

A Survey on Out-of-Distribution Detection in NLP

Anonymous EMNLP submission

Abstract

Out-of-distribution (OOD) detection is essential for the reliable and safe deployment of machine learning systems in the real world. Great progress has been made over the past years. This paper presents the first review of recent advances in OOD detection with a particular focus on natural language processing approaches. First, we provide a formal definition of OOD detection and discuss several related fields. We then categorize recent algorithms into three classes according to the data they used: (1) OOD data available, (2) OOD data unavailable + in-distribution (ID) label available, and (3) OOD data unavailable + ID label unavailable. Third, we introduce datasets, applications, and metrics. Finally, we summarize existing work and present potential future research topics.

1 Introduction

Natural language processing systems deployed in the wild often encounter out-of-distribution (OOD) samples that are not seen in the training phase. A reliable and trustworthy NLP model should not only obtain high performance on samples from seen distributions, i.e., In-distribution (ID) samples, but also accurately detect OOD samples (Amodei et al., 2016; Boulton et al., 2019). For instance, when building task-oriented dialogue systems, it is hard, if not impossible, to cover all possible user intents in the training stage. It is critical for a practical system to detect these OOD intents or classes in the testing phase so that they can be properly handled (Zhan et al., 2021).

However, existing flourishes of neural-based NLP models are built upon the *closed-world assumption*, i.e., the training and testing data are sampled from the same distribution (Vapnik, 1991). This assumption is often violated in practice, where deployed models are generally confronting an *open-world*, i.e., some testing data may come from OOD distributions that are not seen in training (Bendale and Boulton, 2015; Fei and Liu, 2016).

A rich line of work has been proposed to tackle problems introduced by OOD samples. Specifically, distributional shifts in NLP can be broadly divided into two types: 1. semantic shift, i.e., OOD samples may come from unknown categories, and therefore should not be blindly predicted into a known category; 2. non-semantic shift, i.e., OOD samples may come from different domains or styles but share the same semantic with some ID samples (Arora et al., 2021). The detection of OOD samples with semantic shift is the primary focus of this survey, where the label set \mathcal{Y} of ID samples is different from that of OOD samples. Although there already exists surveys on many aspects of OOD, such as OOD generalization (Wang et al., 2022) and OOD detection in computer vision (CV) (Yang et al., 2021), *a comprehensive survey for OOD detection in NLP is still lacking and thus urgently needed for the field*. Concretely, applying OOD detection to NLP tasks requires specific considerations, e.g., tackling discrete input spaces, handling complex output structures, and considering contextual information, which have not been thoroughly discussed. Our key contributions are summarized as follows:

1. We propose a novel taxonomy of OOD detection methods based on the availability of OOD data (Section 3) and discuss their pros and cons for different settings (Section 6.1).

2. We present a survey on OOD detection in NLP and identify various differences between OOD detection in NLP and CV (Section 6.2).

3. We review datasets, applications (Section 4), metrics (Section 5), and future research directions (Section 6.3) of OOD detection in NLP.

2 OOD Detection and Related Areas

Definition 1 (Data distribution). *Let \mathcal{X} denote a nonempty input (non-semantic) space and \mathcal{Y} a label (semantic) space. A data distribution is defined as a joint distribution $P(X, Y)$ over $\mathcal{X} \times \mathcal{Y}$. $P(X)$*

and $P(Y)$ refer to the marginal distributions for inputs and labels, respectively.

In practice, common non-semantic distribution shifts on $P(X)$ include domain shifts (Wang et al., 2022), sub-population shifts (Koh et al., 2021), style changes (Pavlick and Tetreault, 2016), or adversarial examples (Carlini and Wagner, 2017; Rozsa et al., 2017). Typically, the label space \mathcal{Y} remains unchanged in these non-semantic shifts, and sophisticated methods are developed to improve the model’s robustness and generalization performance (Hendrycks et al., 2020). On the contrary, semantic distribution shifts on $P(Y)$ generally lead to a new label space $\tilde{\mathcal{Y}}$ that are different from the one seen in the training phase (Bendale and Boult, 2016). These shifts are usually caused by the occurrence of new classes at the testing stage. In this work, we mainly focus on detecting OOD samples with semantic shifts, the formal definition of which is given as follows:

Definition 2 (OOD detection). *We are given an ID training set $\mathcal{D}_{train} = \{(\mathbf{x}_i, y_i)\}_{i=1}^L \sim P_{train}(X, Y)$, where $\mathbf{x}_i \in \mathcal{X}_{train}$ is a training instance, and $y_i \in \mathcal{Y}_{train} = \{1, 2, \dots, K\}$ is the associated class label. Facing the emergence of unknown classes, we are given a test set $\mathcal{D}_{test} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N \sim P_{test}(X, Y)$, where $\mathbf{x}_i \in \mathcal{X}_{test}$, and $y_i \in \mathcal{Y}_{test} = \{1, \dots, K, K + 1\}$. Note that class $K + 1$ is a group of novel categories representative of OOD samples, which may contain more than one class. The overall goal of OOD detection is to learn a predictive function f from \mathcal{D}_{train} to achieve a minimum expected risk on \mathcal{D}_{test} : $\min_f \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_{test}} \mathbb{I}(y \neq f(\mathbf{x}))$, i.e., not only classify known classes but also detect the unknown categories.*

We briefly describe the related research areas:

Domain generalization (DG) (Wang et al., 2022), or out-of-distribution generalization, aims to learn a model from one or several source domains and expect these learned models generalize well on unseen testing domains (i.e., target domains). DG mainly focuses on the non-semantic drift, i.e., the training and testing tasks share the same label space \mathcal{Y} while they have different distributions over the input space \mathcal{X} . Different from DG, OOD detection handles a different label space during testing.

Domain adaptation (DA) (Blitzer et al., 2006) follows most settings of DG except that DA has access to some unlabeled data from the target domain in the training process (Ramponi and Plank, 2020). Similar to DG, DA also assumes the label space

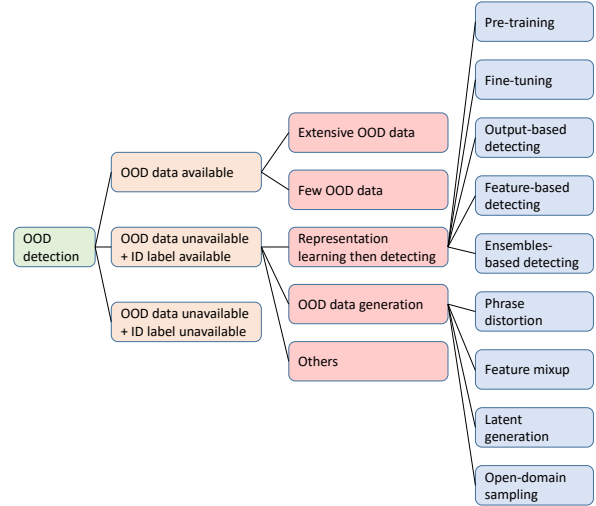


Figure 1: Taxonomy of OOD detection methods.

remains unchanged.

Zero-shot learning (Wang et al., 2019) aims to use learned models to classify samples from unseen classes. However, OOD detection in general aims to classify samples from seen classes while flagging the unseen class in testing.

Meta-learning (Vilalta and Drissi, 2002) aims to learn from the model training process so that models can quickly adapt to new data. Different from meta-learning, achieving strong few-shot performance is not the major focus of OOD detection. Nevertheless, the idea of meta-learning can serve as a strategy for OOD detection (Xu et al., 2019; Li et al., 2021) by simulating the behaviors of predicting unseen classes in the training stage.

Positive-unlabeled Learning (Zhang and Zuo, 2008), or PU learning, aims to train a classifier with only positive and unlabeled examples while being able to distinguish both positive and negative samples in testing. However, OOD detection considers multiple classes in training. PU learning approaches can be applied to tackle the OOD detection problem when only one labeled class exists (Li and Liu, 2003).

3 Methodology

A major challenge of OOD detection is the lack of representative OOD data, which is important for estimating OOD distributions (Zhou et al., 2021b). As shown in Figure 1, we classify existing OOD detection methods into three categories according to the availability of OOD data. Methods covered in our survey are selected following the criteria listed in Appendix A.

