Resonant Anomaly Detection with Multiple Reference Datasets

Mayee F. Chen, a Benjamin Nachman, b; c and Frederic Salad

E-mail: mfchen@stanford.edu, bpnachman@lbl.gov, fredsala@cs.wisc.edu

Abstract: An important class of techniques for resonant anomaly detection in high energy physics builds models that can distinguish between reference and target datasets, where only the latter has appreciable signal. Such techniques, including Classication Without Labels (CWoLa) and Simulation Assisted Likelihood-free Anomaly Detection (SALAD) rely on a single reference dataset. They cannot take advantage of commonly-available multiple datasets and thus cannot fully exploit available information. In this work, we propose generalizations of CWoLa and SALAD for settings where multiple reference datasets are available, building on weak supervision techniques. We demonstrate improved performance in a number of settings with realistic and synthetic data. As an added benet, our generalizations enable us to provide nite-sample guarantees, improving on existing asymptotic analyses.

^aComputer Science Department, Stanford University, Stanford, CA 94305, USA

^bPhysics Division, Lawrence Berkeley National Laboratory, Berkeley, CA 94720, USA

^cBerkeley Institute for Data Science, University of California, Berkeley, CA 94720, USA

^dDepartment of Computer Sciences, University of Wisconsin, Madison, WI 53706, USA

Со	ontents					
1 Introduction						
2	Problem Setup	3				
3	Multi-CWoLa: Learning from Multiple Resonant Features 3.1 Multi-CWoLa Method 3.1.1 Model 3.1.2 Parameter Estimation 3.1.3 Inference and Training 3.2 Theoretical Results 3.3 Empirical Results	3 4 4 5 5 6 7				
4	Multi-SALAD: Learning from Multiple Simulations 4.1 Multi-SALAD Method 4.2 Theoretical Results 4.3 Empirical Results	8 9 10 11				
5	Conclusions and Outlook	13				
Α	Glossary	19				
В	Additional Algorithmic Details B.1 Multi-SALAD Algorithm	19 19				
С	Additional Theoretical Results C.1 The Need for 3 Resonant Features C.2 Rademacher Complexity Bounds C.3 Asymptotic behavior of SALAD's L'S(h; w)	20 20 21 21				
D	Proofs D.1 Proof of Theorem 1 D.2 Proof of Theorem 2	22 22 23				
E	Experiment Details E.1 Multi-CWoLa Experiments E.2 Multi-SALAD Experiments	25 25 25				

1 Introduction

Due to the vast parameter space of Standard Model extensions and to the lack of signicant evidence for new particles or forces of nature, a new model-agnostic search paradigm has emerged. Many of these anomaly detection (AD) strategies are enabled by machine learning (see e.g. Ref. [1{4}]) and the rst results with collision data are now available [5, 6]. One way to characterize AD methods is based on their physics assumption of the new phenomena [2]. Strategies that assume the new physics is \rare" [7] estimate (explicitly or implicitly) the data probability density and focus on events with low density. In contrast, techniques that assume the new physics will manifest as an overdensity in phase space use likelihood ratio methods to compare a reference dataset to a target dataset. The latter approach has been extensively studied in the context of resonant anomaly detection [8], where one resonant feature (usually a mass) is used to create a sideband region (reference dataset) nearly devoid of any anomalous events and a signal region (target dataset) that may contain anomalies. The reference dataset is used to estimate the presence of anomalies in the target dataset via interpolation.

Many existing approaches are dened using one reference dataset and one target dataset [9{18, 18{24}}]. However, in practice one can have access to or construct multiple references. First, there may exist multiple resonant features that can be used to construct sideband and signal regions. For instance, when a particle decays into two new particles, the decay products can be used to construct all three intermediate resonances, a setting present in the LHC Olympics Dataset [3]. Second, there may also exist multiple independent Standard Model simulators available for producing a dataset (e.g. Pythia [25], Herwig [26], or Sherpa [27]). Using multiple reference datasets may improve performance, but it is not clear how to incorporate all of their information when using existing methods designed for a single set.

We explore two generalizations of resonant AD to multiple reference datasets. First, we consider Classication Without Labels (CWoLa) [9, 10, 28], in which the reference is simply the sideband region | a form of weak supervision where the noisy label of \signal" is assigned to events in the signal region and the noisy label of 'background' to events in the sideband region. We propose a new method, Multi-CWoLa, that builds multiple reference datasets by constructing signal and sideband regions along dierent resonant features. We consider a point's membership in each feature's signal region as a noisy vote for anomaly, learn weights on each vote, and aggregate them to produce a higher-quality noisy label. We demonstrate Multi-CWoLa's performance on the LHC Olympics Dataset [3].

Next, we study Simulation Assisted Likelihood-free Anomaly Detection (SALAD) [14]. In this method, a reweighting function between a reference simulation dataset and a target dataset is learned in the sideband conditioned on the resonant feature. The simulated events in the signal region are reweighted by interpolating this function and then are used to distinguish anomalies in the target dataset. We extend this to the case of multiple simulated datasets, each of which may make dierent approximation choices and thus provide complementary accuracy when using SALAD. We introduce Multi-SALAD, which

combines the simulated datasets accordingly and then reweights, with the key nding that combining data helps when each simulator approximates dierent components of the background well. We demonstrate Multi-SALAD's performance on synthetic data.

Finally, we study the nite sample guarantees of our proposed methods. Many resonant AD methods have optimality guarantees in some asymptotic limit, but there is no rst-principles understanding of the methods' performance with nite samples. In particular, approaches like the ones described above that use classiers to distinguish a reference dataset from a target dataset approximate the data-to-background likelihood ratio. When the reference (physics) model is correct, this approach will converge to the optimal Neyman-Pearson likelihood ratio test in the limit of innite statistics, complex enough classier architecture, and exible enough training procedure [15, 29]. However, a nite sample understanding of these approaches is lacking. We draw on results from statistical theory to begin a formal study of resonant AD methods with limited data. Our results lay a foundation for future investigations into the nite sample properties of AD and related methods.

This paper is organized as follows. Section 2 briey set up the resonant AD setting and then Multi-CWoLa and Multi-SALAD are introduced in Secs. 3 and 4, respectively. The paper ends with conclusions and outlook in Sec. 5.

2 Problem Setup

We have an input space of discriminating features $x \in X$ and $x \in X$ and x

3 Multi-CWoLa: Learning from Multiple Resonant Features

We introduce Multi-CWoLa, an approach to anomaly detection that uses multiple reference datasets and is built using principles from the area of weak supervision [30, 31].

Standard CWoLa We have one unlabeled dataset $D = f(x_i; m_i)g_{i=1}^n$ with one resonant feature (k = 1) that we want to use to learn f. We use m to construct the signal and sideband regions, D_{SR} ; D_{SB} D where D_{SR} = D\SR and D_{SB} = D\SB, with distributions p_{SR} and p_{SB} respectively. With the intuition that there are more anomalies in the signal region, we express each distribution as a mixture of signal and background components with weight 0 $_{SR}$; $_{SB}$ 1:

$$p_{SR}(x) = {}_{SR}p(xjy = 1) + (1 {}_{SR})p(xjy = 0)$$
 (3.1)

$$p_{SB}(x) = {}_{SB}p(xjy = 1) + (1 {}_{SB})p(xjy = 0)$$
 (3.2)

Under this construction, the density ratio of the mixtures $\frac{p_{S\,R}\,(x)}{p_{S\,B}(x)}$ can be written in terms of the ratio of the signal and background components, $r(x) = \frac{p(xjy=1)}{p(xjy=0)}$, as $\frac{p_{S\,R}\,(x)}{p_{S\,B}\,(x)} = \frac{s_R\,r(x)+1}{s_B\,r(x)+1}\,s_B\,s_R$. Assuming $s_R > s_B$ (e.g. more signal in the signal region), the mixture ratio is monotonically increasing in r(x). Therefore, we train a classier f to learn $\frac{p_{S\,R}\,(x)}{p_{S\,B}\,(x)}$ by distinguishing between $D_{S\,R}$ and $D_{S\,B}$, and this f provides information about r(x) and can be used for anomaly detection.

3.1 Multi-CWoLa Method

Intuitively, CWoLa uses the resonant feature m as a noisy label that identies the signal versus sideband region and then trains a classier using these. This idea leads to a simple question | if more than one such feature is available (k > 1), how can the multiple noisy labels best be utilized? We tackle this question using principles from weak supervision [30{ 33}].

3.1.1 Model

In our approach, we split D along each resonant feature m^i to produce pairs of datasets D_{SB_i} and D_{SR_i} for each i 2 [k] based on membership in I_{m^i} . A straightforward way to use all datasets $(D_{SB_1}; D_{SR_1}); \ldots; (D_{SB_k}; D_{SR_k})$ is to apply standard CWola k times by training k classiers that we can then ensemble or average. Instead, in Multi-CWola, we construct a binary vector per x consisting of k noisy membership labels, $M(m) = fM_1(m); \ldots; M_k(m)g \ 2 \ f0; 1g^k$, where $M_i(m) = 1$ if $(x; m) \ 2 \ D_{SR_i}$ and $M_i(m) = 0$ if $(x; m) \ 2 \ D_{SB_i}$. We propose to directly aggregate these labels M(m) into an estimate of y, \$, and train a classier on the aggregated \$ along with the discriminative features x. Since each $M_i(m)$'s \vote" can have dierent correlation with the true y, we aim to combine the votes in a weighted fashion. We cannot directly measure each membership label's accuracy since the true y is unknown, so we draw on methods from weak supervision.

