
In or Out? Fixing ImageNet Out-of-Distribution Detection Evaluation

Julian Bitterwolf^{*1} Maximilian Müller^{*1} Matthias Hein¹

Abstract

Out-of-distribution (OOD) detection is the problem of identifying inputs which are unrelated to the in-distribution task. The OOD detection performance when the in-distribution (ID) is ImageNet-1K is commonly being tested on a small range of test OOD datasets. We find that most of the currently used test OOD datasets, including datasets from the open set recognition (OSR) literature, have severe issues: In some cases more than 50% of the dataset contains objects belonging to one of the ID classes. These erroneous samples heavily distort the evaluation of OOD detectors. As a solution, we introduce with NINCO a novel test OOD dataset, each sample checked to be ID free, which with its fine-grained range of OOD classes allows for a detailed analysis of an OOD detector’s strengths and failure modes, particularly when paired with a number of synthetic “OOD unit-tests”. We provide detailed evaluations across a large set of architectures and OOD detection methods on NINCO and the unit-tests, revealing new insights about model weaknesses and the effects of pretraining on OOD detection performance. We provide code and data at <https://github.com/j-cb/NINCO>.

1. Introduction

While deep learning based models have shown impressive performance on many real world tasks, they often exhibit unforeseen behaviour when confronted with unknown situations like receiving an input that is not related to the task it has been trained on. Such samples are regarded as out-of-distribution (OOD) and deep neural network classifiers are known to make very confident predictions that those belong to one of the **in-distribution (ID)** classes (Hendrycks & Gimpel, 2017; Hein et al., 2019). This unwanted behaviour

is a serious obstacle when applying classifiers in real world applications. The purpose of OOD detectors is to reject OOD inputs, which depending on the application can mean requesting human intervention, steering towards a safe state, or simply abstaining from making a prediction, while at the same time letting ID inputs pass through.

Current OOD detection evaluations in image classification rely on the assumption that there is no ID class present in an OOD test image, not even in the background. We follow this definition and consider an input to be **out-of-distribution (OOD)** if it does not contain any of the in-distribution classes. However, we show that this assumption is not fulfilled for most of the current test OOD datasets for ImageNet-1K (IN-1K) of Russakovsky et al. (2015). The closely related task of open set recognition (OSR), which simultaneously demands detection of OOD data and high classification accuracy on the ID data, is evaluated on OOD datasets which have the same requirements as in OOD detection. We also examine the test OOD datasets that have been used in the OSR literature for IN-1K and find similar issues there. We demonstrate that occurrences of objects from ID classes in test OOD datasets are often correctly recognized by state-of-the-art OOD detectors, but as an unwarranted consequence held against them as mistakes in OOD detection evaluations (false “false positive”). Even in cases where current models struggle to identify ID content, e.g. if ID objects are partially occluded or in the background, OOD datasets containing ID objects are not future proof: when evaluating on them, one would not realize if a future model correctly predicts the class of a visible ID object.

The erroneous occurrences of ID objects in existing OOD datasets can be characterized into two failure modes, which we illustrate in Figure 1 and define as follows. **Categorical ID contaminations** show objects from ID classes which already are classes in a base dataset from which the test OOD dataset has been built. Their label coincides with an ID class or semantically designates a subset of an ID class, e.g. the class *hayfield* from the PLACES dataset and the IN-1k class *hay*. **Incidental ID contaminations** on the other hand occur in images which are supposed to belong to an OOD category but which contain an ID object. The object can be in the background or an aspect of the specific instance of the shown main object, e.g. the IN-1k class *plane* in an image of the OOD category *sky*. We show that ID contaminations

^{*}Equal contribution ¹University of Tübingen and Tübingen AI Center. Correspondence to: Julian Bitterwolf <julian.bitterwolf@uni-tuebingen.de>.



Figure 1. Contamination of OOD test sets with ID samples (ImageNet). *Blue*: ImageNet-1K class found in the image. (Brown): Label of the image in the original source dataset. **Top**: Samples from classes of the OOD dataset that by class meaning categorically overlap with ImageNet-1K classes. **Bottom**: Labels alone do not reveal that the images are ID, but incidental ID objects can be found.

strongly impact the conclusions which can be drawn from evaluating OOD detection methods by (1) systematically underestimating the true OOD detection performance and (2) unjustly punishing stronger OOD detectors.

Probing the true performance of OOD detectors for IN-1K requires a range of OOD classes that are challenging, diverse, and most importantly actually OOD. Compiling a test OOD dataset is indeed a challenging task, as the 1000 classes of IN-1K cover a fair portion of the images found in general image datasets. In this paper we introduce the **NINCO (No ImageNet Class Objects) dataset** which contains 5 879 images that we individually checked not to contain any ID object from the classes in IN-1K. These images are ordered into 64 OOD classes, which facilitates a specific analysis of the failure modes of an OOD detector. Additionally, we provide a dataset of “**OOD unit-tests**”, synthetic images which do not resemble real world photos, but are designed to test specific weaknesses that might have impact in real-world applications (e.g. due to a camera failure). We find that surprisingly many OOD detectors struggle to detect these supposedly easy unit-tests, in particular methods that work well on natural test data.

We provide a detailed OOD detection evaluation on NINCO for a range of eleven OOD detection methods across a large number of architectures and training schemes. Surprisingly, it turns out to be difficult for many OOD detectors to improve consistently over the baseline of Maximum Softmax Probability (MSP). While we confirm the observation that pretraining on larger datasets generally helps OOD detectors and particularly methods explicitly using pre-logit feature-information, we find that the type of pretraining has a strong impact.

2. Existing test OOD datasets for ImageNet-1K

First, we give an overview of the datasets that have been used to evaluate OOD detection performance for IN-1K as ID. In the following we use *blue* for the name of an ImageNet class and *brown* for the category name in the source dataset used for the generation of the test OOD dataset.

INATURALIST OOD PLANTS is a subset of 10 000 images curated by Huang & Li (2021) from 110 OOD plant species of iNat2017 (Van Horn et al., 2018) which is sourced from the iNaturalist project. It is frequently used as test OOD dataset (Xia & Bouganis, 2022; Ming et al., 2022).

Table 1. Percentage of ID samples, $p = \frac{ID}{ID+OOD}$, in commonly used test OOD datasets found by visual inspection of 400 random samples per dataset. Unclear samples are ignored (which are at most 6.7% (for PLACES) of the 400 samples).

Dataset	ID samples	Dataset	ID samples
PLACES	59.5%	SPECIES	57.0%
IMAGENET-O	20.2%	TEXTURES	25.6%
INAT. PLANTS	2.5%	TEXTURES43	20.0%
OPENIM.-O	4.9%	IN-1K-OOD	32.1%
SSB-HARD	41.6%	SSB-EASY	53.4%
360OPENSET	26.9%	COOD	38.2%

PLACES is a subset of Places365 (Zhou et al., 2017) curated by Huang & Li (2021) as “50 categories [...] that are not present in IN-1K”. It is used as test OOD dataset in (Huang & Li, 2021; Sun et al., 2021; Ming et al., 2022). The dataset contains 9 822 images from 50 environment classes. We find that several of these classes are either subsets of ID classes, e.g. *hayfield* (*hay*), *cornfield* (*corn*), *lagoon* (*seashore* and *lakeshore*), or contain mostly ID objects, e.g. *underwater* (*coral reef* and *scuba diver*), *ocean* (*seashore*).

TEXTURES (Cimpoi et al., 2014) contains 5640 images of various objects that show one of 47 patterns. It is used as test OOD dataset in (Huang & Li, 2021; Sun et al., 2021; Wang et al., 2021; Xia & Bouganis, 2022; Ming et al., 2022) and others. Wang et al. (2022a) address the issue of overlap with IN-1K and remove four categorically ID textures (*bubbly* (*bubble*), *honeycombed* (*honeycomb*), *cobwebbed* (*spider web*), *spiralled* (*spiral*)). We find that even their version (denoted as TEXTURES43) contains about 20% ID images. SPECIES was proposed in (Hendrycks et al., 2022) as OOD dataset for IN-21K (Deng et al., 2009) and should thus also be OOD for the IN-1K subset. Sourced from iNaturalist, it consists of 700 000 images from 1 316 species which were selected for not being in IN-21K. They sort the species into 10 superclasses. The largest superclass *Fungi* largely coincides with the IN-1K class *mushroom*, and also many of the remaining species are ID. Papers evaluating on SPECIES for IN-1K OOD detection include (Salehi et al., 2021; Yang et al., 2022; Song et al., 2022).

IMAGENET-O (Hendrycks et al., 2021) contains 2 000 images from IN-21K, excluding its subset IN-1K. To make the dataset challenging it was composed from images where a ResNet-50 classifier for a subset of 200 IN-1K classes attains high confidence. The samples being OOD relies on the assumption that IN-21K without IN-1K is OOD for IN-1K. However, this assumption does not hold, due to a significant overlap between ImageNet classes from IN-1K and IN-21K, e.g. *analytical balance/scale* and *pickle/cucumber*, and insufficient filtering for incidental ID objects.

OPENIMAGE-O (Wang et al., 2022a) consists of 17 632 images from the OpenImage-v3 (Krasin et al., 2017) test set which their human labellers categorize as OOD. It is also used in Yang et al. (2022).



Figure 2. A Vision Transformer confidently classifies ID objects in samples from popular OOD datasets (*source label in parentheses*) as the correct IN-1K class, but is marked down with false positives in OOD detection evaluation when using MSP (Max Softmax Prob.) as criterion. The weaker ResNet-50, in contrast, doesn’t recognize the ID objects and hence the MSP is low enough to reject all images wrongly as OOD. This illustrates how a better model (ViT in our case) can be unjustly punished when the test OOD dataset contains ID objects. For both models, the 95%TPR threshold is at a MSP of 38%. Origins of the images: PL=PLACES, SP=SPECIES, OO=OPENIMAGE-O, IO=IMAGENET-O.

360OPENSETCLASSES (Bendale & Boult, 2016) uses those 360 classes (15.000 samples) from ILSVRC2010 which are not part of ILSVRC2012. Like for IMAGENET-O, this leads to large semantic overlap, e.g. the class *organ pipe* coinciding with the ID class *organ*.

SEMANTIC SHIFT BENCHMARK (SSB) (Vaze et al., 2022) contains a *hard* and *easy* OSR benchmark, each consisting of 1000 classes, that were created by regarding the distances between nodes in the WordNet tree. Similar to 360OPENSETCLASSES, we find both categorical and incidental ID contamination, e.g. *rainbow lorikeet/lorikeet*. Papers evaluating on SSB include (Wen et al., 2022).

IMAGENET-1K-OOD (Wang et al., 2022b) contains 50.000 images from 1.000 classes randomly sampled from ImageNet-21K, such that those classes don’t overlap with ImageNet-1K and ImageNet-LT, another dataset introduced by the authors. Categorical examples include *bobwhite quail/quail* and *king vulture/vulture*.

COOD-BENCHMARK (Galil et al., 2023) is a general framework for benchmarking ImageNet-1K OOD detection. Their test set consists of ImageNet-21K samples which were filtered by class. It includes severe contamination, including categorical cases like *orange, orange tree/orange*.

2.1. Prevalence of ID samples in popular OOD datasets

Concerningly, several test OOD datasets for IN-1K that are in use by the community contain a substantial fraction of samples that show ID objects. Figure 1 shows some

typical appearances of ID data in supposedly OOD datasets. The categorical ID failure mode illustrated in the top part is the inclusion of samples from explicitly ID classes of the source dataset from which the OOD dataset has been built. For instance, the class *hayfield* from the PLACES-dataset overlaps with the IN-1K class *hay*. However, also in principally innocuous classes (bottom part), many incidental ID samples can still be found. Here, the occurring failure modes are numerous: some ID objects happen to be in the background, some are a prominent part of the depicted scene, and some happen to realize both the original class and the ID class. For instance, the class *table knife* contains samples which also show a *plate*, and the class *striped* from the TEXTURES-dataset often shows the stripes of a *zebra*.

In order to quantify the severity of ID objects in test OOD datasets, we manually check for ID objects in 400 random samples from each of the most commonly used datasets. For fair treatment, unclear and ambiguous samples, which we would exclude from NINCO introduced below, are ignored in this survey. The results in Table 1 show that for many of these common OOD detection benchmarks, a substantial fraction of samples is actually ID: For both the PLACES and SPECIES datasets, it is more than 50%. Only INATURALIST OOD PLANTS (2.5% of samples ID) and OPENIMAGE-O (4.9% ID) contain comparably few ID images.

2.2. Effect of ID contamination on OOD evaluation

In Figure 2, we show how OOD detection evaluation with incidental ID samples can unrightfully punish strong OOD-detectors: A better model can correctly recognize ID objects with high confidence even if they are in the background of the image, leading to a false “false positive” in the evaluation, while a weaker model not recognizing the ID object and providing a low-confidence prediction is “rewarded” with a false “true negative”. For example, the strong Vision-Transformer (ViT) (Dosovitskiy et al., 2021) identifies the *pole* besides an otherwise empty desert road, and thus has high confidence on the image where the weaker ResNet-50 does not recognize any ID class with high confidence. Similarly, in the second example, the ViT is punished with a false “false positive” for recognizing (above the detection threshold) the *oranges* in the background while ignoring the unknown flying fox (truly OOD), whereas the ResNet-50 even does predict a wrong ID class, namely *squirrel monkey*, but does so with low confidence (below the detection threshold), and is thus rewarded with a false “true negative”.

We quantify the effect of ID contaminations on evaluation results in customary OOD datasets in Figure 3 for the MSP baseline and the Mahalanobis OOD detection method (Lee et al., 2018). For the test OOD datasets which showed a large portion of ID samples in Table 1, we report the FPR at 95% TPR obtained with a ViT when evaluating on the

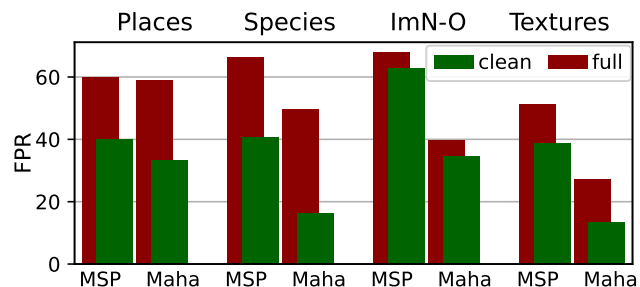


Figure 3. **OOD-detection before and after removing samples with ID-objects:** We show FPR (lower is better) of two OOD detectors (MSP and Mahalanobis distance) for a ViT, evaluated on cleaned and full subsets of four popular OOD datasets.

original 400 samples and our cleaned subsample of it not containing any more ID objects (detailed results for a range of models and methods can be found in Appendix J). We find that ID contaminations strongly impact the conclusions which can be drawn from evaluating OOD detection methods on those datasets. Most clearly, both methods perform substantially *better* after removing the images with ID objects from the OOD datasets, in some cases reducing the FPR by more than 50%. This is unsurprising: If a significant fraction of the dataset is actually ID, this fraction should not be detected as OOD by a well-performing method. Hence, evaluating OOD detection performance with partially ID data leads to a systematic *overestimation* of the true FPR of the OOD detection method and disadvantages better models as they are more likely to detect ID objects as discussed above. Additionally, we observe that the differences between OOD detectors become more pronounced. In Figure 3 it can be seen that for each dataset, the FPR for the Mahalanobis OOD detector decreases more than for the MSP-baseline. The effect is particularly strong for SPECIES (25.6% gain of MSP vs. 33.2% gain of Mahalanobis) and PLACES (19.6% gain vs. 26.3% gain), which are the two datasets we found to contain most ID samples. We further emphasize that due to the presence of large fractions of ID samples in most common benchmarks, even the performance of a perfect detector would saturate significantly above 0% FPR. For example with SPECIES, we find that for a strong current detector already more than 85% of the ‘false positives’ contain ID objects.

