

000
001
002
003
004
005
006
007
008
009
010
011054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

Open-World Learning Without Labels

Anonymous CVPR 2021 submission

Paper ID 2301

Abstract

Open-world learning is a problem where an autonomous agent detects things that it does not know and learns them over time from a non-stationary and never-ending stream of data; in an open-world environment, the training data and objective criteria are never available at once. The agent should grasp new knowledge from learning without forgetting acquired prior knowledge. Researchers proposed a few open-world learning agents for image classification tasks that operate in complex scenarios. However, all prior work on open-world learning has all labeled data to learn the new classes from the stream of images. In scenarios where autonomous agents should respond in near real-time or work in areas with limited communication infrastructure, human labeling of data is not possible. Therefore, supervised open-world learning agents are not scalable solutions for such applications. Herein, we propose a new framework that enables agents to learn new classes from a stream of unlabeled data in an unsupervised manner. Also, we study the robustness and learning speed of such agents with supervised and unsupervised feature representation. We also introduce a new metric for open-world learning without labels. We anticipate our theories and method to be a starting point for developing autonomous true open-world never-ending learning agents.

1 Introduction

Autonomous robots and self-driving vehicles are emerging technologies that are predicted to grow rapidly in quality and quantity in the near future. Vision-based recognition is an important subsystem of such autonomous agents. Visual recognition systems combine a feature extraction (perception) subsystem and inference (decision maker) subsystem. In real-world applications, environments of autonomous robots and self-driving vehicles change over time. This change will often introduce new classes, new attributes, and even a shift in the distribution of existing classes. In an open-world, the agent must detect the new classes/attributes and adapt.

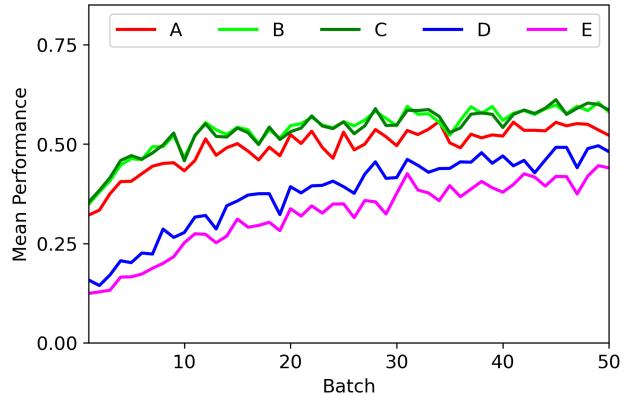


Figure 1: This paper formalizes and explores solutions to open-world learning without labels, including defining a new metric for performance measurement on such problems. Mean performance of five open-world recognition systems as they adapt to a 100 image batches of mixing known and new classes. The performance of open-world learning directly depends on the quality of feature representation, the detection of novelty, and the type of feedback during learning. **Can you determine which of the five curves (A–E) learnt the new classes in each batch using labels and which of them were learning without labels? Which of them used pure supervised feature representation, which used unsupervised features and which fused the two?** Please see experimental section for answers.

Babies can detect novel objects and learn them even if they are not given a semantic label with which to associate them. Similarly, online vision-based systems, autonomous robots, and self-driving vehicles may confront new classes of objects in areas and must learn to deal with them even if they don't know the semantic label to use. These systems should first detect that these objects are new and were not in the training set. Then, they should distinguish between the new classes. Also, they should recognize the new classes when they see them again. Ideally, each of the above steps should be done in an unsupervised manner. To achieve this goal, researchers should address many challenges such as novelty detection, change point detection, feature represen-

108 tation, transfer learning, meta-learning, continual learning,
 109 etc. Here, we investigate open-world unsupervised class
 110 incremental learning of image classifiers for autonomous
 111 agents. Our motivation is to build fundamentals and for-
 112 malize **open-world learning without labels** to be used
 113 along with other theories and solutions in the design of au-
 114 tonomous never-ending learning robots in the future.
 115

116 Computer vision and machine learning has seen a sub-
 117 stantial expansion in the work addressing self-supervised
 118 learning [21, 11, 10, 20, 16], unsupervised learning [8,
 119 23, 42, 36, 31], as well as open-set/out-of-distribution re-
 120 search [14, 13, 15, 40, 45, 32], and incremental learning
 121 [43, 46, 18, 24, 25], and this paper combines results from
 122 these three open topics to address a new problem, the de-
 123 tection and continual learning of new classes in an unsuper-
 124 vised manner – i.e., we formalize the problem of and de-
 125 velop the first class of True Open-World Learning (TOWL)
 126 algorithms to address the problem of open-world learning
 127 without labels.
 128

The contributions of this paper are:

- Formalizing open-world learning without labels prob-
 lem,
- Proposing a new metric to measure the quality of open-
 world learning
- Creating a framework to evaluate autonomous agent’s
 performance in both supervised and unsupervised
 open-world scenarios,
- Enhancing previous open-world image classifier using
 statistical Extreme Value Theory (EVT),
- Designing our TOWL autonomous agents that dis-
 cover, characterize, and learn new classes without la-
 bels from an open-world stream of data, and
- Investigating effect of feature representation in the ro-
 bustness and learning speed of autonomous agents dur-
 ing open-world learning.

2 Background

We cover only the background needed to develop/evaluate our problem and approach, reference to more related work is given in section 7.

2.1 Extreme Value Theory

Extreme Value Theory (EVT) is a branch of statistics that studies the behavior of extreme events on the tails of probability distributions [12, 5, 9]. EVT estimates the probability of events that are more extreme than any of the already observed ones. EVT is an extrapolation from observed samples to unobserved samples. There are two principal parametric approaches to modeling the extremes of a probability distribution: (1) block maxima and (2) threshold exceedance. The Hill Estimator approach is also commonly used which is a non-parametric approach. The block maxima uses Generalized Extreme Value distribution (GEV)

162 and threshold exceedance uses Generalized Pareto Distribu-
 163 tion (GPD). According to Fisher-Tippet asymptotic theorem,
 164 for normalized maxima of blocks of random variables
 165 $M_n = \max(X_1, \dots, X_n)$, there is a non-degenerate distribu-
 166 tion, which is a GEV distribution, which for our case must
 167 follow a Weibull distribution
 168

$$W(x; \mu, \sigma, \xi) = \begin{cases} e^{-(1+\xi(\frac{x-\mu}{\sigma}))^\xi} & , x < \mu - \frac{\sigma}{\xi} \\ 1 & , x \geq \mu - \frac{\sigma}{\xi} \end{cases} \quad (1)$$

