

# The Interplay of Machine Learning–based Resonant Anomaly Detection Methods

---

Tobias Golling,<sup>a</sup> Gregor Kasieczka,<sup>b</sup> Claudius Krause,<sup>c</sup> Radha Mastandrea,<sup>d,e</sup> Benjamin Nachman,<sup>e,f</sup> John Andrew Raine,<sup>a</sup> Debajyoti Sengupta,<sup>a</sup> David Shih,<sup>g</sup> and Manuel Sommerhalder<sup>b</sup>

<sup>a</sup>*Département de physique nucléaire et corpusculaire, Université de Genève, 1211 Genève, Switzerland*

<sup>b</sup>*Institut für Experimentalphysik, Universität Hamburg, 22761 Hamburg, Germany*

<sup>c</sup>*Institut für Theoretische Physik, Universität Heidelberg, 69120 Heidelberg, Germany*

<sup>d</sup>*Department of Physics, University of California, Berkeley, CA 94720, USA*

<sup>e</sup>*Physics Division, Lawrence Berkeley National Laboratory, Berkeley, CA 94720, USA*

<sup>f</sup>*Berkeley Institute for Data Science, University of California, Berkeley, CA 94720, USA*

<sup>g</sup>*NHETC, Dept. of Physics and Astronomy, Rutgers University, Piscataway, NJ 08854, USA*

*E-mail:* [tobias.golling@unige.ch](mailto:tobias.golling@unige.ch), [gregor.kasieczka@uni-hamburg.de](mailto:gregor.kasieczka@uni-hamburg.de),  
[claudius.krause@thphys.uni-heidelberg.de](mailto:claudius.krause@thphys.uni-heidelberg.de), [rmastand@berkeley.edu](mailto:rmastand@berkeley.edu), [bpnachman@lbl.gov](mailto:bpnachman@lbl.gov),  
[john.raine@unige.ch](mailto:john.raine@unige.ch), [debajyoti.sengupta@unige.ch](mailto:debajyoti.sengupta@unige.ch), [shih@physics.rutgers.edu](mailto:shih@physics.rutgers.edu),  
[manuel.sommerhalder@uni-hamburg.de](mailto:manuel.sommerhalder@uni-hamburg.de)

**ABSTRACT:** Machine learning–based anomaly detection (AD) methods are promising tools for extending the coverage of searches for physics beyond the Standard Model (BSM). One class of AD methods that has received significant attention is resonant anomaly detection, where the BSM physics is assumed to be localized in at least one known variable. While there have been many methods proposed to identify such a BSM signal that make use of simulated or detected data in different ways, there has not yet been a study of the methods’ complementarity. To this end, we address two questions. First, in the absence of any signal, do different methods pick the same events as signal-like? If not, then we can significantly reduce the false-positive rate by comparing different methods on the same dataset. Second, if there is a signal, are different methods fully correlated? Even if their maximum performance is the same, since we do not know how much signal is present, it may be beneficial to combine approaches. Using the Large Hadron Collider (LHC) Olympics dataset, we provide quantitative answers to these questions. We find that there are significant gains possible by combining multiple methods, which will strengthen the search program at the LHC and beyond.

---

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Methodology</b>	<b>2</b>
2.1	Overview of resonant anomaly detection	2
2.2	Dataset	4
2.3	Classifier architecture	6
<b>3</b>	<b>Contrasting the synthetic SM samples</b>	<b>7</b>
3.1	Background-only case	7
3.2	Adding in signal	11
<b>4</b>	<b>Combining the samples</b>	<b>12</b>
<b>5</b>	<b>Conclusions</b>	<b>15</b>
<b>A</b>	<b>Robustness of classifier scores against network initialization</b>	<b>20</b>
A.1	Binary classifier initialization	20
A.2	Generative network initialization	20
A.3	Using non-robustness to indicate a breakdown of CWOLA	21
<b>B</b>	<b>Additional plots</b>	<b>25</b>

---

## 1 Introduction

Since the observation of the Higgs Boson in 2012 at the Large Hadron Collider (LHC) [1, 2], no new fundamental particle has been observed. This is not for lack of effort: theoretical models involving supersymmetric particles, dark matter candidates, or heavy matter generations abound, informing past, current, and planned analyses at the LHC [3–9]. Given that such past searches for specific alternatives to the Standard Model (SM) have been unsuccessful, there has been a push to run broader, model-agnostic searches for new physics in parallel. In particular, machine learning (ML) has enabled many new search strategies [10–12].

One of the most popular and well-motivated search strategies for evidence of physics beyond the Standard Model is *resonant anomaly detection*. In such investigations, the new physics signal is expected to take the form of a new particle, i.e. a resonance with respect to a mass-like event variable. The search strategy then involves looking for a localized excess of these new physics events with respect to the SM background.

There now exist many ML methods for resonant anomaly detection (AD)<sup>1</sup> with comparable sensitivities [15–27], some of which have also been applied to data [28–31]. These methods have largely been developed independently of each other, with different strengths and weaknesses. However, there

---

<sup>1</sup>We are not counting generic AD methods applied to the resonant case, see e.g. the recent ATLAS results [13, 14] and method papers they cite.

has not yet been a thorough study of the complementarity of these techniques. In particular, we want to ask the questions: can we improve signal sensitivity by combining these methods? Can we improve robustness in the background-only case by combination? Do these different methods classify the same things as “signal-like” for background and signal events?

In this paper, we evaluate a selection of these resonant AD methods on equal footing, using an identical methodological setup for each. In Sec. 2, we provide a more detailed background of the resonant AD procedure and introduce the four detection methods that we will consider in this paper. In Sec. 3, we consider how similar the detection methods are to each other, gauging whether different methods pick up distinct components of the phase space of resonant anomalies. In Sec. 4, we combine the four sample generation methods with the goal of increasing sensitivity for a resonant AD task. We conclude in Sec. 5, suggesting avenues for further exploration.

As a word of caution: this study is not meant to be an exhaustive summary for machine learning-enhanced anomaly detection across all signals and setups. For illustrative purposes, we focus on one well-studied signal model and signal region and compare our findings with the existing literature. Practitioners should examine different methods in their own region of phase space.

## 2 Methodology

### 2.1 Overview of resonant anomaly detection

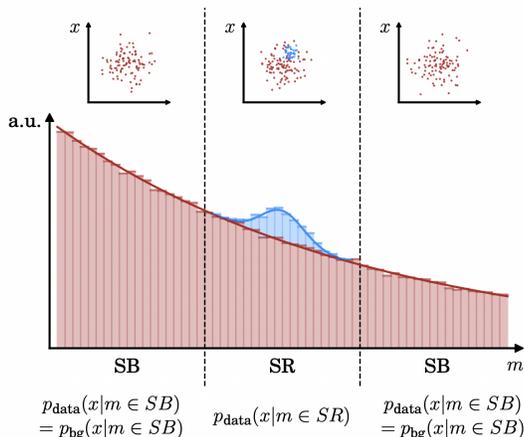
The goal of resonant AD (illustrated schematically in Fig. 1) is to find an excess of beyond-the-Standard Model (BSM) events that are localized in some event variable  $m$  (usually, a mass-like feature). The BSM signal thus corresponds to a new particle with a nonzero  $m$ , expected to be reconstructed within a *signal region*, SR, defined as an interval in  $m$ . In particular, the search makes use of a set of other (i.e. non- $m$ ) features in order to elevate the sensitivity above that of a standard bump hunt. Importantly, the excess events must be observed with respect to a SM background, but this background is nontrivial to construct: using out-of-the-box simulated data is not ideal given the numerous necessary approximations made for the hard-process, showering, and detector simulation steps; using actual data from outside of the signal region (or in *sideband regions*, SB) requires the analysis to only use event features that are statistically independent from the mass variable [15, 16, 19].

An alternative strategy is to construct a set of *synthetic SM samples*, or events that are representative of the SM background process in the same mass space as the BSM events. A binary classifier trained to discriminate the synthetic samples from detected data is then the optimal classifier for discriminating SM background from the new physics (see Ref. [32]), so long as the synthetic samples are indeed a faithful representation of actual SM (i.e. not containing any events derived from a resonant anomaly) events in the probed mass range.

In recent years, there has been much work done on developing procedures to construct such synthetic SM samples. While there now exist many varied methods for sample construction, the vast majority<sup>2</sup> of them can be characterized based on two properties of their construction. First, generation methods can be *data-exclusive* or *simulation-assisted*: data-exclusive methods generate synthetic samples by making use of collider data from the SB mass regions, where the data are far enough from any BSM signal to be treated as representative of SM background; simulation-assisted methods will also use an auxiliary dataset of simulated background-only collisions. Second, generation methods exploit machine learning techniques through either *likelihood(-ratio) learning* or *feature morphing*:

---

<sup>2</sup>The vast majority, but not all. For examples of methods that cannot be so neatly classified on these two characters, see [19, 21, 33, 34].



**Figure 1:** A schematic of the resonant anomaly detection motivation. The goal is to observe an excess of signal (blue) events above a background (red) process. The signal is localized in  $m$  to a signal region (SR), and a model for background can be derived from data in the sidebands (SB) regions. Typically, the signal-background discrimination task makes use of features other than  $m$ . Figure is taken from [22].

methods can either learn the likelihood(-ratio) of an SM-only dataset (this can be either from the auxiliary simulated dataset or the SB data) and interpolate this likelihood into the SR; alternatively, methods can morph features from said background-only regions into the SR.

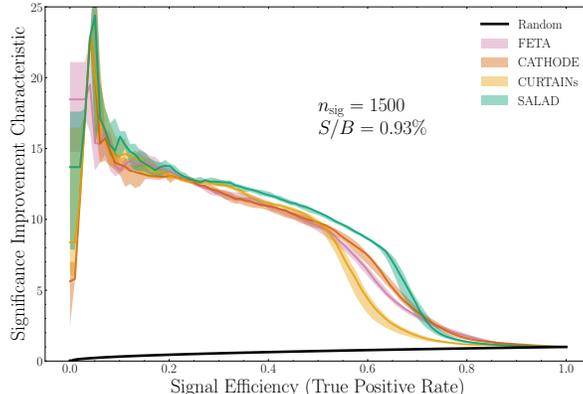
In this paper, we will consider the four methods shown in Table 1, which span this “character space” of methods for resonant AD.

	Simulation-assisted	Data-exclusive
Likelihood learning	SALAD [17]	(LA)CATHODE [22, 24]
Feature morphing	FETA [26]	CURTAINS [27, 35]

**Table 1:** Many methods for constructing Standard Model background templates for resonant anomaly detection can be classified on two axes: on the horizontal, usage of an auxiliary dataset (simulation); on the vertical, how non-signal region Standard Model background processes are morphed into signal region Standard Model template samples.

We provide a brief summary of the four methods considered here.

- SALAD: Simulation Assisted Likelihood-free Anomaly Detection [17] trains a binary classifier to discriminate simulated SM events from detected SM events in the SB (background-only) region, then uses the classifier to reweight simulated SM events in the SR. These reweighted events comprise the synthetic SM samples.
- CATHODE: Classifying Anomalies THrough Outer Density Estimation [22] trains a normalizing flow-based probability density estimator to model detected data in SB, then interpolates the distribution into the SR. A set of events drawn from the interpolated distribution comprises the synthetic SM samples.