3. A new OOD test set for ImageNet-1K

As discussed in Sec. 1, an **OOD input for IN-1K** is an image that does not contain an object from one (or several) of the 1 000 IN-1K classes. These ImageNet classes are based on individual WordNet (Fellbaum, 1998) synsets, each consisting of one or more keywords that are synonymous in some context. During the ImageNet creation process (Deng et al., 2009), images were first collected from the web by using variations of each keyword of a respective class and then verified by humans to fit its synset’s definition.

In or Out? Fixing ImageNet OOD Detection Evaluation

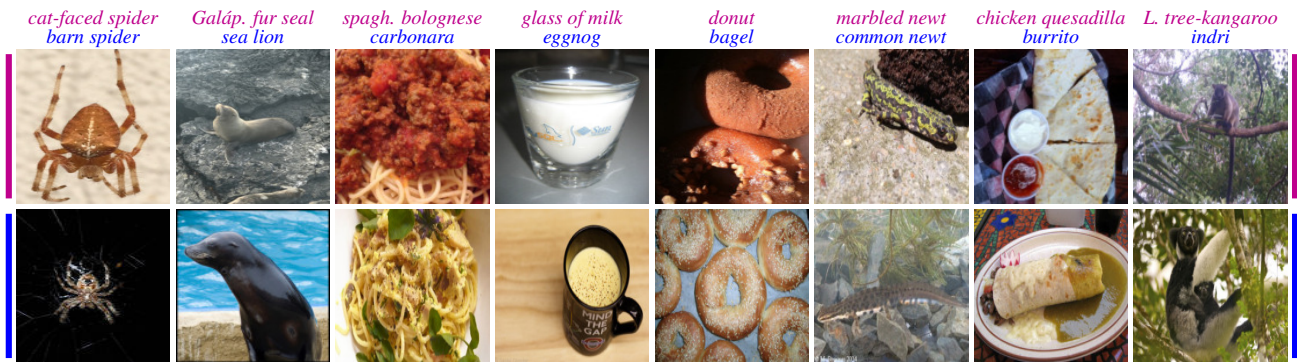


Figure 4. **Difficult OOD classes in NINCO**: Examples of images from some of NINCO’s most difficult (see Table 7) *OOD classes* (first row) and from the *ImageNet-1K class* (second row) which the *OOD class* is most frequently confused for.

Sourcing OOD test samples for ImageNet-1K from ImageNet-21K (or its subsets) based on class-labels has been leading to highly contaminated datasets (5 of the datasets in Table 1 are sourced from ImageNet-21K and all contain between 20% and 53% ID samples and show significant categorical contamination). This is partly due to the class-structure of those datasets: Both ImageNet-1K and ImageNet-21K contain leaf and internal nodes of the WordNet-tree as classes. While the internal nodes of ImageNet-1K are not ancestors to other Imagenet-1K classes, ImageNet-21K internal nodes can be ancestors to ImageNet-1K nodes, and vice versa. Moreover, there are ambiguous class-definitions in WordNet, like e.g. *police dog*, which is not parent or child of another dog class, but mostly shows a *german shepherd*, or an *alley cat* showing one of the many cat classes without being parent or child to other cat classes. Besides, there is significant incidental contamination even for nominally disjoint classes. Since the automation of filtering for challenging OOD data would require a strong detector that already solves the problems that the dataset is meant to pose, we conclude that it is impossible to construct a clean and challenging OOD dataset without manually checking the OOD samples for ID contamination.

In reality, many ImageNet samples fit one but not necessarily *all* keywords of their class label. This means that to make sure that OOD detectors are treated fairly¹, OOD test samples cannot fall into the definition of any keyword of any IN-1K class. For example, photos of the Sumatran orangutan cannot be considered OOD, since they could be included in the IN-1K class (*orangutan, orang, orangutang, Pongo pygmaeus*), even though *Pongo pygmaeus* only refers to the Bornean orangutan. To determine what counts as an ID object, we follow the WordNet glosses² as well as dictionary definitions of keywords and source dataset class labels. For difficult cases, we consult additional sources

¹For fair treatment of previous OOD *datasets*, such unclear samples that don’t fit all keywords were ignored in Table 1.

²One can look up synsets with glosses [here](#).

like Wikipedia. For example, the species *northern elephant seal* does not fall into the ID class *sea lion*, among other biological criteria distinguished by the fact that the former do not have ears while the latter do. An image of an OOD dataset can furthermore not incidentally contain ID objects, to avoid cases as in Figure 1 (bottom) and Figure 2.

3.1. NINCO dataset construction

For each OOD class of our new NINCO dataset, we start by **choosing a base class** which consists of all samples from a named class of an existing or newly scraped dataset. The majority of the NINCO base classes are sourced from SPECIES (Hendrycks et al., 2022), which provides images scraped from iNaturalist. For each base class, we carefully decide, based on WordNet glosses, iNaturalist taxonomy details and Wikipedia, whether it can be included according to the non-permissive interpretation described at the beginning of Section 3. The choice of base classes is not random, since there is no way to randomly sample from the set of concepts that might occur at test time. Rather, we aim for a variety of classes that are challenging, diverse and, most importantly, not actually categorically ID to begin with. Then for each base class, we **individually inspect each image** for ID objects. To help remembering the 1000 ID classes, we display the 5 top ID classes of a ViT’s prediction on each image. If an ID object is at least partially visible, the corresponding sample is removed. In cases where it is ambiguous whether we see an ID object in the image, the sample is not included in the cleaned dataset. As the iNaturalist data (including the SPECIES dataset) has been curated by experts and can be considered very reliable, we generally trust in the main object belonging to the species it is labelled as. For base classes chosen from the other sources, we consider ourselves competent to verify whether a label is correct. In addition to samples showing ID objects, we also remove images where no object from the OOD class is visible, e.g. we exclude pictures of animal traces or remains which frequently appear in iNaturalist. While for most existing datasets, the cleaning has been outsourced to external services like Amazon Me-

chancial Turk or student labellers. By researching all OOD classes and visually inspecting all their samples ourselves, we as authors of NINCO were able to do more in-depth research for each ambiguous case and obtain more coherent decisions, which we are positive leads to a higher quality dataset. Such high data quality is crucial for in-depth evaluations (Vasudevan et al., 2022; Shankar et al., 2021), as only being completely in-distribution free allows understanding a detector’s individual mistakes.

The NINCO (No ImageNet Class Objects) dataset consists of 64 OOD classes with a total of 5879 samples. The base classes which we cleaned to obtain NINCO were sourced from SPECIES (35 classes) (Hendrycks et al., 2022), PLACES (3 classes) (Zhou et al., 2017), which both are discussed in Section 2, as well as from the FOOD-101 dataset (7 classes) (Bossard et al., 2014), CALTECH-101 (4 classes) (Li et al., 2022), MYNURSINGHOME (4 classes) (Ismail et al., 2020), ImageNet-21k (1 class) and newly scraped from iNaturalist.org (2 classes) or other websites like Flickr (8 classes). Details for all NINCO OOD classes are given in Appendix F. We show samples from all NINCO classes in Figures 10 and 11 in Appendix H. In addition to NINCO, we also provide the 2715 OOD images obtained from cleaning 400 samples of eleven test OOD datasets as discussed in Section 2.2. In order to notice ID contaminations potentially biasing the drawn conclusions, we recommend to also evaluate on these cleaned versions when evaluating on those original benchmarks.

3.2. OOD unit-tests

Following common practice (e.g. Hendrycks et al. (2022)), we argue that evaluating an OOD detector on a range of simple, synthetic classes *besides* the variably challenging natural image classes of an OOD dataset can give additional insights about its OOD detection weaknesses. Example images and reproducibility details for all 17 pre-existing and newly proposed OOD unit-tests are included in Appendices G and H. Since these **OOD unit-tests** do not represent a diverse distribution of photos, but different modes of simple, synthetically generated image inputs which any good OOD detector should be expected to detect, we don’t include them in summary metrics or distribution plots. Instead, we suggest to count an OOD unit test as **failed** if a method has an FPR above a user-defined threshold, which we suggest setting at 10%, and to report the number of *failed* OOD unit-tests (which should be 0 for a strong OOD detector) alongside the aggregate results on a test OOD dataset like NINCO. For each OOD unit-test, we provide a set of 400 samples in typical ImageNet format, by mirroring the sizes and file formats of random ImageNet samples. While some OOD unit-tests may appear redundant at first sight, we find that they provide important information as some detectors e.g. mostly pass the *monochrome* test but completely fail

on *black*, which reveals a specific weakness that is very realistic to be encountered in practice.

3.3. OOD detectors and how to evaluate them

An **OOD detector** for inputs from the domain X of possible input images is represented by a score function $S : X \rightarrow \mathbb{R} \cup \{\pm\infty\}$ which is generally supposed to be larger on ID inputs than on OOD inputs. One example is the Maximum Softmax Probability (MSP) or confidence $S_{\text{MSP}}(x) = \max_{k=1,\dots,K} p_k(x)$ of a classifier with output probabilities p for K ID classes. The MSP is the standard baseline OOD detection method (Hendrycks & Gimpel, 2017), since it is intuitively expected to be low on OOD compared to ID inputs. Observing that standard classifiers are frequently overconfident on OOD inputs, OOD detection research aims at finding detectors that improve on this baseline. In Appendix C, we give an overview of a range of OOD detection methods which have been proposed for IN-1K as ID. An OOD detector is usually obtained by combining such an OOD detection method with a concrete classifier model. We analyze OOD detectors in terms of the fraction of falsely accepted OOD inputs at a true positive rate of 95%, short FPR. Detailed definitions can be found in Appendix D.

Different OOD classes (and similarly also different test OOD datasets) represent different probabilistic distributions of inputs that a detector is tested against. An important arising question is how the collective of individual performance measurements can be interpreted and whether they can be aggregated into one number that can be used to make an informed decision on which OOD detector works best. Certainly, the notion of ‘best’ may notably vary depending on the application and situation and we often cannot hope to model a ‘true’ out-distribution, or even be sure that it meaningfully exists. An aggregate number which gives a good overview of an OOD detector’s performance on the class based NINCO dataset is the **mean FPR** of the individual FPR values for each of the 64 OOD classes of NINCO.

However, for many applications it is not possible to model the potential OOD inputs that might be encountered at test time with a fixed probability distribution. Thus a single aggregate number cannot tell the full story, and may hide outliers in the FPR values. For one, some errors might be *less acceptable* than others, e.g. a FPR of 20.0% might be very bad for monochrome inputs, but would lose much significance when subsumed into a mean. For OOD unit tests, where OOD detectors can be expected to be very robust, we therefore propose regarding pass-fail statistics instead of mean FPR. Also, an evaluator might want to be informed about the *concrete failure modes* of the model, e.g. all OOD classes with a particular high FPR. An OOD detector showing consistent improvements on most of the OOD classes (instead of only in terms of the mean) can be seen as strong

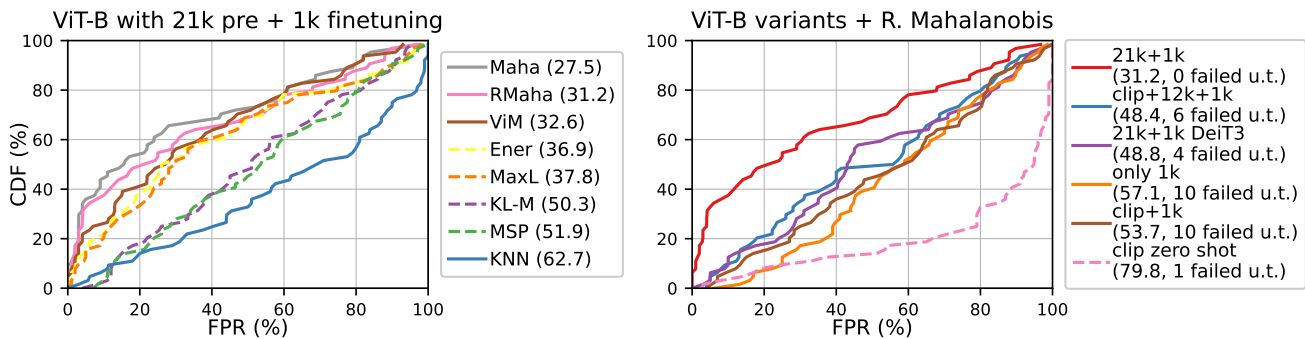


Figure 5. Cumulative distribution of the % of NINCO-classes for which an FPR at least as low as a given x-value is achieved. The area over this curve corresponds to the mean FPR. The further in the top left corner, the better. **The best methods explicitly access pre-logit features (Left):** Different OOD detection methods with a ViT-B pretrained on IN-21k (mean FPR in parentheses, pre-logit feature-accessing methods are solid, others dashed). **Not all pretraining helps (Right):** RMaha applied to ViT-B with different training variants (MCM for CLIP zero-shot is dashed). Only the top model does not fail OOD unit-tests.

evidence for the method yielding actual improvement, as opposed to the detector overfitting to a limited scope of test OOD data, which Wang et al. (2022a) describe as a form of hackability. Due to these considerations, and with the OOD data being organized into *OOD classes* as in NINCO, we suggest evaluations of OOD detectors to always provide the **distribution of results over OOD classes** and additionally to **make the individual results available**, such that the reader can make an informed comparison based on which types of OOD inputs are most relevant to them.