2.2 Extreme Value Machine

The Extreme Value Machine (EVM) [35, 19] is a distance-
 174 based kernel-free non-linear classifier that uses Weibull
 175 families distribution to compute the radial probability of in-
 176 clusion of a point with respect to nearest members of other
 177 classes. For a given point x_i , they fit the Weibull on the
 178 distribution margin distance, half the distance to the nearest
 179 negative samples,

$$m_{i,j} = 0.5 * \|\hat{x}_i - x_j\| \quad (2)$$

181 for the τ closest points x_j from other classes. EVM pro-
 182 vides a compact probabilistic representation of each class’s
 183 decision boundary, characterized in terms of its extreme
 184 vectors. Each extreme vector has a family of Weibull distri-
 185 bution. Probability of a point belonging to each class is de-
 186 fined as the maximum probability of the point belonging to
 187 each extreme vector of the class. EVM uses greedy approx-
 188 imation for Karp’s set cover problem for model size reduc-
 189 tion by deleting redundant extreme vectors. In short, EVM
 190 for each input (point) computes the probability of inclusion
 191 to each class, i.e., the output is a vector of probabilities. The
 192 predicted class is computed by
 193

$$\hat{P}(C_l|x) = \max_k W_{l,k}(x; \mu_{l,k}, \sigma_{l,k}, \xi_{l,k}) \quad (3)$$

194 where $W_{l,k}(x)$ is Weibull probability of x corresponding to
 195 k extreme vector in class l .
 196

197 Source code for EVM is available, and a python version
 198 of EVM can be installed via pip. Our PyTorch enhanced
 199 version will be publicly released with the proposed method
 200 to reproduce the experiments.
 201

2.3 B3 Metric

204 B3 is a fuzzy probabilistic metric that measures the preci-
 205 sion and recall between clustering labels and true labels.
 206 Let’s denote features matrix with X such that each row is a
 207 feature vector that corresponds to a point (sample). Then,
 208 we can show the membership function of the true label with
 209 $\mu_Y(X)$ where the element in row i and column j is mem-
 210 bership of point i belonging to true class label j . Simi-
 211 larly, the membership function of the clustering label can
 212 be shown by $\mu_K(X)$. Let’s represent element-wise multi-
 213 plication by \odot , element-wise multiplication division by \oslash ,
 214 and a vector with all elements equal to one by $\mathbb{1}$. We can
 215

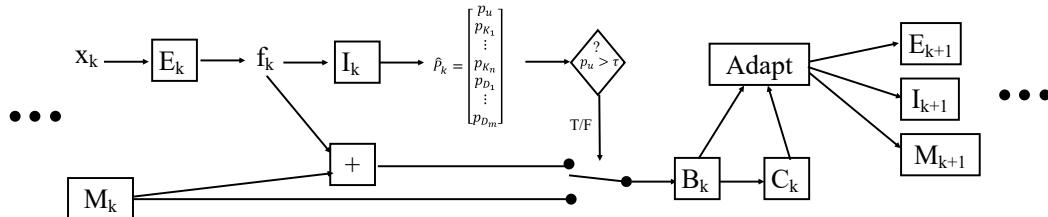


Figure 2: Function blocks diagram of open-world learning. At time step k , agent \mathcal{A}_k can be modeled with a memory M_k , a perception subsystem or feature extractor E_k , and a decision making or inference subsystem I_k . The agent acts on open-world stream \mathcal{S}^O , see Eq. (11). At time step k , feature extractor E_k converts data $x_k \in \mathcal{S}^O$ to feature f_k . Then, inference subsystem I_k predicts probabilities of data belonging to unknowns, knowns and discovered classes, $P_k = [p_u \ p_{K_1} \ p_{K_2} \ \dots \ p_{K_{n-1}} \ p_{K_n} \ p_{D_1} \ p_{D_2} \ \dots \ p_{D_{m_k-1}} \ p_{D_{m_k}}]^T$, where n is the number of known classes in the training set and m_k is the number of discovered classes. If probability of unknown p_u is less than a threshold τ , then, the buffer B_k is equal to the memory M_k , otherwise, the buffer B_k is equal to the concatenation of the feature f_k and the memory M_k . Next, each instance of the buffer B_k , gets a label at function C_K either supervised (human or other agents) or unsupervised via clustering. Finally, the agent \mathcal{A}_k will be updated to \mathcal{A}_{k+1} based on the buffer B_k and the supervised/unsupervised labels C_K . The agent \mathcal{A}_{k+1} will be used in the next time step $k + 1$.

compute B3 metrics by

$$A_{L \times C} = \mu_Y^\top \mu_K \quad M_{L \times C} = A \odot A \quad (4)$$

$$T_{C \times 1} = \sum_L A \quad S_{L \times 1} = \sum_C A \quad (5)$$

$$P_{C \times 1} = (\sum_L M) \oslash (T \odot T) \quad (6)$$

$$R_{L \times 1} = (\sum_C M) \oslash (S \odot S) \quad (7)$$

$$\text{Precision} = \frac{T^\top P}{T^\top \mathbb{1}} \quad \text{Recall} = \frac{S^\top R}{S^\top \mathbb{1}} \quad (8)$$

$$F = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (9)$$

We used the library from paper [3] to compute B3 scores.

3 Open-World Learning Formalizations

While [6] has a formal definition of open-world learning, it is insufficient to characterize open-world learning without labels, so we provide an expanded formalization and new metrics. Fig. 2 demonstrates a cycle of open-world learning. In open-world learning, agents start from an initial (potentially pre-trained) model. The agents confront a continuous stream of data that contains a mixture of known and unknown objects. The agent should (1) distinguish known from unknown, (2) distinguish classes of the unknown from each other, and (3) learn the recognized classes unknown without forgetting previously learned classes.