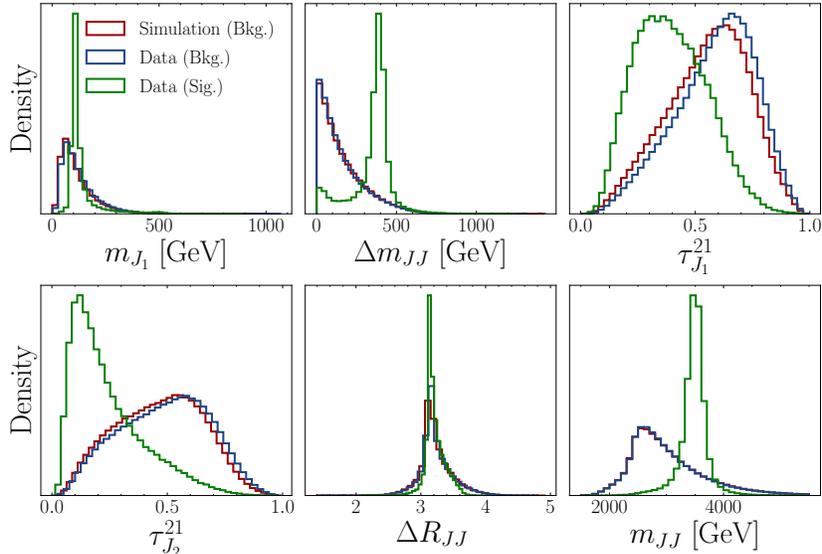
**Figure 2:** Significance Improvement Characteristic for a binary classifier trained to discriminate each synthetic background generation method’s samples from signal-contaminated data. Below this signal injection ( $n_{\text{sig}} = 1500$ ), methods perform less equally. For readability, we withhold a description of the classifier architectures and ensembling choices until Sec. 2.3.

- CURTAINS: Constructing Unobserved Regions by Transforming Adjacent INtervals [27, 35] trains a normalizing flow-based transport function to morph detected data between low- and high-mass SB, then applies the flow to map from SB into the SR. These morphed samples comprise the synthetic SM samples.
- FETA: Flow-Enhanced Transportation for Anomaly detection [26] trains a normalizing flow-based transport function to morph SM simulation in SB to detected data in background-dominated SB, then applies the model to SM simulation in the SR. These morphed samples comprise the synthetic SM samples.

Our goal is then to explore how the synthetic SM samples generated by each of these methods perform in resonant AD tasks, focusing on their relative performances in addition to their absolute performances. In fact, all four methods are comparable at picking up on signal contaminations of  $\sim 0.93\%$  and above: in Fig. 2, we plot the significance improvement characteristic (SIC) as a function of the signal efficiency. Broadly speaking, the SIC corresponds to the multiplicative factor by which a signal significance would improve by making a well-motivated cut on the data; a classifier that is ideally suited to discriminating signal from background would have a high SIC at all signal efficiencies.

## 2.2 Dataset

We use the LHC Olympics 2020 R&D dataset [11, 36], which consists of 1,000,000 background events comprised of QCD dijet production, together with 100,000 signal events from a  $Z'$  resonance at 3.5 TeV, decaying to scalars  $X$  and  $Y$  at 500 GeV and 100 GeV respectively, which then each decay to quark pairs. Since the  $X$  and  $Y$  scalars are highly boosted, their decay products are highly collimated and form large-radius jets. For the main resonant feature, we use the dijet invariant mass ( $m = m_{JJ}$ ) which should reconstruct the  $Z'$  mass for the signal. Events are required to have at least one large-radius anti- $k_T$  [37, 38] jet ( $R = 1$ ) with a  $p_T$  threshold of 1.2 TeV. Each event contains up to 700 particles with three degrees of freedom  $p_T$ ,  $\eta$ ,  $\phi$ . The events are generated with PYTHIA 8.219 [39, 40] and DELPHES 3.4.1 [41]. Also included in the LHC Olympics dataset is a set of 1,000,000 HERWIG++ [42] QCD dijet events generated using the same tunes, which is used for the simulation-assisted approaches



**Figure 3:** The 5-dimensional feature space of dijet collision events used in this resonant AD study. We also show a 6th feature, a mass-like event variable that is used to define a signal region (SR) and sidebands regions (SB).

(SALAD and FETA) as the auxiliary “simulated” data (SIM). We therefore denote the PYTHIA events as “data” (DAT).

In addition to the LHC Olympics dataset, we make use of two auxiliary sets of QCD dijet events generated using the same tunes as the LHC Olympics dataset, but only in the  $m_{JJ}$  region [3.3, 3.7] TeV (our chosen signal region). The first set, consisting of 1,000,000 HERWIG++-based events, is used to generate additional SALAD samples. The second set, consisting of 320k PYTHIA-based events, is used for testing the synthetic samples.

We choose a feature space of six dijet observables  $m_{J_1}$ ,  $\Delta m_{JJ}$ ,  $\tau_{J_1}^{21}$ ,  $\tau_{J_2}^{21}$ ,  $\Delta R_{JJ}$ , and  $m_{JJ}$  (see Figs. 3 and 4). This last feature is our mass-like event variable, which we use to define a SR spanning [3.3, 3.7] TeV. For our sidebands, we use the full amount of available data, down to 1.5 TeV and up to 5.5 TeV.

Synthetic samples are generated using the procedure outlined in each method’s respective paper. However, in an attempt to equalize the quality of the samples as much as possible, all methods are given the same training and validation sets of simulation and data, which is an 80-20 split of the full sidebands data, i.e. all data in the range  $([1.5, 5.5] \setminus [3.3, 3.7])$  TeV. These training and validation event counts, as well as the number of synthetic SM samples generated for each method, are shown in Table 2.

The CATHODE, CURTAINS, and FETA methods all involve training a generative model (i.e. a normalizing flow) to learn an underlying probability distribution. This allows us to reduce statistical uncertainties of the synthetic SM samples by *oversampling* from the generative models. For CATHODE, we can sample from the normalizing flow that has been interpolated into the SR as many times as we want; for CURTAINS (or FETA), we can similarly oversample from the normalizing flows that learn the densities of SB data (or SR simulation). In this study, we use the oversampling factor (also shown in Table 2) that was shown to saturate each model’s performance in its respective paper.

Method	Training data	Validation data	# samples	Oversampling
SALAD	793k SIM, 696k DAT	198K SIM, 174K DAT	1,045k	N/A
CATHODE	696k DAT	174K DAT	400k	3
CURTAINS	373k DAT	93k DAT	1,887k	4
FETA	793k SIM, 696k DAT	198K SIM, 174K DAT	732k	6

**Table 2:** Numerical breakdown of the events used to construct the given number of synthetic SM samples (in the SR) for the four machine learning methods considered in this report. The CURTAINS method uses a slightly narrower SB region of [2.7, 4.5] TeV to avoid transforming events across the  $m_{JJ}$  turn-on region border. The SALAD samples are generated by applying the learned weights to an additional, much larger set of HERWIG++ simulated SM events not contained in the LHC Olympics dataset. Note that CATHODE and CURTAINS are data-exclusive (i.e. fully data-driven), using only the the “detected” (DAT) PYTHIA set, while SALAD and FETA require an auxiliary “simulated” (SIM) HERWIG++ set.

### 2.3 Classifier architecture

As stated earlier, a set of synthetic SM background samples can be used for a resonant AD task by training a binary classifier to discriminate the samples from a set of detected SR data. If the SR data contains a nonzero percentage of BSM events, then the binary classifier should be able to pick up on this difference. In fact, that classifier (in the limit of infinite data and arbitrarily flexible training / architecture) is the optimal one for distinguishing SM background from the BSM signal [32].

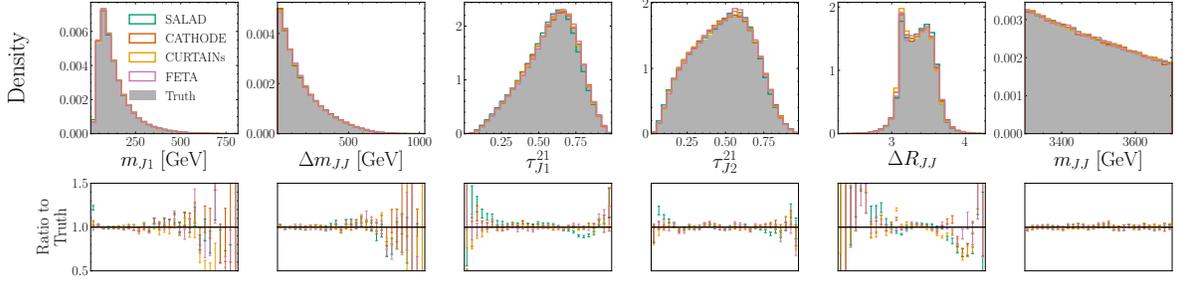
Concretely, the discrimination task in this paper is between a given method’s synthetic SM samples and a set of 121k SR DAT background +  $n_{\text{sig}}$  SR DAT signal ( $Z'$ ) events, where  $n_{\text{sig}}$  is a controllable parameter, the number of injected signal events in the SR.

To evaluate the similarity between two datasets of events, we use a binary classifier network consisting of 3 linear layers of size (64, 32, 1). The network uses rectified linear unit (ReLU) activations and is optimized using ADAM [43] on the Binary Cross Entropy loss. We train using a 5-fold cross validation scheme, keeping the network with the best validation loss. Each fold is trained with a batch size of 128<sup>3</sup> and a learning rate of  $10^{-3}$  for up to 100 epochs with a patience of 5 epochs. All hyperparameters were optimized via manual tuning to give the best possible significance improvement characteristic curves (introduced formally in Sec. 4)<sup>4</sup>. All errorbars and errorbands in plots come from retraining the binary classifier with different random seeds. Note that we do not vary the random seeds for training the architectures that generate each set of synthetic samples, but we find these effects to be small in App. A.

All figures and summary statistics are generated by evaluating the trained binary classifier networks on a “standard test set” of 20,000 signal and 320,000 background dijet events, which are not used at any time during the training procedure. Using a larger number of background events for evaluation allows for smooth summary statistics at low signal efficiencies, which is expected to be the relevant region for resonant AD. Unless explicitly stated otherwise, all plots are score-averaged for each method, i.e. the plots are derived from the average of classifier scores over 10 independent runs.

<sup>3</sup>Interestingly, we found that the choice of batch size significantly affected the classifier, with larger batch sizes leading to sizeable drop-offs in performance.

<sup>4</sup>Note that since the binary classifier is trained on signal region data, the network is no longer signal-agnostic.



**Figure 4:** Distributions of synthetic SM backgrounds generated by each method compared to data with  $n_{\text{sig}} = 0$ .

Method	$m_{J1}$	$\Delta m_{JJ}$	$\tau_{J1}^{21}$	$\tau_{J2}^{21}$	$\Delta R_{JJ}$	$m_{JJ}$
SALAD	0.00775	0.00501	0.02229	0.00610	0.01205	0.00215
CATHODE	0.00405	0.00450	0.00597	0.00534	0.00755	0.00228
CURTAINS	0.00325	0.00255	0.00238	0.00214	0.02122	0.00353
FETA	0.00605	0.00352	0.00588	0.00536	0.00725	0.00386

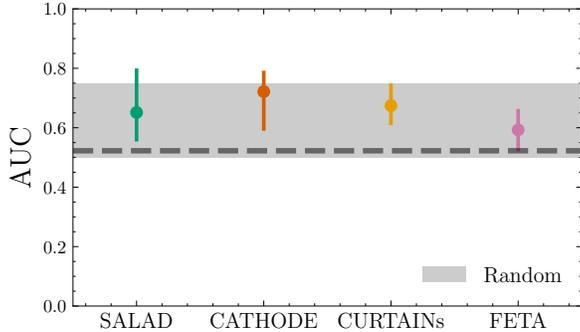
**Table 3:** Kolmogorov-Smirnov test statistics between each method’s marginal distribution and the truth’s marginal distribution. Larger test statistics indicate a greater difference between two distributions, as gauged by the maximum difference in the empirical cumulative distribution functions.

