

# Searching for Exotic Particles in High-Energy Physics with Deep Learning

P. Baldi,<sup>1</sup> P. Sadowski,<sup>1</sup> and D. Whiteson<sup>2</sup>

<sup>1</sup>*Dept. of Computer Science, UC Irvine, Irvine, CA 92617\**

<sup>2</sup>*Dept. of Physics and Astronomy, UC Irvine, Irvine, CA 92617†*

Collisions at high-energy particle colliders are a traditionally fruitful source of exotic particle discoveries. Finding these rare particles requires solving difficult signal-versus-background classification problems, hence machine learning approaches are often used. Standard approaches have relied on ‘shallow’ machine learning models that have a limited capacity to learn complex non-linear functions of the inputs, and rely on a pain-staking search through manually constructed non-linear features. Progress on this problem has slowed, as a variety of techniques have shown equivalent performance. Recent advances in the field of deep learning make it possible to learn more complex functions and better discriminate between signal and background classes. Using benchmark datasets, we show that deep learning methods need no manually constructed inputs and yet improve the classification metric by as much as 8% over the best current approaches. This demonstrates that deep learning approaches can improve the power of collider searches for exotic particles.

## INTRODUCTION

The field of *high energy physics* is devoted to the study of the elementary constituents of matter. By investigating the structure of matter and the laws that govern its interactions, this field strives to discover the fundamental properties of the physical universe. The primary tools of experimental high energy physicists are modern accelerators, which collide protons and/or antiprotons to create exotic particles that occur only at extremely high energy densities. Observing these particles and measuring their properties may yield critical insights about the very nature of matter [1]. Such discoveries require powerful statistical methods, and machine learning tools play a critical role. Given the limited quantity and expensive nature of the data, improvements in analytical tools directly boost particle discovery potential.

To discover a new particle, physicists must isolate a subspace of their high-dimensional data in which the hypothesis of a new particle or force gives a significantly different prediction than the null hypothesis, allowing for an effective statistical test. For this reason, the critical element of the search for new particles and forces in high-energy physics is the computation of the *relative likelihood*, the ratio of the sample likelihood functions in the two considered hypotheses, shown by Neyman and Pearson [2] to be the optimal discriminating quantity. Often this relative likelihood function cannot be expressed analytically, so simulated collision data generated with Monte Carlo methods are used as a basis for approximation of the likelihood function. The high dimensionality of data, referred to as the *feature space*, makes it intractable to generate enough simulated collisions to describe the relative likelihood in the full feature space, and machine learning tools are used for dimensionality reduction. Machine learning classifiers such as neural networks provide a powerful way to solve this learning problem.

The relative likelihood function is a complicated function in a high-dimensional space. While any function

can theoretically be represented by a ‘shallow’ classifier, such as a neural network with a single hidden layer [3], an intractable number of hidden units may be required. Circuit complexity theory tells us that deep neural networks (DN) have the potential to compute complex functions much more efficiently (fewer hidden units), but in practice they are notoriously difficult to train due to the vanishing gradient problem [4, 5]; the adjustments to the weights in the early layers of a deep network rapidly approach zero during training. A common approach is to combine shallow classifiers with high-level features that are derived manually from the raw features. These are generally non-linear functions of the input features that capture physical insights about the data. While helpful, this approach is labor-intensive and not necessarily optimal; a robust machine learning method would obviate the need for this additional step and capture all of the available classification power directly from the raw data.

Recent successes in deep learning – *e.g.* neural networks with multiple hidden layers – have come from alleviating the gradient diffusion problem by a combination of factors, including: 1) speeding up the stochastic gradient descent algorithm with graphics processors; 2) using much larger training sets; 3) using new learning algorithms, including randomized algorithms such as dropout [6, 7]; and 4) pre-training the initial layers of the network with unsupervised learning methods such as autoencoders [8, 9]. With these methods, it is becoming common to train deep networks of five or more layers. These advances in deep learning could have a significant impact on applications in high-energy physics. Construction and operation of the particle accelerators is extremely expensive, so any additional classification power extracted from the collision data is very valuable.

In this paper, we show that the current techniques used in high-energy physics fail to capture all of the available information, even when boosted by manually-constructed physics-inspired features. This effectively reduces the power of the collider to discover new particles.

We demonstrate that recent developments in deep learning tools can overcome these failings, providing significant boosts even without manual assistance.

## RESULTS

The vast majority of particle collisions do not produce exotic particles. For example, though the Large Hadron Collider produces approximately  $10^{11}$  collisions per hour, approximately 300 of these collisions result in a Higgs boson, on average. Therefore, good data analysis depends on distinguishing collisions which produce particles of interest (*signal*) from those producing other particles (*background*).

Even when interesting particles are produced, detecting them poses considerable challenges. They are too small to be directly observed and decay almost immediately into other particles. Though new particles cannot be directly observed, the lighter stable particles to which they decay, called *decay products*, can be observed. Multiple layers of detectors surround the point of collision for this purpose. As each decay product pass through these detectors, it interacts with them in a way that allows its direction and momentum to be measured.

Observable decay products include electrically-charged leptons (electrons or muons, denoted  $\ell$ ), and particle jets (collimated streams of particles originating from quarks or gluons, denoted  $j$ ). In the case of jets we attempt to distinguish between jets from heavy quarks ( $b$ ) and jets from gluons or low-mass quarks; jets consistent with  $b$ -quarks receive a  $b$ -quark *tag*. For each object, the momentum is determined by three measurements: the momentum transverse to the beam direction ( $p_T$ ), and two angles,  $\theta$  (polar) and  $\phi$  (azimuthal). For convenience, at hadron colliders, such as Tevatron and LHC, the *pseudorapidity*, defined as  $\eta = -\ln(\tan(\theta/2))$  is used instead of the polar angle  $\theta$ . Finally, an important quantity is the amount of momentum carried away by the invisible particles. This cannot be directly measured, but can be inferred in the plane transverse to the beam by requiring conservation of momentum. The initial state has zero momentum transverse to the beam axis, therefore any imbalance of transverse momentum (denoted  $\cancel{E}_T$ ) in the final state must be due to production of invisible particles such as neutrinos ( $\nu$ ) or exotic particles. The momentum imbalance in the longitudinal direction along the beam cannot be measured at hadron colliders, as the initial state momentum of the quarks is not known.