4. Evaluation results for OOD Detectors

We evaluate a range of IN-1K models obtained from the public timm-library (Wightman, 2019) and state-of-the-art OOD-detection methods on NINCO. We focus on transformer architectures and convolutional networks, both with and without pretraining. While most pretrained models were initially trained on IN-21K, we also include an EfficientNet trained via noisy student (Xie et al., 2019) on the JFT-300M dataset, and four ViTs with CLIP-pretraining (Radford et al., 2021) and subsequent fine-tuning, as well as a zero-shot CLIP model. A detailed description of all models can be found in Appendix B. We investigate the following commonly used OOD detection methods, which can be grouped into two categories: Max-Softmax (MSP) (Hendrycks & Gimpel, 2017), Max-Logit (Hendrycks et al., 2022), Energy (Liu et al., 2020) and KL-Matching (Hendrycks et al., 2022) derive an OOD-score exclusively from logit outputs, whereas Mahalanobis distance (Maha) (Lee et al., 2018), Virtual Logit Matching (ViM) (Wang et al., 2022a), ReAct (Sun et al., 2021), Relative Mahalanobis distance (RMaha) (Ren et al., 2021), and K-Nearest-Neighbours (KNN) (Sun et al., 2022) also leverage explicit information from the features of the DNN’s penultimate (pre-logit) layer. For the zero-shot evaluation of CLIP, we use Maximum-Concept-Matching (MCM) (Ming et al., 2022) and Cosine-similarity

(Cos) (Galil et al., 2023) to class-specific text-embeddings. Noting that OOD detection based on softmax of a cosine similarity to a specific feature vector has been proposed in different variants (Tack et al. (2020), Techapanurak et al. (2020) and MCM), we find that using it with classifier class means produces reasonable OOD detection results, marked below as relative cosine class similarity (RCos). We call those methods which explicitly access the pre-logit feature layer *feature-based* and provide an overview over all methods in Appendix C.

4.1. Results on NINCO

Comparison of OOD detection Methods. In Figure 5 (left), we illustrate the performance of a single ViT when combined with a range of OOD-methods. Overall, most feature-based methods, like Maha, RMaha and ViM, outperform the MSP-baseline by a clear margin. Notably, MaxLogit and Energy, which do not explicitly access the pre-logit features, are also able to strongly improve over MSP, while KL-Matching performs roughly on par, and KNN much worse. We observe that while Maha, RMaha and ViM improve over MSP in all FPR ranges, this is different for e.g. MaxLogit: For large FPR, it is similar to MSP, indicating that the method brings no advantage over MSP for hard test classes, and its improved mean performance is mainly due to lower FPR for the easier OOD classes. When regarding the mean FPR values of all method-model-combinations shown in Table 3 in Appendix A, we observe that while Maha in combination with a (pretrained) ViT is the single best OOD-detector, this method often performs worse when combined with other models. RMaha, however, yields good results with *all* models, and is together with (Relative) Cosine the only method which can fairly consistently improve over the MSP baseline in terms of mean FPR. For most models, it is either the best-performing method, or close to the best-performing method, which is somewhat surprising, given its relatively poor performance on the unit-tests. We further

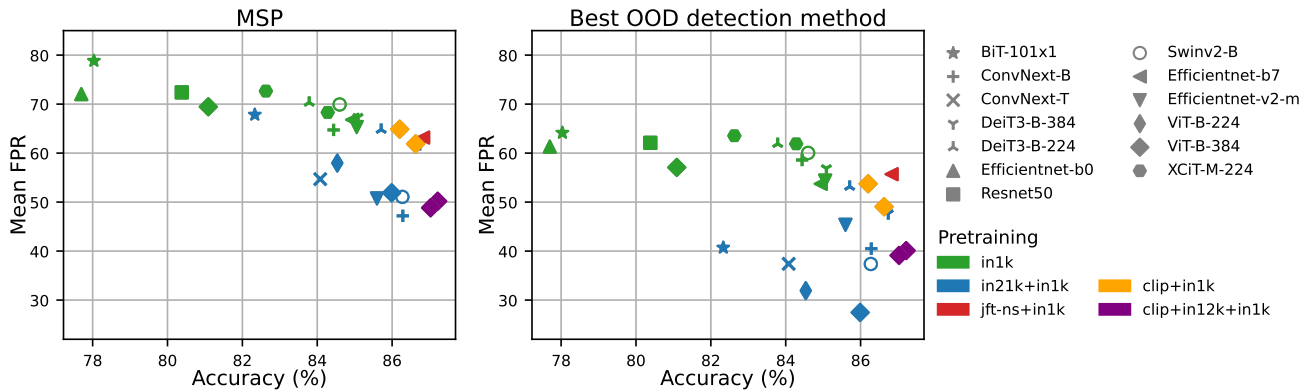


Figure 6. **IN-21K pretraining boosts feature-based OOD detectors on NINCO:** Mean FPR vs. accuracy for MSP and each model’s best detector, which (except for the noisy-student model) always explicitly accesses the pre-logit features. OOD detection strongly improves when using models pretrained on IN-21K. Additional CLIP-pretraining or on JFT can yield higher accuracy, but OOD detection need not be better than with IN-21K pretraining.

note that for all models (except the noisy-student model), the best-performing method always explicitly accesses the pre-logit features, and that in contrast to e.g. KNN, Energy and ReAct, even the adapted methods based on feature space cosine similarity Cos and MCM/RCos fairly consistently improve over the MSP-baseline. Each OOD dataset representing a different out-distribution that can be relevant for certain applications, we find that results vary on the cleaned subsets of eleven previous benchmarks which we evaluate in Appendix J, while the overall conclusions on the methods and models resemble those on NINCO.

Pretraining matters. In Figure 6, we plot the mean FPR on NINCO over the accuracy for all investigated models for both the MSP-baseline (left) and the best-performing OOD detector per model (right). For MSP, the mean FPR decreases roughly linearly with accuracy. Since most pretrained models (blue) have higher accuracy, they typically also show better OOD-detection performance, but also between models of similar accuracy, the pretrained ones achieve better mean FPR. For the best-performing OOD detector, improvements can be observed for models both with and without pretraining. Notably, the linear relation between FPR and accuracy disappears, and all purely 1K models (green) perform roughly on one level. In comparison, the gains for the majority of models pretrained on IN-21K (blue) are larger. In particular ViT and BiT benefit strongly from leveraging their respective best method, which as discussed above is always feature-based. In other words, pretraining helps in two ways: First, it leads to higher ID-performance (accuracy), which benefits methods like the MSP-baseline. Second, it creates better feature-embeddings for this task, which lead to improvements beyond the accuracy-MSP correlation. This is most clearly visible for the pretrained BiT-m, which has comparably low accuracy (82%) and hence no outstanding MSP-performance, but outperforms all 1k-models by a significant margin with features leveraging

ViM. However, as we observe in Figure 5 (right), the benefit of pretraining depends strongly on the specific data and training method: With RMaha, the ViT with ‘traditional’ IN-21K pretraining from (Steiner et al., 2022) clearly outperforms models with the distillation-based training of DeiT3 (Touvron et al., 2022), CLIP-pretraining or even CLIP with interjected IN-12K training. The zero-shot methods for CLIP, despite having shown promising results in (Galil et al., 2023) and (Ming et al., 2022) and performing well on the unit tests, are not competitive to IN-1k classifiers on NINCO. Regarding all methods, the five models trained with different pretraining strategies (EfficientNet-b7 with noisy student and four ViTs with CLIP-pretraining (Radford et al., 2021) and subsequent fine-tuning) show some of the highest accuracies in our survey, yet, their OOD-detection performance is surprisingly poor. Overall, we see strong indication that the precise type of pretraining has a large impact on whether it produces a feature space that is beneficial for feature based methods. In Appendix K we investigate whether IN-21K-pretraining particularly benefits detection of OOD classes that overlap with IN-21K classes, but we notice no substantially different changes between the model with and without pretraining.

Analysis of failure cases. In Figure 7 we plot the individual FPR for each OOD class of NINCO for the combination ViT+Maha, the overall best OOD detector in terms of mean FPR, and contrast it with ConvNext+Maha, which also shows good mean FPR. Performance varies widely between OOD classes, with both models severely struggling for some classes. Where the ViT shows large FPR, the ConvNext rarely performs better, while it also fails to detect certain classes like the *long-tailed silverfish* where the ViT does well. We illustrate samples from hard classes in Figure 4. Both models struggle to detect the *Galápagos fur seal* (98% FPR for the ViT), often confused with the IN-1K class *sea lion*, and *cat-faced spider* (confused with *barn spider*,

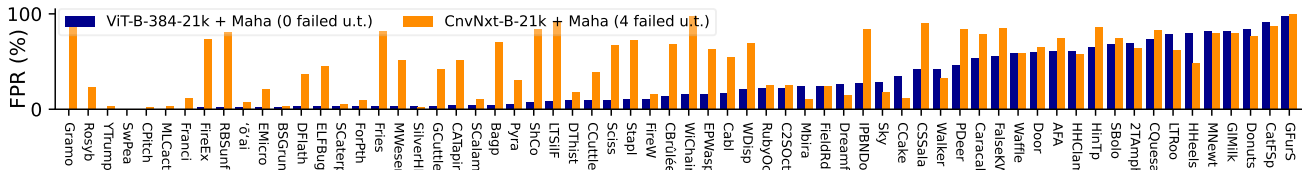


Figure 7. FPR of a pretrained ViT-B and pretrained ConvNext-B for all classes of NINCO.

91% FPR). From a human perspective, those classes are arguably hard to detect. We note, however, that it is possible to tell them apart, as a ViT IN-21K-classifier e.g. identifies the *Galápagos fur seal* as a *fur seal* (IN-21K class) in 92% of samples and misclassifies only 6% of them as a *sea lion*. The networks however also fail for classes more obvious to humans: *donut* (84% FPR ViT, confused with *bagel*), *spaghetti bolognese* (69% FPR, *carbonara*) and *chicken quesadilla* (73% FPR, *burrito*) also confuse both models.

4.2. Results on the OOD unit-tests

Auditing OOD detectors on the OOD unit-tests, we find that surprisingly many combinations of models and OOD detection methods struggle to distinguish supposedly easy inputs from ID-data. While results for all models and methods can be found in Appendix I, we provide some illustrative unit-test results in Table 2 for a ViT pretrained on IN-21k and a ConvNext both with and without IN-21K pretraining. In general, most methods fail fewer unit tests when applied to pre-trained models, however there are still many severely flawed combinations, often involving methods that would otherwise shine based on their detection of natural OOD data discussed above: especially the feature-based methods ViM, Maha and RMaha reveal weaknesses, each failing multiple unit-tests on at least 21 of 26 models. Many tested OOD detectors are vulnerable to *black*, *white* and *grey*, which is concerning as encountering inputs of this kind could occur in many real-world applications due to camera malfunction or occlusion. Here those feature-based methods only provide trustworthy

Table 2. **Some detectors fail OOD unit-tests:** FPR for a ViT and a ConvNext (with and without pretraining) on selected unit-tests. FPR larger than 10% count as failed and are thus marked red. Especially for methods relying on feature representations (like ViM and Maha) the OOD unit-tests reveal difficulties.

	method	<i>bla</i>	<i>whi</i>	<i>gre</i>	<i>hor</i>	<i>SmN</i>	<i>Rad</i>	<i>mon</i>
ViT21k	MSP	0.0	0.0	0.0	0.2	0.5	0.0	0.0
	ViM	0.0	100.0	46.0	0.0	0.0	0.0	0.5
	Maha	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	Cos	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Cnv1k	MSP	0.0	0.0	0.0	60.5	0.8	0.0	0.0
	ViM	100.0	100.0	100.0	98.0	24.5	100.0	100.0
	Maha	100.0	100.0	100.0	87.5	27.5	100.0	100.0
	Cos	0.0	0.0	0.0	27.5	0.0	0.0	0.0
Cnv21k	MSP	0.0	0.0	0.0	13.5	2.2	0.0	0.0
	ViM	100.0	100.0	100.0	0.0	0.0	41.2	0.5
	Maha	100.0	100.0	100.0	0.0	0.0	42.5	2.8
	Cos	0.0	0.0	0.0	0.0	0.0	0.0	0.0

results in combination with ViTs pretrained on IN-21k, the BiT-models and a pretrained EfficientNet-V2. Methods like Cos (7/26 models fail multiple tests) and MCM/RCos (7/26), originally designed for cosine-trained features as in CLIP, achieve remarkably strong OOD-detection performance on the unit-tests across a broad range of models, both with and without CLIP-pretraining. While taking note of these general trends, each OOD detector’s robustness to the OOD unit-tests should be examined individually.

5. Conclusions

We introduce with NINCO a novel, ID-contamination-free and challenging OOD test-dataset for IN-1K with fine-grained class-resolution. We find that many OOD detectors work better than previously thought, when their recorded number of undetected OOD inputs is not inflated by ID contaminations. However, most detection methods cannot reliably be applied with arbitrary classifier models, as even OOD unit-tests are failed by many combinations. We are hopeful for NINCO and the cleaned test OOD subsets to facilitate the more precise development of reliable OOD detectors which do not try to avoid presumed failures which are actually correct decisions.

Acknowledgements

We thank Vaclav Voracek for helpful discussions and suggesting the cdf plots. We acknowledge support from the German Federal Ministry of Education and Research (BMBF) through the Tübingen AI Center (FKZ: 01IS18039A) and from the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany’s Excellence Strategy (EXC number 2064/1, Project number 390727645), as well as from the Carl Zeiss Foundation in the project “Certification and Foundations of Safe Machine Learning Systems in Healthcare”. We also thank the European Laboratory for Learning and Intelligent Systems (ELLIS) for supporting Maximilian Müller.

References

Bendale, A. and Boulton, T. E. Towards open set deep networks. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. URL <https://doi.ieeecomputersociety.org/10.1109/CVPR.2016.173>.