3.1 OOD Data Available

Methods in this category assume access to both labeled ID and OOD data during training. Based on the quantity and diversity of OOD data, we further classify these methods into two subcategories:

3.1.1 Detection with Extensive OOD Data

Some methods assume that we can access extensive OOD data in the training process together with ID data. In this subcategory, one line of work formulates OOD detection as a discriminative classification task, i.e., a special label is allocated in the label space for OOD samples. Fei and Liu (2016); Larson et al. (2019) formed a $(K + 1)$ -way classification problem, where K denoted the number of ID classes and the $(K + 1)^{th}$ class represented OOD samples. Larson et al. (2019); Kamath et al. (2020) regarded OOD detection as a binary classification problem, where the two classes correspond to ID and OOD samples, respectively. Kim and Kim (2018) introduced a neural joint learning model with a multi-class classifier for domain classification and a binary classifier for OOD detection.

Another line of work optimizes an outlier exposure regularization term on these OOD samples to refine the representations and OOD scores learned by the OOD detector. Hendrycks et al. (2018) introduced a generalized outlier exposure (OE) loss to train models on both ID and OOD data. For example, when using the maximum Softmax probability detector (Hendrycks and Gimpel, 2016), the OE loss pushes the predicted distribution of OOD samples to a uniform distribution (Lee et al., 2018a). When the labels of ID data are not available, the OE loss degenerates to a margin ranking loss on the predicted distributions of ID and OOD samples. Zeng et al. (2021b) added an entropy regularization objective to enforce the predicted distributions of OOD samples to have high entropy.

3.1.2 Detection with Few OOD Data

Some methods assume that we can only access a small amount of OOD data besides ID data. This setting is more realistic in practice since it is expensive to annotate large-scale OOD data. Several methods in this subcategory are developed to generate pseudo samples based on a small number of seed OOD data. Chen and Yu (2021) constructed pseudo-labeled OOD candidates using samples from an auxiliary dataset and kept only the most beneficial candidates for training through a novel election-based filtering mechanism. Rather than

directly creating OOD samples in natural language, Zeng et al. (2021b) borrowed the idea of adversarial attack (Goodfellow et al., 2014) to obtain model-agnostic worst-case perturbations in the latent space, where these perturbations or noise can be regarded as augmentations for OOD samples.

3.2 OOD Data Unavailable + ID Label Available

Building OOD detectors using only labeled ID data is the major focus of research communities. We generally classify existing literature into three subcategories based on their learning principles:

3.2.1 Learn Representations Then Detect

Some methods formulize the OOD detector f into two components: a representation extractor g and an OOD scoring function d , i.e., $f(\mathbf{x}) = d(g(\mathbf{x}))$: g aims to capture a representation space \mathcal{H} in which ID and OOD samples are distinct, and d maps each extracted representation into an OOD score so that OOD samples can be detected based on a selected threshold. We provide an overview of methods to enhance these two components:

a. Representation Learning usually involves two stages: (1) a *pre-training* stage leverages massive unlabeled text corpora to extract representations that are suitable for general NLP tasks; (2) a *fine-tuning* stage uses labeled in-domain data to refine representations for specified downstream tasks. An overview of these two stages is given here:

Pre-training Pre-trained transformer models such as BERT (Kenton and Toutanova, 2019) have become the de facto standard to implement text representation extractors. Hendrycks et al. (2020) systematically measured the OOD detection performance on various representation extractors, including bag-of-words models, ConvNets (Gu et al., 2018), LSTMs (Hochreiter and Schmidhuber, 1997), and pre-trained transformer models (Vaswani et al., 2017). Their results show that pre-trained models achieve the best OOD detection performance, while the performances of all other models are often worse than chance. The success of pre-trained models attributes to these diverse corpora and effective self-supervised training losses used in training (Hendrycks et al., 2019).

Moreover, it is observed that better-calibrated models generally produce higher OOD detection performance (Lee et al., 2018a). Desai and Durrett (2020) evaluated the calibration of two pre-trained models, BERT and RoBERTa (Liu et al., 2019), on

different tasks. They found that pre-trained models were better calibrated in out-of-domain settings, where non-pre-trained models like ESIM (Chen et al., 2017) were overconfident. Dan and Roth (2021) also demonstrated that larger pre-trained models are more likely to be better calibrated and thus result in higher OOD detection performance.