### Benchmark Case for Higgs Bosons (HIGGS)

The first benchmark classification task is to distinguish between a signal process where new theoretical Higgs bosons are produced, and a background process with the

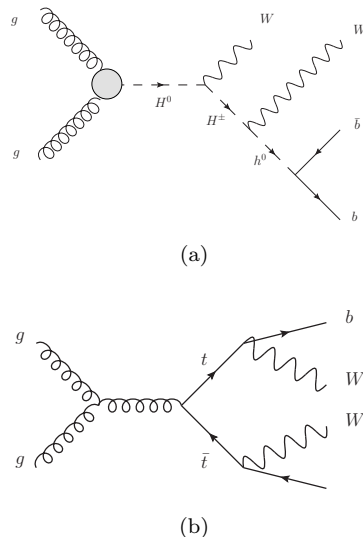


FIG. 1: **Diagrams for Higgs benchmark.** (a) Diagram describing the signal process involving new exotic Higgs bosons  $H^0$  and  $H^\pm$ . (b) Diagram describing the background process involving top-quarks ( $t$ ). In both cases, the resulting particles are two  $W$  bosons and two  $b$ -quarks.

identical decay products but distinct kinematic features. This benchmark task was recently considered by experiments at the LHC [10] and the Tevatron colliders [11].

The signal process is the fusion of two gluons into a heavy electrically-neutral Higgs boson ( $gg \rightarrow H^0$ ), which decays to a heavy electrically-charged Higgs bosons ( $H^\pm$ ) and a  $W$  boson. The  $H^\pm$  boson subsequently decays to a second  $W$  boson and the light Higgs boson,  $h^0$  which has recently been observed by the ATLAS [12] and CMS [13] experiments. The light Higgs boson decays predominantly to a pair of bottom quarks, giving the process:

$$gg \rightarrow H^0 \rightarrow W^\mp H^\pm \rightarrow W^\mp W^\pm h^0 \rightarrow W^\mp W^\pm b\bar{b}, \quad (1)$$

which leads to  $W^\mp W^\pm b\bar{b}$ , see Figure 1. The background process, which mimics  $W^\mp W^\pm b\bar{b}$  without the Higgs boson intermediate state, is the production of a pair of top quarks, each of which decay to  $Wb$ , also giving  $W^\mp W^\pm b\bar{b}$ , see Figure 1.

Simulated events are generated with the MADGRAPH5 [14] event generator assuming 8 TeV collisions of protons as at the latest run of the Large Hadron Collider, with showering and hadronization performed by PYTHIA [15] and detector response simulated by DELPHES [16]. For the benchmark case here,  $m_{H^0} = 425$  GeV and  $m_{H^\pm} = 325$  GeV has been assumed.

We focus on the semi-leptonic decay mode, in which one  $W$  boson decays to a lepton and neutrino ( $\ell\nu$ ) and the other decays to a pair of jets ( $jj$ ), giving decay products  $\ell\nu b jjb$ . We consider events which satisfy the requirements:

- Exactly one electron or muon, with  $p_T > 20$  GeV and  $|\eta| < 2.5$ ;
- at least four jets, each with  $p_T > 20$  GeV and  $|\eta| < 2.5$ ;
- $b$ -tags on at least two of the jets, indicating that they are likely due to  $b$ -quarks rather than gluons or lighter quarks.

Events which satisfy the requirements above are naturally described by a simple set of features which represent the basic measurements made by the particle detector: the momentum of each observed particle. In addition we reconstruct the missing transverse momentum in the event and have  $b$ -tagging information for each jet. Together, these twenty-one features comprise our *low-level feature set*. Figure 2 shows the distribution of a subset of these kinematic features for signal and background processes.

The low-level features show some distinguishing characteristics, but our knowledge of the different intermediate states of the two processes allows us to construct other features which better capture the differences. As the difference in the two hypotheses lies mostly in the existence of new intermediate Higgs boson states, we can distinguish between the two hypotheses by attempting to identify whether the intermediate state existed. This is done by reconstructing its characteristic invariant mass; if a particle A decays into particles B and C, the invariant mass of particle A ( $m_A$ ) can be reconstructed as:

$$m_A^2 = m_{B+C}^2 = (E_B + E_C)^2 - |(\mathbf{p}_B + \mathbf{p}_C)|^2 \quad (2)$$

where  $E$  is the energy and  $\mathbf{p}$  is the three-dimensional momentum of the particle. Similarly, the invariant mass of more than two particles can be defined as the modulus of the particles' Lorentz four-vector sum. In the signal hypothesis we expect that:

- $W \rightarrow \ell\nu$  gives a peak in the  $m_{\ell\nu}$  distribution at the known  $W$  boson mass,  $m_W$ ,
- $W \rightarrow jj$  gives a peak in the  $m_{jj}$  distribution at  $m_W$ ,
- $h^0 \rightarrow b\bar{b}$  gives a peak in the  $m_{b\bar{b}}$  distribution at the known Higgs boson mass,  $m_{h^0}$ ,
- $H^\pm \rightarrow Wh^0$  gives a peak in the  $m_{Wb\bar{b}}$  distribution at the assumed  $H^\pm$  mass,  $m_{H^\pm}$ ,
- $H^0 \rightarrow WH^\pm$  gives a peak in the  $m_{WWb\bar{b}}$  distribution at the assumed  $H^0$  mass, at  $m_{H^0}$ ,

Note that the leptonic  $W$  boson is reconstructed by combining the lepton with neutrino, whose transverse

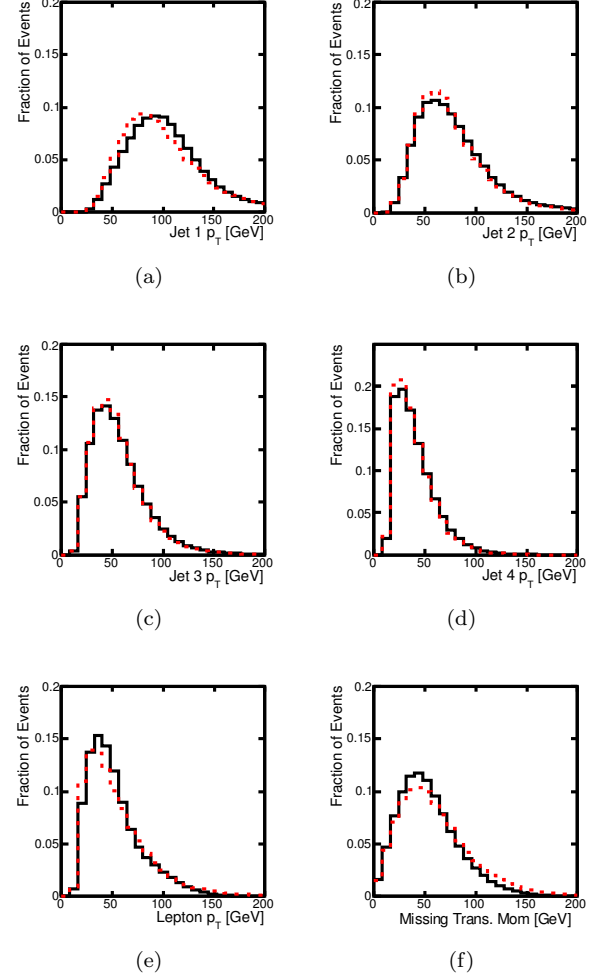


FIG. 2: **Low-level input features for Higgs benchmark.** Distributions in  $\ell\nu jj b\bar{b}$  events for simulated signal (black) and background (red) benchmark events. Shown are the distributions of transverse momenta ( $p_T$ ) of each observed particle (a,b,c,d,e) as well as the imbalance of momentum in the final state (f). Momentum angular information for each observed particle is also available to the network, but is not shown, as the one-dimensional projections have little information.