### 3 Contrasting the synthetic SM samples

#### 3.1 Background-only case

As a first test of the synthetic SM samples created by each of the four generation methods, we compare their distributions to background-only SM data, i.e. with a signal injection  $n_{\text{sig}} = 0$ . These distributions are shown in Fig. 4. We see that at this level all four methods reproduce the true distribution well. In particular, the ratios of marginals for all methods are all close to 1, which indicates that any differences between the sample generation methods and truth cannot be ascribed to a single observable. As a quantitative assessment of the marginal distributions, in Table 3, we provide the Kolmogorov-Smirnov (KS) test statistics for the marginal distributions between each method and the truth. The KS test statistic is defined as the supremum of the differences between two distributions’ empirical cumulative distribution functions, and can therefore provide a gauge of how different two distributions are.

As a next test of the synthetic SM samples created by each of the four generation methods, we analyze classifiers trained to discriminate synthetic SM background from background-only SM data, i.e. with a signal injection  $n_{\text{sig}} = 0$ . Given that there is no BSM signal present in the training data, differences in classifier performances come down to the “nature” of the synthetic SM samples for each method. All of the methods are given the same training and validation data sets, so the statistical fluctuations from the input data should be correlated between methods. Further, all of the binary classifiers are evaluated with the same random seed, so the network initializations should be identical in that respect. However, there are differences stemming from the initialization of the neural networks for each *method*, as well as from the differences of the methods themselves. These differences might be expected to decorrelate the classifier scores.



**Figure 5:** Receiver operating characteristic area-under-the-curves (AUC) for binary classifiers trained to discriminate each method’s synthetic samples from data with  $n_{\text{sig}} = 0$ . The table summarizes 100 classifier runs with different random seeds, with errorbars showing a 68-percentile spread. Also plotted is the spread of 100 random classifiers, with the thick dashed line showing the median of those runs. No score averaging has been done for this plot. To see the AUC spreads corresponding to the combination of samples, see Fig. 18.

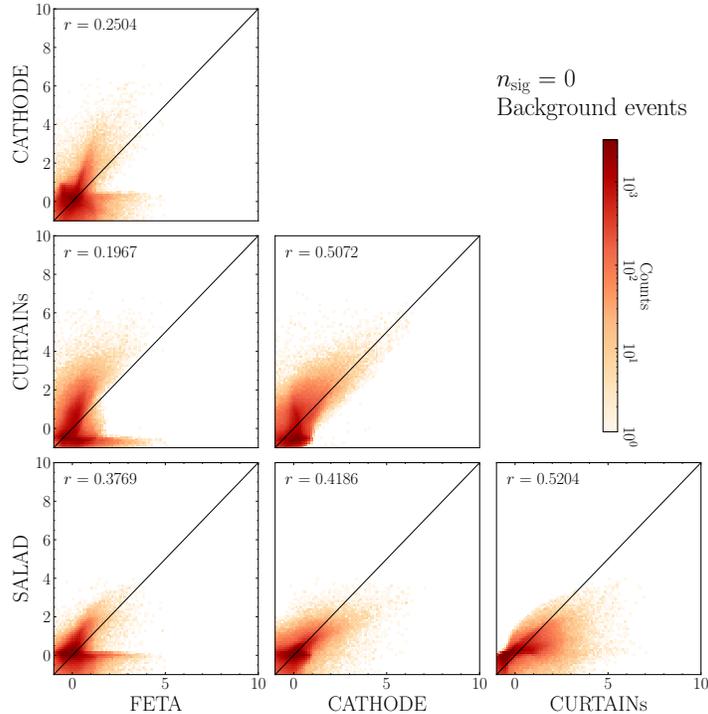
We provide the receiver operating characteristic area-under-the-curves (AUC) for such classifiers in Fig. 5. Also plotted is the AUC spread derived from training a classifier to discriminate truth from truth, which represents the spread of a random classifier given the set of different network initializations and the fact that the network is not infinitely powerful. The ROC spread of each individual method is consistent with that of the spread of a random classifier, again providing evidence that the nature of the synthetic samples is truth- (SM-) like.

Figure 6 plots the classifier scores, averaged over 10 classifier runs, as evaluated on the test set’s background events<sup>5</sup>. Note that we plot the *standardized* scores, since the binary classifier is trained to flag the most anomalous events with the highest scores. We also focus on the first quadrant of the coordinate plane, corresponding to the “anomaly regions” of the plots, or the highest regions in (standardized) score space where the classifier-flagged anomalies are expected to lie. In general, the classifier scores for the networks trained to discriminate detected data from each synthetic sample generation method do not appear to agree across methods: the scores are, for the most part, uncorrelated when evaluated on true background events.

As a next task, we quantify the similarity across sample generation methods of events that are deemed “signal-like” by the binary classifier. As our similarity metric, we consider the background events in the standard test set with classifier scores in the top  $p$  percentile, i.e. the background events classified as the most “signal-like”, or most different to the synthetic background samples. Within the scope of an anomaly detection search, these top  $p$ -percentile events are exactly those that correspond to high-score, likely-to-be-anomalous events.

For each percentile  $p$ , we find the set of the top  $p$  events independently for the classifiers trained on all four methods. We then calculate the overlaps between the sets of selected events between each pair of sample generation methods, then across all four methods. Note that for fully independent event selection methods, we would expect a fraction  $p$  of shared events by chance; therefore for easier visualization, we plot the excess overlap with respect to this random baseline in Fig. 7 (subtracting the correct baseline for the overlap of all four methods). If the excess overlap is 0, then the performance

<sup>5</sup>See App. B for the corresponding plots evaluated on the test set’s signal events.



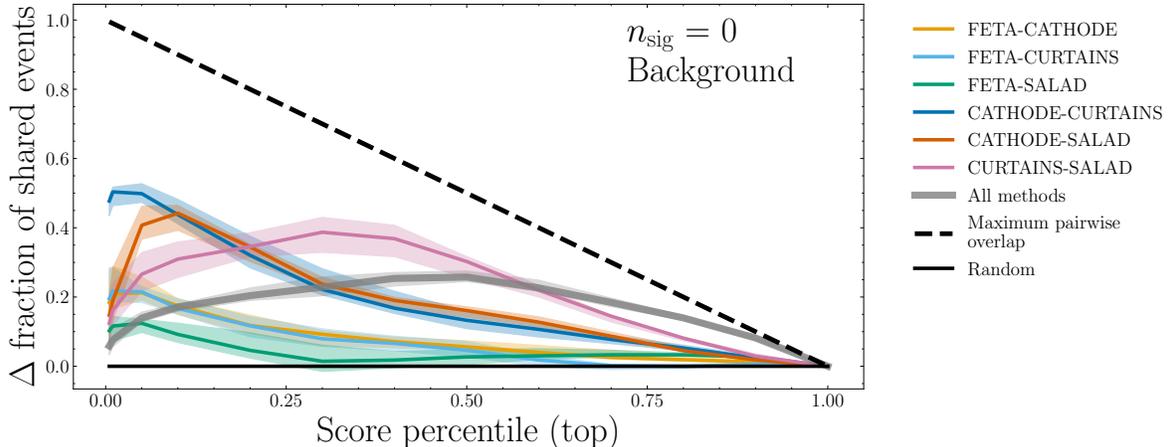
**Figure 6:** Scores for background events for a binary classifier trained to discriminate synthetic SM background from data with  $n_{\text{sig}} = 0$ ). Each axis represents a different method of SM sample generation.  $r$  denotes the Pearson correlation coefficient, computed over all samples (not just the upper right quadrant). The scores are standardized so as to make it clear what events the binary classifier flags as the most anomalous. For each method, scores are averaged over 10 classifier runs.

overlap is no more than expected by random chance.

For most combinations of synthetic sample generation methods, the amount of overlap between any two sample generation methods is slightly greater than what we would expect between two uncorrelated sets of random numbers, especially for larger score percentiles. However, there appears to be a large amount of similarity between CATHODE and SALAD, and between CATHODE and CURTAINS, especially at small percentiles (all of which implies a somewhat smaller amount of similarity between FETA and any other method). There is also large degree of similarity between CURTAINS and SALAD, especially at larger percentiles.

Another way to study the quality of the background-only samples relative to each other utilizes a multiclass classifier. This method of comparing different generative models was first introduced in [44] in the context of up-sampling hydrodynamical galaxy simulations. (See also [45] for a subsequent application to comparing generative models for fast calorimeter simulation.) We use the same classifier architecture as before and only modify the output layer to yield 4 (softmaxed) numbers, which we interpret as the probabilities of the input samples belonging to one of the four methods<sup>6</sup>. We also use a larger batch size of 1000, which we found necessary in order to get repeatable results. We use a subset of 400,000 samples from each of the four methods to train (60%), test (20%), and evaluate

<sup>6</sup>We checked a larger classifier architecture with more nodes and additional dropout, but the averaged results did not change with respect to the ones reported below.



**Figure 7:** Fractional overlap, with respect to a random-choice baseline, of the  $p$ th percentile of true background events classified as the most “signal-like” between different methods of synthetic SM sample generation. Errorbands show a 68-percentile spread across the median and come from 100 repetitions of training 5-fold classifiers on the associated methods with different random seeds and ensembling scores over 10 repetitions. Note that the left-most point corresponds to a percentile of  $p = 0.005$ .

(20%) the classifier. Samples from SALAD get their appropriate sample weight in training, testing, and evaluation. Since the average of these weights is about 0.98, we add an additional class weight of  $1/0.98 = 1.021$  to the SALAD samples to correct for the small imbalance.

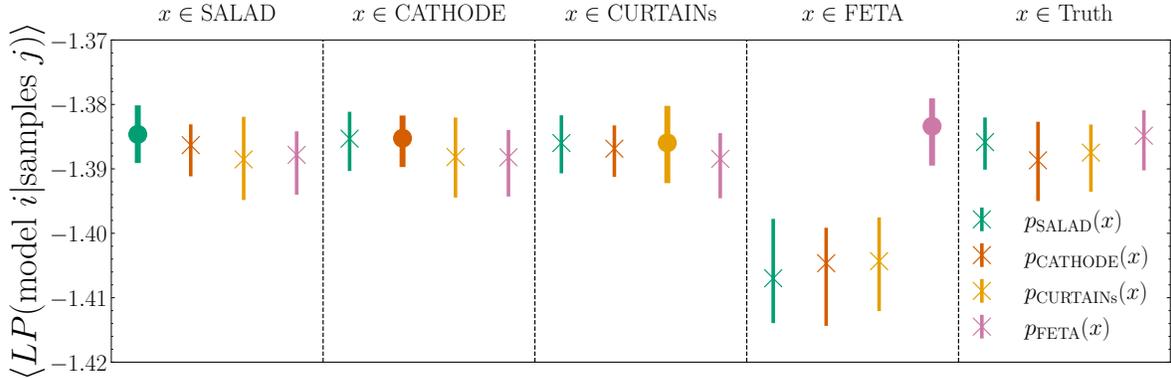
For final evaluation and comparison of the samples, we consider the average of the log posterior [44], which is defined as

$$LP(\text{model } i | \text{samples } j) = \frac{1}{N} \sum_{x_k \in j} \omega_k \log p_{\text{model } i}(x_k), \quad (3.1)$$

evaluated on the held-out test datasets. Here, the sum includes all samples  $x_k$  of the tested model  $j \in$  (SALAD, CATHODE, CURTAINS, FETA),  $\omega_k$  is the sample weight of the sample  $x_k$ , and  $N$  is the number of samples in the set. Since individual runs tend to scatter, we average the log posteriors over 100 independent classifier trainings with different random seeds. A well-trained multiclass classifier should be able to identify the samples belonging to each model, therefore we would expect to have

$$LP(\text{model } i | \text{samples } j = i) > LP(\text{model } i | \text{samples } j \neq i). \quad (3.2)$$

Indeed, this is what we observe in the first four columns of Fig. 8. The probability of belonging to a given model is highest for samples that were generated with that model for SALAD, CATHODE, CURTAINS (albeit only slightly for the latter two), and FETA. These results are consistent with the previous similarity studies: SALAD, CATHODE, and CURTAINS exhibit an above-average degree of similarity with each other, while FETA appears to be more independent. To assess the question which of the methods produces artificial background closest to “truth” (the true SR background SM events), we evaluate the log posterior of Eq. (3.1) for samples from the truth dataset. We see in the right column of Fig. 8 that all four methods are essentially of equivalent quality, with their log posterior scores all well within each other’s error bars. With respect to the truth, the “FETA anomaly” appears to be less pronounced.