We model the distribution p(y; M(m)) as a probabilistic graphical model with the following parametrization:

$$p(y; M(m);) = \frac{1}{Z} exp_{y} ye + \sum_{i=1}^{K} M_{i}(m) ye;$$
(3.3)

where = f_y ; i 8i 2 [k]g are the canonical parameters of the distribution, Z is for normalization, and φ and $M_i(m)$ are y and $M_i(m)$ scaled from f0; 1g to f 1; 1g. Intuitively, i represents the (unobserved) strength of the correlation between $M_i(m)$ and y and thus captures a notion of M_i 's accuracy. This model also implies, for simplicity, that $M_i(m)$? $M_i(m)$ jy; that is, the resonant features are conditionally independent given y.¹

Our goal is to estimate the parameters of the graphical model and use them to perform inference, producing aggregated weak labels γ from the distribution p(y = 1jM(m);) given a vector of noisy labels M(m).

¹We can model some dependencies among resonant features if desired (see [31] for a method and see [34] for how to learn if resonant features are not conditionally independent). However, we need at least three conditionally independent subsets of resonant features in M(m) in order for the estimation method from [31] to recover the correct parameters.

3.1.2 Parameter Estimation

We rst learn the parameters of p(y; M(m);) as dened in (3.3). Of key interest is the accuracy parameter $_i = p(M_i(m) = 1jy = 1) = p(M_i(m) = 0jy = 0)$ of the ith resonant feature, which corresponds to the canonical parameter $_i$ (see [35] for more background on probabilistic graphical models). We estimate the accuracy parameters by adapting the triplet approach from [31]. First, we draw triplets of resonant features a;b;c 2 [k]. If the distribution on y;M(m) follows the graphical model in (3.3), it holds that $yM_a(m)$? $yM_b(m)$ if $M_a(m)$? $M_b(m)$ jy. Then, we have that $E[y_bM_a(m)]E[y_bM_b(m)] = E[M_a(m)M_b(m)]$ since $y^2 = 1$. Writing one such equation for each pair in the triplet (a;b;c), we have that

$$\begin{split} & E[\mathbf{p}\mathbf{\hat{M}}_{a}(m)]E[\mathbf{p}\mathbf{\hat{M}}_{b}(m)] = E[\mathbf{\hat{M}}_{a}(m)\mathbf{\hat{M}}_{b}(m)] \\ & E[\mathbf{p}\mathbf{\hat{M}}_{a}(m)]E[\mathbf{p}\mathbf{\hat{M}}_{c}(m)] = E[\mathbf{\hat{M}}_{a}(m)\mathbf{\hat{M}}_{c}(m)] \\ & E[\mathbf{p}\mathbf{\hat{M}}_{b}(m)]E[\mathbf{p}\mathbf{\hat{M}}_{c}(m)] = E[\mathbf{\hat{M}}_{b}(m)\mathbf{\hat{M}}_{c}(m)] : \end{split}$$

Solving this system, we obtain

$$jE[\gamma_{e}M_{a}(m)]j = t \frac{V_{u}^{v}E[M_{f}(m)M_{f}(m)]E[M_{fa}(m)M_{fc}(m)]}{E[M_{f}(m)M_{fc}(m)]}$$

and similarly for b and c. We assume that each signal region is positively correlated with the true signal, which allows for us to ignore the absolute value and uniquely recover $E[\gamma_e \hat{M}_a(m)]$. Next, we can use $E[\gamma_e \hat{M}_a(m)] = 2p(\gamma_e = {}^fM_a(m))$ 1 to obtain i using properties of the graphical model in (3.3). Note that in practice, all of these quantities are empirical estimates, with terms such as $E[\hat{M}_a(m)\hat{M}_b(m)] = \frac{1}{n} {}^n \hat{M}_a(m_i)\hat{M}_b(m_i)$.

3.1.3 Inference and Training

After we learn the accuracy parameters, we use them to estimate p(y = 1jM(m)) for a given M(m). We use Bayes' rule and the conditional independence among M(m) to write $p(yjM(m)) = \frac{\sum_{i=1}^{m} p(M_i(m)jy=1)p(y=1)}{p(M(m))}$. We assume that the class balance p(y = 1) is known; otherwise, it can be estimated via tensor decomposition [33]. $p(M_i(m)jy = 1)$ is either equal to p(m) = 1 or p(m) = 1 using the estimated accuracies and class balance.

Once p(y = 1jM(m)) is estimated for all M(m) 2 f0; $1g^k$, the aggregated weak label f is drawn from such distribution. With labels f for each f on the weakly labeled dataset $f(x; f)g_{i=1}^n$. This procedure is summarized in Algorithm 1.

 $^{^{2}}$ We assume that k 3. In Lemma 1, we discuss why having k = 1 or k = 2 resonant features does not recover a unique model.

Algorithm 1 Multi-CWoLa

- 1: Input: Dataset D = $f(x_i; m_i)g_{i=1}^n$; thresholds I_{m^i} that split D into signal and sideband regions, D_{SR_i} and D_{SB_i} respectively, for each m^i ; class balance probability of anomaly p(y = 1)
- 2: Construct noisy label $M_i(m) = \begin{pmatrix} 1 & (x; m) \ 2 & D_{SR_i} \\ 0 & (x; m) \ 2 & D_{SB_i} \end{pmatrix}$ for each resonant feature m^i .
- 3: for each triplet a; b; c 2 [k] do

4:

$$\frac{q}{a} := E[\underbrace{M_{a}(m)M_{b}(m)}_{b}[M_{a}(m)M_{c}(m)] = E[M_{b}(m)M_{c}(m)]}_{b} : \underbrace{Q}[M_{a}(m)M_{b}(m)] \underbrace{E[M_{a}(m)M_{c}(m)] = E[M_{b}(m)M_{c}(m)]}_{c} (3.4)$$

$$: \underbrace{Q}[M_{a}(m)M_{b}(m)] \underbrace{E[M_{b}(m)M_{c}(m)] = E[M_{a}(m)M_{c}(m)]}_{c} (3.5)_{c}$$

$$: \underbrace{Q}[M_{a}(m)M_{c}(m)] \underbrace{E[M_{b}(m)M_{c}(m)] = E[M_{a}(m)M_{b}(m)]}_{c};$$

$$(3.6)$$

where f is an empirical estimate of the expectation over D, and f (m) indicates M(m) scaled to f 1;1g.

- 5: end for
- 6: Set accuracy parameter $_{i} = p(M_{i}(m) = 1)y = 1) = p(M_{i}(m) = 0)y = 0) = p(M_{i}(m) = y) = \frac{_{i}+1}{_{2}}$
- 7: Compute estimate $p(y = 1jM(m)) / Q_{i=1}^{m} p(M_i(m)jy = 1)p(y = 1)$.
- 8: Construct $\oint p(y = 1jM(m))$ for each (x; m) 2 D.
- 9: Output: Classier f for anomaly detection trained on $f(x_i; m_i; y) g_{i=1}^n$.

3.2 Theoretical Results

Under (3.3), Multi-CWoLa oers nite-sample generalization guarantees. Suppose the downstream model f^trained on $\$ belongs to class F. Dene a loss function 'C: YY! R and let the expected loss of f be $L_C(f) := E['_C(f(x);y)]$ on true labels. Then, the optimal classier is $f^? = argmin_{f2F} L_C(f)$, which is achieved with unlimited labeled data. Let the empirical loss of f on $\$ be $L_C(f^?) := \frac{1}{n} \int_{i=1}^{n} f_C(f(x_i); \)$. Then, the $f^?$ we learn is constructed from $f^? = argmin_{f2F} \$ $f^?$ C(f), which is learned on nite and noisily labeled data. Note that this construction is dierent from the standard empirical risk minimization (ERM) loss on labeled data, and thus $L^n_C(f)$ does not asymptotically equal $L_C(f)$. We aim to minimize the generalization error $L_C(f^?)$.

We now present our result on an upper bound for $L_{\hat{F}}(f')$ $L_{C}(f^{?})$. Dene the Rademacher complexity of F as $R_{n}('F) = E \sup_{1 \le n} \sup_{1 \le n} \sum_{i=1}^{n} \prod_{j=1}^{n} \prod_{i=1}^{n} \prod_{j=1}^{n} \prod_{i=1}^{n} \prod_{j=1}^{n} \prod_{i=1}^{n} \prod_{j=1}^{n} \prod_{j=1}^{n} \prod_{i=1}^{n} \prod_{j=1}^{n} \prod_{j=1}^{n} \prod_{j=1}^{n} \prod_{i=1}^{n} \prod_{j=1}^{n} \prod_{j=1}^$

Theorem 1. Assume that p(y; M(m)) can be parametrized according to (3.3) and that 'is scaled to be bounded in [0;1]. Assume that the class balance p(y) is known (if not, there are ways to estimate it [33]), and that k 3. Then, with probability at least 1 , the

generalization error of Multi-CWoLa on D is at most

$$L_{C}(f)$$
 $L_{C}(f^{?})$ $4R_{n}('F) + \frac{r}{\log \frac{2}{2}} \frac{c_{1}}{2n} \frac{c_{1}}{e_{min}a_{min}^{5}} \frac{k}{n} + \frac{c_{2}k}{n}$

where c_1 ; c_2 are positive constants.