momentum is deduced from the imbalance of momentum in the final state objects and whose rapidity is set to give  $m_{\ell\nu}$  closest to  $m_W = 80.4$  GeV.

while in the case of the  $t\bar{t}$  background we expect that:

- $W \rightarrow \ell\nu$  gives a peak in  $m_{\ell\nu}$  at  $m_W$ ,
- $W \rightarrow jj$  gives a peak in  $m_{jj}$  at  $m_W$ ,
- $t \rightarrow Wb$  gives a peak in  $m_{j\ell\nu}$  and  $m_{jbb}$  at  $m_t$ .

See Figure 3 for distributions of these high-level features for both signal and background processes. Clearly these contain more discrimination power than the low-level features.

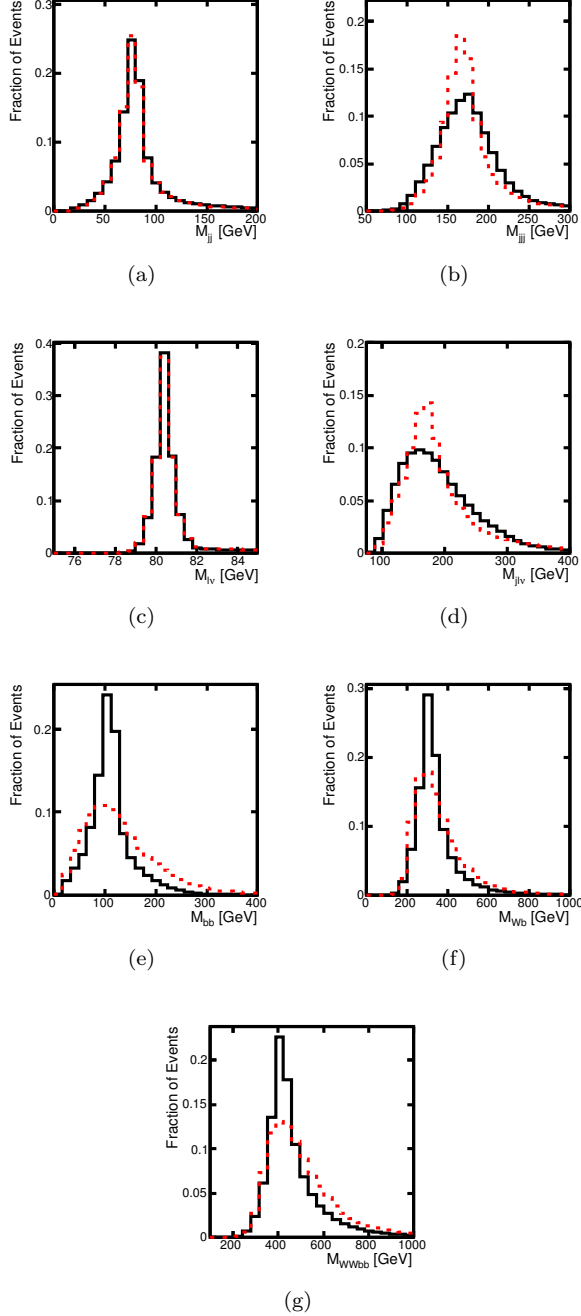


FIG. 3: **High-level input features for Higgs benchmark.** Distributions in simulation of invariant mass calculations in  $\ell\nu jj b\bar{b}$  events for simulated signal (black) and background (red) events.

We have published a dataset containing 11 million simulated collision events for benchmarking machine learning classification algorithms on this task, which can be found in the UCI Machine Learning Repository at [archive.ics.uci.edu/ml/datasets/HIGGS](http://archive.ics.uci.edu/ml/datasets/HIGGS).

### Benchmark Case for Supersymmetry Particles (SUSY)

The second benchmark classification task is to distinguish between a process where new supersymmetric particles are produced, leading to a final state in which some particles are detectable and others are invisible to the experimental apparatus, and a background process with the same detectable particles but fewer invisible particles and distinct kinematic features. This benchmark problem is currently of great interest to the field of high-energy physics, and there is a vigorous effort in the literature [17–20] to build high-level features which can aid in the classification task.

The signal process is the production of electrically-charged supersymmetric particles ( $\chi^\pm$ ), which decay to  $W$  bosons and an electrically-neutral supersymmetric particle  $\chi^0$ , which is invisible to the detector. The  $W$  bosons decay to charged leptons  $\ell$  and invisible neutrinos  $\nu$ , see Fig. 4. The final state in the detector is therefore two charged leptons ( $\ell\ell$ ) and missing momentum carried off by the invisible particles ( $\chi^0\chi^0\nu\nu$ ). The background process is the production of pairs of  $W$  bosons, which decay to charged leptons  $\ell$  and invisible neutrinos  $\nu$ , see Fig. 4. The visible portion of the signal and background final states both contain two leptons ( $\ell\ell$ ) and large amounts of missing momentum due to the invisible particles. The classification task requires distinguishing between these two processes using the measurements of the charged lepton momenta and the missing transverse momentum.

As above, simulated events are generated with the MADGRAPH [14] event generator assuming 8 TeV collisions of protons as at the latest run of the Large Hadron Collider, with showering and hadronization performed by PYTHIA [15] and detector response simulated by DELPHES [16]. The masses are set to  $m_{\chi^\pm} = 200$  GeV and  $m_{\chi^0} = 100$  GeV.

We focus on the fully leptonic decay mode, in which both  $W$  bosons decay to charged leptons and neutrinos,  $\ell\nu\ell\nu$ . We consider events which satisfy the requirements:

- Exactly two electrons or muons, each with  $p_T > 20$  GeV and  $|\eta| < 2.5$ ;
- at least 20 GeV of missing transverse momentum

As above, the basic detector response is used to measure the momentum of each visible particle, in this case the charged leptons. In addition, there may be particle jets induced by radiative processes. A critical quantity is the missing transverse momentum,  $\cancel{E}_T$ . Figure 5 gives distributions of low-level features for signal and background processes.