**Figure 8:** Average log posteriors  $\langle LP(\text{model } i | \text{samples } j) \rangle$  of the multiclass classifier. The circle markers highlight the case  $i = j$ , and the cross markers indicate the cases  $i \neq j$ . “Truth” designates the SR SM background events. Errorbars show a 68-percentile spread of the  $LP$ s of 100 independent retrainings (no score averaging is carried out).

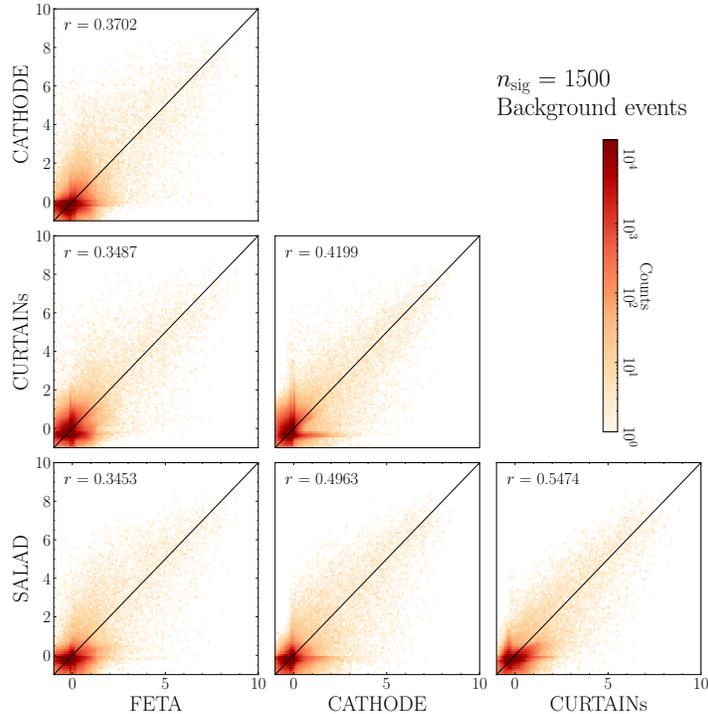
### 3.2 Adding in signal

In Figs. 9 and 10, we provide the scatterplots of the classifier scores evaluated on true background and true signal (respectively) events across different methods, this time for the case with injected signal:  $n_{\text{sig}} = 1500$  ( $S/B = 0.93\%$ ,  $S/\sqrt{B} = 3.24$ ). Note that we fully retrain all classifiers for this new signal injection.

As shown in Fig. 9, for this larger signal injection (as compared with 0 signal injection in Fig. 6), the correlation between classifier scores for background events across synthetic sample generation method is somewhat higher, especially for correlations involving SALAD or FETA. In contrast, Fig. 10 shows that the classifier scores for signal events are highly correlated between any two methods; all methods seem to agree on what anomalous events are. To summarize these results: we see that the classifier scores are rather uncorrelated on background events, but highly correlated on signal events. This might mean that the characteristics (i.e. the 5-dimensional non-mass feature space) of the synthetic background that is created differ non-trivially from method to method; there isn’t overwhelming consensus on how to classify true background. However, all four of the methods produce background that is non-trivially different from true signal, at least different enough that classifiers can reliably distinguish background from signal.

In Fig. 11, we once again plot the overlaps of the top- $p$  percentile most “signal-like” events for a training signal injection of  $n = 1500$ . For background events (Fig. 11a), there is now a sizable amount of overlap between all pairs of methods at  $p \lesssim 0.1$ , though the overlap drops off quickly for larger percentiles. For the signal events (Fig. 11b), there is a significant excess of event overlaps between any two methods down to low-to-mid percentiles. This agrees with intuition: the natures of the synthetic SM samples may differ from method to method, but the hope is that they all differ significantly from a BSM resonance such that they can be used as a suitable background against which to discriminate the resonance. Importantly, there is an excess in event overlaps above random chance between all four methods across the board, at all percentiles.

In Fig. 12, we consider a slightly different view of the percentile overlaps by fixing the percentile of the most “signal-like” events and plotting this as a function of  $n_{\text{sig}}$  in the classifier training set. For the top 5 percentile of the most signal-like true background events, there appears to be slightly increasing



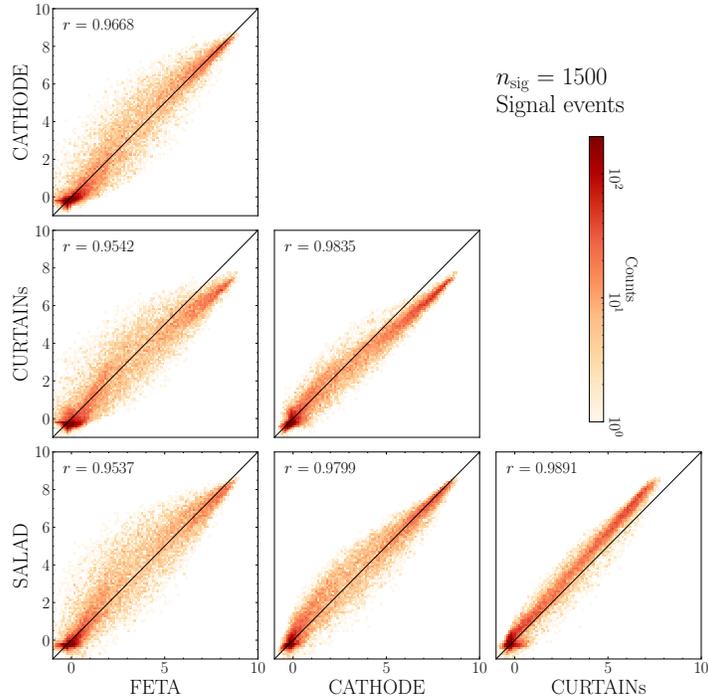
**Figure 9:** Scores for background events for a binary classifier trained to discriminate synthetic SM background from data with  $n_{\text{sig}} = 1500$ ,  $S/\sqrt{B} = 3.24$ ). Each axis represents a different method of SM sample generation.

similarity with the number of injected signal events  $n_{\text{sig}}$  across all four methods, but not between any two methods. For the top 5 percentile of the most signal-like true signal events, the agreement increases with  $n_{\text{sig}}$ , leveling out at  $n_{\text{sig}} \approx 1200$ . Put another way, the four methods considered here agree on what the 5% most anomalous events are when trained to discriminate their own synthetic SM samples from a dataset containing signal injections as low as 0.62%.

## 4 Combining the samples

In this section, we investigate the extent to which *combining* the synthetic samples can provide a more faithful approximation for SM background than taking samples from any of the individual generation methods alone. We have seen previously that classifiers trained to discriminate an individual method's synthetic samples from data tend to agree on what anomalous, signal-like events are more often than random. However, the agreement is not absolute. This might indicate that the events that each method's classifier are flagging as anomalous occupy slightly different parts of phase space. Therefore by combining the synthetic samples, we could hope to be more broadly sensitive to a larger phase space.

There are numerous ways to combine the sample generation methods, as the combination can in principle be done at one of many stages of an analysis. Additionally, one could imagine combining methods in a way that prioritizes one method over the other three. In this section, we will investigate



**Figure 10:** Scores for signal events for a binary classifier trained to discriminate synthetic SM background from data with  $n_{\text{sig}} = 1500$ ,  $S/\sqrt{B} = 3.24$ ). Each axis represents a different method of SM sample generation.

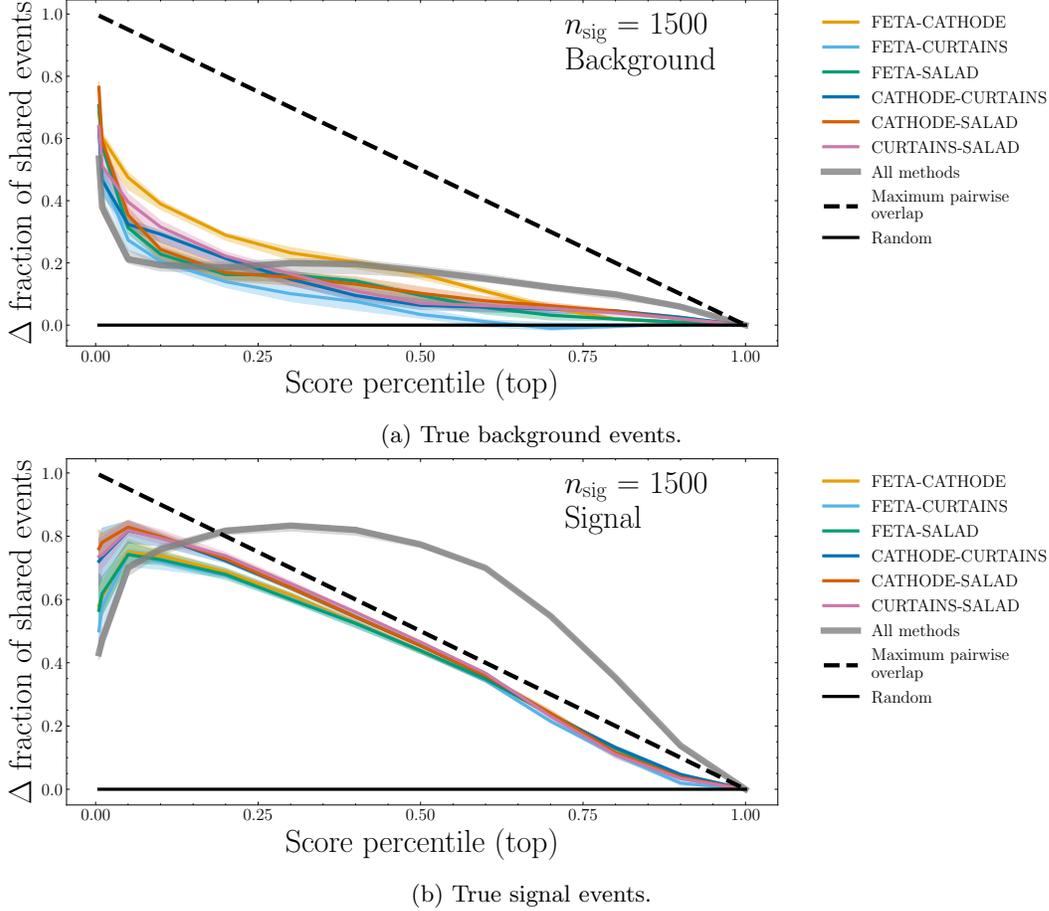
the two most straightforward combinations that weight all four methods in the same way: first at the sample level, and second at the (classifier) score level.

To combine generation methods at the sample level, we take 250k samples from each method so that each contributes equally. We then train a binary classifier to discriminate the combined synthetic samples from the SR data, varying  $n_{\text{sig}}$  from 0 to 1500 (corresponding to  $S/B$  in the SR of 0% to 0.93%). We evaluate each classifier on the standard test set. To combine different methods at the score level, we simply average the score attributed to a given test set event over each of the four methods.

To aggregate classifier runs, we train 100 such binary classifiers, average scores across ensembles of 10 runs, and generate classifier metric curves using the ensembled scores. This has the effect of tightening the errorbands for Figs. 13 and 14, making them easier to parse. For all methods, we apply a further level of aggregation by ensembling over the *generator seed*. In other words, we create three instantiations of each generative ML model, repeat the analysis outlined in this and the previous paragraph, and amass all the classifier metric curves across the instantiations<sup>7</sup>.

