested. Two unsupervised  
 280 transformations of the input variables were tested to try to improve the performance of  
 281 the nominal BDT. Combining the unsupervised clustering with the supervised structure  
 282 of the BDT converts the unsupervised features into a semi-supervised algorithm whose  
 283 performance can be compared to other supervised methods.

284 The first transformation was to pass the nominal kinematic inputs through a k-means  
 285 clustering stage before training the BDT. The second transformation involved performing a  
 286 principal component analysis (PCA) decomposition on the nominal kinematic inputs before  
 287 passing through the clustering step and finally the BDT. PCA is a technique for finding  
 288 an orthogonal basis of input data that minimizes the variance along each new axis. No  
 289 transformation was found to improve the performance of the nominal configuration, and  
 290 the results are shown in Table 5.

Method	$S/\sqrt{B}$	$N_{sig}$	$N_{bkg}$
Nominal BDT	$1.84 \pm 0.09$	986.3	$2.9 \cdot 10^5$
15 Clusters	$1.29 \pm 0.02$	2100.2	$2.7 \cdot 10^6$
15 Clusters + PCA	$1.25 \pm 0.02$	2189.5	$3.1 \cdot 10^6$
20 Clusters	$1.30 \pm 0.02$	2260.6	$3.0 \cdot 10^6$
20 Clusters + PCA	$1.27 \pm 0.03$	21756.4	$1.9 \cdot 10^6$
40 Clusters	$1.44 \pm 0.03$	1704.6	$1.4 \cdot 10^6$
40 Clusters + PCA	$1.34 \pm 0.02$	2144.5	$2.0 \cdot 10^6$

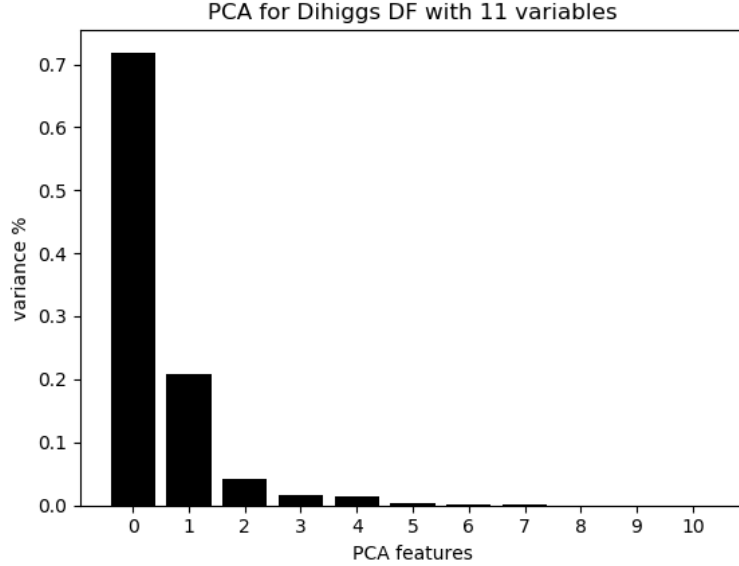
**Table 5.** Significance and yields showing BDT performance for the nominal kinematic inputs, clustered kinematic inputs, and clustered inputs from a PCA decomposition. All yields are normalized to full HL-LHC dataset of  $3000 \text{ fb}^{-1}$ .

## 291 4.2 Autoencoder

292 An autoencoder (AE) is an unsupervised machine learning architecture used for detecting  
 293 anomalies that differ significantly from the data used to train the network. The structure of  
 294 the AE compresses the input information into a lower-dimensional representation called the  
 295 latent space. In principle this compression encodes the most important features of the train-  
 296 ing data, and the second half of the network ‘decodes’ this latent representation back into  
 297 a representation approaching the original inputs. This construction fundamentally changes  
 298 the meaning of the loss function; rather than computing the loss between a prediction and  
 299 a target, the AE loss is a measure of how well the network reproduces the original inputs  
 300 after encoding and decoding. Inputs that differ significantly from the data used to train the  
 301 AE will not be properly reconstructed, and anomalies can be identified by selected events  
 302 with large losses. Monte Carlo simulations enable a semi-supervised cross-check on AE