We observe that there are three quantities controlling the above bound:

- The Rademacher complexity of F: this term describes the model's expressivity. Smaller Rademacher complexity means that the model is easier to learn and that our f will be closer to the best model in F. This quantity can be readily computed for a variety of function classes F, such as decision trees, linear models, and two-layer feedforward networks, which makes our bound in Theorem 1 tractable. See Appendix C.2 for exact values.
- Using n nite samples: as the amount of data increases, the error decreases in $O(n^{-1-2})$.
- Using noisy labels \$\forall instead of y: for our weak supervision algorithm and graphical model, using \$\forall rather than y contributes an additional O(n 1=2) error. Asymptotically, our approach thus does no worse than training with labeled data.

By contrast, the standard CWoLa approach with k = 1 does not utilize any aggregation or weak supervision, which requires k = 3. For standard CWoLa, the second term in the generalization error is irreducible due to the fact that using any single resonant feature in place of y is biased. On the other hand, Multi-CWoLa corrects for some of this bias; the second term asymptotically approaches 0 with more data.

3.3 Empirical Results

In Figure 1, we compare Multi-CWoLa with standard CWoLa as well as two other baselines. We use simulation data from the LHC Olympics Dataset [3]; in particular from Pythia 8 [25], where the signal is boson decay and the background is generic 2! 2 parton scattering. This dataset contains 5 features; in the standard CWoLa setup, we use one thresholded resonant feature (k = 1) and use 4 discriminative features as x. For Multi-CWoLa, we have generated k = 3 mixtures by varying how the 3 resonant features (the jet masses in addition to the dijet mass) are thresholded and use 2 discriminative features as x. We have three other baselines that utilize 3 resonant features:

- CWoLa + intersect denes the signal region as the intersection of the resonant features' signal regions, e.g. $SR = SR_1 \setminus SR_2 \setminus SR_3$, but this can be overly conservative.
- CWoLa +x thresholding has one resonant feature as the noisy label ŷ = M₁(m), and includes the remaining thresholded features as discriminative fM₂(m); M₃(m); xg.
- CWoLa + average runs standard CWoLa three times, once per resonant feature and with the 2 discriminative features. The three model scores are averaged to produce the nal output.

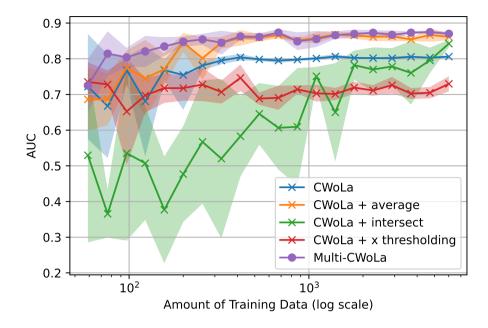


Figure 1. Comparison between CWoLa and Multi-CWoLa. Using multiple mixed samples helps performance across a range of dataset sizes. Access to multiple weak sources enables better AUC and lower variance compared to the single-feature version.

We vary the number of samples available on a logarithmic scale from n = 59 to 6003 and plot the AUC averaged over 5 runs per sample size in 1. We nd that Multi-CWoLa oers a higher AUC and lower variance, especially when there is limited data. We also plot the SI curves averaged over 5 runs for n = 59;530;6003 in 2.

4 Multi-SALAD: Learning from Multiple Simulations

We often have access to a(n approximate) simulation of the background process. We rst provide an overview of SALAD, which reweighs samples from the simulation to better assist with classication on the real dataset. Then, we present Multi-SALAD, a variant of SALAD that uses multiple simulations.

Standard SALAD We have a background simulation dataset $D^{sim} = f(x_i; m_i)g_{i=1}^{n_{sim}}$ with $y_i = 0$ for all i in addition to one true dataset $D = f(x_i; m_i)g_{1}^n D^{sim}$ is drawn from some distribution P_{sim} with density p_{sim} . While CWoLA learns the likelihood ratio between the signal and sideband regions of D alone, SALAD utilizes D^{sim} as well. Note that if p_{sim} is equal to p(jy = 0), we could directly train a model to distinguish between D and D^{sim} in the signal region to get a classier that could detect anomalies. However, since D^{sim} may not match the true background data, we instead rst need to learn a reweighting function that captures the dierences between D^{sim} and D's background data, and then we train a model to distinguish between D and the reweighted D^{sim} in the signal region. Formally, given xed SR and SB for both datasets, the method can be broken into two steps:

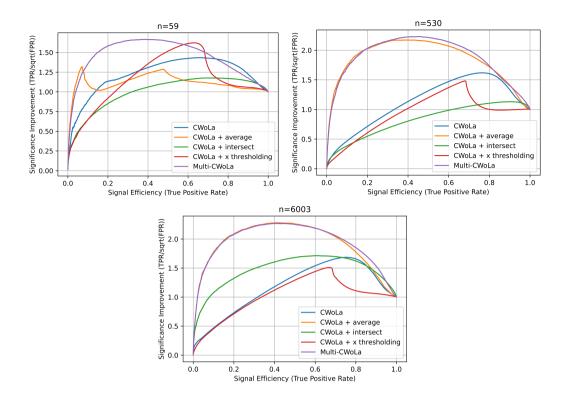


Figure 2. Signicance Improvement (SI) curve for Multi-CWoLa at sizes n = 59;530; and 6003.

- 1. Reweighting: a classier g is trained to distinguish between $D_{SB}^{sim} = D^{sim} \setminus SB$ and D_{SB} . Assuming that the sideband region has no anomalies, this g is able to produce an estimate of the weight ratio³ $w(x;m) = \frac{p}{p} \frac{p(x;m)y=0}{\sqrt{\frac{x}{x}mjy=0}} \frac{p(x;m)}{\sqrt{\frac{g(x;m)y}{x}mjy}}$ assuming that the datasets are the same size $(jD^{sim}j = \frac{p}{SB}jD_{SB}j)$.
- 2. Detection: Using a loss function L_S with estimated $\Psi(x;m)$ applied to $D_{SR}^{sim} = D^{sim} \setminus SR$, a classier h's trained to distinguish between D_{SR} and D_{SR}^{sim}

If the estimate $\psi(x;m)$ is exactly equal to $\psi(x;m)$ (e.g. g is Bayes-optimal), then the second step will be equivalent in expectation to learning the ratio $\frac{p(x)}{p(x)y=0)}$ (see Lemma 2 in Appendix C.3), from which one can detect anomalies.

4.1 Multi-SALAD Method

Now, we have multiple simulation datasets D_1^{sim} ;:::; D_k^{sim} . One approach would be to maintain distinctions among simulations by reweighing each pair to learn k weight functions $w_i(x;m)$, and then using one overall loss function that weights points from each $D_{SR;i}^{sim}$ with w_i . However, it has been shown that importance reweighting, despite working in expectation, can be highly unstable and result in poor performance of tasks on the target data D [40]. To understand why, Ref. [41] showed that the generalization error of an empirical loss function with importance weights w depends on the magnitude of w. Applied to our

³This is with the binary cross entropy loss function (also works for other functions [36]). This likelihoodratio trick is well-known (see e.g. Ref. [37, 38]), also in high-energy physics (see e.g. Ref. [39]).

setting, it suggests that the more inaccurate the simulation is, the less the reweighted loss recovers the true $\frac{p(x)}{p(x)y=0)}$, and the model may instead pick up on dierences between D_{SR} and the reweighted D_{SR} that are noise rather than the anomaly. As a result, aggregating individual SALAD outputs can be equivalent to ensembling many poor classiers.

Given these observations, Multi-SALAD uses multiple simulation datasets in a very simple yet theoretically principled way: control the magnitude of the overall w by combining all the D_i^{sim} to produce one large simulation dataset \mathcal{B}^{sim} whose distribution best approximates the true background p(xjy=0), and then use standard SALAD with \mathcal{B}^{sim} and D. Note that this approach both improves sample complexity and can \suppress" a simulation that on its own has high w, while the approach of learning k weight functions would not oer such improvements. In Algorithm 2 and Appendix B.1, we write this procedure out where we simply concatenate all D_i^{sim} together. However, with domain knowledge on the strengths and weaknesses of each simulation across features, one could produce D_i^{sim} by sampling accordingly from each. We leave this direction for future work.

4.2 Theoretical Results

We now present a nite sample generalization error bound on Multi-SALAD that also applies to SALAD. To measure the generalization error, recall $w(x;m) = \frac{p}{p} \frac{x;mjy=0}{(x;mjy=0)} dn$ let w be the classier g's estimate. We denote h as the reweighted classier. Let h? = $argmin_{h2H} L_S(h; w)$ and let h = $argmin_{h2H} L_S(h; w)$. We aim to bound $L_S(h, w)$ $L_S(h, w)$.