Definition 1 Open-World Stream

Let us define $\mathcal{K} = \cup_i K_i$ for known classes seen in training, as well as $\mathcal{U} = \cup_j U_j$, classes unseen in training. The world set is defined as $\mathcal{W} = \mathcal{K} \cup \mathcal{U}$. Let x_n be a sample drawn from \mathcal{W} at time step n . The closed-set stream is a time series

$$\mathcal{S}^C = \{x_n \in \mathcal{K} \quad \forall n \in \mathbb{N}\} \quad (10)$$

The open-world stream is a time series

$$\mathcal{S}^O = \{x_n \in \mathcal{W} \quad \forall n \in \mathbb{N} \mid (\exists x_i \in \mathcal{K}) \wedge (\exists x_j \in \mathcal{U})\} \quad (11)$$

Definition 2 Open-World Learner

Let the classifier of the agent \mathcal{A} at the time step be $f_n : \mathcal{W} \mapsto \mathbb{R}^{k_n+u_n}$, which maps an input $x_n \in \mathcal{W}$ to a vector of probabilities of x_n belonging to one of the currently known k_n classes $K_1 \dots K_{k_n}$ or one of the hypothesized u_n unknown classes $U_1 \dots U_{u_n}$ where we allow the number of known classes to expand as new labels are provided, and the number of hypothesized unknown classes to expand as new data is processed and determined to form a new class. We further break down the agent into its d dimensional feature representation extractor ($R(x_n) : \mathcal{W} \mapsto \mathbb{R}^d$), and its classification engine $C(x) : \mathbb{R}^d \mapsto \mathbb{R}^{k_n+u_n}$. The agent \mathcal{A} is an open-world learner if it acts on open-world streams and discovers and learns new classes $U_j \in \mathcal{U}$ in the stream after confronted with sufficient but bounded inputs drawn from each class. Learning a new class means predicting with a probability equal or greater than 0.5 for already seen instances in the class.

If supervised labels are provided to an open-world learner for instances in its unknown class U_j then we have a new known class $k_{n+1} = k_n + 1$ and $K_{k_n+1} = U_j$. This supervised model of open-world learning, converting unknown classes into known classes, is what is considered in prior work such as [6, 35].

However, we note that an agent may continue to function with many identified unknown classes that have only unsupervised pseudo-labels. Such a system may continue to improve its representation of that class even without labels as well as distinguish it from new unknown classes. This leads to the new definition for *open-world learning without labels*:

324
Definition 3 *Open-World Unsupervised Learner*
 325 *The agent \mathcal{A} is an open-world unsupervised learner if it is*
 326 *an open-world learner, and it learns the new classes without*
 327 *using labeled data from humans or other agents.*

328
 329 Full open-world learning agents may update their feature representation subsystems $R(x)$ based on the increasing stream of data.

330
Definition 4 *Open-World Class Incremental Learner*
 331 *If the agent only updates the inference subsystem $C(x)$ and*
 332 *keeps the feature representation $R(x)$ during learning, we*
 333 *call it an open-world class incremental learner.*

334 The latter two definitions can be combined, yielding open-world unsupervised class incremental learners, which is the focus of the remainder of the paper.

335 **Metric for Open-World Learning**

336 Because open-world learning mixes recognition of known and unknown classes, directly applying traditional metrics designed for either supervised or unsupervised learning does not necessarily work well.

337 Accuracy and balanced accuracy are the most popular metrics in supervised learning research. Unfortunately, accuracy cannot be defined when we do not have labels and hence cannot be applied to the unknowns. Even if we have ground truth labels for the data that goes into the unknowns used in testing since no label is provided, the unsupervised learning may split class or merge them, and hence we need unsupervised metrics, a.k.a clustering metrics. B3 (section 2.3) and Normalized Mutual Information (NMI) are two most widely used metrics in clustering research [2]. B3 and NMI are good metrics when the number of samples is large enough to represent the probability distribution of each class. In early versions of this work (see supplemental material), we were using just B3 or NMI on batches of data and eventually discovered that they were not well suited to open-world learning where we may have a large number of classes but only a small number of samples. None of them captures misclassifications of the unknown into an otherwise empty "known" class or the splitting of a known class into a mix of known plus unknown classes, e.g., breaking novel views into new classes. Therefore, we are proposing a new metric to overcome the issue of accuracy, B3, and NMI in open-world learning without labels. We call this the "Open-World Metric."

338 **Definition 5** *Open-World Metric*

339 Let N be the number of items to be evaluated in data X .
 340 Let Acc be accuracy for known data and $B3$ be the $B3$ metric (Eq. 9) for unknown data. Let us use subscripts ground truth and predicted categories of known and unknown such that known predicted as known is KK , known data which

341 was (incorrectly) predicted as unknown by classifier with
 342 KU , unknown data that (incorrectly) predicted as known as
 343 UK , and unknown data that predicted unknown by classifier
 344 with UU . For correct known predictions, we can use accu-
 345 racy and for correct unknown predictions, we can use $B3$,
 346 and we use incorrect predictions only in normalizing, then
 347 the Open-World Metric (OWM) score is computed by
 348
$$OWM = \frac{N_{UU} Acc(X_{KK}) + N_{UU} B3(X_{UU})}{N_{KK} + N_{KU} + N_{UK} + N_{UU}} \quad (12)$$

349 While we prefer $B3$, this measure can be generalized to
 350 combine other supervised or unsupervised metrics, e.g.,
 351 $OWM_{F1, NMI}$ would use the above definition with macro- $F1$
 352 instead of accuracy and NMI instead of $B3$.

353 **4 Evaluation Framework**

354 Prior evaluations of open-world learning in [6, 35], were
 355 fundamentally flawed because they used feature extractors
 356 that were trained on ImageNet 2012, but then they artifi-
 357 cially defined subsets of the 1000 classes as the base of
 358 knowns and incrementally tried to detect other ImageNet
 359 2012 classes as the unknowns. Thus, their feature space
 360 was trained using the "unknowns" as known and hence
 361 not a meaningful framework for proper open-world evalua-
 362 tion, even in a supervised setting. Therefore, we require a
 363 new evaluation framework, even for supervised open-world
 364 learning agents, and we do not reproduce data/tables from
 365 those prior works.

366 To evaluate and compare the performance of open-world
 367 learning algorithm in the task of image classification, (1)
 368 we use all 1000 classes of ImageNet 2012 train data set
 369 for training autonomous agents, (2) we use combinations of
 370 validation data set of ImageNet 2012 (known classes) and
 371 166 classes of ImageNet 2010 train data set that do not over-
 372 lap with ImageNet 2012 (unknown classes). We define four
 373 levels of tests: varying the number of instances per class
 374 and the number of unknown classes. Each test consists of
 375 50 batches, where the batch size was 100 images. Regard-
 376 less of the test level, each test has 100 classes of known, and
 377 each class of known has 25 images. So, each test has 2500
 378 known images. Test U10 has 10 unknown classes, where
 379 each class has 250 images. Test U25 uses 25 unknown
 380 classes, where each class has 100 images. Test U50 uses
 381 50 unknown classes, where each class has 50 images. And
 382 U100 tests have 100 unknown classes, where each class has
 383 25 images. Known classes, unknown classes, and images in
 384 each class are selected randomly. All images, known and
 385 unknown, were distributed randomly across each test, and
 386 we run each test 5 times, and we report the average OWM.
 387 (standard deviation is shown in supplemental). The base-
 388 line for this evaluation is an EVM model which does not
 389 adapt during open-world learning but just classifies items
 390 as known or unknown, placing all unknowns into one group
 391 instead of evaluation.