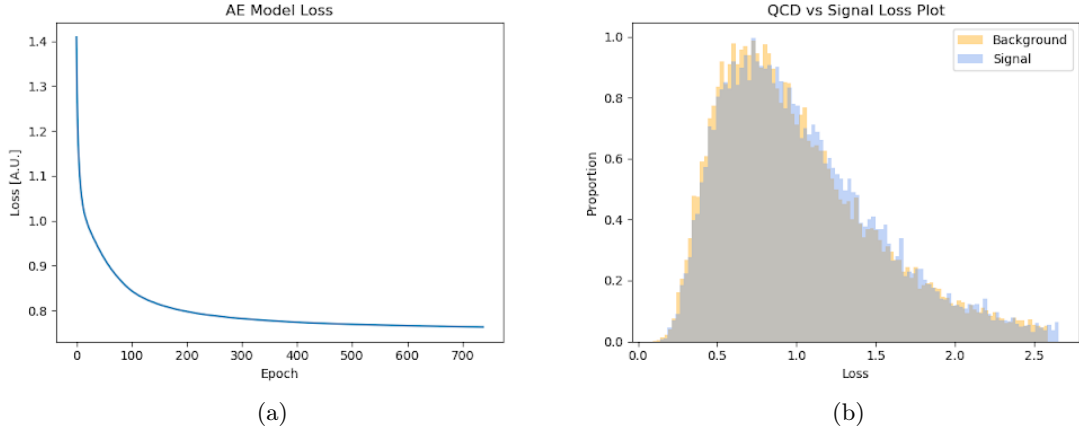
performance, meaning AEs may have exciting implications for signal detection in particle physics.

Because AE anomaly detection relies on a well-modeled understanding of all background processes, the network was first trained using only QCD events. Additional models were trained by substituting the pure-QCD training sample with training samples consisting of QCD+dihiggs mixtures in order to test the stability of the method. No significant deterioration was observed for reasonable levels of contamination. The AE used for dihiggs detection was built using the Keras package [11] and consists of an input layer, a single hidden layer, and an output layer. Eleven reconstructed variables were selected for use in the AE. The number of hidden layers and the number of hidden nodes per layer are hyperparameters that were optimized. A PCA analysis shown in Figure 13 was used to determine that a latent space of three nodes was optimal.



**Figure 13.** PCA performed on the selected eleven kinematic inputs. The x-axis indicates the number of PCA features, and the y-axis indicates the variance. Choosing the optimal size for the latent space requires identifying the point of diminishing variance return.

The output layer is a mirror of the input layer and therefore has 11 nodes. The hidden layer and output layer use ReLU and sigmoid activation functions respectively. An L2 regularization term was added to the hidden layer to avoid overfitting. Because training an autoencoder is an unsupervised and unlabelled process, there is no prediction of whether a given event is signal or background. Instead of using a prediction to evaluate the efficacy of the AE approach, the loss output from a testing dataset is used as a proxy. Since dihiggs events should be relatively anomalous compared to the QCD training set, signal events should have a relatively larger average loss value. Cutting on the loss function allows for a significance to be calculated for comparison with other supervised methods.



**Figure 14.** (Left) The loss of the AE during QCD training/reconstruction converged after 700 epochs. (Right) The loss distribution generated by the AE when being tested on QCD and diHiggs event data separately.

Training for several hundred epochs leads the model to converge, reaching a loss near 0.8 (see Fig.14). Requiring the loss to be larger than 0.1 yields a best  $S/\sqrt{B}$  of  $0.67 \pm 0.01$ . This significance value may be somewhat misleading since the signal and background loss distributions have little separation. The highest significance result effectively is a cut that keeps nearly all events.

## 5 Results

The methods covered in this paper are by no means exhaustive of the ML landscape available to high energy physics, but a wide range of techniques and philosophies are covered. A weak observed trend is that the more complex the inputs are allowed to remain, the more information sophisticated algorithms can wring out of the data. This is less of a hard rule and more of a loose heuristic. Table 6 provides a summary of the methods described in the previous sections. The results of a traditional 1-D sequential cut technique is shown for comparison though it has not been previously described.

An important caveat to keep in mind is that all results discussed here were determined in conditions with zero pileup. In higher pileup environments like those expected at the HL-LHC, reconstruction algorithms see serious reductions in correct combinatoric matching. This effect will certainly degrade the expected performance of techniques that rely on explicit event reconstruction. Methods that do not rely on event reconstruction (CNN, PFN) might be more robust to these effects, and this should be studied in further work. The unsupervised AE technique performed the worst among all the methods tested, but this is likely a reflection of a mismatch between problem and proposed solution rather than a statement on the use of unsupervised techniques in general.