- Bitterwolf, J., Meinke, A., and Hein, M. Certifiably adversarially robust detection of out-of-distribution data. In *NeurIPS*, 2020.
- Bitterwolf, J., Meinke, A., Augustin, M., and Hein, M. Breaking down out-of-distribution detection: Many methods based on ood training data estimate a combination of the same core quantities. In *ICML*, 2022.
- Bossard, L., Guillaumin, M., and Van Gool, L. Food-101 – mining discriminative components with random forests. In Fleet, D., Pajdla, T., Schiele, B., and Tuytelaars, T. (eds.), *Computer Vision – ECCV 2014*, pp. 446–461, Cham, 2014. Springer International Publishing. ISBN 978-3-319-10599-4.
- Cimpoi, M., Maji, S., Kokkinos, I., Mohamed, S., and Vedaldi, A. Describing textures in the wild. In *CVPR*, 2014.
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., and Hounsby, N. An image is worth 16x16 words: Transformers for image recognition at scale. *ICLR*, 2021.
- Fellbaum, C. Wordnet: An electronic lexical database. *Cambridge, MA: MIT Press*, 1998.
- Fort, S., Ren, J., and Lakshminarayanan, B. Exploring the limits of out-of-distribution detection. In *NeurIPS*, 2021. URL <https://openreview.net/forum?id=j5NrN8ffXC>.
- Galil, I., Dabbah, M., and El-Yaniv, R. A framework for benchmarking class-out-of-distribution detection and its application to imagenet. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=Iuubb9W6Jtk>.
- Hein, M., Andriushchenko, M., and Bitterwolf, J. Why ReLU networks yield high-confidence predictions far away from the training data and how to mitigate the problem. In *CVPR*, 2019.
- Hendrycks, D. and Gimpel, K. A baseline for detecting misclassified and out-of-distribution examples in neural networks. In *ICLR*, 2017. URL <https://openreview.net/forum?id=Hkg4TI9xl>.
- Hendrycks, D., Mazeika, M., and Dietterich, T. Deep anomaly detection with outlier exposure. In *ICLR*, 2019. <https://github.com/hendrycks/outlier-exposure>.
- Hendrycks, D., Zhao, K., Basart, S., Steinhardt, J., and Song, D. Natural adversarial examples. In *CVPR*, 2021.
- Hendrycks, D., Basart, S., Mazeika, M., Zou, A., Kwon, J., Mostajabi, M., Steinhardt, J., and Song, D. Scaling out-of-distribution detection for real-world settings. In *ICML*, 2022.
- Huang, R. and Li, Y. Mos: Towards scaling out-of-distribution detection for large semantic space. In *CVPR*, 2021.
- Ismail, A., Ahmad, S. A., Che Soh, A., Hassan, M. K., and Harith, H. H. Mynursinghome: A fully-labelled image dataset for indoor object classification. *Data in Brief*, 2020. doi: <https://doi.org/10.1016/j.dib.2020.106268>. URL <https://www.sciencedirect.com/science/article/pii/S2352340920311628>.
- Kolesnikov, A., Beyer, L., Zhai, X., Puigcerver, J., Yung, J., Gelly, S., and Hounsby, N. Big transfer (bit): General visual representation learning. In *ECCV*, 2020. URL https://doi.org/10.1007/978-3-030-58558-7_29.
- Koner, R., Sinhamahapatra, P., Roscher, K., Günnemann, S., and Tresp, V. Oodformer: Out-of-distribution detection transformer. *arXiv preprint arXiv:2107.08976*, 2021.
- Krasin, I., Duerig, T., Alldrin, N., Ferrari, V., Abu-El-Haija, S., Kuznetsova, A., Rom, H., Uijlings, J., Popov, S., Kamali, S., Mallocci, M., Pont-Tuset, J., Veit, A., Belongie, S., Gomes, V., Gupta, A., Sun, C., Chechik, G., Cai, D., Feng, Z., Narayanan, D., and Murphy, K. Openimages: A public dataset for large-scale multi-label and multi-class image classification. *Dataset available from storage.googleapis.com/openimages/web/index.html*, 2017.
- Krizhevsky, A. and Hinton, G. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.
- Lee, K., Lee, K., Lee, H., and Shin, J. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. In *NeurIPS*, 2018.
- Li, F.-F., Andreeto, M., Ranzato, M., and Perona, P. Caltech 101, 2022.
- Liu, W., Wang, X., Owens, J., and Li, Y. Energy-based out-of-distribution detection. *Advances in Neural Information Processing Systems*, 2020.
- Meinke, A., Bitterwolf, J., and Hein, M. Provably adversarially robust detection of out-of-distribution data (almost) for free. *NeurIPS*, 2022.

- Ming, Y., Cai, Z., Gu, J., Sun, Y., Li, W., and Li, Y. Delving into out-of-distribution detection with vision-language representations. In *NeurIPS*, 2022. URL <https://openreview.net/forum?id=KnCS9390Va>.
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., and Sutskever, I. Learning transferable visual models from natural language supervision. In *ICML*, 2021.
- Ren, J., Fort, S., Liu, J., Roy, A. G., Padhy, S., and Lakshminarayanan, B. A simple fix to mahalanobis distance for improving near-ood detection, 2021. URL <https://arxiv.org/abs/2106.09022>.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A. C., and Fei-Fei, L. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 2015.
- Salehi, M., Mirzaei, H., Hendrycks, D., Li, Y., Rohban, M. H., and Sabokrou, M. A unified survey on anomaly, novelty, open-set, and out-of-distribution detection: Solutions and future challenges. *arXiv preprint arXiv:2110.14051*, 2021.
- Shankar, V., Roelofs, R., Mania, H., Fang, A., Recht, B., and Schmidt, L. Evaluating machine accuracy on imagenet. In *NeurIPS 2021 Workshop on ImageNet: Past, Present, and Future*, 2021. URL <https://openreview.net/forum?id=Q3R088Eftng>.
- Song, Y., Sebe, N., and Wang, W. Rankfeat: Rank-1 feature removal for out-of-distribution detection. In *NeurIPS*, 2022. URL <https://openreview.net/forum?id=-deKNiSOXLG>.
- Steiner, A. P., Kolesnikov, A., Zhai, X., Wightman, R., Uszkoreit, J., and Beyer, L. How to train your vit? data, augmentation, and regularization in vision transformers. *TMLR*, 2022. URL <https://openreview.net/forum?id=4nPswr1KcP>.
- Sun, Y., Guo, C., and Li, Y. React: Out-of-distribution detection with rectified activations. *NeurIPS*, 2021.
- Sun, Y., Ming, Y., Zhu, X., and Li, Y. Out-of-distribution detection with deep nearest neighbors. *ICML*, 2022.
- Tack, J., Mo, S., Jeong, J., and Shin, J. Csi: Novelty detection via contrastive learning on distributionally shifted instances. In *NeurIPS*, 2020.
- Techapanurak, E., Sukanuma, M., and Okatani, T. Hyperparameter-free out-of-distribution detection using cosine similarity. In *Proceedings of the Asian Conference on Computer Vision*, 2020.
- Touvron, H., Cord, M., and Jegou, H. Deit iii: Revenge of the vit. *ECCV*, 2022.
- Van Horn, G., Mac Aodha, O., Song, Y., Cui, Y., Sun, C., Shepard, A., Adam, H., Perona, P., and Belongie, S. The inaturalist species classification and detection dataset. In *CVPR*, 2018.
- Vasudevan, V., Caine, B., Gontijo-Lopes, R., Fridovich-Keil, S., and Roelofs, R. When does dough become a bagel? analyzing the remaining mistakes on imagenet. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho, K. (eds.), *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=mowt1WNhTC7>.
- Vaze, S., Han, K., Vedaldi, A., and Zisserman, A. Open-set recognition: a good closed-set classifier is all you need? In *International Conference on Learning Representations*, 2022.
- Wang, H., Liu, W., Bocchieri, A., and Li, Y. Can multi-label classification networks know what they don't know? *NeurIPS*, 2021.
- Wang, H., Li, Z., Feng, L., and Zhang, W. Vim: Out-of-distribution with virtual-logit matching. In *CVPR*, 2022a.
- Wang, H., Zhang, A., Zhu, Y., Zheng, S., Li, M., Smola, A., and Wang, Z. Partial and asymmetric contrastive learning for out-of-distribution detection in long-tailed recognition. In *ICML 2022*, 2022b.
- Wen, X., Zhao, B., and Qi, X. A simple parametric classification baseline for generalized category discovery. *ArXiv*, abs/2211.11727, 2022.
- Wightman, R. Pytorch image models. <https://github.com/rwightman/pytorch-image-models>, 2019.
- Xia, G. and Bouganis, C.-S. On the usefulness of deep ensemble diversity for out-of-distribution detection. *arXiv preprint arXiv:2207.07517*, 2022.
- Xie, Q., Luong, M.-T., Hovy, E., and Le, Q. V. Self-training with noisy student improves imagenet classification. *arXiv preprint arXiv:1911.04252*, 2019.
- Yang, J., Wang, P., Zou, D., Zhou, Z., Ding, K., Peng, W., Wang, H., Chen, G., Li, B., Sun, Y., et al. Openood: Benchmarking generalized out-of-distribution detection. *arXiv preprint arXiv:2210.07242*, 2022.
- Zhou, B., Lapedriza, A., Khosla, A., Oliva, A., and Torralba, A. Places: A 10 million image database for scene recognition. *IEEE transactions on pattern analysis and machine intelligence*, 40(6):1452–1464, 2017.

A. Detailed results on NINCO

A detailed overview over the results on the NINCO benchmark is presented in Table 3, where we show the mean FPR for all models and methods across the dataset’s OOD classes. Tables 4-6 show AUROC, AUPR-S and AUPR-E with the same conclusions. The best method per model is marked bold, and the difference to the MSP-baseline is shown in green where a model outperforms the MSP-baseline and in red if it performs worse than MSP. It is clearly visible that there is no one-fits-all method. Instead, different models synergize with different methods. Overall, the two ViT models pretrained only on IN-21K in combination with Mahalanobis distance outperform other models and methods by a clear margin. This is in line with the observations of previous works (Koner et al., 2021; Fort et al., 2021; Galil et al., 2023), which also found the ViTs to perform exceptionally well. In terms of MSP, the ViTs are not better than e.g. the ConvNext, indicating that their improved OOD detection capabilities stem from a favourably structured feature-space. It is further interesting to see that for models without pretraining, out of all methods only Relative Mahalanobis and the cosine-based methods improve over the MSP-baseline fairly consistently. Apart from KL-Matching and KNN, most methods improve over the MSP-baseline for most pretrained models and the CLIP-methods Cosine and Rcos perform comparably well, yielding their best results with models pretrained *both* on CLIP and IN-12k. Since CLIP models are trained with cosine-similarity, it is likely that the structure of the feature space after finetuning remains favorable to cosine-based methods, while it might harm the performance of other feature-based methods like Mahanobis compared to models pretrained *only* on IN-21k.

It has been remarked (Hendrycks et al., 2022) that the advantage of models pretrained with IN-21K in the OOD detection task CIFAR-10 vs. CIFAR-100 (Krizhevsky & Hinton, 2009) might partially be explained by the CIFAR-100 classes not truly being unseen at train time, as they have a large overlap with IN-21K classes. We checked each NINCO class for overlap with the 21 843 classes of IN-21K with the help of a ViT classifier for IN-21K, see Table 9. This allows us to test whether the pretrained models have a larger advantage over purely IN-1K-trained models when trying to detect those classes with overlap compared to the classes without overlap. In Appendix K notice no substantially different changes between the models with and without pretraining. We remark, however, that even for several models without pretraining, the subselections of classes show quite different results.

In Figure 8 we contrast the results on NINCO with the results from previously used datasets. We show all methods for a pretrained ViT-B-384 and all models for the MSP-baseline. In both cases we observe several ranking changes: For the ViT, the best-performing method changes from ViM to Mahalanobis, and Relative Mahalanobis improves from sixth to second place. For the MSP-baseline, the clip-pretrained ViTs were the strongest OOD detectors on the previously used datasets, but are outperformed by the ConvNext-B on NINCO.

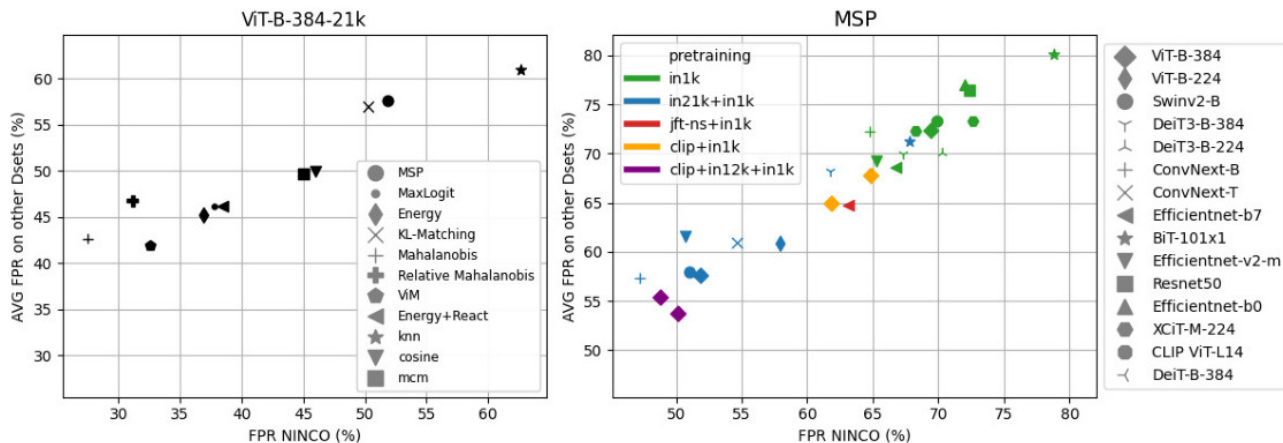


Figure 8. Mean FPR on NINCO vs. mean-FPR on previously used contaminated datasets with fixed model (left) and fixed method (right). We observe several ranking changes, including the best-performing method and model.

B. Models

In Table 8 we give an overview over the evaluated models. All model implementation and model weights were taken from the publicly available timm-repository (Wightman, 2019), except for the BiT-s weights, which can be obtained via the [github repository](#) of (Kolesnikov et al., 2020), and the zero-shot CLIP models, which are also available via [github](#). For the ViTs finetuned from CLIP and the ViT without pretraining we used the timm-version *0.8.0dev0*, for all other models version *0.6.12*. IN-12k ([description](#) and [defining synsets](#)) is a subset of IN-21k, for which the classes with few samples are excluded, leading to an overlap of roughly 85%.