432

Algorithm 1: Image classification with EVM

Input: single image (optionally a batch of images)
Output: probabilities of all classes and top-1 predicted label

```

433   x ← normalize image to range [-1, +1]
434   f ← CNN(x)           // Deep feature
435   q ← EVM(X)          // Equation 3
436   m ← max (q)         // Maximum probability
437   u ← 1 - m            // Uncertainty
438   v ← concatenate ( u , q )
439   s ← ∑ v
440   p ←  $\frac{v}{s}$           // Estimated probabilities
441   y ← argmax (p)      // Predicted label
442   return p and y
443

```

448

To better understand the different aspects of the system, in evaluation, we consider three phases: closed-world where only data from known classes are present in the stream; Open-set, where unknowns are present but the system is not allowed to adapt; and open-world, where the system is allowed to learn from the data. The open-world stage has no access to the unknowns from the open-set stage. One should expect, and experimental data confirms, a drop in performance moving from closed-set to open-set, and then some level of recovery during the open-world stage.

459

5 Method

460

Our true open-world learning algorithm is summarized in Alg. 2 with three main elements: deep feature, enhanced EVM-based incremental-learning classifier for classification and detection of novel inputs, and clustering of detected novel inputs to form the basis of new classes.

466

The EVM has three important parameters: cover threshold, tail size (τ), and distance multiplier. We use 0.7 and 33998 for cover threshold and tail size, the same as the original EVM paper [35]. However, the original EVM formulation with its margin theorem concept using Eq. 2, is somewhat problematic for true open-worlds. The intuition behind the margin is that EVM is claiming half the space to the nearest other known class. That is fine for well-separated known classes, but it can easily be taking over too much open space for open-world learning as the assumption implies there are no classes in between the class being fitted and the nearest known classes. Because the original EVM experiments were tested on using pre-trained features that already separated all classes, this oversight may not have been apparent. Also, we find that margin is poorly defined in a highly imbalanced setting where a new class may have only a few samples. In such settings, we might need greater generalization from the few samples. Again this was not a problem in their experiments as they used balanced samples of well-separated classes, real open-world learning cannot

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

Algorithm 2: True Open-world Learner

Input: Single image and EVM model

Initialize: Empty clustered and residual sets

Config: $\delta = 0.001$, minimum number of images to start learning ψ , minimum number of cluster to start learning $\gamma = 2$, minimum cluster size to create a new class ρ , pre-trained features Ω

Output: new EVM

$f, p \leftarrow$ run Algorithm 1

// f : extracted feature

// p : class probabilities

$\phi \leftarrow$ first element of p // Unknown prob.

if $\phi > \delta$ **then**

 Insert f in Residual

if $\text{size}(\text{Residual}) > \psi$ **then**

$L, M \leftarrow \text{Clustering}(\text{Residual})$

// L : cluster labels

// M : Number of clusters

if $M > \gamma$ **then**

foreach cluster K **do**

if $\text{size}(K) > \rho$ **then**

$R^- \leftarrow \text{Residual} - K$

$N \leftarrow \text{concatenate} (\text{Clustered}, \Omega,$

$R^-)$

 Insert new class to EVM with K as + and N as -

 Insert K to Clustered

 Delete covered clusters from Residual

return EVM

presume these situations.

To address these issues, our enhanced EVM includes the idea of a distance multiplier d_m , which replaces the multiplier of 0.5 in Eq. 2 with a free parameter. If $d_m < 0.5$, then the model is smaller (more specialized), leaving some room between it and the nearest other known class. If we choose a higher value for distance multiplier $d_m > 0.5$ during incremental class addition, we can expand the class generalizing more. We tested on a range of values of the distance multiplier using a held-out validation data, which is used for optimizing open-set classification accuracy. Among them, 0.45 demonstrates the best separation between known validation and unknown validation sets of ImageNet—slightly less generalization than the original EVM paper. With $d_m = 0.45$, we leave some room for classes between two known classes.

For features representations, we used EfficientNet-B3 [39] from Timm library [41] as feature extractor. The EfficientNet-B3 was trained on all 1000 classes of ImageNet 2012 train data set. Then, we extracted features from the last layer of the network before logit, which has 1536 dimensions. Then, an extreme value machine was trained on the extracted features (frozen features). Algorithm 1, shows

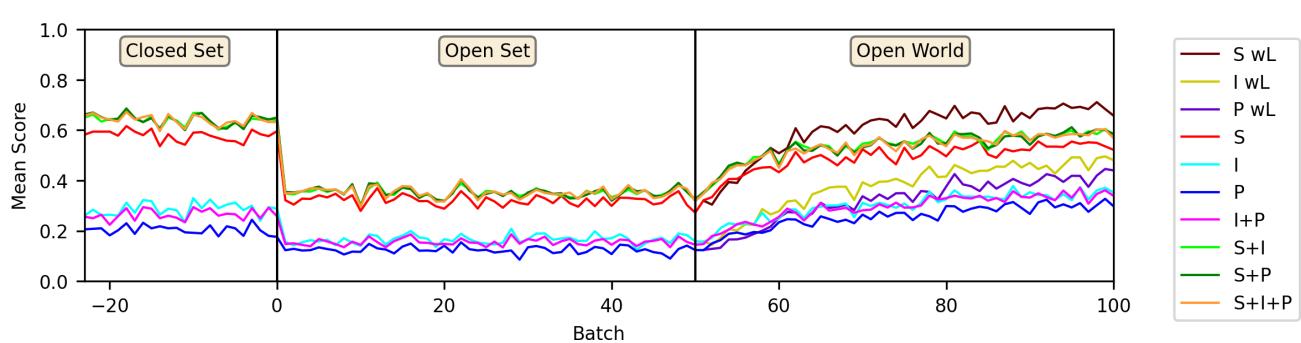


Figure 3: Average open-world metric scores over 5 runs of the for variations of TOWL algorithm when there are 100 unknown classes. All models start in the closed-set setting trained with supervised labels to build their base EVM. While supervised features updating with supervised labels (S_{wl}) is eventually the best open-world learner, during the open-set testing, it is not and the various fusion of supervised + unsupervised features ($S+I$, $S+P$, $S+I+P$) are better. Interestingly, when doing open-world learning without labels using the fused features ($S+I$, $S+P$, $S+I+P$), all outperform using just supervised features (S). Pure unsupervised features (I , P , $I+P$) are consistently worse, even when they are provided labels during the open-world phase. See Table 1 for the full feature/label combinations in the legend.

the details of EVM-based image classifier for the proposed autonomous agent.