Method	0PU		
	Best $S/\sqrt{B}$	$N_{\text{Signal}}$	$N_{\text{Background}}$
Autoencoder	$0.67 \pm 0.01$	5849.4	$7.7 \cdot 10^7$
1D-Rectangular Cuts	$0.82 \pm 0.XX$	3621.0	$1.97 \cdot 10^7$
k-Means Clustering	$1.44 \pm 0.02$	1703.6	$1.4 \cdot 10^6$
Particle Flow Network	$1.62 \pm 0.01$	$1.8 \cdot 10^4$	$1.2 \cdot 10^8$
Boosted Decision Tree	$1.84 \pm 0.09$	986.3	$2.8 \cdot 10^5$
Lorentz Boost Network	$1.87 \pm 0.08$	1123.3	$3.6 \cdot 10^5$
Random Forest	$2.04 \pm 0.07$	1154.7	$3.2 \cdot 10^5$
Feed-Forward NN	$2.20 \pm 0.07$	1520.1	$4.8 \cdot 10^5$
Convolutional NN	$2.85 \pm 0.02$	-	-

**Table 6.** Comparison of method significance and signal/background yields normalized to full HL-LHC dataset of  $3000 \text{ fb}^{-1}$ .

## 346 6 Conclusions

347 Measuring the rate of diHiggs production will be a problem facing the high energy physics  
348 community through the end of the HL-LHC era. The techniques explored in this paper  
349 show the power of machine learning techniques to identify diHiggs signals from amongst  
350 overwhelming QCD signals. The significance obtained for many of these methods is im-  
351 pressive given the simplicity of the evaluation metric and bodes well for current and future  
352 measurements of double Higgs production at the LHC.

## References

- [1] V. M. Abazov, B. Abbott, M. Abolins, B. S. Acharya, M. Adams, T. Adams, E. Aguilo, M. Ahsan, G. D. Alexeev, G. Alkhazov, and et al., “Observation of single top-quark production,” *Physical Review Letters*, vol. 103, Aug 2009.
- [2] K. A. et al., “Machine learning in high energy physics community white paper,” 2018.
- [3] R. Brun and F. Rademakers, “ROOT: An object oriented data analysis framework,” *Nucl. Instrum. Meth. A*, vol. 389, pp. 81–86, 1997.
- [4] J. Alwall, R. Frederix, S. Frixione, V. Hirschi, F. Maltoni, O. Mattelaer, H. S. Shao, T. Stelzer, P. Torrielli, and M. Zaro, “The automated computation of tree-level and next-to-leading order differential cross sections, and their matching to parton shower simulations,” *JHEP*, vol. 07, p. 079, 2014.
- [5] T. Sjöstrand, S. Ask, J. R. Christiansen, R. Corke, N. Desai, P. Ilten, S. Mrenna, S. Prestel, C. O. Rasmussen, and P. Z. Skands, “An introduction to pythia 8.2,” *Computer Physics Communications*, vol. 191, p. 159–177, Jun 2015.
- [6] J. de Favereau, C. Delaere, P. Demin, A. Giammanco, V. Lemaître, A. Mertens, and M. Selvaggi, “Delphes 3: a modular framework for fast simulation of a generic collider experiment,” *Journal of High Energy Physics*, vol. 2014, Feb 2014.
- [7] B. T. et al., “dihiggsmlproject.” <https://github.com/neu-physics/dihiggsMLProject>, 2020.
- [8] G. Aad, T. Abajyan, B. Abbott, J. Abdallah, S. Abdel Khalek, A. Abdelalim, O. Abdinov, R. Aben, B. Abi, M. Abolins, and et al., “Observation of a new particle in the search for the standard model higgs boson with the atlas detector at the lhc,” *Physics Letters B*, vol. 716, p. 1–29, Sep 2012.
- [9] S. Chatrchyan, V. Khachatryan, A. Sirunyan, A. Tumasyan, W. Adam, E. Aguilo, T. Bergauer, M. Dragicevic, J. Erö, C. Fabjan, and et al., “Observation of a new boson at a mass of 125 gev with the cms experiment at the lhc,” *Physics Letters B*, vol. 716, p. 30–61, Sep 2012.
- [10] T. Chen and C. Guestrin, “Xgboost: A scalable tree boosting system,” *CoRR*, vol. abs/1603.02754, 2016.
- [11] F. Chollet, “Keras.” <https://github.com/fchollet/keras>, 2015.
- [12] M. A. et al., “Tensorflow: Large-scale machine learning on heterogeneous distributed systems,” *CoRR*, vol. abs/1603.04467, 2016.
- [13] J. Alison, S. An, P. Bryant, B. Burkle, S. Gleyzer, M. Narain, M. Paulini, B. Poczoz, and E. Usai, “End-to-end particle and event identification at the Large Hadron Collider with CMS Open Data,” in *Meeting of the Division of Particles and Fields of the American Physical Society*, 10 2019.