Table 8. Overview over the evaluated models.

model	pretraining	top-1 acc.	params	timm name
ViT-B-384-12b-12k	laion2b + IN-12k	87.2	87M	vit_base_patch16_clip_384.laion2b_ft_in12k_in1k
ViT-B-384-oai-12k	openai + IN-12k	87.0	87M	vit_base_patch16_clip_384.openai_ft_in12k_in1k
ViT-B-384-12b	laion2b	86.6	87M	vit_base_patch16_clip_384.laion2b_ft_in1k
ViT-B-384-oai	openai	86.2	87M	vit_base_patch16_clip_384.openai_ft_in1k
ViT-B-384-21k	IN-21k	86.0	87M	vit_base_patch16_384
ViT-B-224-21k	IN-21k	84.5	87M	vit_base_patch16_224
Swinv2-B-256-21k	IN-21k	86.3	88M	swinv2_base_window12to16_192to256_22kft1k
DeiT3-B-384-21k	IN-21k	86.7	87M	deit3_base_patch16_384_in21ft1k
DeiT3-B-224-21k	IN-21k	85.7	87M	deit3_base_patch16_224_in21ft1k
CnvNxt-B-21k	IN-21k	86.3	89M	convnext_base_in22ft1k
CnvNxt-T-21k	IN-21k	84.1	29M	convnext_tiny_384_in22ft1k
BiT-m	IN-21k	82.3	45M	resnetv2_101x1_bitm
EffNetv2-M-21k	IN-21k	85.6	54M	tf_efficientnetv2_m_in21ft1k
EffNetb7-ns	JFT - noisy student	86.8	66M	tf_efficientnet_b7_ns
ViT-B-384	—	81.1	87M	vit_base_patch16_384.augreg_in1k
Swinv2-B-256	—	84.6	88M	swinv2_base_window16_256
DeiT3-B-384	—	85.1	87M	deit3_base_patch16_384
DeiT3-B-224	—	83.8	87M	deit3_base_patch16_224
XCiT-M-224	—	82.6	84M	xcit_medium_24_p16_224
XCiT-M-224-d	—	84.3	84M	xcit_medium_24_p16_224_dist
CnvNxt-B	—	84.4	89M	convnext_base
BiT-s	—	78.0	45M	resnetv2_101x1_bitm
EffNetv2-M	—	85.0	54M	tf_efficientnetv2_m
EffNetb7	—	84.9	66M	tf_efficientnet_b7
EffNet-B0	—	77.7	5M	efficientnet_b0
ResNet50	—	80.4	26M	resnet50
CLIP-ViT-B16	openai	66.6	150M	—
CLIP-ViT-B16	openai	74.2	428M	—

C. Methods

Here we give an overview over the evaluated OOD detection methods. For clarity, we denote vectors in bold and lowercase letters and matrices in bold uppercase letters. We write neural networks as functions n , which are parametrized by weights θ , take an input sample \mathbf{x} and produce an output vector \mathbf{o} of size C , where C is typically the number of classes in a classification task (1000 in the case of IN-1K). We refer to \mathbf{o} as the logits of \mathbf{x} , which can be transformed to a probability vector \mathbf{p} (also of size C) via the softmax function: $p_i = \exp(o_i) / \sum_c \exp(o_c)$. The network n can be decomposed into a feature extractor h and the networks last layer g :

$$\mathbf{o} = n(\mathbf{x}) = g(h(\mathbf{x})),$$

where g is a fully connected, linear layer, i.e. $g(\mathbf{h}) = \mathbf{W}^T \mathbf{h} + \mathbf{b}$ with weight \mathbf{W} and bias \mathbf{b} . We refer to $\mathbf{h} = h(\mathbf{x})$ as the *features* or the *embeddings* of \mathbf{x} w.r.t. the network n . As presented in Section 4, for each sample \mathbf{x} , a method returns an OOD-score $s = f(\mathbf{x})$, a scalar value which is supposed to be larger for ID data and smaller for OOD data. Methods accessing $h(\mathbf{x})$ directly in order to compute the OOD-score are referred to as feature-based methods, in contrast to methods that derive their OOD-score from the logits \mathbf{o} (even though obviously the logits implicitly also depend on these features). In the following, we will describe how each method computes the score s for a test input \mathbf{x} .

MSP (Hendrycks & Gimpel, 2017): The most popular OOD-detection baseline uses the confidence, i.e. the max softmax probability of a models probability output vector:

$$s = \max_c(p_c)$$

Max-Logit (Hendrycks et al., 2022): Similar to MSP, Max-Logit returns the largest entry of the logit-vector \mathbf{o} , i.e.

$$s = \max_c(o_c)$$

Energy (Liu et al., 2020): The Energy based OOD detection method uses the denominator of the softmax-function as OOD-score:

$$s = \log \sum_c \exp(o_c)$$

KL-Matching (Hendrycks et al., 2022): KL-Matching computes a mean probability vector \mathbf{d}_c for each of the C classes. For a test input, the KL-distances of all \mathbf{d}_c vectors to its probability vector \mathbf{p} are computed, and the OOD-score is the negative of the smallest of those distances:

$$s = -\min_c \text{KL}[\mathbf{p} \parallel \mathbf{d}_c]$$

In the original paper by (Hendrycks et al., 2022), the average for \mathbf{d}_c is computed over an additional validation set. Since none of the other methods leverages extra data and we are interested in fair comparison, we deploy KL-Matching like in (Wang et al., 2022a; Yang et al., 2022), where the average is computed over the train set.

KNN (Sun et al., 2022): KNN is a non-parametric method that computes distances in the feature-space. Specifically, the feature vector of a test input is normalized to $\mathbf{z} = \mathbf{h} / \|\mathbf{h}\|_2$ and the pairwise distances $r_i(\mathbf{z}) = \|\mathbf{z} - \mathbf{z}_i\|_2$ to the normalized features $\mathcal{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_N\}$ of all samples of the training set are computed. The distances $r_i(\mathbf{z})$ are then sorted according to their magnitude and the K^{th} smallest distance, denoted $r^K(\mathbf{z})$ is used as negative OOD-score:

$$s = -r^K(\mathbf{z})$$

Like suggested in (Sun et al., 2022), we use $K = 1000$.

Mahalanobis distance (Lee et al., 2018): This popular method fits a class-conditional Gaussian with shared covariance matrix to the train set, i.e. computes

$$\hat{\mu}_c = \frac{1}{N_c} \sum_{i:y_i=c} \mathbf{h}_i, \quad \hat{\Sigma} = \frac{1}{N} \sum_c \sum_{i:y_i=c} (\mathbf{h}_i - \hat{\mu}_c)(\mathbf{h}_i - \hat{\mu}_c)^T$$

where N_c is the number of train samples in class c and N is the total number of train samples. The OOD-score of a test sample is then the Mahalanobis distance induced by $\hat{\Sigma}$ between its feature \mathbf{h} and the closest class mean:

$$s = -\min_c (\mathbf{h} - \hat{\mu}_c) \hat{\Sigma}^{-1} (\mathbf{h} - \hat{\mu}_c)^T$$

Relative Mahalanobis distance (Ren et al., 2021): A modification of the Mahalanobis distance method, thought to improve near-OOD detection, is to additionally fit a global Gaussian distribution to the train set without taking class-information into account:

$$\hat{\mu}_{\text{global}} = \frac{1}{N} \sum_i \mathbf{h}_i, \quad \hat{\Sigma}_{\text{global}} = \frac{1}{N} \sum_i (\mathbf{h}_i - \hat{\mu}_{\text{global}})(\mathbf{h}_i - \hat{\mu}_{\text{global}})^T$$

The OOD-score is then defined as the difference between the original Mahalanobis distance and the Mahalanobis distance w.r.t. the global Gaussian distribution:

$$s = -\min_c \left((\mathbf{h} - \hat{\mu}_c) \hat{\Sigma}^{-1} (\mathbf{h} - \hat{\mu}_c)^T - (\mathbf{h} - \hat{\mu}_{\text{global}}) \hat{\Sigma}_{\text{global}}^{-1} (\mathbf{h} - \hat{\mu}_{\text{global}})^T \right)$$

ReAct (Sun et al., 2021): The authors propose to perform a truncation of the feature vector, $\bar{\mathbf{h}} = \min(\mathbf{h}, r)$, where the min operation is to be understood element-wise and r is the truncation threshold. The truncated features can then be converted to so-called rectified logits via $\bar{\mathbf{o}} = g(\bar{\mathbf{h}}) = \mathbf{W}^T \bar{\mathbf{h}} + \mathbf{b}$. While the rectified logits can now be used with a variety of existing detection methods, we follow (Sun et al., 2021) and use the rectified Energy as OOD-score:

$$s = \log \sum_c \exp(\bar{o}_c)$$

As suggested in (Wang et al., 2022a), we set the threshold r such that 1% of the activations from the train set would be truncated.

Virtual Logit Matching (Wang et al., 2022a): The idea behind ViM is that meaningful features are thought to lie in a low-dimensional manifold, called the principal space P , whereas features from OOD-samples should also lie in P^\perp , the space orthogonal to P . P is the D -dimensional subspace spanned by the eigenvectors with the largest D eigenvalues of the matrix $\mathbf{F}^T \mathbf{F}$, where \mathbf{F} is the matrix of all train features offsetted by $\mathbf{u} = -(\mathbf{W}^T)^+ \mathbf{b}$ (+ denotes the Moore-Penrose inverse). A sample with feature vector \mathbf{h} is then also offset to $\tilde{\mathbf{h}} = \mathbf{h} - \mathbf{u}$ and can be decomposed into $\tilde{\mathbf{h}} = \tilde{\mathbf{h}}^P + \tilde{\mathbf{h}}^{P^\perp}$, and $\tilde{\mathbf{h}}^{P^\perp}$ is referred to as the *Residual* of \mathbf{h} . ViM leverages the Residual and converts it to a virtual logit $o_0 = \alpha \|\tilde{\mathbf{h}}^{P^\perp}\|_2$, where

$$\alpha = \frac{\sum_{i=1}^N \max_c o_i^c}{\sum_{i=1}^N \|\mathbf{h}_i^{P^\perp}\|_2}$$

is designed to match the scale of the virtual logit to the scale of the real train logits. The virtual logit is then appended to the original logits of the test sample, i.e. to \mathbf{o} , and a new probability vector is computed via the softmax function. The probability corresponding to the virtual logit is then the final OOD-score:

$$s = -\frac{\exp(o_0)}{\sum_{c=1}^C \exp(o_c) + \exp(o_0)}$$

Like suggested in (Wang et al., 2022a), we use $D = 1000$ if the dimensionality of the feature space d is $d \geq 2048$, $D = 512$ if $2048 \geq d \geq 768$, and $D = d/2$ rounded to integers otherwise.

Cosine (Tack et al., 2020; Galil et al., 2023): This method computes the maximum cosine-similarity between the features of a test-sample and embedding vectors $\tilde{\mathbf{u}}_c$ (sometimes also called concept-vector):

$$s = \max_c \tilde{\mathbf{u}}_c^T \mathbf{h} / \|\tilde{\mathbf{u}}_c\|_2 \quad (1)$$

For zero-shot CLIP, $\tilde{\mathbf{u}}_c$ can be obtained by creating text-embeddings from the ImageNet class names. Encoding 'A photo of a ...' yields an embedding from the corresponding class. For classifiers, we use the class-wise train means $\hat{\mu}_c$, that are also used for Mahalanobis distance.

MCM/RCos (Ming et al., 2022; Techapanurak et al., 2020): Maximum-Concept-Matching was recently introduced as a zero-shot OOD detection method for CLIP and applies additional softmax-scaling to the cosine-similarities of the *Cosine* method, potentially with a temperature scaling (which we omit, following (Ming et al., 2022)). Again, we extend this method to work with conventional classifiers by using the class-means $\hat{\mu}_c$ like they are used for Mahalanobis distance as embedding/concept vectors. We then refer to it as relative cosine (short: MCM/RCos or just RCoS) in order to distinguish it from CLIPs zero-shot method.

D. Definitions of OOD detection metrics

The performance of OOD detectors is commonly reported in terms of the *false positive rate at a fixed true positive rate* Q , denoted as **FPR@TPRQ**, short **FPR**. This means that the detector is interpreted as making the decision to *accept* an unknown input x if $S(x) \geq \tau$, for a threshold τ that is chosen such that $Q\%$ of ID inputs are accepted, and *rejecting* the input as OOD if $S(x) < \tau$. The FPR@TPRQ counts the fraction of falsely accepted OOD inputs under this decision scheme. This means the *lower* the FPR@TPRQ, the *better* the OOD detection performance. In the OOD detection literature, the most commonly used value for Q is 95%, which we too use throughout this paper. We also report results in terms of the mean *area under the receiver-operator characteristic curve*, short **AUROC** in Table 4. It represents the probability that an ID input receives a higher score (equal scores counted half) than an OOD input when both are drawn randomly from their respective evaluation datasets (Bitterwolf et al., 2022). Like for the FPR, the mean AUROC corresponds to first uniformly drawing an OOD class and then drawing a sample from that class.

E. Illustrative examples from the cleaning process

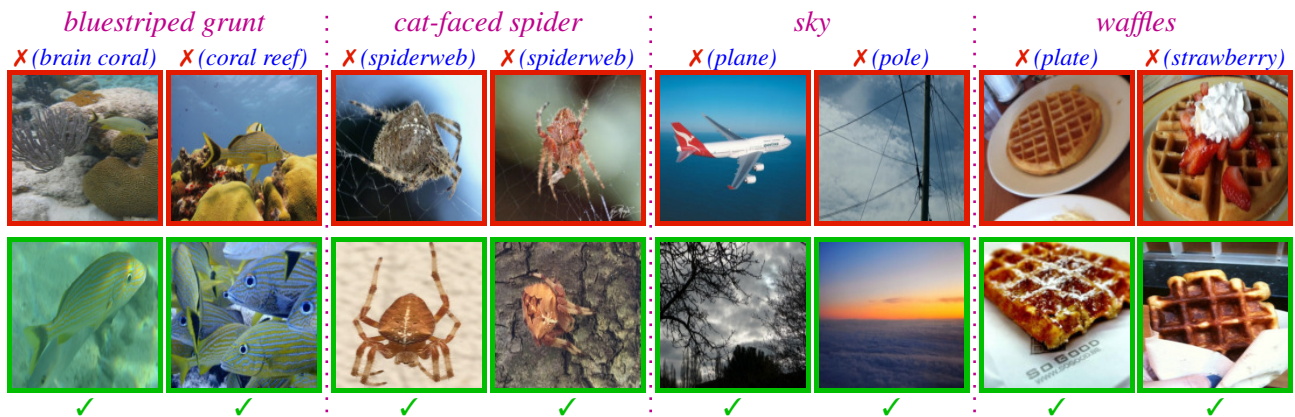


Figure 9. **Cleaning the OOD classes.** **Top:** Samples that were excluded due to overlap with ID classes. **Bottom:** Samples from the same OOD class that were included in the cleaned datasets.

F. Details of the NINCO dataset.

Table 9. Detailed information for each OOD class. For determining overlap with classes of IN-21K, we checked the 8 most common predictions of a ViT classifier for IN-21K on the NINCO OOD class.