To detect a novel instance, we threshold on the enhanced EVM probability. The Weibull family distribution often converges to zero rather quickly, and EVM generates a very sharp boundary. Thus, we declare an image as novel (and hence nominate it to create a new class) if the probability of the class of unknown of EVM (the first class, i.e., the class with label zero) is above a very small threshold ($\delta = 0.001$).

To add class incrementally from a stream of images, we determine which images should be combined to create a new class. All prior work pursued supervised open-world learning, getting labels for each of the detected novel images, and used to update the model. To develop our true open-world learner (TOWL), we proposed to collect nominated novel images into a residual set, and when the size of the set becomes greater than a threshold, we group them together to create new classes. While one might consider classic clustering, such as K-means, we don't have any prior expectations on the number of new classes. Automatically discovering related groups of data in unsupervised data, without parameters, is an important and still unsolved problem. There are only a few published clustering methods that are appropriate. In this paper, we use the Finch algorithm [36] for clustering, which, while it is formally parameter-free, still does not provide fully automatic operation since it produces multiple potential partitions among which we must choose. We choose the smallest partition as we don't expect a lot of classes. Then, for each cluster with sufficient points for EVM fitting threshold, the agent creates a class and adds it to the current EVM. Algorithm 2 shows a summary of the proposed TOWL approach combining enhanced EVM and Finch.

6 Experimental Results and Discussion

We trained three EfficientNet-B3 networks: (1) supervised learning on ImageNet 2012, (2) unsupervised learning on ImageNet 2012 data set using MoCo V2, and (3) unsupervised learning on Places2 data set using MoCo V2. Then, we extracted features from a layer before Logit and froze them. Next, we trained an EVM for each frozen feature set. Then, we trained four EVM models with the concatenation of the different features sets. In testing, we consider both open-world supervised learners with labels and open-world unsupervised learners, see table 1 for the 10 combinations tests.

We repeated each test 5 times with different random selections. Fig. 3 shows the mean of open-world scores of 10 configurations when the number of unknown classes is 100. Fig. 4 in the supplementary illustrates the performances for different number of unknown classes. Table 2 states the open-world score of the last 1000 samples.

We first discuss open-world learning without labels from the point of the view of the open-world metric; the supplemental material contains tables and discusses simple K+1 class open-set accuracy, where not surprisingly, there is not much gain from trying learn new unseen classes when only being measured with an open-set measure.

Figure 3 shows performance plots for the 10 systems when tested with 100 unknown classes; plots for 10, 25, and 50 unknown classes are in the supplemental material. Table 2, summarizes performance compared to the non-learning baseline for each of the 10, 25, 50 and 100 class experiments. From the plots, we see a significant drop from closed-set to open-set performance, but the drop is more dramatic for systems using at least some supervised features. Once open-world learning starts, they all improve, both with and without labels. The best performance, our

Table 1: Primary feature/label combinations used for experiments/plots

648	S	Supervised features	702
649	I	MoCo V2 on ImageNet 2012 features	704
650	P	MoCo V2 on Places 2 features	705
651	I+P	Concatenation of MoCo V2 on ImageNet 2012 and MoCo V2 on Places 2 features	706
652	S+I	Concatenation of supervised and MoCo V2 on ImageNet 2012s 2 features	707
653	S+P	Concatenation of supervised and MoCo V2 on Places 2 features	708
654	S+I+P	Concatenation of supervised, MoCo V2 on ImageNet 2012, and MoCo V2 on Places 2 features	709
655	wL	with label, i.e., supervised open-world learning	710
656			711
657			712
658			713
659			714
660			715
661			716
662			717
663			718
664			719
665			720
666			721
667			722
668			723
669			724
670			725
671			726
672			727
673			728
674			729
675			730
676			731
677			732
678			733
679			734
680			735
681			736
682			737
683			738
684			739
685			740
686			741
687			742
688			743
689			744
690			745
691			746
692			747
693			748
694			749
695			750
696			751
697			752
698			753
699			754
700			755
701			756

Table 2: Average on 5 tests, open-world scores of last 1000 images.

# Unknown classes	10		25		50		100	
	Feature extractor	Base	Algorithm	Base	Algorithm	Base	Algorithm	Base
S	0.3638	0.3694	0.333	0.3599	0.3109	0.3561	0.3031	0.3716
I	0.1677	0.1878	0.158	0.1918	0.1429	0.2012	0.1495	0.2206
P	0.1304	0.1547	0.1205	0.166	0.1112	0.1643	0.109	0.1844
I+P	0.1635	0.1814	0.1567	0.1968	0.1415	0.1945	0.1404	0.2095
S+I	0.3966	0.422	0.3633	0.4211	0.3415	0.4106	0.3342	0.4177
S+P	0.3934	0.4192	0.3619	0.4173	0.3361	0.4056	0.3372	0.4239
S+I+P	0.3953	0.4188	0.3597	0.414	0.336	0.4055	0.3352	0.4174
SwL	0.3641	0.62	0.3338	0.5749	0.3103	0.5553	0.3047	0.5893
IwL	0.1681	0.3737	0.1582	0.3346	0.1437	0.3204	0.15	0.3744
PwL	0.1304	0.2988	0.1208	0.2806	0.1112	0.2708	0.1098	0.3336

computational upper bound, is given by using supervised features and supervised open-world learning. However, we see that when using fused features (S+I, S+P, or S+I+P), unsupervised open-world learning comes very close to the upper bound and is superior to using either just supervised features for unsupervised open-world learning or using just unsupervised features with supervised open-world learning (Iwl, Pwl). We see the learning rate (improvement rate) of supervised open-world learning can be higher than unsupervised learning. However, the improved learning rate is easier to have when starting at lower performance.