Fine-tuning With the help of labeled ID data, various approaches are developed to fine-tune the representation extractor to widen margins between ID and OOD samples. Lin and Xu (2019) proposed a large margin cosine loss (LMCL) to maximize the decision margin in the latent space. LMCL simultaneously maximizes inter-class variances and minimizes intra-class variances. Yan et al. (2020) introduced a semantic-enhanced Gaussian mixture model to enforce ball-like dense clusters in the feature space, which injects semantic information of class labels into the Gaussian mixture distribution.

Zeng et al. (2021a); Zhou et al. (2021b) proposed a contrastive learning framework (Chen et al., 2020) to increase the discrepancy for representations extracted from different classes. They hypothesized that increasing inter-class discrepancies helps the model learn discriminative features for ID and OOD samples and therefore improves OOD detection performances. Concretely, a supervised contrastive loss (Khosla et al., 2020; Gunel et al., 2020) and a margin-based contrastive loss was investigated. Zeng et al. (2021b) proposed a self-supervised contrastive learning framework to extract discriminative representations of OOD and ID samples from unlabeled data. In this framework, positive pairs are constructed using the back-translation scheme. Zhou et al. (2022) applied KNN-based contrastive learning losses to OOD detectors and Wu et al. (2022) used a reassigned contrastive learning scheme to alleviate the over-confidence issue in OOD detection.

Moreover, there are some regularized fine-tuning schemes to tackle the over-confidence issue of neural-based OOD detectors. Kong et al. (2020) addressed this issue by introducing an off-manifold regularization term to encourage producing uniform distributions for pseudo off-manifold samples. Shen et al. (2021) designed a novel domain-regularized module that is probabilistically motivated and empirically led to a better generalization in both ID classification and OOD detection.

b. OOD Scoring processes usually involve a scoring function d to map the representations of

input samples to OOD detection scores. A higher OOD score indicates that the input sample is more likely to be OOD. The implementation of d can be generally categorized into three types: (1) *output-based detecting*, (2) *feature-based detecting*, and (3) *ensembles-based detecting*:

Output-based Detecting compute the OOD score based on the predicted probabilities. Hendrycks and Gimpel (2016) used the maximum Softmax probability as the detection score, and Liang et al. (2018) improved this scheme with the temperature scaling approach. Shu et al. (2017) employed K 1-vs-rest Sigmoid classifiers for K predefined ID classes and used the maximum probabilities from these classifiers as the detection score. Liu et al. (2020) proposed an energy score for better distinguishing ID/OOD samples. The energy score is theoretically aligned with the probability density of the inputs.

Feature-based Detecting leverages features derived from intermediate layers of the model to implement density-based and distance-based scoring functions. Gu et al. (2019) proposed a nearest-neighbor based method with a distance-to-measure metric. Breunig et al. (2000) used a local outlier factor as the detection score, in which the concept “local” measured how isolated an object was with respect to surrounding neighborhoods. Lee et al. (2018b); Podolskiy et al. (2021) obtained the class-conditioned Gaussian distributions with respect to features of the deep models under Gaussian discriminant analysis. This scheme resulted in a confidence score based on the Mahalanobis distance. While Mahalanobis imposes a strong distributional assumption on the feature space, Sun et al. (2022) demonstrated the efficacy of non-parametric nearest neighbor distance for OOD detection. Zhang et al. (2021) proposed a post-processing method to learn an adaptive decision boundary (ADB) for each ID class. Specifically, the ADB is learned by balancing both the empirical and open space risks (Scheirer et al., 2014). Recently, Ren et al. (2022) proposed to detect OOD samples for conditional language generation tasks (such as abstractive summarization and translation) by calculating the distance between testing input/output and a corresponding background model in the feature space.

Ensembles-based Detecting uses predictive uncertainty of a collection of supporting models to compute OOD scores. Specifically, an input sample is regarded as an OOD sample if the variance of

these models' predictions is high. Gal and Ghahramani (2016) modeled uncertainties by applying dropouts to neural-based models. This scheme approximates Bayesian inference in deep Gaussian processes. Lakshminarayanan et al. (2017) used deep ensembles for uncertainty quantification, where multiple models with the same architecture were trained in parallel with different initialization. Lukovnikov et al. (2021) further proposed a heterogeneous ensemble of models with different architectures to detect compositional OOD samples for semantic parsing.