OOD class name	shortname	# samples	source dataset	ImageNet-21K overlap
<i>AFA (cyanobacterium)</i>	<i>AFA</i>	46	SPECIES	<i>microorganism</i>
<i>bagpipe</i>	<i>Bagp</i>	97	Imagenet-21k	<i>bagpipe</i>
<i>bluestriped grunt</i>	<i>BSGrunt</i>	96	SPECIES	<i>grunt</i>
<i>cable</i>	<i>Cabl</i>	88	scraped	<i>scraped television</i>
<i>California pitcher plant</i>	<i>CPitch</i>	100	SPECIES	<i>pitcher plant</i>
<i>California slender salamander</i>	<i>CSSala</i>	100	SPECIES	<i>slender salamander</i>
<i>California two-spot octopus</i>	<i>C2SOct</i>	100	SPECIES	<i>octopus</i>
<i>caracal</i>	<i>Caracal</i>	100	iNat. Download	<i>caracal</i>
<i>cat-faced spider</i>	<i>CatFSp</i>	100	SPECIES	<i>unclear/very broad class</i>
<i>Central American tapir</i>	<i>CATapir</i>	100	SPECIES	<i>tapir</i>
<i>chicken quesadilla</i>	<i>CQuesa</i>	100	FOOD-101	-
<i>common cuttlefish</i>	<i>CCuttle</i>	100	SPECIES	<i>cuttlefish</i>
<i>crème brûlée</i>	<i>CBrûlée</i>	99	FOOD-101	<i>creme brulee</i>
<i>cupcakes</i>	<i>CCake</i>	80	FOOD-101	-
<i>donuts</i>	<i>Donuts</i>	100	FOOD-101	<i>doughnut</i>
<i>door</i>	<i>Door</i>	100	MyNursingHome	<i>interior door</i>
<i>dreamfish</i>	<i>Dreamf</i>	100	SPECIES	<i>sea bream</i>
<i>dune thistle</i>	<i>DThist</i>	100	SPECIES	<i>creme brulee</i>
<i>dusky flathead (fish)</i>	<i>DFlath</i>	100	SPECIES	<i>flathead</i>
<i>E. micromeris (cactus)</i>	<i>EMicro</i>	100	SPECIES	-
<i>Eastern leaf-footed bug</i>	<i>ELFBug</i>	100	SPECIES	<i>leaf-footed bug</i>
<i>European paper wasp</i>	<i>EPWasp</i>	100	SPECIES	<i>paper wasp</i>
<i>false killer whale</i>	<i>FalseKW</i>	67	SPECIES	<i>unclear/very broad class</i>
<i>field road</i>	<i>FieldRd</i>	96	PLACES	<i>byway</i>
<i>fire extinguisher</i>	<i>FireEx</i>	106	MyNursingHome	<i>fire extinguisher</i>
<i>fireworks</i>	<i>FireW</i>	100	scraped	-
<i>forest path</i>	<i>ForPth</i>	100	PLACES	<i>unclear/very broad class</i>
<i>Franciscan wallflower</i>	<i>Franci</i>	100	SPECIES	<i>wallflower</i>
<i>French fries</i>	<i>Fries</i>	100	FOOD-101	<i>french fries</i>
<i>Galápagos fur seal</i>	<i>GFurS</i>	91	SPECIES	<i>arcella</i>
<i>giant cuttlefish</i>	<i>GCuttle</i>	99	SPECIES	<i>cuttlefish</i>
<i>glass of milk</i>	<i>GMilk</i>	89	scraped	<i>milk</i>
<i>gramophone</i>	<i>Gramo</i>	56	scraped	<i>gramophone</i>
<i>high heels</i>	<i>HHeels</i>	99	scraped	-
<i>Hindu temple</i>	<i>HinTp</i>	51	scraped	<i>unclear/very broad class</i>
<i>Horse Hoof clam</i>	<i>HHClam</i>	31	SPECIES	<i>seashell</i>
<i>Indo-Pacific bottlenose dolphin</i>	<i>IPBNDol</i>	100	SPECIES	<i>dolphin</i>
<i>long-tailed silverfish</i>	<i>LTSilF</i>	100	SPECIES	<i>silverfish</i>
<i>Lumholtz's tree-kangaroo</i>	<i>LTRoo</i>	100	SPECIES	<i>tree wallaby</i>
<i>M. wesenbergii (cyanobacterium)</i>	<i>MWesen</i>	33	SPECIES	<i>microorganism</i>
<i>marbled newt</i>	<i>MNewt</i>	100	SPECIES	<i>newt</i>
<i>mbira</i>	<i>Mbira</i>	67	scraped	-
<i>Mexican lime cactus</i>	<i>MLCact</i>	100	SPECIES	<i>barrel cactus</i>
<i>Pampas deer</i>	<i>PDeer</i>	82	SPECIES	<i>buck</i>
<i>pyramid</i>	<i>Pyra</i>	100	caltech-101	<i>Cheops</i>
<i>redbreast sunfish</i>	<i>RBSunf</i>	100	SPECIES	<i>sunfish</i>
<i>rosybell (flowering plant)</i>	<i>Rosyb</i>	100	SPECIES	-
<i>ruby octopus</i>	<i>RubyOct</i>	100	SPECIES	<i>octopus</i>
<i>scissors</i>	<i>Sciss</i>	100	caltech-101	<i>scissors</i>
<i>shuttlecock</i>	<i>ShCo</i>	67	scraped	<i>shuttlecock</i>
<i>silver-haired bat</i>	<i>SilverHB</i>	99	SPECIES	<i>bat</i>
<i>skipper caterpillar</i>	<i>SCaterp</i>	100	iNat. Download	<i>caterpillar</i>
<i>sky</i>	<i>Sky</i>	68	PLACES	<i>sky</i>
<i>southern calamari</i>	<i>SCalam</i>	99	SPECIES	<i>squid</i>
<i>spaghetti bolognese</i>	<i>SBolo</i>	67	FOOD-101	<i>spaghetti</i>
<i>stapler</i>	<i>Stapl</i>	100	caltech-101	<i>stapler</i>
<i>sweet pea</i>	<i>SwPea</i>	100	SPECIES	<i>unclear/very broad class</i>
<i>two-toed amphiuma (salamander)</i>	<i>2TAmph</i>	176	SPECIES	<i>amphiuma</i>
<i>waffles</i>	<i>Waffle</i>	61	FOOD-101	-
<i>walker</i>	<i>Walker</i>	99	MyNursingHome	<i>walker</i>
<i>water dispenser (jugless)</i>	<i>WDisp</i>	100	MyNursingHome	<i>water cooler</i>
<i>Windsor chair</i>	<i>WiChair</i>	71	caltech-101	<i>Windsor chair</i>
<i>yellow trumpets</i>	<i>YTrump</i>	100	SPECIES	<i>yellow trumpet</i>
<i>'ōhelo 'ai (flowering plant)</i>	<i>'ō'ai</i>	100	SPECIES	-

G. Details and recipes for OOD unit-tests

We provide 400 samples for each of 17 OOD unit-tests, mirroring the sizes and file formats of random ImageNet samples. Their reproducible definitions are given as follows:

- **uniform noise (Hendrycks & Gimpel, 2017):** Each RGB colour channel of each pixel is independently sampled uniformly between 0.0 or 1.0.
- **Gaussian noise (Hendrycks & Gimpel, 2017):** For each image, first σ is chosen randomly between (0.05, 0.075, 0.1, 0.15, 0.2, 0.3, 0.5). Then each RGB colour channel of each pixel is independently sampled from $\mathcal{N}(0.5, \sigma)$.
- **Rademacher noise (Hendrycks et al., 2019):** Then each RGB colour channel of each pixel is independently set to 0.0 or 1.0 with 50% probability.
- **IN pixel permutations (Hein et al., 2019):** We choose a random IN-1K validation image and randomly shuffle its pixels (no remixing of colours).
- **black:** All colour channels are set to 0.0.
- **white:** All colour channels are set to 1.0.
- **shades of grey:** All colour channels are set to the same value, sampled uniformly between 0.0 or 1.0.
- **monochrome:** All pixels are set to a uniformly random RGB-colour (sampled uniformly from $[0.0, 1.0]^3$).
- **tricolour:** The image is split into three stripes of equal size, vertically or horizontally with probability 50%. Each stripe is set to an independent uniformly random RGB-colour.
- **primary tricolour:** The image is split into three stripes of equal size, vertically or horizontally with probability 50%. Each stripe is set to a colour where each RGB-channel value is chosen randomly as either 0.0 or 1.0.
- **horizontal stripes:** The image is split into a random number chosen between (4, 5, 7, 10, 15, 20) of horizontal stripes of equal size. Each stripe is set to an independent uniformly random RGB-colour.
- **vertical stripes:** The image is split into a random number chosen between (4, 5, 7, 10, 15, 20) of vertical stripes of equal size. Each stripe is set to an independent uniformly random RGB-colour.
- **smooth noise (Hein et al., 2019; Bitterwolf et al., 2020; Meinke et al., 2022):** For each image, first σ is chosen randomly between (10, 15, 25, 40, 60, 85). A uniform noise image is sampled. Then we apply a Gaussian filter with a size of σ pixels. Finally, the pixel values are scaled linearly such that the minimum brightness over all channels and pixels is 0.0 and the maximum is 1.0.
- **smooth noise+:** For each image, first σ is chosen randomly between (10, 15, 25, 40, 60, 85). A uniform noise image is sampled. Then we apply a Gaussian filter with a size of σ pixels. Finally, each RGB channel is scaled linearly such that its minimum brightness over all pixels is 0.0 and the maximum is 1.0.
- **smooth color:** For each image, first σ is chosen randomly between (10, 15, 25, 40, 60, 85), δ uniformly between 0.1 and 0.3, and a uniformly random RGB-colour c . A uniform noise image is sampled. Then we apply a Gaussian filter with a size of σ pixels. Finally, each RGB channel is scaled linearly such that $c - \delta$ is the 2.5th quantile of its values and $c + \delta$ the 97.5th.

- **smooth IN pixel permutations (Hein et al., 2019):** For each image, first σ is chosen randomly between (1, 1.5, 2, 3, 4, 6, 8).
An IN pixel permutations image is sampled.
Then we apply a Gaussian filter with a size of σ pixels.
- **blobs (Hendrycks et al., 2019):** For each image, first σ is chosen randomly between (1.5, 2, 2.5, 3, 3.5, 4).
Each RGB colour channel of each pixel is independently set to 1.0 with 70% probability or 0.0 with 30%.
Then we apply a Gaussian filter with a size of σ pixels.
Finally, all channel values below 0.75 are set to 0.0.

Where necessary, the resulting channel values are clipped to $[0, 1]$. We show samples of each unit-test in the following Appendix H in Figure 13.

H. Examples images from each OOD class in NINCO and from OOD unit-tests

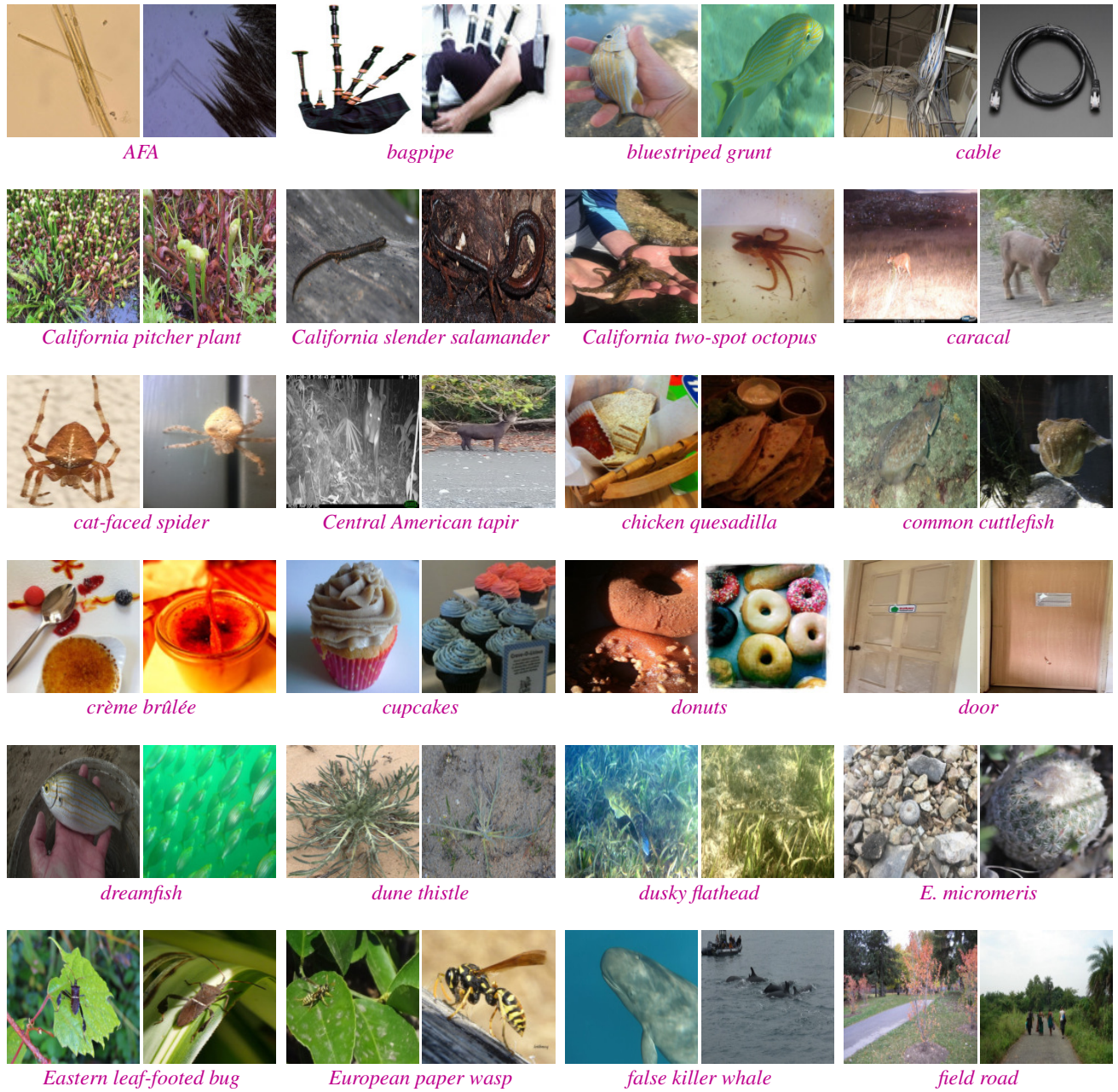


Figure 10. Samples of each class of the NINCO dataset (1/3).