The proposed TOWL algorithm builds on extended EVM, and its model uses its extreme vectors during evaluation and during the update; thus, it does not face catastrophic forgetting. While not done in these experiments, if the number of new classes grows significantly, we expect that simple updating of EVM might not be sufficient, and the feature extractor also should be updated, which might raise issues of catastrophic forgetting, which are beyond the scope of this paper.

In this paper, we use EVM with a distance multiplier of 0.45, and we select a threshold of 0.001 to be considered as a novel class to be updated. These parameters are fixed and not varied during testing. While these values are good for EfficientNet-B3 to be evaluated on ImageNet, for other data sets, they should be reevaluated on validation data.

Another parameter is how many detected novel class points are needed to begin clustering. Here, we choose

threshold 50 to start clustering. If we choose a higher threshold, the quality of clusters will increase, and the speed of learning decreases. Therefore, there is a trade-off between the quality of learning and the speed of learning. The minimum value for these thresholds depends on the clustering algorithm and quality of the feature extractor.

We choose threshold 5 new points in a cluster to use it to instantiate a new class – so few but not one-shot unsupervised learning. If the number of required samples is larger, the new class is better defined and generalizes better; however, again, the learning speed decrease. Thus, there is a trade-off between the quality of learning and the speed of learning. Informally, we observed that a smaller class size threshold required a larger value of distance multiplier to be effective. If the class size threshold is rather large (very low speed), the distance multiplier should be equal to that for known (training) classes. Future works should investigate an optimal policy to adapt these thresholds based on data.

In this paper, we used the Finch clustering algorithm [36] to cluster instances that are predicted as unknown. The Finch generates several partitions. Producing too few clusters is dangerous as it will cause two classes to merge, and once merged, the current approach does not have the ability to separate; thus, the confusion is permanent. Over clustering will cause the new class of EVM not to generalize to the full semantic concept of the original class. Therefore, selecting a partition with the proper number of clusters is necessary. In our tests, we used the Finch partition with the

756 minimum number of clusters because the threshold 50 to
757 clustering is small; future work should evaluate this choice
758 and ideally develop a fully automatic algorithm.
759

760 7 Related Works

761 This is the first paper to tackle Open-World Learning Without
762 Labels (OWLWL) problem. We deferred discussion of
763 related work until here, so our new problem and solution approach
764 was well defined, as a result, related work is given in context.
765 We now discuss related works but again note that none of them has directly addressed the OWLWL problem.
766

767 Unsupervised incremental Learning

768 Unsupervised incremental learning has been used in many
769 applications such as the prediction of musical audio signals
770 [30], hand shape and pose estimation [22], Financial Fraud
771 Detection [29], road traffic congestion detection [4], etc. In
772 the following, we briefly describe the closest works.
773

774 In [34], authors used a mixture of Gaussian latent space,
775 which uses dynamic expansion and mixture generative re-
776 play to minimize catastrophic forgetting in continual un-
777 supervised learning. They did experiments on MNIST
778 and Omniglot data sets. In [28], they designed a per-
779 son re-identification algorithm based on pedestrian Spatial-
780 temporal patterns in the target domain that consists of a
781 feature extractor (CNN) and a matching model (Bayesian).
782 Temporal patterns are not accessible in many image clas-
783 sifier agents. In [1], researchers proposed spike-timing-
784 dependent plasticity for spiking neural networks to learn
785 digits 0 to 9 incrementally. All of these three approaches
786 had very limited experiments and may not work in more
787 complex data such as ImageNet or Places2.
788

789 In excellent research [33], VGGface has used feature ex-
790 traction in videos. Then, a modified version of the nearest
791 neighbor to learn new faces (classes) incrementally. Also,
792 they designed a feature forgetting strategy to control mem-
793 ory size in the long run. The results in [35] shows that EVM
794 has better performance than the nearest neighbor. Thus, in
795 this paper, we do not use the nearest neighbor classifier. In
796 [27], they compared Support Vector Machines (SVM) with
797 Extreme Learning Machines (ELM) in the task of incremen-
798 tal learning for face verification in video surveillance. They
799 found that ELM is slightly better than SVM.

800 Continual Recognition Inspired by Babies (CRIB) [38]
801 is an unsupervised incremental object learning environment
802 that can produce data that models visual imagery produced
803 by object exploration in early infancy. They reported that
804 single exposure yields catastrophic forgetting. The algo-
805 rithm's accuracy stays constant or decreases with a greater
806 number of objects. Also, a smaller learning exposure length
807 results in lower final accuracy. Their algorithm exploits 3D
808 models, so it could not be used as a basis of comparison in
809 this paper.

810 Open-world learning

811 Obviously, the most related work is open-world learning,
812 which was first formalized [6] with subsequent work in
813 [35, 26, 7, 37, 44, 17]. *In these prior works, supervised*
814 *labels were provided to support their incremental learning*
815 *in the open-world.* With the exception of [6, 35] these were
816 done using textual-data, or small images and hence are not
817 suitable to be used for comparison, for which we use the
818 state-of-the-art for ImageNet scale problems: the Extreme
819 Value Machine (EVM) [35].
820

821 8 Conclusions

822 In open-world learning, an autonomous agent learns con-
823 tinuously, discovering new classes in its non-stationary and
824 never-ending stream of data. In an open-world environment,
825 labels for data are often unavailable, and hence supervised
826 open-world learning is not a scalable solution for online or
827 real-time applications.
828

829 Here, we formalized the unsupervised open-world learn-
830 ing problem. Then we created a framework to evaluate au-
831 tonomous agent performance in open-world scenarios and a
832 new open-world metric suited for the evaluation of unsuper-
833 vised open-world learning. Also, we extended a prior open-
834 world image classifier using statistical extreme value theory
835 to better handle open-world learning. Then, we designed
836 a new true open-world learner (TOWL) and autonomous
837 agents that discover, characterize, and learn new classes
838 without labels from an open-world stream of data. TOWL
839 combines state-of-the-art clustering (Finch) with an exten-
840 sion of the extreme value machine to provide the first solu-
841 tion to unsupervised open-world learning.
842

843 We also investigated the effect of feature representation
844 on the robustness and learning speed of autonomous agents
845 during both supervised and unsupervised open-world learn-
846 ing. The learning speed was fastest with supervised open-
847 world learning. We found that combining supervised and
848 unsupervised trained features with our TOWL produced the
849 best results for unsupervised open-world learning, coming
850 close to a fully supervised system using supervised features
851 used with supervised learning. We conclude that having a
852 feature representation that is not overly tuned to the known
853 classes provides improved robustness in an open-set setting
854 and improves learning in an unsupervised open-world learn-
855 ing.
856