In Fig. 13, we provide summary plots across the range of tested  $n_{\text{sig}}$  values. In Fig. 13a, we calculate the classifier significance improvement characteristic (SIC) as a function of the signal efficiency, then take the maximum of the SIC. The  $\text{max}(\text{SIC})$  gives the best multiplicative improvement to signal significance, corresponding to the best-motivated cut (which we do not know a priori). In Fig. 13b, we plot the significance at a background rejection of  $10^3$ , which is less sensitive to the low-signal effi-

<sup>7</sup>This combination of generator seeds was found to give more robust results at the smallest signal injections. We explore this more in App. A.



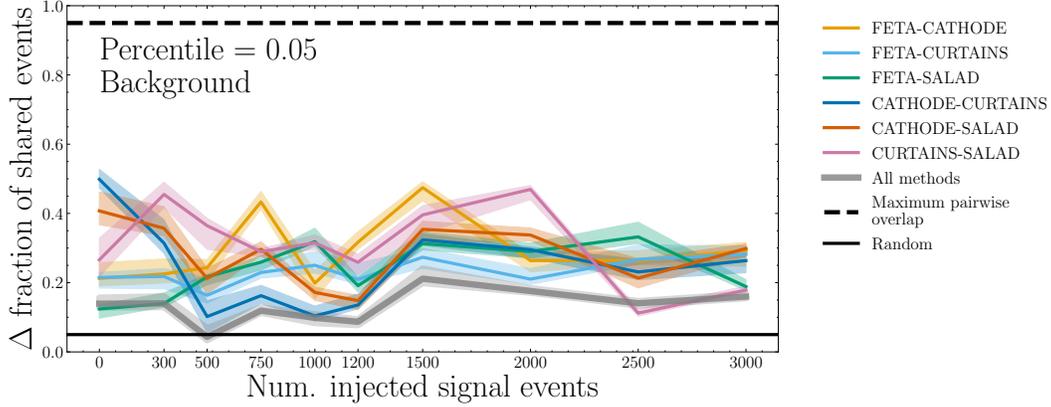
**Figure 11:** Fractional overlap, with respect to a random-choice baseline, of the  $p$ th percentile of events classified as the most “signal-like” between different methods of synthetic SM sample generation.

ciency fluctuations of the max(SIC). Based on these metrics, the median performance of the combined synthetic samples is competitive with, but not necessarily better than, any of the individual sample generation methods; however, the spread of the combined samples is much tighter, implying greater stability.

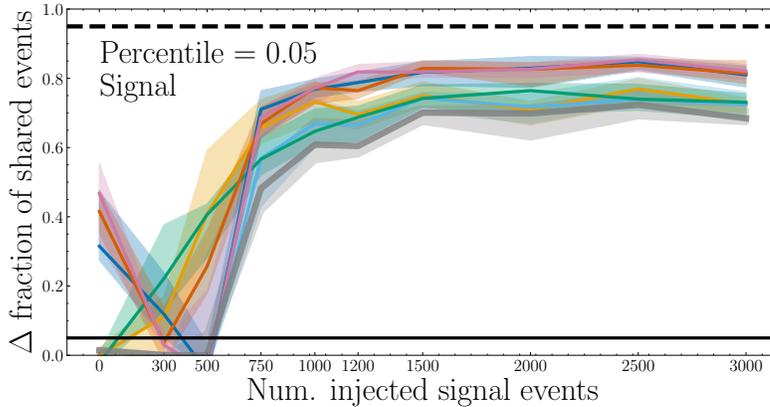
The summary statistics alone may not be the most helpful gauge for performance in an AD task since we do not necessarily know the cut value corresponding to the max(SIC). In Fig. 14, we provide additional summary plots for the lowest signal injection  $n_{\text{sig}} = 750$  ( $S/B = 0.47\%$ ) where using the combined samples leads to an improvement over using any individual method. While the two combination methods (i.e. at the sample and score levels) are comparable, the sample-level combination does appear to give better performance at most signal efficiencies<sup>8</sup>.

Based on these plots, using the combined samples as the SM background leads to a classifier that is uniformly better (with respect to signal efficiency) at detecting the small amount of signal. This implies that when the score cutoff corresponding to the max(SIC) is unknown — as it is in virtually

<sup>8</sup>We provide equivalents to these plots computed at  $n_{\text{sig}} = 500$  in Fig. 21 in App. B, which shows that all methods fail to pick up on the signal.



(a) True background events.



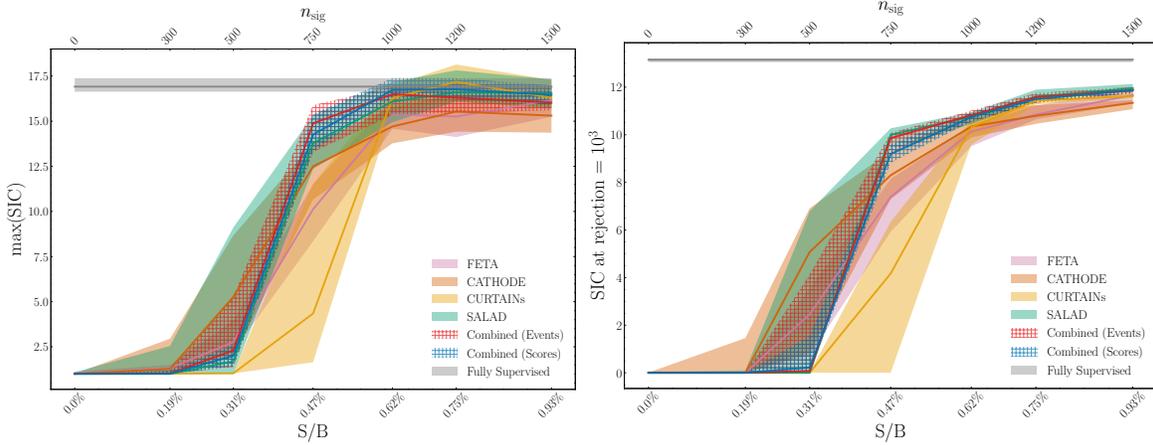
(b) True signal events. For small signal injections ( $n_{\text{sig}} < 750$ ), the negative values are likely a statistical fluctuation due to the small amount of signal present – within errorbands, the values are consistent with 0.

**Figure 12:** Fractional overlap, with respect to a random-choice baseline, of the 5th percentile of events classified as the most “signal-like” between different methods of synthetic SM sample generation, scanning over  $n_{\text{sig}}$ .

all AD tasks — *combining* synthetic SM samples is the optimal strategy.

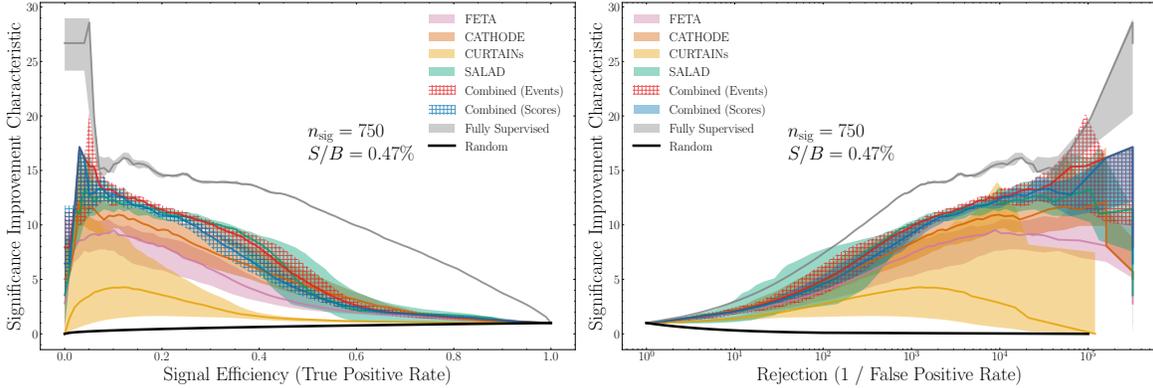
## 5 Conclusions

In this paper, we have explored four conceptually different methods of generating synthetic Standard Model (SM) background samples to be used for resonant anomaly detection (AD) tasks: SALAD, CATHODE, CURTAINS, and FETA. Each method uses a different means of generating a set of synthetic Standard Model samples to be used as a background set for resonant anomaly detection, but all use the same meta-format: shift in some way a sample of background-only events (pulled from an auxiliary dataset or background-only regions in data) to a signal region of interest, and search for a resonant anomaly within that region, such as by estimating an anomaly score based on the data-to-background likelihood ratio in non- $m$  features and cutting on it to enhance on the standard bump hunt procedure.



(a) Maximum of the Significance Improvement Characteristic. (b) Significance Improvement Characteristic at a classifier rejection of  $10^3$ .

**Figure 13:** Various metrics for a classifier trained to discriminate a combination of FETA, CATHODE, and CURTAINS synthetic SM samples from data over a range of  $n_{\text{sig}}$  values. Errorbands show a 68-percentile spread across the median and come from training a 5-fold classifier 100 times with different random seeds, over 3 independent generative model seeds, and ensembling scores over 10 runs.



(a) Significance Improvement Characteristic plotted against the signal efficiency. (b) Significance Improvement Characteristic plotted against classifier rejection.

**Figure 14:** Various classifier metrics for a classifier trained to discriminate a combination of FETA, CATHODE, and CURTAINS synthetic SM samples from data with  $n_{\text{sig}} = 750$ . Errorbands show a 68-percentile spread across the median and come from training a 5-fold classifier 100 times with different random seeds, over 3 independent generative model seeds, ensembling scores over 10 runs, and averaging classifier metrics over the ensembles.

In general, the four construction methods produce synthetic SM samples that perform similarly when used for resonant AD tasks. Binary classifiers trained to discriminate SM samples from SALAD, CATHODE, CURTAINS, and FETA against data (SM background + injected signal) assign scores

to the signal events that are generally correlated. Furthermore, the four methods agree on what the top  $p$  percentile of the most “signal-like” events are. While all methods perform similarly on their own with respect to being able to detect evidence of anomalous events as quantified by their  $\max(\text{SICs})$ , *combining* the four methods allows for a more sensitive AD tool at any given signal efficiency. This is especially useful in practice, when we do not know the optimal score cutoff corresponding to the  $\max(\text{SIC})$ . We find that there is enough evidence to recommend that future AD tasks make use of this combined strategy for generating synthetic SM background samples.

With an eye towards future work: the LHC Olympics dataset is used as a benchmarking tool in the majority of resonant AD R&D, but it represents just one resonant anomaly type out of a vast landscape. It is very possible that the results found in this report do not perfectly extend to other BSM particles, and therefore it would be useful to carry out similar tests of the SM generation methods on vastly different types of signal models. In that respect, it is interesting to remember the background (SM) -only studies in Sec. 3.1, which showed that the four methods considered in this work do seem to produce samples that cover non-overlapping regions of phase space. On a related vein, it would be worthwhile for future studies to explore how to make the anomaly-detecting CWoLa binary classifier signal-agnostic: it is standard in the field (especially for benchmarking studies) to optimize that classifier manually, which adds a degree of model-dependence into the anomaly detection procedure.

Finally, it would also be useful to consider other means of combining the synthetic samples, perhaps at the level of classifier metrics other than at the event-level or score-level, or in ways that prioritize one method in particular.

## Acknowledgements

BN and RM are supported by the U.S. Department of Energy (DOE), Office of Science under contract DE-AC02-05CH11231. RM is additionally supported by the National Science Foundation Graduate Research Fellowship Program under Grant No. DGE 2146752; any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation. The work of DSh was supported by DOE grant DOE-SC0010008. TG, JAR and DSe acknowledge funding through the SNSF Sinergia grant “Robust Deep Density Models for High-Energy Particle Physics and Solar Flare Analysis (RODEM)” with funding number CRSII5\_193716, and the SNSF project grant 200020.212127 “At the two upgrade frontiers: machine learning and the ITk Pixel detector”. GK and MS acknowledge support by the Deutsche Forschungsgemeinschaft under Germany’s Excellence Strategy – EXC 2121 Quantum Universe – 390833306. The work of MS was supported by BMBF grant 05H21GUCC1. CK would like to thank the Baden-Württemberg-Stiftung for financing through the program *Internationale Spitzenforschung*, project *Uncertainties – Teaching AI its Limits* (BWST\_IF2020-010).