We rst set up some denitions. Dene n^{SR} as the number of points from D and D^{sim} belonging to the signal region, and n^{SB} as the number of points belonging to the sideband. Let n^{SR} be the number of points in D^{sim} belonging to the signal region. Let $g(x) = [g_{min}; g_{max}]$ and $g(x) = [g_{min}; g_{max}]$, where $g(x) = g_{max}$, where $g(x) = g_{max}$ is the optimal classier. Let $g(x) = g_{max}$ be the Rademacher complexity of the overall loss $g(x) = g_{max}$ be the Rademacher complexity of the overall loss $g(x) = g_{max}$ be the Rademacher complexity of the overall loss $g(x) = g_{max}$ be the Rademacher the simulation and true background. Let $g(x) = g_{max}$ log $g(x) = g_{max}$ log

Theorem 2. With probability at least 1 $\,$, there exists a constant c>0 such that the generalization error of Multi-SALAD on D^{sim} and D is at most

$$L_{S}(H; w) \quad L_{S}(h^{?}; w) \quad 2R_{S_{R}^{R}}('_{S} fH; Gg) + (1 + WB_{1}) \\ + \frac{n_{S_{Im}}^{SR}}{(1 - g_{max})(1 - g_{max}^{?})} n^{SR} \quad 4cR_{n}^{SB}(', G) + 2c \frac{r_{log}^{4=}}{2n^{SB}} + B_{2} \frac{s}{2n^{SR}} \\ \vdots$$

$$(4.1)$$

We make several observations about this bound:

- The bound scales in (n^{SB}) ¹⁼² and (n^{SR}_{sim}) ¹⁼², where the former comes from the initial reweighting step while the latter comes from the weighted classication step.
- The bound is also dependent on the Rademacher complexities of both classiers g and h used.
- The bound depends on the dierence between the simulation and data distributions through quantities W, B₁; B₂;; g_{max}; g_{max}. If the distributions have very dierent densities, these quantities will all be large, increasing the generalization error.

We comment how this bound is dierent when instantiated for SALAD versus Multi-SALAD. The following example shows how SALAD with one simulation can result in a large W (and other large constants), while Multi-SALAD with two simulations combined can reduce W in the bound.

Example 1. Let $P_{sim}(xjy = 0) = N(;^2),^2P_{sim}(xjy = 0) = N(;^2)$ be Gaussian distributions on x with ; 2 2 R, and let the true background distribution P(jy = 0) be a mixture of the Gaussians on x, $P(xjy = 0) = \frac{1}{2}P_{sim}^1 + \frac{1}{2}P_{sim}^2$. Let $P_{sim}^1; P_{sim}^2;$ and P have the same marginal distribution over m with x ? mjy. Then, if we only use one simulation P_{sim}^1 ,

$$w(x; m) = \frac{p(x; mjy = 0)}{p_{sim}^{1}(x; mjy = 0)} = \frac{p(xjy = 0)}{p_{sim}^{1}(xjy = 0)}$$

$$= \frac{\frac{\frac{1}{2} - \exp(\frac{(x -)^{2}}{2^{2}} + \frac{1}{2} - \exp(\frac{(x +)^{2}}{2^{2}})}{-\frac{1}{2} - \exp(\frac{(x -)^{2}}{2^{2}})}$$

$$= \frac{1}{2} + \frac{1}{2} \exp(\frac{(x -)^{2}}{2^{2}} - \frac{(x -)^{2}}{2^{2}} - \frac{1}{2} + \frac{1}{2} \exp(\frac{x}{2} - \frac{x}{2})$$

Therefore, as $x \mid 1$, $W \mid 1$. However, if we dene P_{sim} as the distribution of the two simulation datasets concatenated, we have that $p_{sim}(xjy=0)=p(xjy=0)$, and as a result, $W \mid 1$, making the generalization error bound smaller.

From this example, we can see that signicantly diering simulation and data distributions can result in large, unbounded weight ratios, which are correlated with poor performance.⁴ This concretely motivates our algorithmic objective to combine multiple simulation datasets as to closely approximate the true data.

4.3 Empirical Results

To demonstrate how Multi-SALAD can improve over using only one simulation and over using simulations separately, we consider a synthetic experiment with two simulation

⁴The bound in Theorem 2 is meant to provide a general understanding of SALAD's performance. It can be made tighter by replacing terms that are maxima like M and B₂ with terms that are based on the overall data distributions (e.g. variance, as in Ref. [41]). Variance-based bounds are less likely to be vacuous, but will still demonstrate how performance is dependent on the intrinsic dierences between the two distributions.

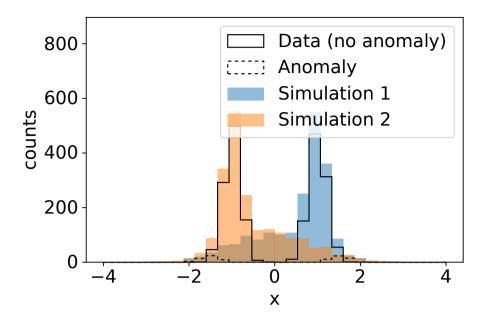


Figure 3. Synthetic data for evaluating Multi-SALAD.

datasets⁵. The true background is $P(jy = 0) = {}_{2}N(1;0:2) + {}_{2}N(1;0:2)$, and the anomaly is $P(jy = 1) = {}_{2}N\{2;0:2) + {}_{2}N(2;0:2)$. Simulation 1 is $P_{sim} = {}_{2}N(1;0:2) + {}_{2}N(1;0:2) + {}_{2}N(1;0:2)$, and simulation 2 is $P_{sim} \not= {}_{2}N(1;0:2) + {$

Intuitively, the anomaly is only slightly dierent from the background data, which makes it important to learn a good reweighting function from the simulations. Because each simulation alone diverges greatly from the data for one mode, each individual reweighting may not approximate the true P(jy=1) well. On the other hand, if we combine both simulation datasets together, the aggregate distribution has smaller weights with lower variance, which can allow for more accurate reweighting. This is demonstrated in Figure 4, which depicts the reweighting in the sideband region. Figure 5 depicts the reweighting's interpolation into the signal region, where we introduce an additional baseline SALADSwitch, which uses k separate weight functions $w_i(x;m)$ and switches among them in the reweighted loss function L_S . In all but the bottom right subgure in both gures, the reweighted simulation data poorly approximates the true background data. As a result, a classier trained to distinguish between the high-variance reweighted simulation and the true background data plus some small anomaly will more likely learn the distinctions coming from poor approximation, rather than anomaly. In particular, note that SALAD-

⁵We nd that the dierences between the simulations in the LHC Olympics are not enough to see a noticeable gain from Multi-SALAD over SALAD.

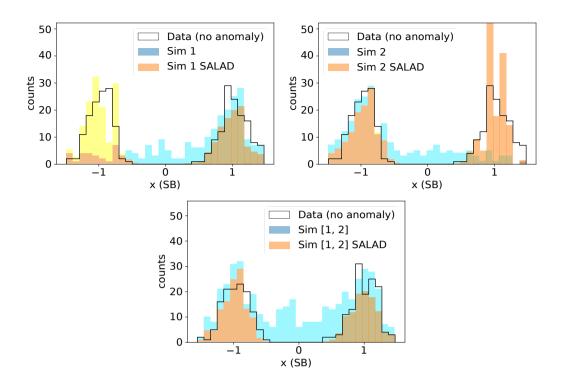


Figure 4. Top left: SALAD reweighting using simulation 1 on sideband region. Top right: reweighting using simulation 2. Bottom: reweighting using simulation 1 and 2 combined.

Switch results in signicant overweighting in Figure 5.

With these observations, we present the signal eciency to rejection rate of each method in Figure 6, where we compare Multi-SALAD against SALAD using simulation 1 only, SALAD using simulation 2 only, and SALAD-Switch. Table 1 contains the accuracy and AUC scores for each method. Averaged over 10 random seeds, Multi-SALAD outperforms other methods. The signal eciency to rejection rate for each of the 10 runs is available in Appendix E.

	Simulation 1		Simulation 2		Simulation 1 and 2		
Method	None	SALAD	None	SALAD	None	SALAD-Switch	Multi-SALAD
Accuracy	43:82:2	62:5 _{8:8}	42:73:6	64:3 _{12:3}	50:0 _{0:0}	54:3 _{6:2} 74:7 _{17:0}	64:89:3 90:810:2
AUC	28:54.2	80:714.5	27:44.5	78:718:2	15:45:3		

Table 1. Accuracy and AUC scores (%) for Multi-SALAD on two simulation datasets. We compare to SALAD-Switch (dierent reweighting), as well as standard SALAD on individual simulations and no reweighting. Performance is averaged over 10 random runs with one standard deviation reported.

5 Conclusions and Outlook

We extend two resonant AD approaches to incorporate multiple reference datasets. For Multi-CWoLa, we draw from weak supervision models to handle multiple resonant features. For Multi-SALAD, we combine multiple simulation datasets to best approximate

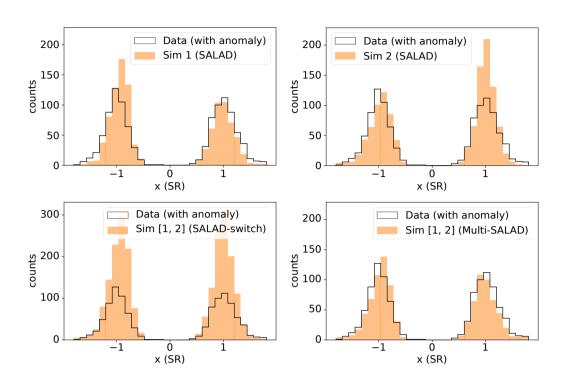


Figure 5. Top left: SALAD reweighting using simulation 1 on signal region. Top right: reweighting using simulation 2. Bottom left: using both simulation 1 and 2 weights separately. Bottom right: reweighting using simulation 1 and 2 combined.