3.2.2 Generate Pseudo OOD Samples

A scheme to tackle the problem of lacking OOD training samples is to generate pseudo OOD samples during training (Lang et al., 2022). With these generated pseudo OOD samples, OOD detectors can be solved by methods designed for using both labeled ID and OOD data. There are mainly four types of approaches to generate pseudo OOD samples: (1) *phrase distortion*, (2) *feature mixup*, (3) *latent generation*, and (4) *open-domain sampling*:

Phrase Distortion approaches generate pseudo OOD samples for NLP tasks by selectively replacing text phrases in ID samples. Ouyang et al. (2021) proposed a data manipulation framework to generate pseudo OOD utterances with importance weights. Choi et al. (2021) proposed Out-Flip, which revised a white-box adversarial attack method HotFlip to generate OOD samples. Shu et al. (2021) created OOD instances from ID examples with the help of a pre-trained language model.

Feature Mixup strategy (Zhang et al., 2018) is also a popular technique for pseudo data generation. Zhan et al. (2021) generated OOD samples by performing linear interpolations between ID samples from different classes in the representation space. Zhou et al. (2021a) leveraged the manifold Mixup scheme (Verma et al., 2019) for pseudo OOD sample generation. Intermediate layer representations of two samples from different classes are mixed using scalar weights sampled from the Beta distribution. These feature-mixup-based methods achieved promising performance while remaining conceptually and computationally straightforward.

Latent Generation approaches considered to use generative adversarial networks (GAN) (Goodfellow et al., 2020) to produce high-quality pseudo OOD samples. Lee et al. (2018a) proposed to generate boundary samples in the low-density area of the ID distribution as pseudo-OOD samples. Ryu

et al. (2018) built a GAN on ID data and used the discriminator to generate OOD samples in the continuous feature space. Zheng et al. (2020) generated pseudo OOD samples using an auto-encoder with adversarial training in the discrete text space. Marek et al. (2021) proposed OodGAN, in which a sequential generative adversarial network (SeqGAN) (Yu et al., 2017) was used for OOD sample generation. This model follows the idea of Zheng et al. (2020) but works directly on texts and hence eliminates the need to include an auto-encoder.

Open-domain Sampling approaches directly uses sentences from other corpora as pseudo OOD samples (Zhan et al., 2021).

3.2.3 Other Methods

We also review some representative methods that do not belong to the above two categories. Vyas et al. (2018) proposed to use an ensemble of classifiers to detect OOD, where each classifier was trained in a self-supervised manner by leaving out a random subset of training data as OOD data. Li et al. (2021) proposed k Folden, which included k classifiers for k class labels. Each classifier was trained on a subset with $k - 1$ classes while leaving one class unknown. Tan et al. (2019) tackled the problem of OOD detection with limited labeled ID training data and proposed an OOD-resistant Prototypical Network to build the OOD detector. Ren et al. (2019); Gangal et al. (2020) used the likelihood ratio produced by generative models to detect OOD samples. The likelihood ratio effectively corrects confounding background statistics for OOD detection. Ryu et al. (2017) employed the reconstruction error as the detection score.

3.3 OOD data unavailable + ID label unavailable

OOD detection using only unlabeled ID data can be used for non-classification tasks. In fact, when ID labels are unavailable, our problem setting falls back to the classic anomaly detection problem, which is developed with a rich set of literature (Pang et al., 2021; Chalapathy and Chawla, 2019). However, this problem setting is rarely investigated in NLP studies. We keep this category here for the completeness of our survey while leaning most of our focus on NLP-related works.

Methods in this category mainly focus on extracting more robust features and making a more accurate estimation for the data distribution. Zong et al. (2018) proposed a DAGMM model for un-

supervised OOD detection, which utilized a deep auto-encoder to generate low-dimensional representations to estimate OOD scores. Xu et al. (2021) transformed the feature extracted from each layer of a pre-trained transformer model into one low-dimension representation based on the Mahalanobis distance, and then optimized an OC-SVM for detection. Some works also use language models (Nourbakhsh and Bang, 2019) and word representations (Bertero et al., 2017) to detect OOD inputs on various tasks such as log analysis (Yadav et al., 2020) and data mining (Agrawal and Agrawal, 2015).