## Code availability

All analysis code can be found at [https://github.com/rmastand/synthetic\\_SM\\_AD/tree/main](https://github.com/rmastand/synthetic_SM_AD/tree/main).

## References

- [1] ATLAS collaboration, G. Aad et al., *Observation of a new particle in the search for the Standard Model Higgs boson with the ATLAS detector at the LHC*, *Phys. Lett. B* **716** (2012) 1–29, [[1207.7214](#)].

- [2] CMS collaboration, S. Chatrchyan et al., *Observation of a New Boson at a Mass of 125 GeV with the CMS Experiment at the LHC*, *Phys. Lett. B* **716** (2012) 30–61, [1207.7235].
- [3] ATLAS Collaboration, *Exotic Physics Searches*, 2019.  
<https://twiki.cern.ch/twiki/bin/view/AtlasPublic/ExoticsPublicResults>.
- [4] ATLAS Collaboration, *Supersymmetry searches*, 2019.  
<https://twiki.cern.ch/twiki/bin/view/AtlasPublic/SupersymmetryPublicResults>.
- [5] ATLAS Collaboration, *Higgs and Diboson Searches*, 2019.  
<https://twiki.cern.ch/twiki/bin/view/AtlasPublic/HDBSPublicResults>.
- [6] CMS Collaboration, *CMS Exotica Public Physics Results*, 2019.  
<https://twiki.cern.ch/twiki/bin/view/CMSPublic/PhysicsResultsEXO>.
- [7] CMS Collaboration, *CMS Supersymmetry Physics Results*, 2019.  
<https://twiki.cern.ch/twiki/bin/view/CMSPublic/PhysicsResultsSUS>.
- [8] CMS Collaboration, *CMS Beyond-two-generations (B2G) Public Physics Results*, 2019.  
<https://twiki.cern.ch/twiki/bin/view/CMSPublic/PhysicsResultsB2G>.
- [9] LHCb Collaboration, *Publications of the QCD, Electroweak and Exotica Working Group*, 2019. [http://lhcbproject.web.cern.ch/lhcbproject/Publications/LHCbProjectPublic/Summary\\_QEE.html](http://lhcbproject.web.cern.ch/lhcbproject/Publications/LHCbProjectPublic/Summary_QEE.html).
- [10] G. Karagiorgi, G. Kasieczka, S. Kravitz, B. Nachman and D. Shih, *Machine Learning in the Search for New Fundamental Physics*, 2112.03769.
- [11] G. Kasieczka et al., *The LHC Olympics 2020 a community challenge for anomaly detection in high energy physics*, *Rept. Prog. Phys.* **84** (2021) 124201, [2101.08320].
- [12] T. Aarrestad et al., *The Dark Machines Anomaly Score Challenge: Benchmark Data and Model Independent Event Classification for the Large Hadron Collider*, *SciPost Phys.* **12** (2022) 043, [2105.14027].
- [13] ATLAS collaboration, G. Aad et al., *Anomaly detection search for new resonances decaying into a Higgs boson and a generic new particle X in hadronic final states using  $\sqrt{s} = 13$  TeV pp collisions with the ATLAS detector*, 2306.03637.
- [14] ATLAS collaboration, *Search for new phenomena in two-body invariant mass distributions using unsupervised machine learning for anomaly detection at  $\sqrt{s} = 13$  TeV with the ATLAS detector*, tech. rep., CERN, Geneva, 2023.
- [15] J. H. Collins, K. Howe and B. Nachman, *Anomaly Detection for Resonant New Physics with Machine Learning*, *Phys. Rev. Lett.* **121** (2018) 241803, [1805.02664].
- [16] J. H. Collins, K. Howe and B. Nachman, *Extending the search for new resonances with machine learning*, *Phys. Rev.* **D99** (2019) 014038, [1902.02634].
- [17] A. Andreassen, B. Nachman and D. Shih, *Simulation Assisted Likelihood-free Anomaly Detection*, *Phys. Rev. D* **101** (2020) 095004, [2001.05001].
- [18] B. Nachman and D. Shih, *Anomaly Detection with Density Estimation*, *Phys. Rev. D* **101** (2020) 075042, [2001.04990].
- [19] K. Benkendorfer, L. L. Pottier and B. Nachman, *Simulation-assisted decorrelation for resonant anomaly detection*, *Phys. Rev. D* **104** (2021) 035003, [2009.02205].
- [20] G. Stein, U. Seljak and B. Dai, *Unsupervised in-distribution anomaly detection of new physics through conditional density estimation*, in *34th Conference on Neural Information Processing Systems*, 12, 2020, 2012.11638.

- [21] O. Amram and C. M. Suarez, *Tag n' train: a technique to train improved classifiers on unlabeled data*, *Journal of High Energy Physics* **2021** (jan, 2021) .
- [22] A. Hallin, J. Isaacson, G. Kasieczka, C. Krause, B. Nachman, T. Quadfasel et al., *Classifying anomalies through outer density estimation*, *Phys. Rev. D* **106** (2022) 055006, [[2109.00546](#)].
- [23] J. F. Kamenik and M. Szewc, *Null hypothesis test for anomaly detection*, *Phys. Lett. B* **840** (2023) 137836, [[2210.02226](#)].
- [24] A. Hallin, G. Kasieczka, T. Quadfasel, D. Shih and M. Sommerhalder, *Resonant anomaly detection without background sculpting*, [2210.14924](#).
- [25] M. F. Chen, B. Nachman and F. Sala, *Resonant Anomaly Detection with Multiple Reference Datasets*, [2212.10579](#).
- [26] T. Golling, S. Klein, R. Mastandrea and B. Nachman, *Flow-enhanced transportation for anomaly detection*, *Phys. Rev. D* **107** (2023) 096025, [[2212.11285](#)].
- [27] D. Sengupta, S. Klein, J. A. Raine and T. Golling, *CURTAINs Flows For Flows: Constructing Unobserved Regions with Maximum Likelihood Estimation*, [2305.04646](#).
- [28] ATLAS collaboration, G. Aad et al., *Dijet resonance search with weak supervision using  $\sqrt{s} = 13$  TeV pp collisions in the ATLAS detector*, *Phys. Rev. Lett.* **125** (2020) 131801, [[2005.02983](#)].
- [29] D. Shih, M. R. Buckley, L. Necib and J. Tamanas, *Via Machinae: Searching for Stellar Streams using Unsupervised Machine Learning*, [2104.12789](#).
- [30] D. Shih, M. R. Buckley and L. Necib, *Via Machinae 2.0: Full-Sky, Model-Agnostic Search for Stellar Streams in Gaia DR2*, [2303.01529](#).
- [31] M. Pettee, S. Thanvantri, B. Nachman, D. Shih, M. R. Buckley and J. H. Collins, *Weakly-Supervised Anomaly Detection in the Milky Way*, [2305.03761](#).
- [32] E. M. Metodiev, B. Nachman and J. Thaler, *Classification without labels: Learning from mixed samples in high energy physics*, *JHEP* **10** (2017) 174, [[1708.02949](#)].
- [33] S. Choi, J. Lim and H. Oh, *Data-driven Estimation of Background Distribution through Neural Autoregressive Flows*, [2008.03636](#).
- [34] CMS collaboration, *Evidence for four-top quark production in proton-proton collisions at  $\sqrt{s} = 13$  TeV*, [2303.03864](#).
- [35] J. A. Raine, S. Klein, D. Sengupta and T. Golling, *CURTAINs for your Sliding Window: Constructing Unobserved Regions by Transforming Adjacent Intervals*, [2203.09470](#).
- [36] G. Kasieczka, B. Nachman and D. Shih, *Official Datasets for LHC Olympics 2020 Anomaly Detection Challenge (Version v6) [Data set]*, 2019. <https://doi.org/10.5281/zenodo.4536624>.
- [37] M. Cacciari, G. P. Salam and G. Soyez, *FastJet User Manual*, *Eur. Phys. J. C* **72** (2012) 1896, [[1111.6097](#)].
- [38] M. Cacciari, G. P. Salam and G. Soyez, *The anti- $k_t$  jet clustering algorithm*, *JHEP* **04** (2008) 063, [[0802.1189](#)].
- [39] T. Sjöstrand, S. Mrenna and P. Z. Skands, *PYTHIA 6.4 Physics and Manual*, *JHEP* **05** (2006) 026, [[hep-ph/0603175](#)].
- [40] T. Sjöstrand, S. Ask, J. R. Christiansen, R. Corke, N. Desai, P. Ilten et al., *An introduction to PYTHIA 8.2*, *Comput. Phys. Commun.* **191** (2015) 159–177, [[1410.3012](#)].

- [41] DELPHES 3 collaboration, J. de Favereau, C. Delaere, P. Demin, A. Giammanco, V. Lemaitre, A. Mertens et al., *DELPHES 3, A modular framework for fast simulation of a generic collider experiment*, *JHEP* **02** (2014) 057, [[1307.6346](#)].
- [42] M. Bähr, S. Gieseke, M. A. Gigg, D. Grellscheid, K. Hamilton, O. Latunde-Dada et al., *Herwig++ physics and manual*, *The European Physical Journal C* **58** (nov, 2008) 639–707.
- [43] D. P. Kingma and J. Ba, *Adam: A method for stochastic optimization*, 2014. 10.48550/ARXIV.1412.6980.
- [44] S. H. Lim, K. A. Raman, M. R. Buckley and D. Shih, *GalaxyFlow: Upsampling Hydrodynamical Simulations for Realistic Gaia Mock Catalogs*, [2211.11765](#).
- [45] S. Diefenbacher, E. Eren, F. Gaede, G. Kasieczka, C. Krause, I. Shekhzadeh et al., *L2LFlows: Generating High-Fidelity 3D Calorimeter Images*, [2302.11594](#).

## A Robustness of classifier scores against network initialization

One might hope that the nature of the synthetic samples produced by each of the four generation methods considered in this report is robust in the sense that the samples produce the same results despite variances in initialization of neural networks. To this end, there are two sources of randomness: initialization of the binary classifier architectures, and initialization of the generator-level architectures, i.e. the networks used to generate the synthetic samples.

### A.1 Binary classifier initialization

We first gauge the robustness of scores against the binary classifier seed. In Fig. 15, we plot the classifier scores derived from samples of the *same* method, but for classifiers trained starting with a different random seed, on a dataset with  $n_{\text{sig}} = 1500$  injected signal events.

These plots illustrate that the scores output by the random classifier are relatively robust with respect to the training procedure when trained on a single sample generation method: the correlation between different random seeds is greater than 90% for all methods for  $n_{\text{sig}} = 1500$ . For smaller signal injections, the score correlation is reduced (especially below  $n_{\text{sig}} = 300$ ), but does remain positive.

In the main text of this paper, we elect to show results derived from scores averaged over 10 classifier trials (i.e. with different random seeds). This choice helps to stabilize the results, such that we are not comparing an uncharacteristically good classifier run for one method with an uncharacteristically bad classifier run for another.

### A.2 Generative network initialization

The question of score correlation across generative model seeds is an interesting one to ask particularly in the zero-signal injection case, when the binary classifiers can not leverage the characteristic observable distributions of signal events to make up for the differences between the different types of background that each synthetic sample generation method creates. In this case, it is useful to know if each sample generation method consistently produces similarly-functioning synthetic background samples. In Fig. 16, we plot the classifier scores derived from samples of the *same* method, but for generative models initialized starting with a different random seed, on a dataset with  $n_{\text{sig}} = 0$  injected signal events. Note that we average over 10 instantiations of binary classifiers initialized with different seeds.