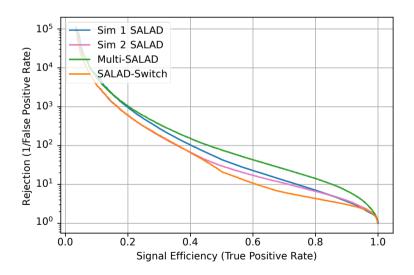


Figure 6. Signal eciency to rejection of Multi-SALAD versus other baselines (weighted and unweighted).

the background process. Future work includes 1) exploring Multi-SALAD's applicability on real data and algorithms for sampling from simulation datasets 2) extending Multi-

CWoLa to model more complex relationships among resonant features and 3) using such approaches together over multiple simulations and resonant features, eectively utilizing as much information as possible.

Acknowledgments

We thank David Shih and Jesse Thaler for useful discussions and comments about the manuscript. BN was supported by the Department of Energy, Oce of Science under contract number DE-AC02-05CH11231. FS is grateful for the support of the NSF under CCF2106707 and the Wisconsin Alumni Research Foundation (WARF). We gratefully acknowledge the support of NIH under No. U54EB020405 (Mobilize), NSF under Nos. CCF1763315 (Beyond Sparsity), CCF1563078 (Volume to Velocity), and 1937301 (RTML); ARL under No. W911NF-21-2-0251 (Interactive Human-Al Teaming); ONR under No. N000141712266 (Unifying Weak Supervision); ONR N00014-20-1-2480: Understanding and Applying Non-Euclidean Geometry in Machine Learning; N000142012275 (NEPTUNE); NXP, Xilinx, LETI-CEA, Intel, IBM, Microsoft, NEC, Toshiba, TSMC, ARM, Hitachi, BASF, Accenture, Ericsson, Qualcomm, Analog Devices, Google Cloud, Salesforce, Total, the HAI-GCP Cloud Credits for Research program, the Stanford Data Science Initiative (SDSI), and members of the Stanford DAWN project: Facebook, Google, and VMWare. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright notation thereon. Any opinions, ndings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reect the views, policies, or endorsements, either expressed or implied, of NIH, ONR, or the U.S. Government.

References

- [1] HEP ML Community, \A Living Review of Machine Learning for Particle Physics."
- [2] G. Karagiorgi, G. Kasieczka, S. Kravitz, B. Nachman and D. Shih, Machine Learning in the Search for New Fundamental Physics, 2112.03769.
- [3] G. Kasieczka et al., The LHC Olympics 2020: A Community Challenge for Anomaly Detection in High Energy Physics, 2101.08320.
- [4] T. Aarrestad et al., The Dark Machines Anomaly Score Challenge: Benchmark Data and Model Independent Event Classication for the Large Hadron Collider, 2105.14027.
- [5] ATLAS Collaboration, Dijet resonance search with weak supervision using 13 TeV pp collisions in the ATLAS detector, 2005.02983.
- [6] ATLAS Collaboration collaboration, Anomaly detection search for new resonances decaying into a Higgs boson and a generic new particle X in hadronic nal states using 13 TeV pp collisions with the ATLAS detector, tech. rep., CERN, Geneva, 2022.
- [7] G. Kasieczka, R. Mastandrea, V. Mikuni, B. Nachman, M. Pettee and D. Shih, Anomaly Detection under Coordinate Transformations, 2209.06225.
- [8] G. Kasieczka, B. Nachman and D. Shih, New Methods and Datasets for Group Anomaly Detection From Fundamental Physics, ANDEA (2021), [2107.02821].

- [9] 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].
- [10] J. H. Collins, K. Howe and B. Nachman, Extending the search for new resonances with machine learning, Phys. Rev. D99 (2019) 014038, [1902.02634].
- [11] R. T. D'Agnolo and A. Wulzer, Learning New Physics from a Machine, Phys. Rev. D99 (2019) 015014, [1806.02350].
- [12] R. T. D'Agnolo, G. Grosso, M. Pierini, A. Wulzer and M. Zanetti, Learning Multivariate New Physics, 1912.12155.
- [13] K. Benkendorfer, L. L. Pottier and B. Nachman, Simulation-Assisted Decorrelation for Resonant Anomaly Detection, 2009.02205.
- [14] A. Andreassen, B. Nachman and D. Shih, Simulation Assisted Likelihood-free Anomaly Detection, Phys. Rev. D 101 (2020) 095004, [2001.05001].
- [15] B. Nachman and D. Shih, Anomaly Detection with Density Estimation, Phys. Rev. D 101 (2020) 075042, [2001.04990].
- [16] O. Amram and C. M. Suarez, Tag N' Train: A Technique to Train Improved Classiers on Unlabeled Data, 2002.12376.
- [17] A. Hallin, J. Isaacson, G. Kasieczka, C. Krause, B. Nachman, T. Quadfasel et al., Classifying Anomalies Through Outer Density Estimation (CATHODE), 2109.00546.
- [18] J. A. Raine, S. Klein, D. Sengupta and T. Golling, CURTAINs for your Sliding Window: Constructing Unobserved Regions by Transforming Adjacent Intervals, 2203.09470.
- [19] R. T. d'Agnolo, G. Grosso, M. Pierini, A. Wulzer and M. Zanetti, Learning New Physics from an Imperfect Machine, 2111.13633.
- [20] P. Chakravarti, M. Kuusela, J. Lei and L. Wasserman, Model-Independent Detection of New Physics Signals Using Interpretable Semi-Supervised Classier Tests, 2102.07679.
- [21] B. M. Dillon, R. Mastandrea and B. Nachman, Self-supervised Anomaly Detection for New Physics, 2205.10380.
- [22] M. Letizia, G. Losapio, M. Rando, G. Grosso, A. Wulzer, M. Pierini et al., Learning new physics eciently with nonparametric methods, 2204.02317.
- [23] K. Krzyzanska and B. Nachman, Simulation-based Anomaly Detection for Multileptons at the LHC, 2203.09601.
- [24] S. Alvi, C. Bauer and B. Nachman, Quantum Anomaly Detection for Collider Physics, 2206.08391.
- [25] T. Sjostrand, S. Mrenna and P. Z. Skands, A Brief Introduction to PYTHIA 8.1, Comput. Phys. Commun. 178 (2008) 852{867, [0710.3820].
- [26] J. Bellm et al., Herwig 7.0/Herwig++ 3.0 release note, Eur. Phys. J. C 76 (2016) 196, [1512.01178].
- [27] Sherpa collaboration, E. Bothmann et al., Event Generation with Sherpa 2.2, SciPost Phys. 7 (2019) 034, [1905.09127].
- [28] E. M. Metodiev, B. Nachman and J. Thaler, Classication without labels: Learning from mixed samples in high energy physics, JHEP 10 (2017) 174, [1708.02949].

- [29] J. Neyman and E. S. Pearson, On the problem of the most ecient tests of statistical hypotheses, Phil. Trans. R. Soc. Lond. A 231 (1933) 289.
- [30] A. Ratner, S. H. Bach, H. Ehrenberg, J. Fries, S. Wu and C. Re, Snorkel: Rapid training data creation with weak supervision, in Proceedings of the 44th International Conference on Very Large Data Bases (VLDB), (Rio de Janeiro, Brazil), 2018.
- [31] D. Y. Fu, M. F. Chen, F. Sala, S. M. Hooper, K. Fatahalian and C. R/'e, Fast and three-rious: Speeding up weak supervision with triplet methods, in International Conference on Machine Learning, 2020, https://arxiv.org/pdf/2002.11955.pdf.
- [32] A. Ratner, C. D. Sa, S. Wu, D. Selsam and C. Re, Data programming: Creating large training sets, quickly, in Proceedings of the 30th International Conference on Neural Information Processing Systems, NIPS'16, (Red Hook, NY, USA), p. 3574{3582, Curran Associates Inc., 2016.
- [33] A. Ratner, B. Hancock, J. Dunnmon, F. Sala, S. Pandey and C. Re, Training complex models with multi-task weak supervision, in Proceedings of the AAAI Conference on Articial Intelligence, Jul, 2019.
- [34] P. Varma, F. Sala, A. He, A. Ratner and C. Re, Learning dependency structures for weak supervision models, in Proceedings of the 36th International Conference on Machine Learning, 2019.
- [35] M. J. Wainwright, M. I. Jordan et al., Graphical models, exponential families, and variational inference, Foundations and Trends® in Machine Learning 1 (2008) 1{305.
- [36] B. Nachman and J. Thaler, E Pluribus Unum Ex Machina: Learning from Many Collider Events at Once, 2101.07263.
- [37] T. Hastie, R. Tibshirani and J. Friedman, The Elements of Statistical Learning. Springer Series in Statistics. Springer New York Inc., New York, NY, USA, 2001.
- [38] M. Sugiyama, T. Suzuki and T. Kanamori, Density Ratio Estimation in Machine Learning. Cambridge University Press, 2012, 10.1017/CBO9781139035613.
- [39] K. Cranmer, J. Pavez and G. Louppe, Approximating Likelihood Ratios with Calibrated Discriminative Classiers, 1506.02169.
- [40] S. Dasgupta and P. M. Long, Boosting with diverse base classiers, in Learning Theory and Kernel Machines (B. Schelkopf and M. K. Warmuth, eds.), (Berlin, Heidelberg), pp. 273{287, Springer Berlin Heidelberg, 2003.
- [41] C. Cortes, Y. Mansour and M. Mohri, Learning bounds for importance weighting, in Advances in Neural Information Processing Systems (J. Laerty, C. Williams, J. Shawe-Taylor, R. Zemel and A. Culotta, eds.), vol. 23, Curran Associates, Inc., 2010, https://proceedings.neurips.cc/paper/2010/le/59c33016884a62116be975a9bb8257e3-Paper.pdf.
- [42] T. Ma, Lecture notes for machine learning theory (cs229m/stats214), June, 2022.
- [43] P. L. Bartlett and S. Mendelson, Rademacher and gaussian complexities: Risk bounds and structural results, Journal of Machine Learning Research 3 (2002) 463{482.
- [44] M. Mohri, A. Rostamizadeh and A. Talwalkar, Foundations of machine learning. MIT press, 2018.