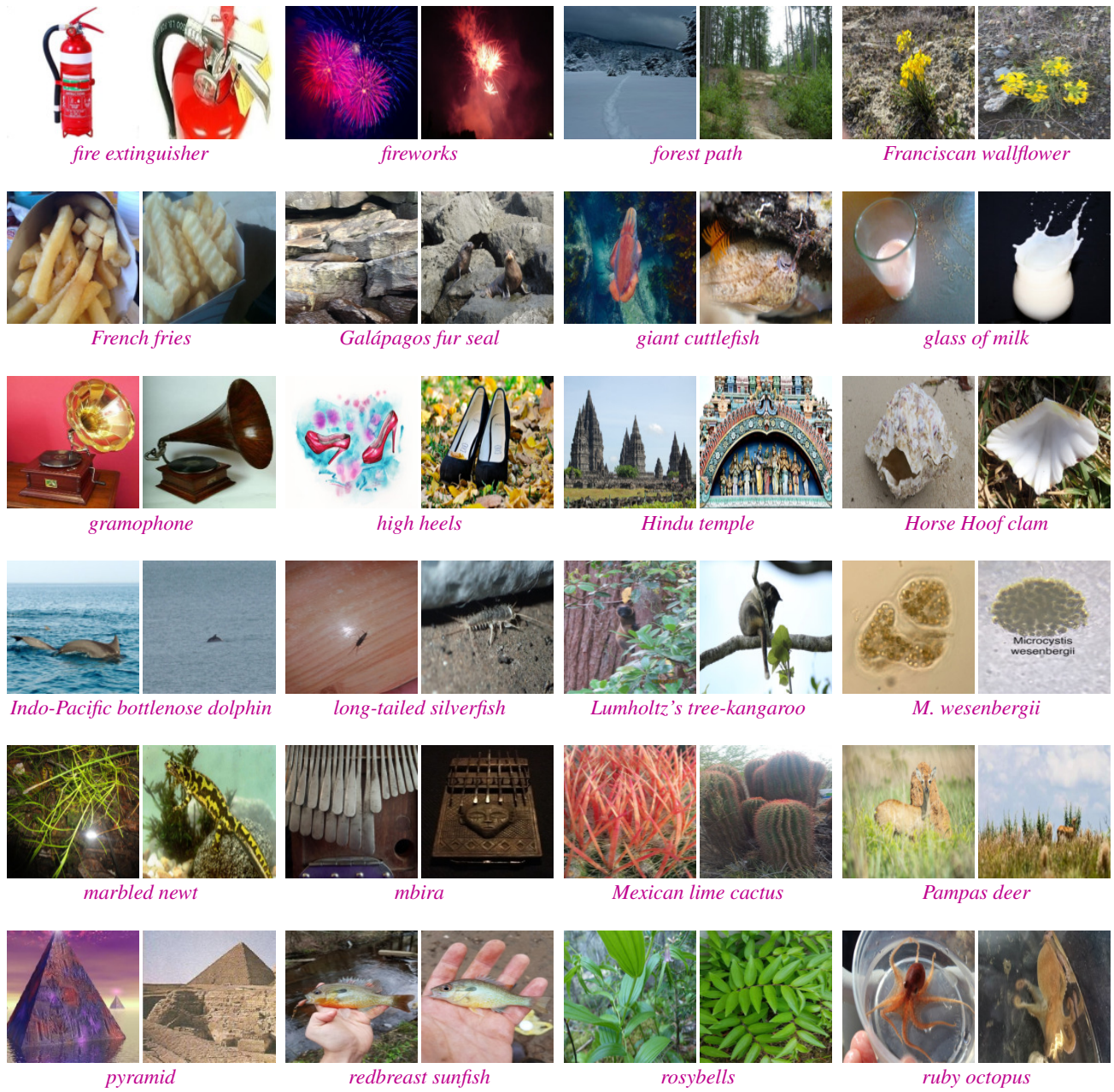


Figure 11. Samples of each class of the NINCO dataset (2/3).

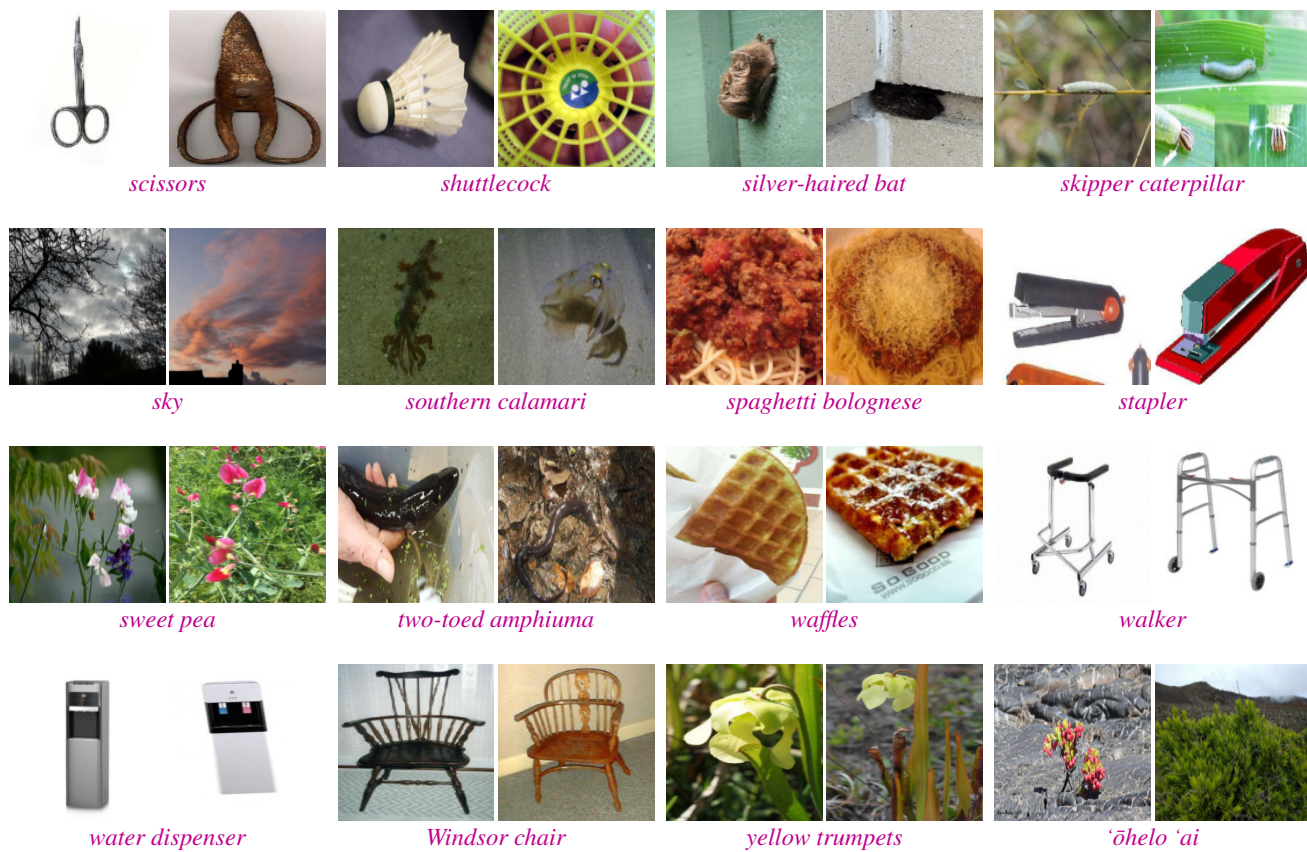


Figure 12. Samples of each class of the NINCO dataset (3/3).

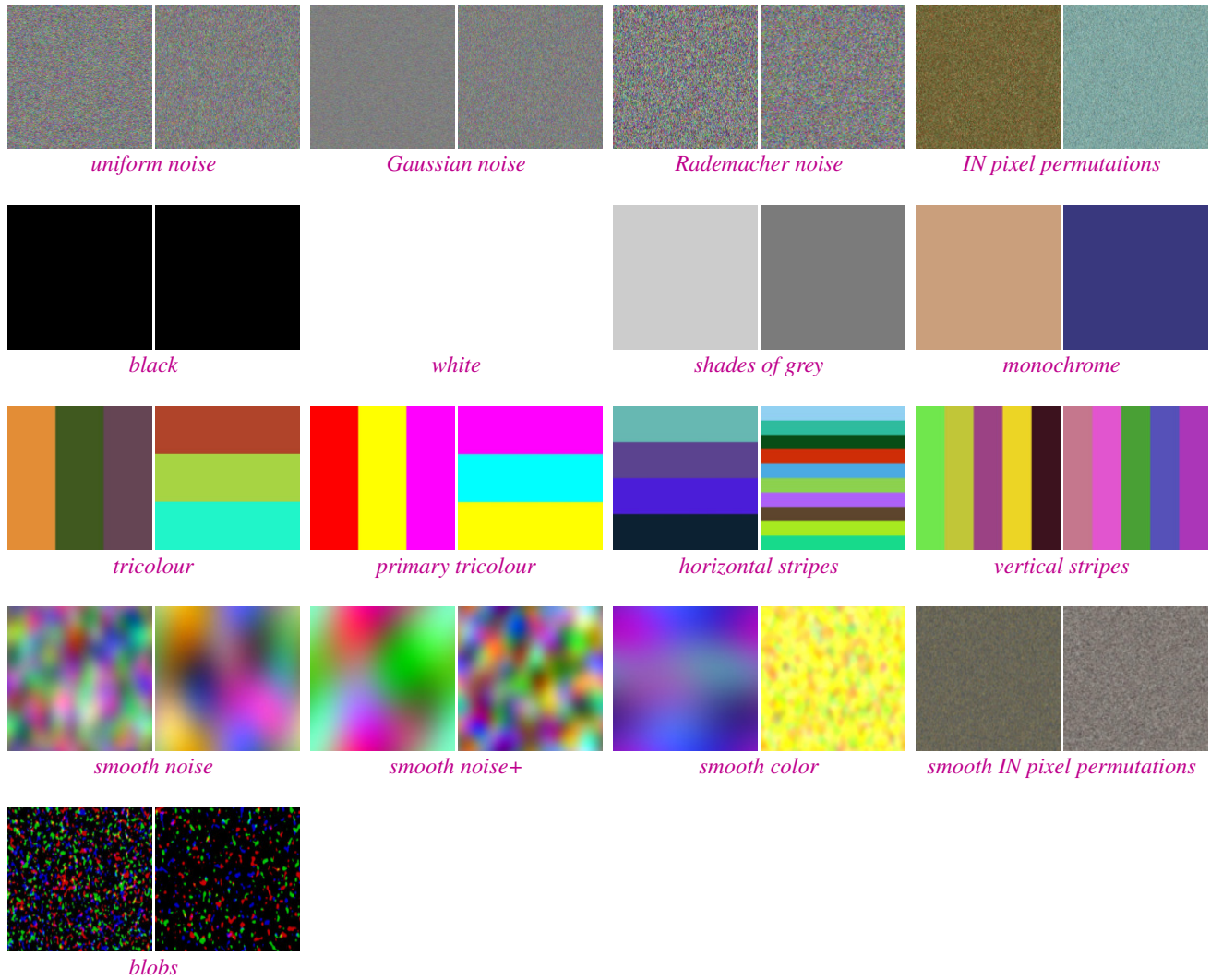


Figure 13. Samples of each OOD unit-test.

J. Effect of ID contamination on all models

In Table 14 we show the FPR values averaged across the cleaned subsampled datasets on which Table 1 in the main paper is based.

Detailed results on the individual datasets are shown in Tables 15-18. There, we show results on the uncleaned full (-f) and cleaned subsampled (-c) datasets: PLACES (Pl), SPECIES (SpC), IMAGENET-O (IN), TEXTURES (txt) & TEXTURES43, OPENIMAGE-O (OpO), INATURALIST OOD PLANTS (iNat), IMAGENET-1K-OOD (IN1K), 360OPENSETCLASSES (OS), SEMANTIC SHIFT BENCHMARK EASY (SBe) & HARD (SBh), and COOD (CO).

Since TEXTURES and INATURALIST are fairly easy test OOD datasets, the FPR values of most models in Table 14 are lower than on NINCO. In general, the results allow similar conclusions: Feature-based methods outperform methods not explicitly accessing pre-logit feature-information, yet still fail for some models, and pretraining only on IN-21k yields the best OOD-detectors. Again, Cosine and MCM/RCos improve fairly consistently over MSP, and are in some cases even the best-performing method.

Table 14. Mean FPR on subsampled datasets (averaged).

pre	acc.	model	MSP	MaxL	Ener	KL-M	Maha	RMaha	ViM	E+R	KNN	Cos	MCM/RCos	
21k	86.0	ViT-B-384	39.7	27.0 -13	25.7 -14	38.4 -1	22.4 -17	25.5 -14	22.4 -17	27.5 -12	48.2 +8	30.6 -9	30.4 -9	
	84.5	ViT-B-224	43.3	30.8 -13	29.3 -14	42.7 -1	23.8 -19	28.2 -15	24.7 -19	32.6 -11	53.3 +10	37.0 -6	36.1 -7	
	86.3	Swinv2-B-256	41.9	32.3 -10	31.5 -10	46.4 +4	47.4 +5	40.4 -2	37.5 -4	27.8 -14	43.1 +1	35.5 -6	34.2 -8	
	86.7	Deit3-B-384	53.4	45.4 -8	46.4 -7	52.5 -1	40.8 -13	37.8 -16	41.2 -12	39.9 -13	40.1 -13	36.3 -17	36.0 -17	
	85.7	Deit3-B-224	55.1	46.9 -8	47.2 -8	56.1 +1	46.6 -9	42.6 -12	47.5 -8	42.0 -13	45.1 -10	41.4 -14	40.4 -15	
	86.3	CnvNxt-B	38.6	32.9 -6	35.3 -3	43.6 +5	36.3 -2	30.5 -8	29.9 -9	31.1 -8	37.0 -2	30.0 -9	29.5 -9	
	84.1	CnvNxt-T	44.1	37.6 -7	35.7 -8	50.7 +7	36.2 -8	37.0 -7	27.7 -16	34.0 -10	44.1 -0	40.2 -4	38.9 -5	
	82.3	BiT-m	59.9	52.0 -8	52.6 -7	55.3 -5	30.9 -29	32.7 -27	26.9 -33	46.3 -14	37.2 -23	32.9 -27	38.2 -22	
	85.6	EffNetv2-M	43.4	42.5 -1	49.7 +6	46.3 +3	43.7 +0	41.1 -2	37.0 -6	89.0 +46	50.2 +7	32.4 -11	38.5 -5	
none	81.1	ViT-B-384	63.5	59.4 -4	58.8 -5	59.6 -4	49.1 -14	48.2 -15	61.4 -2	55.4 -8	64.0 +0	59.1 -4	60.9 -3	
	84.6	Swinv2-B-256	63.5	63.0 -1	68.6 +5	60.9 -3	49.4 -14	46.0 -17	52.0 -11	60.5 -3	57.0 -6	52.1 -11	50.4 -13	
	85.1	Deit3-B-384	60.0	64.8 +5	83.2 +23	57.8 -2	51.2 -9	48.5 -11	44.9 -15	89.2 +29	65.6 +6	57.2 -3	43.8 -16	
	83.8	Deit3-B-224	60.4	62.2 +2	76.1 +16	58.9 -1	57.6 -3	52.8 -8	48.9 -11	80.4 +20	73.7 +13	64.4 +4	49.5 -11	
	82.6	XCiT-M-224	65.8	65.2 -1	71.4 +6	65.4 -0	58.3 -7	55.7 -10	55.4 -10	66.9 +1	63.1 -3	57.3 -8	56.4 -9	
	84.3	XCiT-M-224-d	63.9	61.6 -2	69.9 +6	61.0 -3	55.4 -8	52.8 -11	50.4 -13	66.4 +3	59.5 -4	53.6 -10	52.3 -12	
	84.4	CnvNxt-B	63.1	72.3 +9	92.1 +29	62.8 -0	55.5 -8	52.1 -11	53.7 -9	88.7 +26	60.8 -2	53.6 -9	50.6 -12	
	78.0	BiT-s	75.3	77.7 +2	79.8 +5	59.8 -15	68.9 -6	51.2 -24	60.1 -15	65.8 -10	71.2 -4	56.0 -19	84.0 +9	
	85.1	EffNetv2-M	59.0	59.4 +0	70.5 +12	56.8 -2	48.2 -11	42.9 -16	57.4 -2	59.9 +1	54.7 -4	50.2 -9	43.5 -16	
	84.9	EffNetb7	60.1	63.4 +3	75.7 +16	56.0 -4	57.9 -2	47.6 -13	63.4 +3	66.6 +6	58.4 -2	52.7 -7	44.4 -16	
JFT	77.7	EffNet-B0	69.3	69.9 +1	77.3 +8	68.2 -1	75.9 +7	68.6 -1	65.7 -4	67.8 -1	77.0 +8	51.7 -18	63.8 -6	
	80.4	ResNet50	68.3	70.0 +2	76.6 +8	64.5 -4	81.0 +13	75.9 +8	73.0 +5	97.6 +29	65.9 -2	51.7 -17	55.0 -13	
	86.8	EffNetb7-ns	53.8	49.9 -4	62.5 +9	52.7 -1	79.5 +26	53.2 -1	82.4 +29	57.0 +3	55.0 +1	47.0 -7	46.5 -7	
	clip	87.2	ViT-B-384-12b	37.3	33.7 -4	35.6 -2	40.5 +3	43.6 +6	36.9 -0	36.7 -1	31.6 -6	35.0 -2	29.5 -8	29.3 -8
		+12k	87.0	ViT-B-384-oai	38.7	33.1 -6	32.9 -6	40.7 +2	45.9 +7	37.4 -1	38.4 -0	31.2 -8	33.7 -5	29.2 -9
	clip	86.6	ViT-B-384-12b	54.2	52.5 -2	57.2 +3	51.0 -3	40.2 -14	40.4 -14	38.4 -16	54.0 -0	44.0 -10	40.0 -14	39.5 -15
		86.2	ViT-B-384-oai	56.7	55.0 -2	59.0 +2	53.9 -3	40.6 -16	40.8 -16	41.4 -15	56.0 -1	45.6 -11	41.3 -15	40.3 -16
	clip	74.3	clip-ViT-L-336	---	---	---	---	---	---	---	---	---	64.4	51.8
	z. shot	66.6	clip-ViT-B-224	---	---	---	---	---	---	---	---	---	71.4	60.0