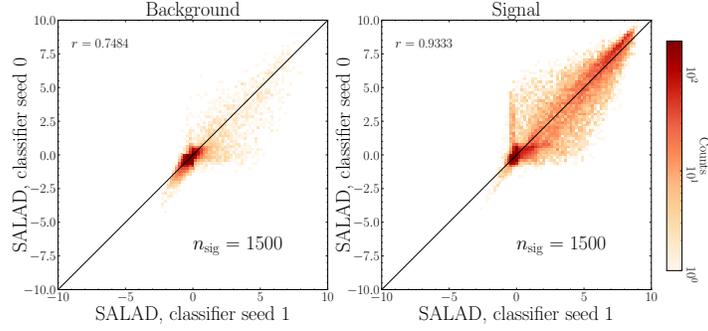
We see that for no injected signal, there is a good deal of correlation for the classifier scores of the true signal events between the scores derived from differently seeded generator models. However, the correlation for the scores of the true background events is only mild for CATHODE samples, and is negative for SALAD samples (note that the correlations jump to 74% and 72%, respectively, for the case where  $n_{\text{sig}} = 1500$ ). Therefore a maximally robust study should make use of averaging over generator initializations in addition to over binary classifier initializations, or perhaps by modifying the synthetic sample generator training procedures to reduce sensitivity to network instantiation.

### A.3 Using non-robustness to indicate a breakdown of CWoLa

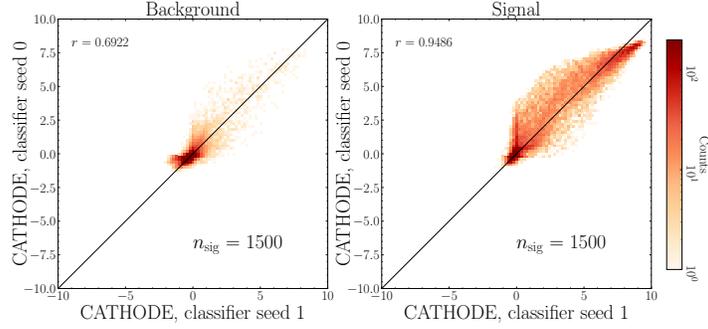
While the previous section provides evidence in favor of robustness of network scores, there is a limit to this correlation with respect to the anomaly detection procedure. When the signal injection is below a certain threshold (which is likely to be model-dependent), the CWoLa procedure will break down, and even an idealized anomaly detection classifier will fail to pick up on the signal event. This can be seen in two ways:

1. When the full sample combination study (corresponding to Sec. 4) is rerun, regenerating all of the synthetic samples for each method by retraining the generative architectures with a different initialization, the performance of the individual methods is highly variable below  $n_{\text{sig}} = 750$ , but is stable above that point. In particular, if we do not average of generator initializations (as shown in the rows of Fig. 17), then for one instantiation, the SALAD method appears to win out at  $n_{\text{sig}} = 500$ ; for another, the CATHODE method performs best at that signal injection; for a third, CATHODE and SALAD both exhibit low- $n_{\text{sig}}$  fluctuations.
2. From Fig. 12 (the fractional overlap between different sample generation methods of the 5th percentile of events classified as the most “signal-like”), the agreement of the most anomalous signal events is relatively stable with respect to  $n_{\text{sig}}$  above 750 events, but fluctuates highly below that value. This is again a sign that there is too little signal in the dataset for any of the individual learning methods to effectively recognize it.

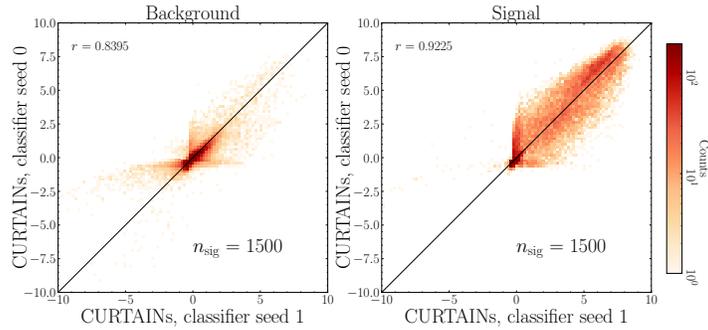
Given these findings, another benefit of combination emerges: the combined methods break down below  $n_{\text{sig}} = 750$ , in a sense flagging the low signal statistics and indicating a poor regime to use CWoLa procedure. In this situation, using an individual sample generation method might give an untrustworthy result.



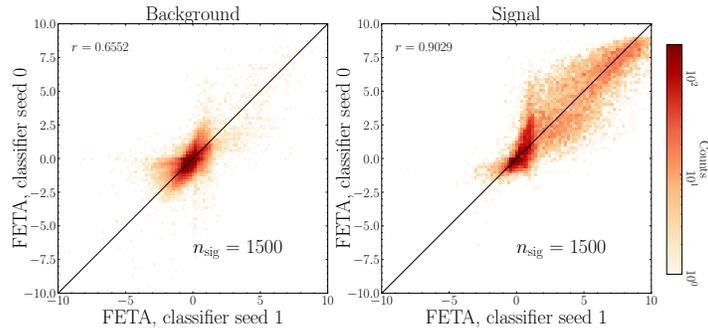
(a) SALAD against SALAD.



(b) CATHODE against CATHODE.

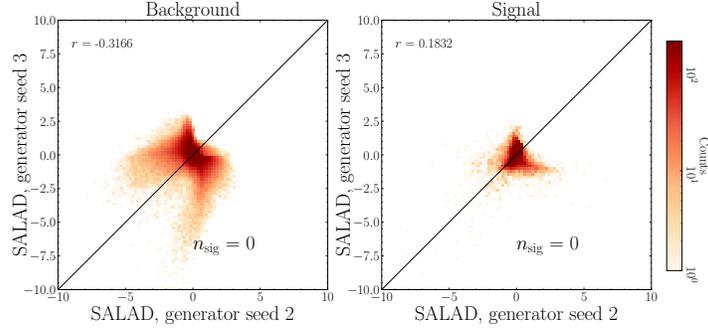


(c) CURTAINS against CURTAINS.

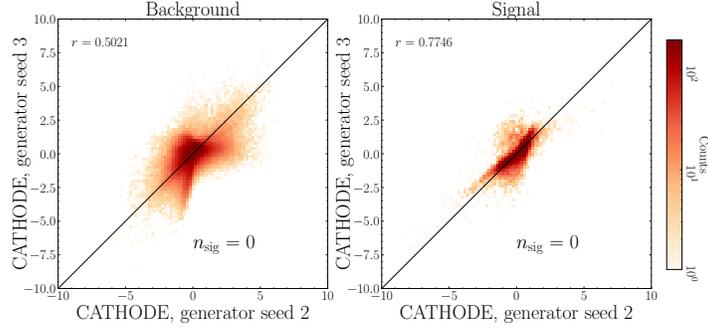


(d) FETA against FETA.

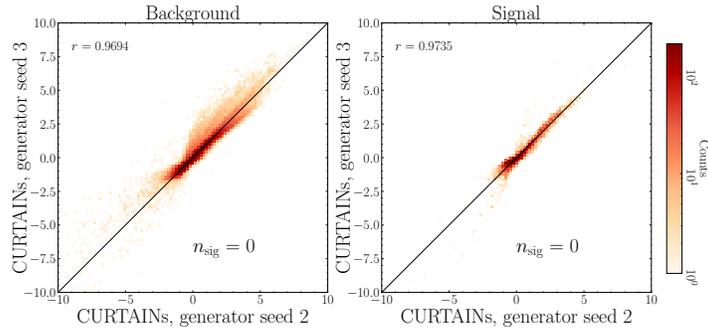
**Figure 15:** Classifier scores for a binary classifier trained to discriminate synthetic SM background from data with  $n_{\text{sig}} = 1500$ ). Each axis represents a different binary classifier random seed for the same sample generation method.



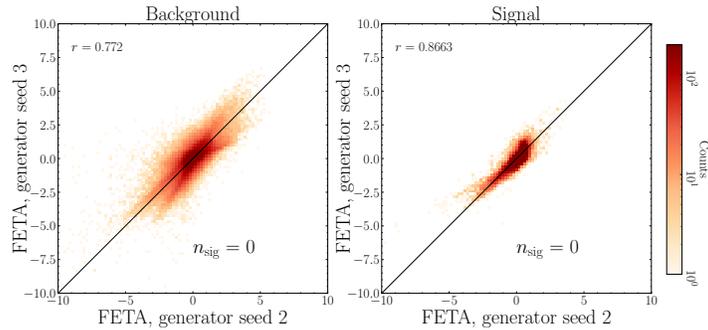
(a) SALAD against SALAD.



(b) CATHODE against CATHODE.

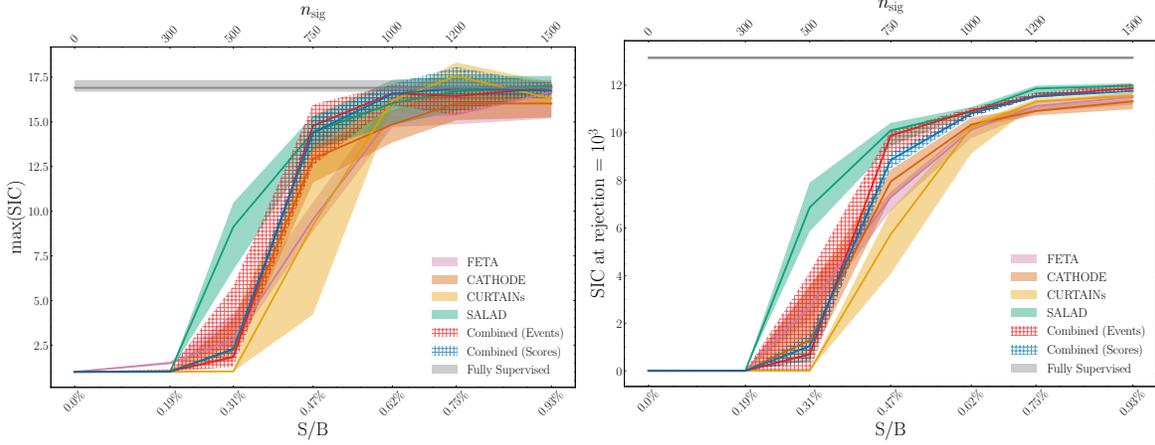


(c) CURTAINS against CURTAINS.

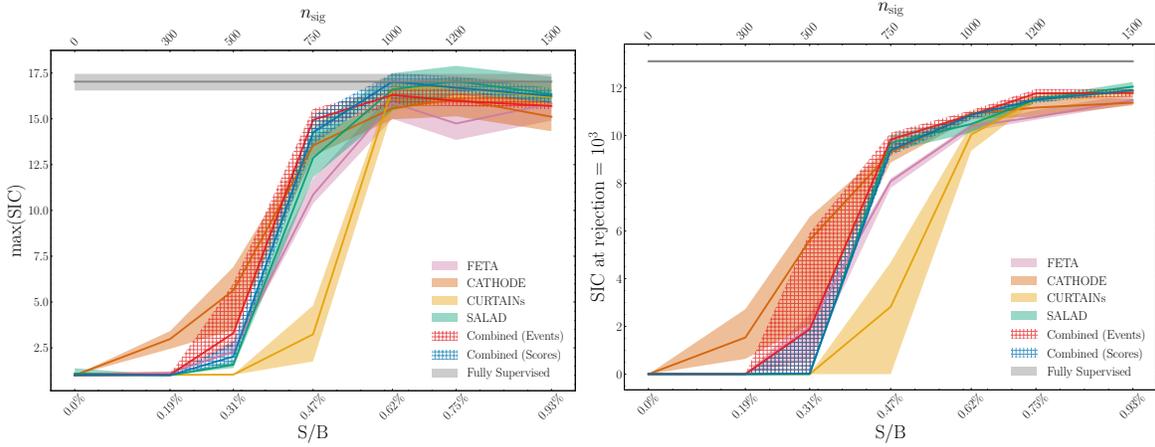


(d) FETA against FETA.