The search for supersymmetric particles is a central piece of the scientific mission of the Large Hadron Collider. The strategy we applied to the Higgs boson bench-

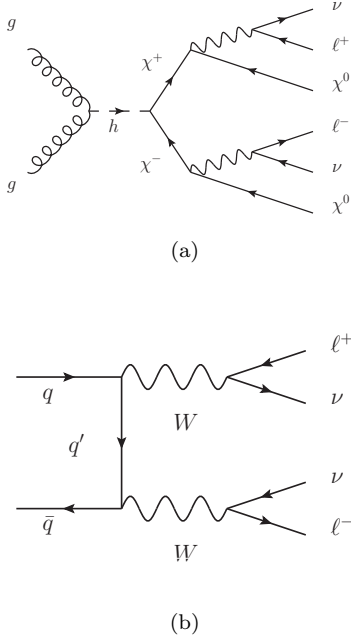


FIG. 4: **Diagrams for SUSY benchmark.** Example diagrams describing the signal process involving hypothetical supersymmetric particles  $\chi^\pm$  and  $\chi^0$  along with charged leptons  $\ell^\pm$  and neutrinos  $\nu$  (a) and the background process involving  $W$  bosons (b). In both cases, the resulting observed particles are two charged leptons, as neutrinos and  $\chi^0$  escape undetected.

mark, of reconstructing the invariant mass of the intermediate state, is not feasible here, as there is too much information carried away by the escaping neutrinos (two neutrinos in this case, compared to one for the Higgs case). Instead, a great deal of intellectual energy has been spent in attempting to devise features which give additional classification power. These include high-level features such as:

- Axial  $\cancel{E}_T$ : missing transverse energy along the vector defined by the charged leptons,
- stransverse mass  $M_{T2}$ : estimating the mass of particles produced in pairs and decaying semi-invisibly [17, 18],
- $\cancel{E}_T^{Rel}$ :  $\cancel{E}_T$  if  $\Delta\phi \geq \pi/2$ ,  $\cancel{E}_T \sin(\Delta\phi)$  if  $\Delta\phi < \pi/2$ , where  $\Delta\phi$  is the minimum angle between  $\cancel{E}_T$  and a jet or lepton,
- razor quantities  $\beta, R$ , and  $M_R$  [19],
- super-razor quantities  $\beta_{R+1}$ ,  $\cos(\theta_{R+1})$ ,  $\Delta\phi_R^\beta$ ,  $M_\Delta^R$ ,  $M_R^T$ , and  $\sqrt{s_R}$  [20].

See Figure 6 for distributions of these high-level features for both signal and background processes.

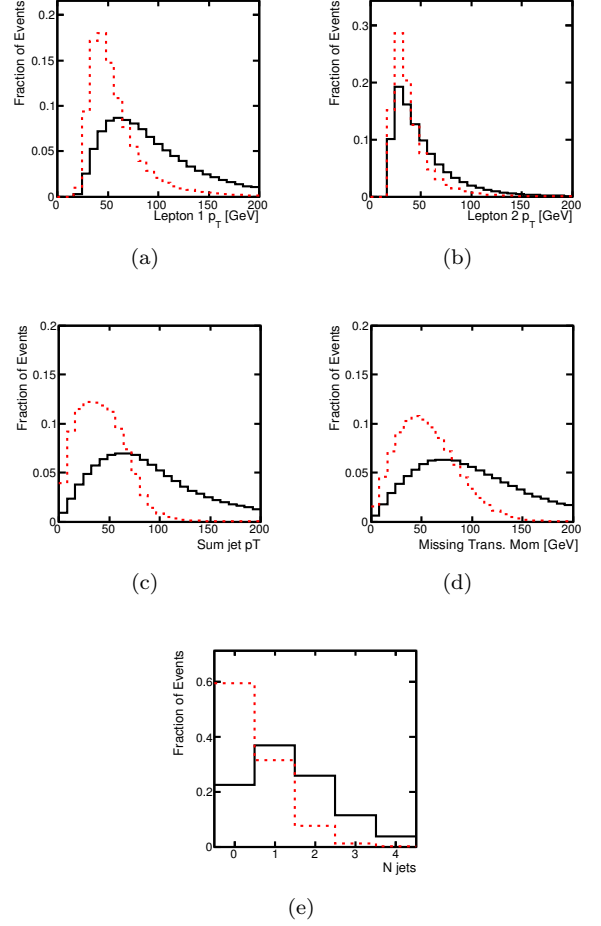


FIG. 5: **Low-level input features for SUSY benchmark.** Distribution of low-level features in simulated samples for the SUSY signal (black) and background (red) benchmark processes.

A dataset containing five million simulated collision events are available for download at [archive.ics.uci.edu/ml/datasets/SUSY](http://archive.ics.uci.edu/ml/datasets/SUSY).

### Current Approach

Standard techniques in high-energy physics data analyses include feed-forward neural networks with a single hidden layer and boosted-decision trees. We use the widely-used TMVA package [21], which provides a standardized implementation of common multi-variate learning techniques and an excellent performance baseline.

### Deep Learning

We explored the use of deep neural networks as a practical tool for applications in high-energy physics. Hyper-



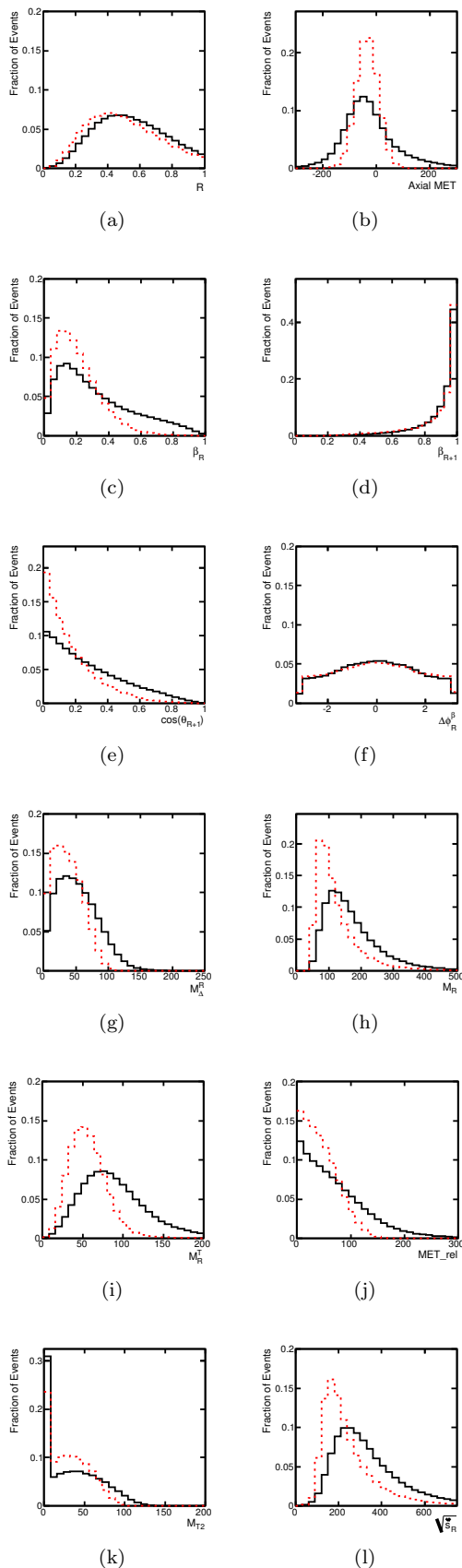


FIG. 6: **High-level input features for SUSY benchmark.** Distribution of high-level features in simulated samples for the SUSY signal (black) and background (red) benchmark processes.