- [45] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan et al., Pytorch: An imperative style, high-performance deep learning library, in Advances in Neural Information Processing Systems 32, pp. 8024{8035. Curran Associates, Inc., 2019.
- [46] F. Chollet et al., Keras, 2015.

Appendix

We provide a glossary of notation in A. We provide algorithmic details for Multi-SALAD in Section B. We present additional theoretical results on Rademacher complexities and the asymptotic behavior of SALAD in Section C. In section D, we provide proofs for our theoretical results. In section E, we provide additional experimental details.

A Glossary

The glossary is given in Table 2.

B Additional Algorithmic Details

B.1 Multi-SALAD Algorithm

Multi-SALAD is described in Algorithm 2. We have simulation datasets D_1^{sim} ;::: D_k^{sim} , where $D_i^{sim} = f(x_j; m_j)g_{j=1}^{n_{sim}}$ and all points belong to the background (y = 0). As discussed in Section 4, we propose using these simulation datasets by aggregating them into a single simulation dataset D^{sim} (whether it be with uniform or stratied sampling, etc.) Then the rest of this section proceeds as follows and is a review of the standard SALAD method.

Reweighting First, we learn weights to correct for the bias of the simulated background data. We split the both simulation and true data along m to produce sets $D_{S\,R}^{sim}$; $D_{S\,B}^{sim}$ and $D_{S\,R}$ and $D_{S\,B}$. We train a classier over $D_{S\,B}^{sim}$ to distinguish between simulation and real data in the sideband region. That is, we train a binary classier g over points g over points g in the sideband where g is either from g over g or g over points g over points g in the sideband where g is either from g over g or g or g or g or g over points g over points g over points g over points g in the sideband where g is either from g over g or g or g or g over points g o

$$\begin{split} \hat{w}(x;m) &= \frac{g(x;m)}{1} \quad \frac{q(z=1jx;m)}{q(z=0jx;m)} = \frac{q(x;mjz=1)}{q(x;mjz=0)} \quad \frac{q(z=1)}{q(z=0)} \\ &= \frac{q(x;mjz=1)}{q(x;mjz=0)} = \frac{p(x;mjy=0)}{p_{sim}(x;mjy=0)} \\ \vdots \end{split}$$

Here, we assume that q(z = 1) = q(z = 0) (i.e. balanced simulation and real dataset, which we can always ensure by generating more or less simulation data). Equality is obtained in the expression above when g is Bayes-optimal.

Training The above $\psi(x;m)$ is dened on the sideband region. Next, we interpolate and correct the bias of the simulation in the signal region. Let $D_{S\,R}^{sim}$ be the set of simulation data in the signal region of size $n_{sim}^{s\,R}$, and let $D_{S\,R}$ be the set of true data in the signal region of size $n_{data}^{s\,R}$, for a total of $n^{s\,R}$ points. We train a classier h to distinguish between the reweighted simulated data, which approximates true background data, and the true data. In particular, the loss function used is

$$L_{XS}(h; w) = \underbrace{\frac{1}{R}}_{n} \times 2D_{SR} \times 2D_{SR} + X \quad w(x; m) \log(1 \quad h(x; m)) : \qquad (B.1)$$

Algorithm 2 Multi-SALAD

- 1: Input: Simulation datasets D_1^{sim} ;:::; D_k^{sim} and real dataset D.
- 2: Construct overall simulation dataset $D^{sim} = S_{i=1}^{k} D_{i}^{sim}$.
- 3: Split each dataset into signal region and sideband region using resonant feature m to get fD_{SR}^{sim} ; D_{SB}^{sim} g and fD_{SR} ; D_{SB} g.
- 4: Learn weight $\psi(x; m) = \frac{g(x; m)}{1 g(x; m)}$, where g is a classier that distinguishes data D_{SB} from simulation D_{SB}^{sim} in the sideband region.
- 5: Train a new classier h'on the signal region to distinguish between points in D_{SR}^{sim} reweighted by ψ , using the following loss:

$$L_{S}(h; w) = \int_{\Lambda}^{R} X \log h(x; m) + X w(x; m) \log(1 h(x; m)):$$

$$\frac{1}{n} \times 2D_{SR} \times 2D_{Sim}$$
(B.2)

6: Output: Classier output h(x; m), which yields a score that is thresholded for anomaly detection.

In expectation with an optimal w, we can see that minimizing this loss is equivalent to minimizing the cross-entropy loss on a task that distinguishes between points drawn from p and points drawn from p(jy = 0) in the signal region. Therefore, h can be used for anomaly detection. The procedure is summarized in Algorithm 2.

C Additional Theoretical Results

C.1 The Need for 3 Resonant Features

We show that to identify the model (3.3), we need at least k = 3 resonant features.

Lemma 1. If k = 1 or k = 2 in model (3.3), the parameters $_1$ and $_2$ cannot be recovered from the observable quantities.

Proof. The strategy we use to show that the model cannot be identied for k=1 or k=2 is to prove that the observable distributions $P(M_k^f(m); :::; M_k^f(m))$ are consistent with multiple values of . We do so by direct calculation.

First, consider the case of k = 1. Set y = 0 for simplicity. Then, the model is $z = \exp(M_1 f_m) \varphi$. Then $Z = 2 \exp(1 + 2 \exp(1))$, and

$$P(fM_1(m) = 1) = \frac{exp() + exp()}{2 exp() + 2 exp()} = \frac{1}{2}$$
:

Thus, any value produces the same observable distribution, so that we cannot identify .

Next, we consider k = 2. Again, set y = 0. The model is now $z = \exp({}_{1}M_{1}f(m)y + 2M_{2}(m)y = 0$. We similarly compute

$$Z = 2(exp(_1 + _2) + exp(_1 + _2) + exp(_1 _2) + exp(_1 _2))$$
:

The observable distribution is now $P(M_1(m); M_2(m))$. We have that

$$P(fM_1(m) = 1; fM_2(m) = 1) = \frac{1}{Z}(exp(_1 + _2) + exp(_1$$
 _2));

and

$$P(fM_1(m) = 1; fM_2(m) = 1) = \frac{1}{7}(exp(_1 _2) + exp(_1 + _2)):$$

Note that we have $P(\mathring{M}_1(m) = 1; \mathring{M}_2(m) = 1) = P(\mathring{M}_1(m) = 1; \mathring{M}_2(m) = 1)$ and $P(\mathring{M}_1(m) = 1; \mathring{M}_2(m) = 1) = P(\mathring{M}_1(m) = 1; \mathring{M}_2(m) = 1)$.

As a result, we have the same distribution $P(M_1^f(m); M_2^f(m))$ for the parameters $_1;_2 = a;_b$ and for $_1;_2 = b;_a$, where $a;_b$ are some non-negative values. If a = b, we end up with at least two solutions that cannot be distinguished, completing the proof.

C.2 Rademacher Complexity Bounds

We present bounds on the Rademacher complexity $R_n(F)$ of various models F. For all of the F below, we obtain $R_n('F)$ by computing $R_n(F)$. These two Rademacher complex-ities are equal when we assume that ' is 1-Lipschitz and apply Talagrand's lemma.

- Linear models: We dene f(x) = x with kk_2 B and $E[kxk^2]$ C^2 , $R_n(F)$ $\frac{BC}{n}$ [42, Theorem 5.5].
- Two-layer feed-forward neural networks (MLPs): We dene f(x) where = (U; w) are the parameters for the weights for the two layers of an MLP. Here U 2 R^{md} and w 2 R^{m} . Suppose ReLU is the activation function, kwk₂ B_{w} , ku_ik₂ B_{u} for all 1 i m, and that $E[kxk_{2} \ C^{2}$. Then, $R_{n}(F) \ 2B_{w}B_{u}C$ $\frac{p}{n}$ [42, Theorem 5.9].
- Kernels: Let $k: X \times Y = R$ be a continuous symmetric function so that for $x_1; \ldots; x_n$, the matrix given by $K_{ij} = k(x_i; x_j)$ is positive semidenite. The class of kernel estimators consists of functions $f(x) = \prod_{i=1}^n q^i \frac{k(X_i; x)}{k(X_i; x)}$. Suppose that $K_{ij} = K_{ij} = K_{ij}$

C.3 Asymptotic behavior of SALAD's $\mathcal{L}_{S}^{\Lambda}(h; w)$

Lemma 2. Assume that the reweighting function is Bayes-optimal, meaning that $\psi(x; m) = \psi(x; m)$. Then,

where $L_{CE}(h) = E_{x;m;z_0^0=1}[\log h(x;m)] + E_{x;m;z_0^0=0}[\log (1-h(x;m))]$ is the cross entropy loss on label $z_0 = \begin{pmatrix} 1 & x;m & P \\ 0 & x;m & p(jy=0) \end{pmatrix}$.