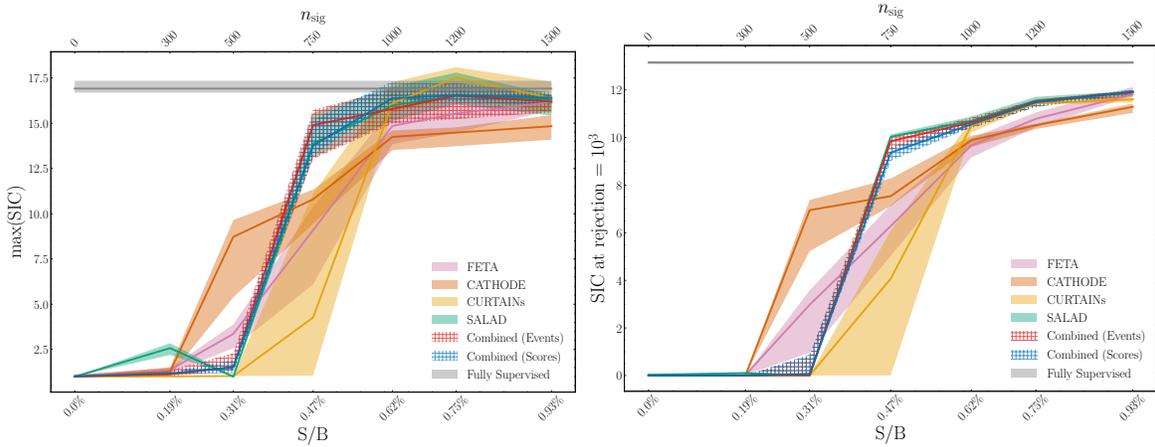
**Figure 16:** Classifier scores for a binary classifier trained to discriminate synthetic SM background from data with  $n_{\text{sig}} = 0$ ). Each axis represents a different generator architecture random seed for the given sample generation method. Scores are averaged over 10 different binary classifier instantiations.



(a) Maximum of the SIC, all generator seeds set to “2”. (b) SIC at a classifier rejection of  $10^3$ , all generator seeds set to “2”.



(c) Maximum of the SIC, all generator seeds set to “3”. (d) SIC at a classifier rejection of  $10^3$ , all generator seeds set to “3”.



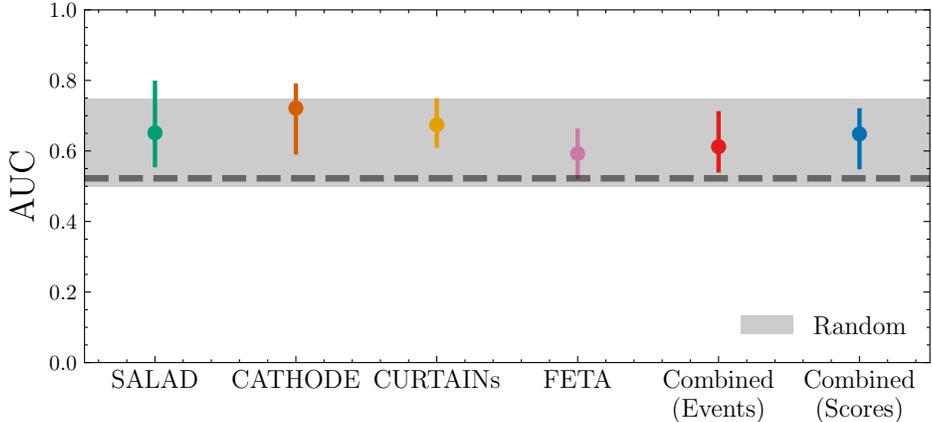
(e) Maximum of the SIC, all generator seeds set to “4”. (f) SIC at a classifier rejection of  $10^3$ , all generator seeds set to “4”.

**Figure 17:** Various metrics for a classifier trained to discriminate synthetic SM samples from data over a range of  $n_{\text{sig}}$  values. Equivalent to Fig. 13, but without averaging over generative network initializations.

## B Additional plots

In this section, we provide a small number of supplementary plots to complement the main text figures.

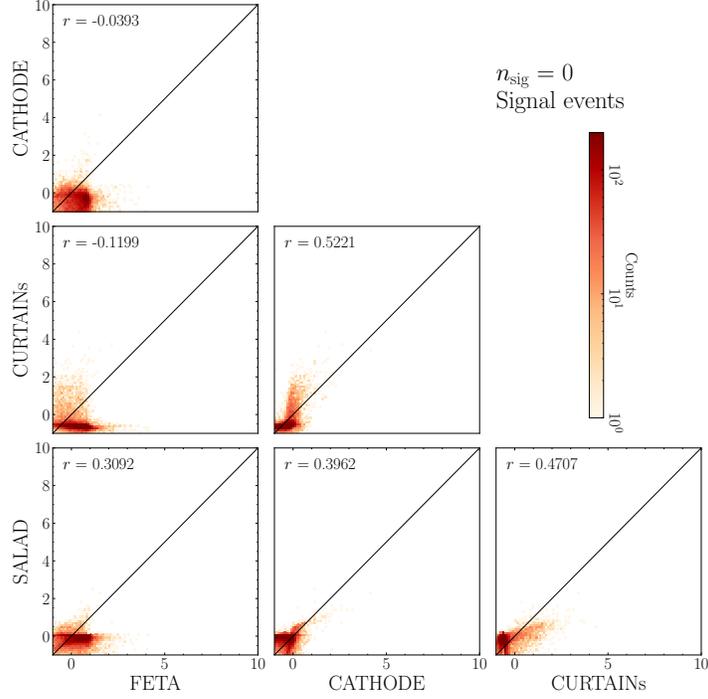
In Fig. 18, we plot a companion plot to Fig. 5, illustrating the spread of ROC AUCs for a binary classifier trained to discriminate sets of synthetic samples (or their combination) from data with  $n_{\text{sig}} = 0$ . Sample combination at both the event and the score level appears to more closely reproduce the spread coming from the random classifier, which indicates that the combination provides a more representative sample of SM background-like events.



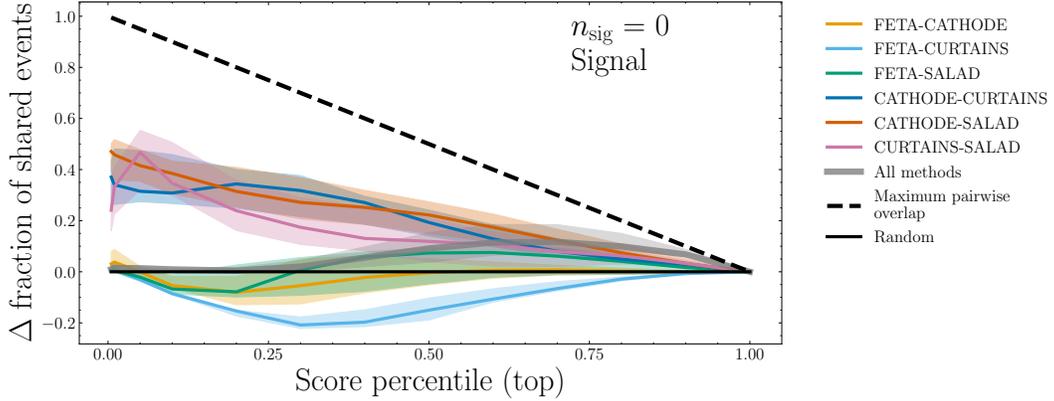
**Figure 18:** Receiver operating characteristic area-under-the-curves (AUC) for binary classifiers trained to discriminate each method’s synthetic samples from data with  $n_{\text{sig}} = 0$ . The table summarizes 100 classifier runs with different random seeds, with errorbars showing a 68-percentile spread.

In Fig. 19, we plot the standardized scores for true signal events, as calculated by a binary classifier trained to discriminate the synthetic samples from data with  $n_{\text{sig}} = 0$ . This is a companion plot to Fig. 6. In Fig. 20, we plot the fractional overlap, with respect to a random-choice baseline, of the  $p$ th percentile of true signal events classified as the most “signal-like” between different methods of synthetic SM samples. This is a companion plot to Fig. 7.

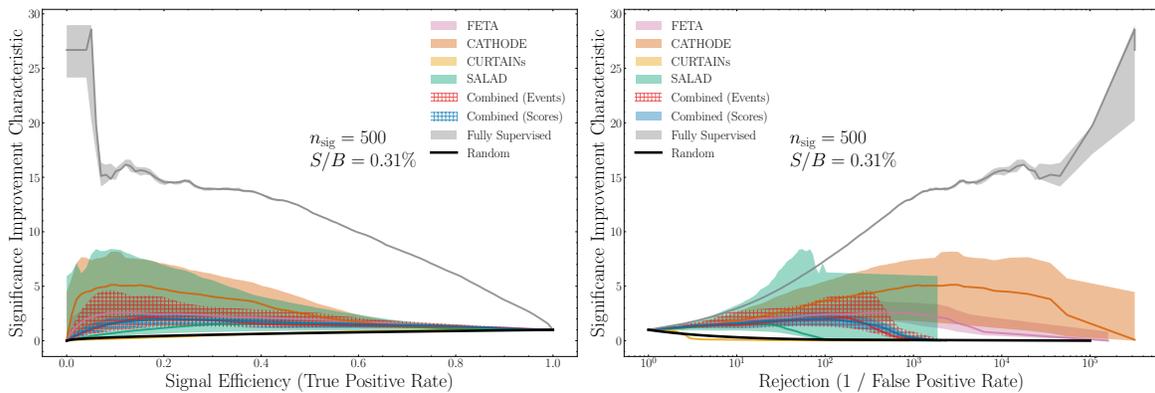
In Fig. 21, we plot the Significance Improvement Characteristic, as a function of the signal efficiency and the rejection, for an ensemble of classifiers trained to discriminate a combination of FETA, CATHODE, and CURTAINS synthetic SM samples from data with  $n_{\text{sig}} = 500$ . At this point, all methods fail to reliably pick up on this low signal injection. This is a companion plot to Fig. 14.



**Figure 19:** Scores for signal events for a binary classifier trained to discriminate synthetic SM background from data with  $n_{\text{sig}} = 0$ ). Each axis represents a different method of SM sample generation.  $r$  denotes the Pearson correlation coefficient. The scores are standardized so as to make it clear what events the binary classifier flags as the most anomalous. For each method, scores are averaged over 10 classifier runs. This is a companion plot to Fig. 6



**Figure 20:** Fractional overlap, with respect to a random-choice baseline, of the  $p$ th percentile of true signal events classified as the most “signal-like” between different methods of synthetic SM sample generation. Errorbands show a 68-percentile spread across the median and come from 100 repetitions of training 5-fold classifiers on the associated methods with different random seeds and ensembling scores over 10 repetitions. This is a companion plot to Fig. 7.



(a) Significance Improvement Characteristic plotted against the signal efficiency.

(b) Significance Improvement Characteristic plotted against classifier rejection.

**Figure 21:** Various classifier metrics for a classifier trained to discriminate a combination of FETA, CATHODE, and CURTAINS synthetic SM samples from data with  $n_{\text{sig}} = 750$ . Errorbands show a 68-percentile spread across the median and come from training a 5-fold classifier 100 times with different random seeds, over 3 independent generative model seeds, ensembling scores over 10 runs, and averaging classifier metrics over the ensembles. This is a companion plot to Fig. 14.