parameters were chosen using a subset of the HIGGS data consisting of 2.6 million training examples and 100,000 validation examples. Due to computational costs, this optimization was not thorough, but included combinations of the pre-training methods, network architectures, initial learning rates, and regularization methods shown in Supplementary Table 1. We selected a five-layer neural network with 300 hidden units in each layer, a learning rate of 0.05, and a weight decay coefficient of  $1 \times 10^{-5}$ . Pre-training, extra hidden units, and additional hidden layers significantly increased training time without noticeably increasing performance. To facilitate comparison, shallow neural networks were trained with the same hyper-parameters and the same number of units per hidden layer. Additional training details are provided in the Methods section below.

The hyper-parameter optimization was performed using the full set of HIGGS features. To investigate whether the neural networks were able to learn the discriminative information contained in the high-level features, we trained separate classifiers for each of the three feature sets described above: low-level, high-level, and combined feature sets. For the SUSY benchmark, the networks were trained with the same hyper-parameters chosen for the HIGGS, as the datasets have similar characteristics and the hyper-parameter search is computationally expensive.

## Performance

Classifiers were tested on 500,000 simulated examples generated from the same Monte Carlo procedures as the training sets. We produced Receiver Operating Characteristic (ROC) curves to illustrate the performance of the classifiers. Our primary metric for comparison is the area under the ROC curve (AUC), with larger AUC values indicating higher classification accuracy across a range of threshold choices.

This metric is insightful, as it is directly connected to classification accuracy, which is the quantity optimized for in training. In practice, physicists may be interested in other metrics, such as signal efficiency at some fixed background rejection, or discovery significance as calculated by  $p$ -value in the null hypothesis. We choose AUC as it is a standard in machine learning, and is closely correlated with the other metrics. In addition, we calculate discovery significance – the standard metric in high-energy physics – to demonstrate that small increases in AUC can represent significant enhancement in discovery significance.

Note, however, that in some applications the determining factor in the sensitivity to new exotic particles is determined not only by the discriminating power of the selection, but by the uncertainties in the background model itself. Some portions of the background model may be

TABLE I: **Performance for Higgs benchmark.** Comparison of the performance of several learning techniques: boosted decision trees (BDT), shallow neural networks (NN), and deep neural networks (DN) for three sets of input features: low-level features, high-level features and the complete set of features. Each neural network was trained five times with different random initializations. The table displays the mean Area Under the Curve (AUC) of the signal-rejection curve in Figure 7, with standard deviations in parentheses as well as the expected significance of a discovery (in units of Gaussian  $\sigma$ ) for 100 signal events and  $1000 \pm 50$  background events.

Technique	AUC		
	Low-level	High-level	Complete
BDT	0.73 (0.01)	0.78 (0.01)	0.81 (0.01)
NN	0.733 (0.007)	0.777 (0.001)	0.816 (0.004)
DN	0.880 (0.001)	0.800 ( $< 0.001$ )	0.885 (0.002)

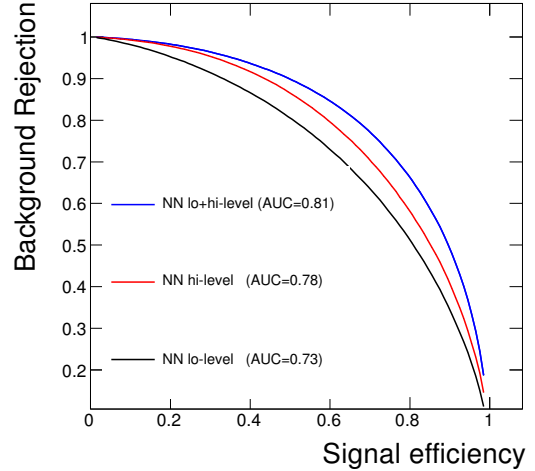
Technique	Discovery significance		
	Low-level	High-level	Complete
NN	$2.5\sigma$	$3.1\sigma$	$3.7\sigma$
DN	$4.9\sigma$	$3.6\sigma$	$5.0\sigma$

better understood than others, so that some simulated background collisions have larger associated systematic uncertainties than other collisions. This can transform the problem into one of reinforcement learning, where per-collision truth labels no longer indicate the ideal network output target. This is beyond the scope of this study, but see Refs. [22, 23] for stochastic optimization strategies for such problems.

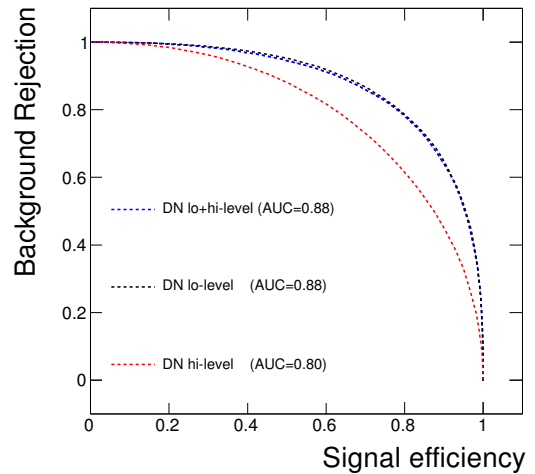
Figure 7 and Table I show the signal efficiency and background rejection for varying thresholds on the output of the neural network (NN) or boosted decision tree (BDT).

A shallow NN or BDT trained using only the low-level features performs significantly worse than one trained with only the high-level features. This implies that the shallow NN and BDT are not succeeding in independently discovering the discriminating power of the high-level features. This is a well-known problem with shallow learning methods, and motivates the calculation of high-level features.