Proof. Let n_{data}^{SR} be the number of points from D that belong to the signal region. Under our assumptions, the empirical loss function can be written as

$$\begin{split} \text{L'}(h; \, \text{w}) \, / & \quad \frac{n_{data}^{S\,R}}{n^{S\,R}} \, \frac{1}{n_{data}^{S\,R}} \, \begin{array}{c} X \\ \text{log } h(x; \, m) \\ \text{x2D}_{S\,R} \\ \\ n^{S\,R} \, \frac{1}{n_{sim}^{S\,R}} \, \frac{1}{n_{sim}^{S\,R}} \, \begin{array}{c} X \\ \text{x2D}_{S\,R} \\ \text{x2D}_{S\,R}^{s\,i\,m} \end{array} \, \frac{p(x; \, mjy \, = \, 0)}{p_{sim}(x; \, mjy \, = \, 0)} \, log(1 \, h(x; \, m)) ; \end{split}$$

As n^{SR} ! 1, the rst term approaches $Pr(z^0 = 1) E_{x;mP} [log h(x;m)] = Pr(z^0 = 1) E_{x;mP} [log h(x;m)]$ 1) $E_{x;mjz^0=1}$ [log h(x; m)]. For the second term, we can construct $n_{data}^{SR;0}$, the amount of data where x is from p(jy = 0), to be equal to n_{sim}^{SR} such that the expression asymptotically approaches $Pr(z^0=0)$ $E_{x;mP_{sim}}$ $p_{p(x;mjy=0)}$ expression log(1 expression h(x; m)). Performing a change of expectation, this is equal to expression expressitogether, we have that

$$\lim_{n^{SR}!} \hat{L}(h; \hat{w}) / Pr(z^0 = 1) E_{x;mjz^0=1} [\log h(x; m)] Pr(z^0 = 0) E_{x;mjz^0=0} [\log(1 h(x; m))]$$

$$= L_{CE}(h):$$

Proofs

Proof of Theorem 1

Proof. From Theorem 3 of [31], we have that $L_C(f^\circ)$ $L_C(f^\circ)$ is bounded by the traditional ERM generalization gap of $L_C(f)$ $L_C(f^\circ)$, where $f = argmin_{f^\circ} f^\circ \int_{i=1}^{n^1} f^\circ f(x_i; m_i); y_i)$ is the classier learned on labeled data, plus the term $\frac{c_1}{e^{\min} a_{\min}^5} \frac{q^\circ \int_{i=1}^{n^1} f(x_i; m_i); y_i)}{n}$. We can apply standard learning theory bounds on $L_C(f)$ $L_C(f^\circ)$. In particular, this

quantity is equal to

$$\begin{split} L_{C}(f) & L_{C}(f^{?}) = (L_{C}(f) \quad \mathring{C}_{C}(f)) + (\mathring{C}_{C}(f^{?})) + (\mathring{C}_{C}(f^{?})) + (\mathring{C}_{C}(f^{?})) \\ & L_{C}(f) \quad L_{C}(f^{?}) + L_{C}(f^{?}) \quad L_{C}(f^{?}) \\ & 2 \sup_{f \geq F} L_{C}(f) \qquad \mathring{C}_{C}(f)j; \end{split}$$

where we have used the fact that $L_C^{\uparrow}(f) = L_C^{\uparrow}(f^2)$. Then, using uniform convergence bounds, such as Theorem 3.3 of [44], we have

$$L_{c}(f)$$
 $L_{c}(f^{?})$ 22R_n('F) + $\frac{r}{\log 2}$ $\frac{1}{2}$

This gives us our desired result.

D.2 Proof of Theorem 2

Proof. We dene the true (cross-entropy) loss as

$$L_{S}(h; w) = Pr(z^{0} = 1)E_{z^{0}=1}[log h(x; m)] Pr(z^{0} = 0)E_{x; m 2P_{sim}^{SR}}[w(x; m) log(1 h(x; m))];$$

where $z_0=1$ for x; m P and 0 for x; m P(jy = 0). Next, dene w(x; m) = $\frac{q(x;mjz=1)}{q(x;mjz=0)}$ and let ${\mathfrak W}$ be the weight ratio learned by our model. Let ${\mathfrak h}= \underset{h \geq 0}{\operatorname{argmin}} h_{1} {\mathfrak L}(h; {\mathfrak W})$, and let ${\mathfrak h}^2= \underset{h \geq 0}{\operatorname{argmin}} h_{2} {\mathfrak H}(h; {\mathfrak W}^2)$. Intuitively, ${\mathfrak K}$ corresponds to the true dierence between P_{SR} and P_{SR} and P

$$L_{S}(\hat{h}; \hat{w}) \quad L_{S}(h^{?}; \hat{w}) = [L_{S}(\hat{h}; \hat{w}) \quad \hat{L}_{S}(\hat{h}; \hat{w})] + [\hat{L}_{S}(\hat{h}; \hat{w}) \quad \hat{L}_{S}(\hat{h}^{?}; \hat{w})] + [\hat{L}_{S}(\hat{h}^{?}; \hat{w})] + (\hat{L}_{S}(\hat{h}^{?}; \hat$$

We know that $L_S(h, w)$ $L_S(h, w)$, so

We rst bound $\sup_{h;w} jL_S(h;w)$ $flact{L}_S(h;w) j$. For notation, we rewrite $llact{L}_S(h;w)$ as $llact{L}_S(h;g)$, where $llact{W}(x;m) = \frac{g(x;m)}{1-g(x;m)}$ and $llact{g}$ belongs to some function class $llact{G}$. Then, using $llact{g}$ heorem 3.3 from [44], we get that $llact{S}$ suph; $llact{W}$ $llact{L}_S(h;w)$ $llact{L}_S$

density ratio, and let $B_1 = \max_{x;m} f \log h^?(x;m); \log(1 h^?(x;m))g$. Assume that $B_1 < 1$. We can apply standard concentration inequalities here (Hoeding) to get that f(x) = f(x) + f(x) +

Finally, we bound $L_S^{\circ}(h^?; w) = L_S^{\circ}(h^?; w)$. We can write $L_S^{\circ}(h^?; w) = L_S^{\circ}(h^?; w)$ as

$$I_{NS}(h^{?}; W) = I_{NS}(h^{?}; w) = \frac{1}{n} X (W(x; m) - W(x; m)) (-\log(1 - h^{?}(x; m))): (D.3)$$

Dene = $max(log(1 h^?(x; m)))$ 0 for x; m 2 D_{SR}^{sim} , which is small as long as $h^?(x; m)$ suciently classies x and is hence a property of how separated the reweighted simulation and true data is. Then,

$$j \downarrow_{NS}(h^{?}; w) \qquad \downarrow_{NS}(h^{?}; w) j \qquad \sum_{SR} X \qquad j w(x; m) \qquad w(x; m) j:$$
 (D.4)

Recall that $\psi(x;m) = \frac{g(x;m)}{1-g(x;m)}$ and $w(x;m) = \frac{g^2(x;m)}{1-g^2(x;m)}$ where $g^2(x;m) = \Pr(z = 1jx;m)$, so $j\psi(x;m) = \psi(x;m) = \frac{jg(x;m)}{(1-g(x;m))(1-g^2(x;m))}$. This denominator is greater than

$$(1 \quad g_{max}^{n})(1 \quad g_{max}^{n}). \text{ Then,}$$

$$j\hat{C}_{S}(h^{?}; w) \quad \hat{C}_{S}(h^{?}; w)j \quad \frac{X}{(1 \quad g_{max}^{n})(1 \quad g_{max}^{n})n^{SR}} \int_{x; m2D_{SR}^{sim}}^{y} jg(x; m) \quad g^{?}(x; m)j: \quad (D.5)$$

We now look at the classier for training g. The per-point cross entropy loss for (x;m;z) is $'(g(x;m);z) = \log g(x;m)$ for z=1 and $\log(1-g(x;m))$ for z=0. WLOG, assume for some x and m, $g^?(x;m) > g(x;m)$. Then $j'(g^?(x;m);1) - (g(x;m);1)j = \log \frac{g^?(x;m)}{g(x;m)} = \log 1 + \frac{g^?(x;m)}{g^?(x;m)} = 1 - \frac{g^?(x;m) - g(x;m)}{g^?(x;m) - g(x;m)} = \frac{g^?(x;m) - g(x;m)}{g^?(x;m)} = \frac{g^?(x;m) - g($

$$j\overset{\wedge}{L_{S}}(h^{?}; w) \overset{\wedge}{L_{S}}(h^{?}; w)j \underbrace{(1 \overset{\circ}{g_{max}})(1 \overset{\circ}{g_{max}} n^{S} \overset{\circ}{n^{S}} j'(g(x; m); z) }_{x; m 2S} (g^{?}(x; m); z))j} \\ \underbrace{(1 \overset{\circ}{g_{max}})(1 \overset{\circ}{im} g^{?})}_{max} n^{S} \overset{\circ}{lm} [j'(g(x; m); z))} (g^{?}(x; m); z)j] + B_{2} \overset{\circ}{log} \frac{1}{2} \frac{1}{2} \frac{1}{n^{S}} \frac{1}{n^{S}}$$