857 9 References

- [1] Jason M Allred and Kaushik Roy. Unsupervised incremental
858 stdp learning using forced firing of dormant or idle neurons.
859 In *2016 International Joint Conference on Neural Networks*
(IJCNN), pages 2492–2499. IEEE, 2016. 8
- [2] Enrique Amigó, Julio Gonzalo, Javier Artiles, and Felisa
860 Verdejo. A comparison of extrinsic clustering evaluation
861 metrics based on formal constraints. *Information retrieval*,
862 12(4):461–486, 2009. 4

- 864 [3] Breck Baldwin, Tom Morton, Amit Bagga, Jason Baldridge,
865 Raman Chandrasekar, Alexis Dimitriadis, Kieran Snyder,
866 and Magdalena Wolska. Description of the upenn camp
867 system as used for coreference. In *Seventh Message Under-*
868 *standing Conference (MUC-7): Proceedings of a Conference Held in Fairfax, Virginia, April 29-May 1, 1998*, 1998. 3
- 869 [4] Tharindu Bandaragoda, Daswin De Silva, Denis Kleyko,
870 Evgeny Osipov, Urban Wiklund, and Damminda Alahakoon.
871 Trajectory clustering of road traffic in urban environments
872 using incremental machine learning in combination with
873 hyperdimensional computing. In *2019 IEEE intelligent*
874 *transportation systems conference (ITSC)*, pages 1664–1670.
875 IEEE, 2019. 8
- 876 [5] Jan Beirlant, Yuri Goegebeur, Johan Segers, and Jozef L
877 Teugels. *Statistics of extremes: theory and applications*.
878 John Wiley & Sons, 2006. 2
- 879 [6] Abhijit Bendale and Terrance Boult. Towards open world
880 recognition. In *IEEE CVPR*, pages 1893–1902, 2015. 3, 4, 8
- 881 [7] Terrance E Boult, Steve Cruz, Akshay Raj Dhamija, M Gun-
882 ther, James Henrydoss, and Walter J Scheirer. Learning and
883 the unknown: Surveying steps toward open world recogni-
884 tion. In *Proceedings of the AAAI Conference on Artificial*
885 *Intelligence*, volume 33, pages 9801–9807, 2019. 8
- 886 [8] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Pi-
887 otter Bojanowski, and Armand Joulin. Unsupervised learning
888 of visual features by contrasting cluster assignments. *Ad-*
889 *vances in Neural Information Processing Systems*, 33, 2020.
890 2
- 891 [9] Enrique Castillo. *Extreme value theory in engineering*. El-
892 sevier, 2012. 2
- 893 [10] Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad
894 Norouzi, and Geoffrey E Hinton. Big self-supervised mod-
895 els are strong semi-supervised learners. *Advances in Neural*
896 *Information Processing Systems*, 33, 2020. 2
- 897 [11] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He.
898 Improved baselines with momentum contrastive learning.
899 *arXiv preprint arXiv:2003.04297*, 2020. 2
- 900 [12] Stuart Coles, Joanna Bawa, Lesley Trenner, and Pat Dorazio.
901 *An introduction to statistical modeling of extreme values*,
902 volume 208. Springer, 2001. 2
- 903 [13] Chuanxing Geng and Songcan Chen. Collective decision for
904 open set recognition. *IEEE Transactions on Knowledge and*
905 *Data Engineering*, 2020. 2
- 906 [14] Chuanxing Geng, Sheng-jun Huang, and Songcan Chen. Re-
907 cent advances in open set recognition: A survey. *IEEE Trans-*
908 *actions on Pattern Analysis and Machine Intelligence*, 2020.
909 2
- 910 [15] Chuanxing Geng, Lue Tao, and Songcan Chen. Guided
911 cnn for generalized zero-shot and open-set recognition us-
912 ing visual and semantic prototypes. *Pattern Recognition*,
913 102:107263, 2020. 2
- 914 [16] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin
915 Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch,
916 Bernardo Avila Pires, Zhaohan Guo, Mohammad Ghesh-
917 laghi Azar, et al. Bootstrap your own latent-a new approach
918 to self-supervised learning. *Advances in Neural Information*
919 *Processing Systems*, 33, 2020. 2
- 920 [17] Xiaojie Guo, Amir Alipour-Fanid, Lingfei Wu, Hemant
921 Purohit, Xiang Chen, Kai Zeng, and Liang Zhao. Multi-
922 stage deep classifier cascades for open world recognition. In
923 *Proceedings of the 28th ACM International Conference on*
924 *Information and Knowledge Management*, pages 179–188,
925 2019. 8
- 926 [18] Jianguo He, Runyu Mao, Zeman Shao, and Fengqing Zhu.
927 Incremental learning in online scenario. In *Proceedings of*
928 *the IEEE/CVF Conference on Computer Vision and Pattern*
929 *Recognition*, pages 13926–13935, 2020. 2
- 930 [19] James Henrydoss, Steve Cruz, Ethan M Rudd, Manuel Gun-
931 ther, and Terrance E Boult. Incremental open set intrusion
932 recognition using extreme value machine. In *2017 16th IEEE*
933 *International Conference on Machine Learning and Applica-*
934 *tions (ICMLA)*, pages 1089–1093. IEEE, 2017. 2
- 935 [20] Simon Jenni, Hailin Jin, and Paolo Favaro. Steering self-
936 supervised feature learning beyond local pixel statistics. In
937 *Proceedings of the IEEE/CVF Conference on Computer Vi-*
938 *sion and Pattern Recognition*, pages 6408–6417, 2020. 2
- 939 [21] Longlong Jing and Yingli Tian. Self-supervised visual fea-
940 ture learning with deep neural networks: A survey. *IEEE*
941 *Transactions on Pattern Analysis and Machine Intelligence*,
942 2020. 2
- 943 [22] Pratik Kalshetti and Parag Chaudhuri. Unsupervised incre-
944 mental learning for hand shape and pose estimation. In *ACM*
945 *SIGGRAPH 2019 Posters*, pages 1–2. ACM, 2019. 8
- 946 [23] Artúr István Károly, Róbert Fullér, and Péter Galambos. Un-
947 supervised clustering for deep learning: A tutorial survey.
948 *Acta Polytechnica Hungarica*, 15(8):29–53, 2018. 2
- 949 [24] Xialei Liu, Chenshen Wu, Mikel Menta, Luis Herranz, Bog-
950 dan Raducanu, Andrew D Bagdanov, Shangling Jui, and
951 Joost van de Weijer. Generative feature replay for class-
952 incremental learning. In *Proceedings of the IEEE/CVF Con-*
953 *ference on Computer Vision and Pattern Recognition Work-*
954 *shops*, pages 226–227, 2020. 2
- 955 [25] Yaoyao Liu, Yuting Su, An-An Liu, Bernt Schiele, and
956 Qianru Sun. Mnemonics training: Multi-class incremental
957 learning without forgetting. In *Proceedings of the IEEE/CVF*
958 *Conference on Computer Vision and Pattern Recognition*,
959 pages 12245–12254, 2020. 2
- 960 [26] Ziwei Liu, Zhongqi Miao, Xiaohang Zhan, Jiayun Wang,
961 Boqing Gong, and Stella X Yu. Large-scale long-tailed
962 recognition in an open world. In *Proceedings of the IEEE*
963 *Conference on Computer Vision and Pattern Recognition*,
964 pages 2537–2546, 2019. 8
- 965 [27] Eric Lopez-Lopez, Carlos V Regueiro, Xosé M Pardo, An-
966 analisa Franco, and Alessandra Lumini. Incremental learning
967 techniques within a self-updating approach for face verifi-
968 cation in video-surveillance. In *Iberian Conference on Pat-*
969 *tern Recognition and Image Analysis*, pages 25–37. Springer,
970 2019. 8
- 971 [28] Jianming Lv, Weihang Chen, Qing Li, and Can Yang. Un-
972 supervised cross-dataset person re-identification by transfer
973 learning of spatial-temporal patterns. In *Proceedings of the*
974 *IEEE Conference on Computer Vision and Pattern Recog-*
975 *nition*, pages 7948–7956, 2018. 8
- 976 [29] Tian Ma, Shiyou Qian, Jian Cao, Guangtao Xue, Jiadi Yu,
977 Yanmin Zhu, and Minglu Li. An unsupervised incremental