K. Results on NINCO classes with and without overlap with IN-21K

Since the classes of NINCO can be distinguished by whether they belong to an IN-21k class or not, we present results on both of these groups here. We note that they should be taken with care, since the groups differ both in size (9 vs. 55 classes) and difficulty of the individual classes. Most models and methods perform better on the classes *with* IN-21k overlap, and ViT+Maha is the best OOD-detector in both cases. While RMaha and (Relative) Cosine yield the most consistent improvements over MSP in both cases, ViM performs comparably better on the classes without overlap. Pretraining *only* on IN-21k yields the best OOD-detectors in both cases.

Table 19. Mean FPR for classes without 21k overlap.

pre	acc.	model	MSP	MaxL	Ener	KL-M	Maha	RMaha	ViM	E+R	KNN	Cos	MCM/RCos
21k	86.0	ViT-B-384	56.5	41.8 -15	39.6 -17	51.7 -5	31.7 -25	36.9 -20	32.2 -24	40.9 -16	67.3 $+11$	46.7 -10	42.2 -14
	84.5	ViT-B-224	64.8	50.6 -14	48.3 -17	60.2 -5	34.1 -31	43.7 -21	34.8 -30	50.2 -15	68.5 $+4$	54.8 -10	53.5 -11
	86.3	Swinv2-B-256	66.3	58.7 -8	59.1 -7	62.0 -4	40.1 -26	42.7 -24	34.3 -32	50.3 -16	54.8 -11	47.5 -19	47.3 -19
	86.7	DeiT3-B-384	72.9	71.1 -2	73.3 $+0$	68.6 -4	43.0 -30	43.6 -29	44.1 -29	64.4 -9	49.3 -24	47.2 -26	46.8 -26
	85.7	DeiT3-B-224	75.1	72.8 -2	72.6 -3	69.5 -6	47.7 -27	48.7 -26	47.1 -28	67.5 -8	56.3 -19	52.9 -22	53.5 -22
	86.3	CnvNxt-B	61.4	60.0 -1	67.0 $+6$	57.6 -4	31.0 -30	37.4 -24	27.5 -34	61.6 $+0$	47.0 -14	40.6 -21	39.7 -22
	84.1	CnvNxt-T	62.9	57.2 -6	54.4 -9	61.6 -1	34.7 -28	42.2 -21	30.6 -32	52.9 -10	53.3 -10	49.1 -14	46.2 -17
	82.3	BiT-m	69.7	62.2 -7	63.9 -6	62.6 -7	40.9 -29	42.1 -28	31.5 -38	60.2 -10	39.1 -31	36.0 -34	42.1 -28
	85.6	EffNetv2-M	55.9	51.8 -4	56.3 $+0$	55.7 -0	48.6 -7	46.9 -9	40.9 -15	96.5 $+41$	55.3 -1	33.8 -22	42.4 -14
none	81.1	ViT-B-384	70.0	64.5 -5	61.1 -9	65.0 -5	56.6 -13	56.2 -14	62.8 -7	59.7 -10	66.3 -4	63.0 -7	63.5 -6
	84.6	Swinv2-B-256	72.4	67.7 -5	68.2 -4	68.2 -4	58.9 -13	56.9 -15	57.6 -15	65.8 -7	67.8 -5	62.2 -10	60.5 -12
	85.1	DeiT3-B-384	70.4	75.1 $+5$	85.4 $+15$	64.4 -6	59.3 -11	57.4 -13	51.5 -19	91.2 $+21$	70.7 $+0$	65.1 -5	49.2 -21
	83.8	DeiT3-B-224	76.4	77.1 $+1$	83.3 $+7$	69.5 -7	62.3 -14	60.0 -16	57.9 -18	83.9 $+8$	75.7 -1	69.4 -7	55.8 -21
	82.6	XCiT-M-224	79.5	79.1 -0	82.4 $+3$	76.1 -3	71.6 -8	69.7 -10	69.2 -10	78.5 -1	76.6 -3	73.3 -6	73.0 -7
	84.3	XCiT-M-224-d	72.6	71.7 -1	78.8 $+6$	66.6 -6	63.4 -9	60.8 -12	60.0 -13	75.3 $+3$	69.6 -3	62.7 -10	60.9 -12
	84.4	CnvNxt-B	74.1	82.3 $+8$	94.5 $+20$	63.9 -10	59.3 -15	56.8 -17	56.2 -18	90.8 $+17$	65.7 -8	59.2 -15	58.0 -16
	78.0	BiT-s	74.2	74.5 $+0$	76.5 $+2$	58.2 -16	83.2 $+9$	56.8 -17	64.4 -10	71.3 -3	81.3 $+7$	66.8 -7	77.2 $+3$
	85.1	EffNetv2-M	70.0	69.5 -1	74.4 $+4$	65.3 -5	52.1 -18	51.4 -19	59.6 -10	61.7 -8	60.3 -10	56.6 -13	53.0 -17
	84.9	EffNetb7	69.0	70.5 $+2$	81.3 $+12$	62.5 -7	55.5 -14	50.4 -19	59.2 -10	71.0 $+2$	61.7 -7	58.0 -11	50.4 -19
	77.7	EffNet-B0	75.0	75.9 $+1$	84.0 $+9$	68.7 -6	71.0 -4	66.8 -8	62.2 -13	75.0 $+0$	85.8 $+11$	58.7 -16	62.8 -12
80.4	ResNet50	76.0	76.6 $+1$	77.5 $+1$	69.0 -7	77.0 $+1$	66.4 -10	75.1 -1	94.8 $+19$	64.0 -12	57.6 -18	56.6 -19	
JFT	86.8	EffNetb7-ns	71.3	64.8 -7	67.5 -4	66.5 -5	83.7 $+12$	72.0 $+1$	85.2 $+14$	65.8 -6	70.3 -1	64.2 -7	63.8 -7
clip +12k	87.2	ViT-B-384-l2b	53.7	51.1 -3	55.9 $+2$	52.7 -1	37.8 -16	40.2 -14	31.7 -22	47.3 -6	41.1 -13	37.3 -16	37.0 -17
	87.0	ViT-B-384-oai	56.0	51.8 -4	54.6 -1	53.6 -2	40.9 -15	39.8 -16	36.9 -19	50.4 -6	36.6 -19	33.8 -22	34.1 -22
clip	86.6	ViT-B-384-l2b	65.8	63.5 -2	62.5 -3	59.0 -7	49.6 -16	50.4 -15	46.1 -20	61.0 -5	53.8 -12	49.5 -16	48.3 -18
	86.2	ViT-B-384-oai	65.8	64.1 -2	67.7 $+2$	62.4 -3	52.4 -13	54.7 -11	48.1 -18	65.4 -0	57.1 -9	53.9 -12	53.4 -12
clip	74.3	clip-ViT-L-336	—	—	—	—	—	—	—	—	—	55.7	55.8
z. shot	66.6	clip-ViT-B-224	—	—	—	—	—	—	—	—	—	56.9	62.8

Table 20. Mean FPR for classes with 21k overlap.

pre	acc.	model	MSP	MaxL	Ener	KL-M	Maha	RMaha	ViM	E+R	KNN	Cos	MCM/RCos
21k	86.0	ViT-B-384	51.1	37.2 -14	36.5 -15	50.1 -1	26.8 -24	30.2 -21	32.7 -18	38.1 -13	61.9 +11	45.9 -5	45.5 -6
	84.5	ViT-B-224	56.8	45.8 -11	45.7 -11	56.7 -0	31.6 -25	35.7 -21	39.0 -18	49.3 -8	68.9 +12	54.6 -2	54.4 -2
	86.3	Swinv2-B-256	48.6	38.2 -10	36.9 -12	55.0 +6	66.5 +18	55.7 +7	58.2 +10	35.3 -13	63.1 +15	52.0 +3	48.3 -0
	86.7	Deit3-B-384	60.0	53.5 -6	53.6 -6	59.0 -1	55.7 -4	49.6 -10	59.0 -1	49.5 -10	54.1 -6	48.6 -11	47.8 -12
	85.7	Deit3-B-224	63.1	57.0 -6	55.8 -7	64.5 +1	62.0 -1	54.7 -8	65.0 +2	53.1 -10	59.1 -4	54.4 -9	53.1 -10
	86.3	CnvNxt-B	44.9	38.0 -7	39.4 -5	54.4 +10	52.6 +8	43.2 -2	43.8 -1	37.0 -8	52.6 +8	44.8 -0	43.0 -2
	84.1	CnvNxt-T	53.3	46.4 -7	44.0 -9	60.6 +7	48.9 -4	46.4 -7	38.5 -15	42.7 -11	57.1 +4	51.5 -2	49.7 -4
	82.3	BiT-m	67.5	62.0 -6	63.0 -4	65.3 -2	51.5 -16	45.6 -22	42.2 -25	56.6 -11	61.1 -6	54.2 -13	56.5 -11
	85.6	EffNetv2-M	49.9	47.8 -2	53.8 +4	54.4 +4	65.3 +15	52.4 +2	55.6 +6	88.7 +39	69.5 +20	47.3 -3	51.9 +2
none	81.1	ViT-B-384	69.4	68.2 -1	69.3 -0	67.0 -2	60.6 -9	57.2 -12	70.4 +1	66.8 -3	74.8 +5	69.6 +0	70.8 +1
	84.6	Swinv2-B-256	69.5	67.6 -2	72.9 +3	67.4 -2	64.7 -5	60.6 -9	67.9 -2	69.3 -0	69.5 -0	63.7 -6	62.3 -7
	85.1	Deit3-B-384	66.8	72.4 +6	87.9 +21	64.6 -2	64.8 -2	59.7 -7	61.3 -5	90.0 +23	75.0 +8	67.5 +1	58.1 -9
	83.8	Deit3-B-224	69.3	71.1 +2	82.1 +13	68.2 -1	70.1 +1	64.9 -4	64.4 -5	83.0 +14	81.2 +12	73.6 +4	62.9 -6
	82.6	XCiT-M-224	71.5	72.3 +1	78.6 +7	71.1 -0	65.4 -6	62.5 -9	64.2 -7	76.0 +4	71.1 -0	66.1 -5	64.9 -7
	84.3	XCiT-M-224-d	67.6	65.2 -2	72.2 +5	66.9 -1	66.9 -1	62.1 -6	62.7 -5	72.0 +4	70.6 +3	64.9 -3	62.9 -5
	84.4	CnvNxt-B	63.2	69.7 +7	88.2 +25	68.7 +6	66.8 +4	61.2 -2	67.0 +4	85.1 +22	71.2 +8	61.7 -2	58.7 -5
	78.0	BiT-s	79.6	82.3 +3	83.9 +4	70.0 -10	83.6 +4	65.3 -14	75.0 -5	78.9 -1	83.5 +4	73.0 -7	85.3 +6
	85.1	EffNetv2-M	64.5	64.6 +0	74.6 +10	62.4 -2	64.2 -0	55.5 -9	74.7 +10	70.9 +6	65.1 +1	60.0 -4	54.6 -10
	84.9	EffNetb7	66.4	68.7 +2	81.6 +15	62.7 -4	70.2 +4	55.3 -11	74.9 +8	77.2 +11	67.7 +1	61.0 -5	54.3 -12
	77.7	EffNetB0	71.6	71.9 +0	78.9 +7	72.8 +1	85.3 +14	75.2 +4	77.3 +6	75.1 +4	87.1 +16	61.8 -10	70.9 -1
80.4	ResNet50	71.8	73.9 +2	78.0 +6	69.0 -3	87.3 +16	70.0 -2	79.2 +7	97.9 +26	80.2 +8	63.8 -8	63.0 -9	
JFT	86.8	EffNetb7-ns	61.8	54.2 -8	60.5 -1	64.1 +2	88.0 +26	68.2 +6	89.8 +28	61.1 -1	74.3 +13	65.4 +4	63.7 +2
clip	87.2	ViT-B-384-12b	49.6	46.8 -3	49.4 -0	52.1 +3	55.0 +5	48.5 -1	48.1 -2	44.5 -5	46.2 -3	40.6 -9	40.7 -9
+12k	87.0	ViT-B-384-oai	47.7	42.3 -5	42.3 -5	48.9 +1	60.4 +13	49.8 +2	55.1 +7	40.9 -7	46.3 -1	40.2 -7	39.9 -8
clip	86.6	ViT-B-384-12b	61.2	61.3 +0	66.4 +5	57.3 -4	53.2 -8	50.5 -11	52.6 -9	63.6 +2	57.5 -4	51.0 -10	49.2 -12
clip	86.2	ViT-B-384-oai	64.7	65.1 +0	70.1 +5	61.7 -3	56.3 -8	53.6 -11	58.3 -6	67.7 +3	62.0 -3	57.0 -8	54.4 -10
clip	74.3	clip-ViT-L-336	---	---	---	---	---	---	---	---	---	75.2	68.9
z. shot	66.6	clip-ViT-B-224	---	---	---	---	---	---	---	---	---	82.8	82.6