where $B_2 = \max_{x;y} f'(g'(x;m);z); '(g'(x;m);z)g = \log(\min f''_{min}; g''_{min}g)$. We assume that B_2 is nite, so there exists a constant c such that

$$j\hat{\mathcal{L}}_{S}(h^{?}; \hat{w}) \quad \hat{\mathcal{L}}_{S}(h^{?}; w)j \quad \frac{n \frac{SR}{sim}}{g_{max}^{2}(1 \frac{g_{max}^{2}}{g_{max}^{2}})n^{SR}} cjL(g^{2}) \quad L(g^{?})j + B_{2} \quad \frac{1}{og} \frac{1}{2n_{sim}^{SR}} cjL(g^{2})$$

where L(g) = $E_{x;m2SR}$ ['(g(x; m); z)]. Since g?(x; m) is Bayes optimal, jL(g) L(g?) = L(g) L(g?) = L(g) L(g?) + L(g?) L(g?) + L(g?) L(g?) L(g?) 2 sup_{g2G} jL(g) L(g)j. From Theorem 3.3 in [44], this is bounded by $2R_{n^{SB}}$ ('G) + $\frac{\log 1}{2n^{SB}}$ with probability at least 1. Then, applying a union bound, with probability 1, we have

$$j \hat{C}_{S}(h^{?}; \hat{w}) = \hat{C}_{S}(h^{?}; \hat{w})j \\ (1 - \frac{1}{g} - \frac{1}{max}) \frac{\$_{R}}{max} n_{max} n$$

Putting everything together with another union bound, with probability $\mathbf{1}$, the generalization error is at most

$$L_{S}(H'; w) \quad L_{S}(h^{?}; w) \quad 2R_{S_{R}^{R}}('_{S} fH; Gg) + (1 + WB_{1}) = \frac{r}{\log \frac{8}{8}}$$

$$+ \frac{n_{Sim}^{SR}}{(1 - g_{max})(1 - g_{max}^{?})} n_{SR}^{SR} 4cR_{nSB}(' G) + 2c_{2n}^{T} \frac{4g_{+B}}{S_{B}} \frac{s}{s} = \frac{s}{2n_{Sim}^{SR}}$$

$$(D.6)$$

E Experiment Details

E.1 Multi-CWoLa Experiments

For the Multi-CWoLa experiment, we used the anomaly and simulation data from the Pythia 8 simulations in the LHC Olympics Dataset to create an unlabeled dataset we want to perform anomaly detection on [3]. We have k=3, and construct $M_i(m)$ based on the thresholds [[3:3;3:7]; [0:09;0:13]; [0:3;0:35]] on the rst three features. For standard CWoLa, only the third feature is regarded as the resonant feature, and it is thresholded with the interval [0:3;0:35]. We constructed training datasets of varying sizes with class balance Pr(y=1)=0:149. We used one test dataset with 65755 randomly sampled anomaly points and 161658 randomly sampled background points.

All methods were trained using scikit-learn's MLPClassier with $max_iter=5000$. For Multi-CWoLa's weak supervision step, we learn the parameters of the graphical model using SGD and PyTorch [45] with class balance Pr(y = 1) = 0.25, 30000 epochs, and learning rate = 1e 6.

E.2 Multi-SALAD Experiments

Setup We use MLPs from Keras [46], each with 3 hidden layers of dimension 32, ReLu activation, and trained with cross-entropy loss and the Adam optimizer. We train for 50 epochs, batch size 200, and default parameters otherwise. Finally, we evaluate our approach on a new test set containing 200000 background points and 200000 anomaly points. This test set is used to produce the signal eciency to rejection rate. All experiments were run on a personal laptop.

Additional Results In Figure 7, we show our results on individual runs. This is because computing the condence intervals of these curves averaged across the 10 random runs is too noisy due to the magnitude of the reciprocal 1/FPR.

Symbol	Used for
х	Discriminative feature x 2 X.
m	Resonant feature vector of length k, $m = [m^1;; m^k] 2 R^k$.
У	True unknown label y 2 Y = $f0$; +1g, where 0 is background and 1 is signal.
P; p	Distribution and density of data (x; m; y).
I_{mi}	Interval along which ith resonant feature m ⁱ is thresholded to produce
	signal region and sideband.
SR; SB	Signal region and sideband. For an interval I_{m^i} , $SR_i = f(x; m) : m^i 2 I_{m^i}g$
	and $SB_i = f(x; m) : m^i \ge I_{m^i} g$.
f	Classier f: X! Y used for anomaly detection.
D	Unlabeled dataset D = $f(x_i; m_i)g_{i=1}^n$ of discriminative and resonant features.
D_{SR} ; D_{SB}	Signal region and sideband of D, $D_{SR} = D \setminus SR$, $D_{SB} = D \setminus SB$.
SR; SB	Mixture weights corresponding to $p(y = 1jx 2 SR)$ and $p(y = 1jx 2 SB)$.
	It is assumed that $SR > SB$.
$M_i(m)$	Noisy membership label for the ith resonant feature, equal to 0 if x 2 $D_{SB_{i}}$
	and 1 if x 2 D_{SR_i} . $M(m) = M_1(m); :::; M_k(m)$.
ý y; i	Weak label drawn from estimated distribution on $p(yjM(m))$.
	Canonical parameters of graphical model on y ; $M(m)$ in (3.3).
Z	$_{y}$ scales with the class balance of y and $_{i}$ scales with the accuracy of $M_{i}(m)$.
ye;fM(m)	Partition function used for normalizing distribution $p(y; M(m))$ in (3.3).
i	y and M(m) scaled from f0; 1g to f 1; 1g.
	Accuracy parameter $i = p(M_i(m) = 1jy = 1)$ for the membership label
'c	of the ith resonant feature.
L _C (f)	Loss function 'c : Y Y! R for training classier f. Expected
•	loss on labeled data using f, $L_C(f) = E['_C(f(x); y)]$.
C (f) f	Optimal classier trained on innite labeled data, $f^? = argmin{f2F} L_C(f)$.
	Empirical loss on D with weak labels using f, $L_C(f) = \frac{1}{n} \int_{i=1}^{P} (c(f(x_i); \gamma_i)).$
D ^{sim}	Classier learned using Multi-CWoLa, $f \triangleq \operatorname{argmin}_{f2F} L_C^{\Lambda}(f)$.
nsim nsim	Simulation dataset used in standard SALAD, $D^{sim} = f(x_i; m_i)g_{i=1}^{n_{sim}}$.
Dsim, Dsim	Has distribution P_{sim} and density $p_{sim}()$.
w(x; m)	$D_{SB}^{sim} = D^{sim} \setminus SB$, $D_{SR}^{sim} = D^{sim} \setminus SR$.
	Density ratio between D_{SB}^{sim} and D_{SB} used for reweighting,
ĝ	$w(x; m) = \frac{p(x; mjy=0)}{p_{sim}(x; mjy=0)}.$
	Classier trained to classify D_{SB}^{sim} vs D_{SB} , used for approximating $w(x; m)$
L _S (h; w)	when $jD_{SB}^{sim}j = jD_{SB}j$.
ĥ	Cross-entropy loss function used to classify D_{SR}^{sim} reweighted with w vs D_{SR} .
$D_1^{sim}; :::; D_k^{sim}$	Classier trained using L _S .
ß sim	k multiple simulation datasets used in Multi-SALAD.
n ^{SR}	Dataset aggregated from D_1^{sim} ;; D_k^{sim} .
n ^{SB}	$n^{SR} = jD_{SR}j.$
n ^{S, R}	$n^{SB} = jD_{SB}j.$ $n^{SB} = iD_{SB}m^{SB}$
h [?]	$n_{sim}^{SR} = jD_{SR}^{sim}j$.
W	The optimal classier h^2 = argmin _{h2H} L _S (h; w).
	The maximum ratio between the simulation and true background,
	$W = \max_{x,m} w(x; m).$

Table 2. Glossary of variables and symbols used in this paper.

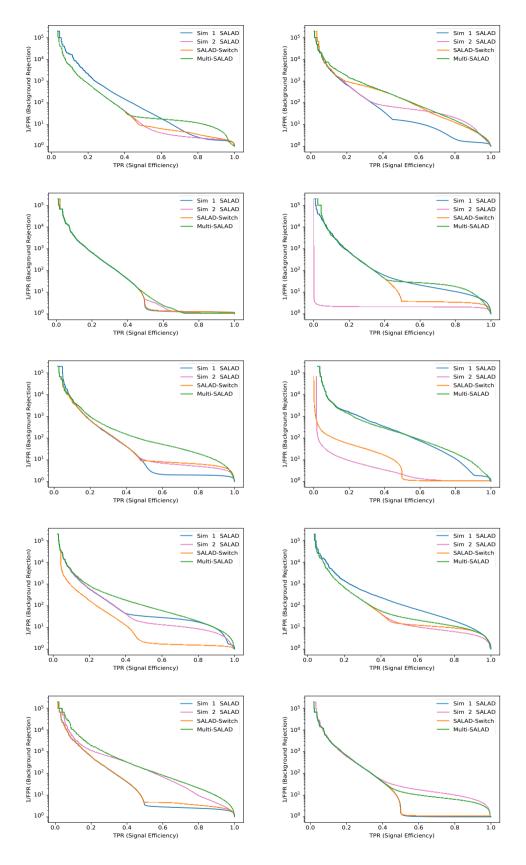


Figure 7. Results on individual runs.