Methods trained with only the high-level features, however, have a weaker performance than those trained with the full suite of features, which suggests that despite the insight represented by the high-level features, they do not capture all of the information contained in the low-level features. The deep learning techniques show nearly equivalent performance using the low-level features and the complete features, suggesting that they are *automatically discovering the insight contained in the high-level features*. Finally, the deep learning technique finds additional separation power beyond what is contained in the high-level features, demonstrated by the superior perfor-



(a)

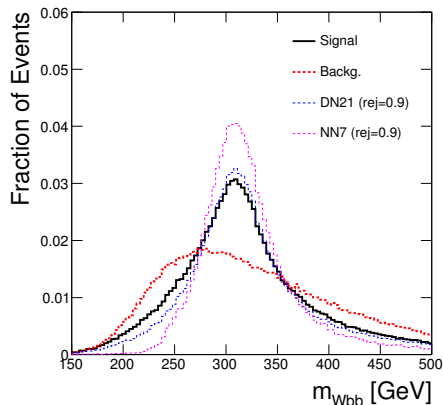


(b)

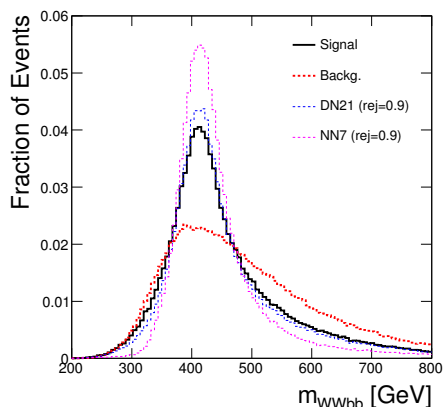
FIG. 7: **Performance for Higgs benchmark.** For the Higgs benchmark, comparison of background rejection versus signal efficiency for the traditional learning method (a) and the deep learning method (b) using the low-level features, the high-level features and the complete set of features.

mance of the deep network with low-level features to the traditional network using high-level features. These results demonstrate the advantage to using deep learning techniques for this type of problem.

The internal representation of a NN is notoriously difficult to reverse engineer. To gain some insight into the mechanism by which the deep network (DN) is improving upon the discrimination in the high-level physics features, we compare the distribution of simulated events selected by a minimum threshold on the NN or DN output, chosen to give equivalent rejection of 90% of the background



(a)



(b)

FIG. 8: **Performance comparisons.** Distribution of events for two rescaled input features: (a)  $m_{Wbb}$  and (b)  $m_{WWbb}$ . Shown are pure signal and background distributions, as well as events which pass a threshold requirement which gives a background rejection of 90% for a deep network with 21 low-level inputs (DN21) and a shallow network with 7 high-level inputs (NN7).

events. Figure 8 shows events selected by such thresholds in a mixture of 50% signal and 50% background collisions, compared to pure distributions of signal and background. The NN preferentially selects events with values of the features close to the characteristic signal values and away from background-dominated values. The DN, which has a higher efficiency for the equivalent rejection, selects events near the same signal values, but *also retains events away from the signal-dominated region*. The likely explanation is that the DN has discovered the same signal-rich region identified by the physics features, but has in addition found avenues to carve into the background-dominated region.

In the case of the SUSY benchmark, the deep neural networks again perform better than the shallow networks.

TABLE II: **Performance comparison for the SUSY benchmark.** Each model was trained five times with different weight initializations. The mean AUC is shown with standard deviations in parentheses as well as the expected significance of a discovery (in units of Gaussian  $\sigma$ ) for 100 signal events and  $1000 \pm 50$  background events.

Technique	AUC		
	Low-level	High-level	Complete
BDT	0.850 (0.003)	0.835 (0.003)	0.863 (0.003)
NN	0.867 (0.002)	0.863 (0.001)	0.875 ( $< 0.001$ )
NN <sub>dropout</sub>	0.856 ( $< 0.001$ )	0.859 ( $< 0.001$ )	0.873 ( $< 0.001$ )
DN	0.872 (0.001)	0.865 (0.001)	0.876 ( $< 0.001$ )
DN <sub>dropout</sub>	0.876 ( $< 0.001$ )	0.869 ( $< 0.001$ )	0.879 ( $< 0.001$ )
Technique	Discovery significance		
	Low-level	High-level	Complete
NN	$6.5\sigma$	$6.2\sigma$	$6.9\sigma$
DN	$7.5\sigma$	$7.3\sigma$	$7.6\sigma$

The improvement is less dramatic, though statistically significant.

An additional boost in performance is obtained by using the dropout training algorithm, in which we stochastically drop neurons in the top hidden layer with 50% probability during training. For deep networks trained with dropout, we achieve an AUC of 0.88 on both the low-level and complete feature sets. Table II, Supplementary Fig. 2 and Supplementary Fig. 3 compare the performance of shallow and deep networks for each of the three sets of input features.

In this SUSY case, neither the high-level features nor the deep network finds dramatic gains over the shallow network of low-level features. The power of the deep network to automatically find non-linear features reveals something about the nature of the classification problem in this case: it suggests that there may be little gain from further attempts to manually construct high-level features.

## Analysis

To highlight the advantage of deep networks over shallow networks with a similar number of parameters, we performed a thorough hyper-parameter optimization for the class of single-layer neural networks over the hyperparameters specified in Supplementary Table 2 on the HIGGS benchmark. The largest shallow network had 300,001 parameters, slightly more than the 279,901 parameters in the largest deep network, but these additional hidden units did very little to increase performance over a shallow network with only 30,001 parameters. Supplementary Table 3 compares the performance of the best shallow networks of each size with deep networks of varying depth.



While the primary advantage of deep neural networks is their ability to *automatically* learn high-level features from the data, one can imagine facilitating this process by pre-training a neural network to compute a particular set of high-level features. As a proof of concept, we demonstrate how deep neural networks can be trained to compute the high-level HIGGS features with a high degree of accuracy (Supplementary Table 4). Note that such a network could be used as a module within a larger neural network classifier.

## DISCUSSION

It is widely accepted in experimental high-energy physics that machine learning techniques can provide powerful boosts to searches for exotic particles. Until now, physicists have reluctantly accepted the limitations of the shallow networks employed to date; in an attempt to circumvent these limitations, physicists manually construct helpful non-linear feature combinations to guide the shallow networks.

Our analysis shows that recent advances in deep learning techniques may lift these limitations by **automatically discovering powerful non-linear feature combinations and providing better discrimination power than current classifiers** – even when aided by manually-constructed features. This appears to be the first such demonstration in a semi-realistic case.

We suspect that the novel environment of high energy physics, with high volumes of relatively low-dimensional data containing rare signals hiding under enormous backgrounds, can inspire new developments in machine learning tools. Beyond these simple benchmarks, deep learning methods may be able to tackle thornier problems with multiple backgrounds, or lower-level tasks such as identifying the decay products from the high-dimensional raw detector output.

## METHODS

### Neural Network Training

In training the neural networks, the following hyperparameters were predetermined without optimization. Hidden units all used the *tanh* activation function. Weights were initialized from a normal distribution with zero mean and standard deviation 0.1 in the first layer, 0.001 in the output layer, and 0.05 all other hidden layers. Gradient computations were made on mini-batches of size 100. A momentum term increased linearly over the first 200 epochs from 0.9 to 0.99, at which point it remained constant. The learning rate decayed by a factor of 1.0000002 every batch update until it reached a minimum of  $10^{-6}$ . Training ended when the momentum had

reached its maximum value and the minimum error on the validation set (500,000 examples) had not decreased by more than a factor of 0.00001 over 10 epochs. This early stopping prevented overfitting and resulted in each neural network being trained for 200-1000 epochs.