- 972 virtual learning method for financial fraud detection. In *2019 IEEE/ACS 16th International Conference on Computer Systems and Applications (AICCSA)*, pages 1–6. IEEE, 2019. 8
- 973
- 974
- 975 [30] Ricard Marxer and Hendrik Purwins. Unsupervised incremental online learning and prediction of musical audio signals. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 24(5):863–874, 2016. 8
- 976
- 977
- 978
- 979 [31] Bhaskar Mukhoty, Ruchir Gupta, K Lakshmanan, and Mayank Kumar. A parameter-free affinity based clustering. *Applied Intelligence*, pages 1–14, 2020. 2
- 980
- 981
- 982 [32] Poojan Oza and Vishal M Patel. C2ae: Class conditioned auto-encoder for open-set recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2307–2316, 2019. 2
- 983
- 984
- 985 [33] Federico Pernici and Alberto Del Bimbo. Unsupervised incremental learning of deep descriptors from video streams. In *2017 IEEE International Conference on Multimedia & Expo Workshops (ICMEW)*, pages 477–482. IEEE, 2017. 8
- 986
- 987
- 988
- 989 [34] Dushyant Rao, Francesco Visin, Andrei Rusu, Razvan Pascanu, Yee Whye Teh, and Raia Hadsell. Continual unsupervised representation learning. In *Advances in Neural Information Processing Systems*, pages 7647–7657, 2019. 8
- 990
- 991
- 992 [35] Ethan M Rudd, Lalit P Jain, Walter J Scheirer, and Terrence E Boult. The extreme value machine. *IEEE transactions on pattern analysis and machine intelligence*, 40(3):762–768, 2017. 2, 3, 4, 5, 8
- 993
- 994
- 995 [36] Saquib Sarfraz, Vivek Sharma, and Rainer Stiefelhagen. Efficient parameter-free clustering using first neighbor relations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8934–8943, 2019. 2, 6, 7
- 996
- 997
- 998 [37] Vikash Sehwag, Arjun Nitin Bhagoji, Liwei Song, Chawin Sitawarin, Daniel Cullina, Mung Chiang, and Prateek Mittal. Analyzing the robustness of open-world machine learning. In *Proceedings of the 12th ACM Workshop on Artificial Intelligence and Security*, pages 105–116, 2019. 8
- 999
- 1000
- 1001 [38] Stefan Stojanov, Samarth Mishra, Ngoc Anh Thai, Nikhil Dhanda, Ahmad Humayun, Chen Yu, Linda B Smith, and James M Rehg. Incremental object learning from contiguous views. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8777–8786, 2019. 8
- 1002
- 1003
- 1004 [39] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International Conference on Machine Learning*, pages 6105–6114, 2019. 5
- 1005
- 1006
- 1007 [40] Edoardo Vignotto and Sebastian Engelke. Extreme value theory for anomaly detection—the gpd classifier. *Extremes*, pages 1–20, 2020. 2
- 1008
- 1009
- 1010 [41] Ross Wightman. *PyTorch image models*, 2020. 5
- 1011
- 1012
- 1013 [42] Garrett Wilson and Diane J Cook. A survey of unsupervised deep domain adaptation. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 11(5):1–46, 2020. 2
- 1014
- 1015
- 1016 [43] Yue Wu, Yinpeng Chen, Lijuan Wang, Yuancheng Ye, Zicheng Liu, Yandong Guo, and Yun Fu. Large scale incremental learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 374–382, 2019. 2
- 1017
- 1018
- 1019
- 1020
- 1021
- 1022
- 1023
- 1024
- 1025
- 1026
- 1027
- 1028
- 1029
- 1030
- 1031
- 1032
- 1033
- 1034
- 1035
- 1036
- 1037
- 1038
- 1039
- 1040
- 1041
- 1042
- 1043
- 1044
- 1045
- 1046
- 1047
- 1048
- 1049
- 1050
- 1051
- 1052
- 1053
- 1054
- 1055
- 1056
- 1057
- 1058
- 1059
- 1060
- 1061
- 1062
- 1063
- 1064
- 1065
- 1066
- 1067
- 1068
- 1069
- 1070
- 1071
- 1072
- 1073
- 1074
- 1075
- 1076
- 1077
- 1078
- 1079