Autoencoder pretraining was performed by training a stack of single-hidden-layer autoencoder networks as in [9], then fine-tuning the full network using the class labels. Each autoencoder in the stack used tanh hidden units and linear outputs, and was trained with the same initialization scheme, learning algorithm, and stopping parameters as in the fine-tuning stage. When training with dropout, we increased the learning rate decay factor to 1.0000003, and only ended training when the momentum had reached its maximum value and the error on the validation set had not decreased for 40 epochs.

### Datasets

The data sets were nearly balanced, with 53% positive examples in the HIGGS data set, and 46% positive examples in the SUSY data set. Input features were standardized over the entire train/test set with mean zero and standard deviation one, except for those features with values strictly greater than zero – these we scaled so that the mean value was one.

### Computation

Computations were performed using machines with 16 Intel Xeon cores, an NVIDIA Tesla C2070 graphics processor, and 64 GB memory. All neural networks were trained using the GPU-accelerated Theano and Pylearn2 software libraries [24, 25]. Our code is available at <https://github.com/uci-igb/higgs-susy>.

## ACKNOWLEDGMENTS

We are grateful to Kyle Cranmer, Chris Hays, Chase Shimmin, Davide Gerbaudo, Bo Jayatilaka, Jahred Adelman, and Shimon Whiteson for their insightful comments. We wish to acknowledge a hardware grant from NVIDIA.

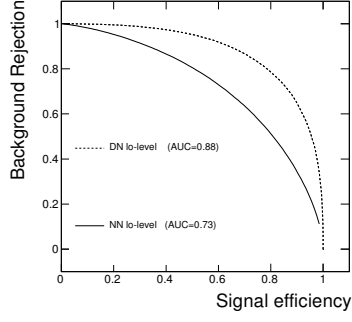
## AUTHOR CONTRIBUTIONS

PB conceived the idea of applying deep learning methods to high energy physics and particle detection. DW chose the benchmarks and generated the data and structured the problem. PB and PS designed the architectures

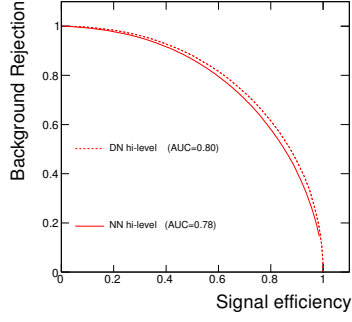
and the algorithms. PS implemented the code and performed the experiments. All authors analyzed the results and contributed to writing up the manuscript.

## REFERENCES

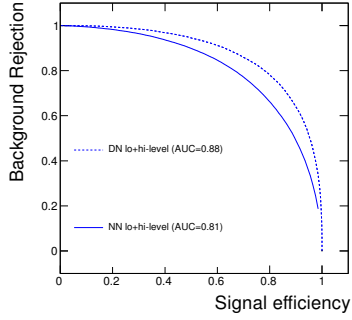
- 
- \* Electronic address: pbaldi@uci.edu  
† Electronic address: daniel@uci.edu
- [1] Dawson, S. *et al.* Higgs Working Group Report of the Snowmass 2013 Community Planning Study, Preprint at <http://arxiv.org/abs/1310.8361> (2013).
  - [2] Neyman, J. & Pearson, E. *Philosophical Transactions of the Royal Society* **231**, 694–706 (1933).
  - [3] Hornik, K., Stinchcombe, M. & White, H. Multilayer feedforward networks are universal approximators. *Neural Netw.* **2**, 359366 (1989). URL [http://dx.doi.org/10.1016/0893-6080\(89\)90020-8](http://dx.doi.org/10.1016/0893-6080(89)90020-8).
  - [4] Hochreiter, S. *Recurrent Neural Net Learning and Vanishing Gradient* (1998).
  - [5] Bengio, Y., Simard, P. & Frasconi, P. Learning long-term dependencies with gradient descent is difficult. *Neural Networks, IEEE Transactions on* **5**, 157–166 (1994).
  - [6] Hinton, G., Srivastava, N., Krizhevsky, A., Sutskever, I. & Salakhutdinov, R. R. Improving neural networks by preventing co-adaptation of feature detectors. <http://arxiv.org/abs/1207.0580> (2012).
  - [7] Baldi, P. & Sadowski, P. The dropout learning algorithm. *Artificial Intelligence* **210**, 78–122 (2014). URL <http://www.sciencedirect.com/science/article/pii/S0004370214000216>.
  - [8] Hinton, G. E., Osindero, S. & Teh, Y.-W. A fast learning algorithm for deep belief nets. *Neural Computation* **18**, 1527–1554 (2006). URL <http://dx.doi.org/10.1162/neco.2006.18.7.1527>.
  - [9] Bengio, Y. *et al.* Greedy layer-wise training of deep networks. In *Advances in Neural Information Processing Systems 19* (MIT Press, 2007).
  - [10] Aad, G. *et al.* Search for a Multi-Higgs Boson Cascade in  $W^+Wb\bar{b}$  events with the ATLAS detector in pp collisions at  $s = 8$  TeV (2013).
  - [11] Aaltonen, T. *et al.* Search for a two-Higgs-boson doublet using a simplified model in  $p\bar{p}$  collisions at  $\sqrt{s} = 1.96$  TeV. *Phys.Rev.Lett.* **110**, 121801 (2013).
  - [12] ATLAS Collaboration. Observation of a new particle in the search for the Standard Model Higgs boson with the ATLAS detector at the LHC. *Phys.Lett.* **B716**, 1–29 (2012).
  - [13] CMS Collaboration. Observation of a new boson at a mass of 125 GeV with the CMS experiment at the LHC. *Phys.Lett.* **B716**, 30–61 (2012).
  - [14] Alwall, J. *et al.* MadGraph 5 : Going Beyond. *JHEP* **1106**, 128 (2011).
  - [15] Sjostrand, T. *et al.* PYTHIA 6.4 physics and manual. *JHEP* **05**, 026 (2006).
  - [16] Ovyn, S., Roubey, X. & Lemaitre, V. DELPHES, a framework for fast simulation of a generic collider experiment (2009).
  - [17] Cheng, H.-C. & Han, Z. Minimal Kinematic Constraints and  $m(T_2)$ . *JHEP* **0812**, 063 (2008).
  - [18] Barr, A., Lester, C. & Stephens, P.  $m(T_2)$ : The Truth behind the glamour. *J.Phys.* **G29**, 2343–2363 (2003).
  - [19] Rogan, C. Kinematical variables towards new dynamics at the LHC (2010).
  - [20] Buckley, M. R., Lykken, J. D., Rogan, C. & Spiropulu, M. Super-Razor and Searches for Sleptons and Charginos at the LHC (2013).
  - [21] Hocker, A. *et al.* TMVA - Toolkit for Multivariate Data Analysis. *PoS ACAT*, 040 (2007).
  - [22] Whiteson, S. & Whiteson, D. Stochastic optimization for collision selection in high energy physics. In *IAAI 2007: Proceedings of the Nineteenth Annual Innovative Applications of Artificial Intelligence Conference*, 1819–1825 (2007).
  - [23] Aaltonen, T. *et al.* Measurement of the top quark mass with dilepton events selected using neuroevolution at CDF. *Phys.Rev.Lett.* **102**, 152001 (2009).
  - [24] Bergstra, J. *et al.* Theano: a CPU and GPU math expression compiler. In *Proceedings of the Python for Scientific Computing Conference (SciPy)* (Austin, TX, 2010). Oral Presentation.
  - [25] Goodfellow, I. J. *et al.* Pylearn2: a machine learning research library. *arXiv preprint arXiv:1308.4214* (2013).



(a)

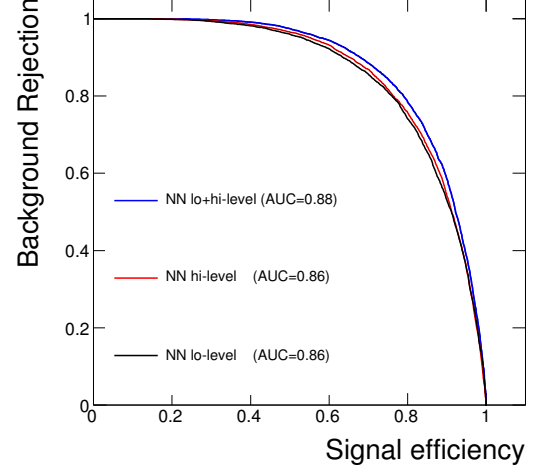


(b)

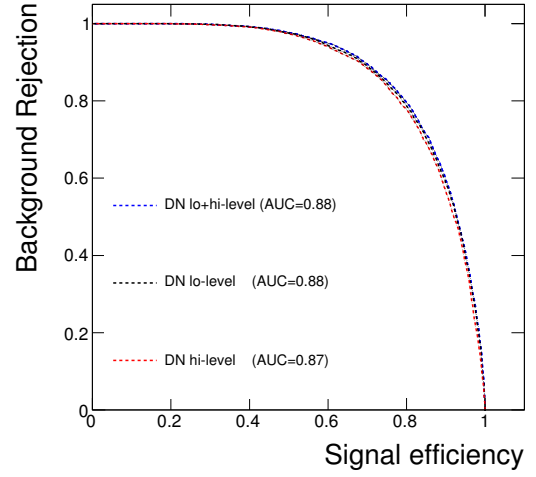


(c)

Supplementary Figure 1: **Details of performance for Higgs benchmark.** For the HIGGS benchmark, comparison of background rejection versus signal efficiency for the low-level features (a), high-level features (b) and complete set of features (c) using traditional and deep learning methods.

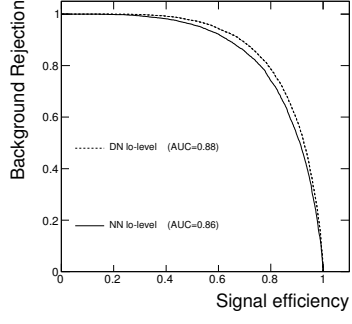


(a)

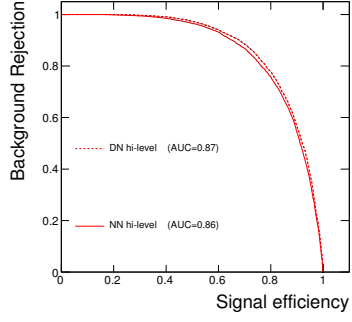


(b)

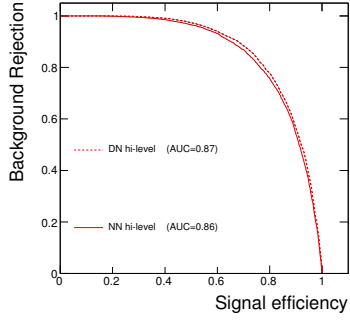
Supplementary Figure 2: **Performance for SUSY benchmark.** Comparison of background rejection versus signal efficiency for the traditional learning method (a) and the deep learning method (b) using the low-level features, the high-level features and the complete set of features.



(a)



(b)



(c)

Supplementary Figure 3: **Details of performance for SUSY benchmark.** Comparison of background rejection versus signal efficiency for the low-level features (a), high-level features (b) and complete set of features (c) using traditional and deep learning methods.

Supplementary Table 1: **Hyper-parameter choices for deep networks.** Shown are the considered value for each hyper-parameter.

Hyper-parameters	Choices
Depth	2, 3, 4, 5, 6 layers
Hidden units per layer	100, 200, 300, 500
Initial learning rate	0.01, 0.05
Weight decay	0, $10^{-5}$
Pre-training	none, autoencoder

Supplementary Table 2: **Hyper-parameter choices for shallow networks.** Shown are the considered value for each hyper-parameter.

Hyper-parameters	Choices
Hidden units	300, 1000, 2000, 10000
Initial learning rate	0.05, 0.005, 0.0005
Weight decay	0, $10^{-5}$

Supplementary Table 3: **Study of network size and depth.** Comparison of shallow networks with different numbers of hidden units (single hidden layer), and deep networks with varying hidden layers in terms of the Area Under the ROC Curve (AUC) for the HIGGS benchmark. The deep networks have 300 units in each hidden layer.

Technique	AUC		
	Low-level	High-level	Complete
NN 300-hidden	0.733	0.777	0.816
NN 1000-hidden	0.788	0.783	0.841
NN 2000-hidden	0.787	0.788	0.842
NN 10000-hidden	0.790	0.789	0.841
DN 3 layers	0.836	0.791	0.850
DN 4 layers	0.868	0.797	0.872
DN 5 layers	0.880	0.800	0.885
DN 6 layers	0.888	0.799	0.893



Supplementary Table 4: **Error Analysis.** Mean Squared Error (MSE) of networks trained to compute the seven high-level features from the 21 low-level features on HIGGS. All networks were trained using the same hyperparameters as the deep networks in the previous experiments, except with linear output units, no momentum, and a smaller initial learning rate of 0.005.

Technique	Feature Regression MSE
Linear Regression	0.1468
NN	0.0885
DN 3 layers	0.0821
DN 4 layers	0.0818
DN 5 layers	0.0815
DN 6 layers	0.0812