4 Datasets and Applications

In this section, we briefly discuss representative datasets and applications for OOD detection. We classify existing OOD detection datasets into three categories according to the construction schemes of OOD samples in the testing stage:

(1) Annotate OOD Samples: This category of datasets contains OOD samples that are manually annotated by crowd-source workers. Specifically, CLINIC150 (Larson et al., 2019) is a manually labeled single-turn dialogue dataset that consists of 150 ID intent classes and 1,200 out-of-scope queries. STAR (Mosig et al., 2020) is a multi-turn dialogue dataset with annotated turn-level intents, in which OOD samples are labeled as “out_of_scope”, “custom”, or “ambiguous”. ROSTD (Gangal et al., 2020) is constructed by annotating about 4,000 OOD samples on the basis of the dataset constructed by Schuster et al. (2019).

(2) Curate OOD samples using existing classes: This category of datasets curate OOD examples by holding out a subset of classes in a given corpus (Zhang et al., 2021). Any text classification datasets can be adopted in this process.

(3) Curate OOD samples using other corpora: This category of datasets curates OOD samples using samples extracted from other datasets (Hendrycks et al., 2020; Zhou et al., 2021b), i.e., samples from other corpora are regarded as OOD samples. In this way, different NLP corpora can be combined to construct OOD detection tasks.

OOD detection tasks have also been widely applied in various NLP applications. We generally divide these applications into three types:

(1) Classification Tasks are natural applications for OOD detectors. Almost every text classifier built in the closed-world assumption needs the

OOD detection ability before deploying to production. Specifically, intent classification for dialogue systems is the most common application for OOD detection (Larson et al., 2019; Lin and Xu, 2019). Other popular application scenarios involve general text classification (Zhou et al., 2021b; Li et al., 2021), sentiment analysis (Shu et al., 2017), and topic prediction (Rawat et al., 2021).

(2) Conditional Language Generation Tasks aim to auto-regressively generate sequences of tokens. Specifically, tokens in each time step are predicted by a classification process over the vocabulary. Some studies explore the OOD detection problem on these sequence generation tasks, such as semantic parsing (Lukovnikov et al., 2021) and translation (Ren et al., 2022).

(3) Selective Prediction Tasks predict higher-quality outputs while abstaining on uncertain ones (Geifman and El-Yaniv, 2017; Varshney et al., 2022). This setting can be combined naturally with OOD detection techniques. A few studies use OOD detection approaches for selective prediction in question answering, semantic equivalence judgments, and entailment classification (Kamath et al., 2020; Xin et al., 2021).

5 Metrics

The main purposes of OOD detectors are separating OOD and ID input samples, which is essentially a binary classification process. Most methods mentioned above try to compute an *OOD score* for this problem. Therefore, threshold-free metrics that are generally used to evaluate binary classifiers are commonly used to evaluate OOD detectors:

AUROC: Area Under the Receiver Operating Characteristic curve (Davis and Goadrich, 2006). The Receiver Operating Characteristic curve is a plot showing the true positive rate $TPR = \frac{TP}{TP+FN}$ and the false positive rate $FPR = \frac{FP}{FP+TN}$ against each other, in which TP, TN, FP, FN denotes true positive, true negative, false positive, false negative, respectively. For OOD detection tasks, ID samples are usually regarded as positive. Specifically, a random OOD detector yields an AUROC score of 50% while a “perfect” OOD detector pushes this score up to 100%.

AUPR: Area Under the Precision-Recall curve (Manning and Schütze, 1999). The Precision-Recall curve plots the precision $\frac{TP}{TP+FP}$ and recall $\frac{TP}{TP+FN}$ against each other. The metric AUPR is used when the positive and negative classes in the

testing phase are severely imbalanced because the metric AUROC is biased in this situation. Generally, two kinds of AUPR scores are reported: 1) **AUPR-IN** where ID samples are specified as positive; 2) **AUPR-OUT** where OOD samples are specified as positive.

Besides these threshold-free metrics, we are also interested in the performance of OOD detectors after the deployment, i.e., when a specific threshold is selected. The following metric is usually used to measure this performance:

FPR@ N : The value of FPR when TPR is $N\%$ (Liang et al., 2018; Lee et al., 2018a). This metric measures the probability that an OOD sample is misclassified as ID when the TPR is at least $N\%$. Generally, we set $N = 95$ or $N = 90$ to ensure high performance on ID samples. This metric is important for a deployed OOD detector since obtaining a low FPR score while achieving high ID performance is important for practical systems.

In addition to the ability to detect OOD samples, some OOD detectors are also combined with downstream ID classifiers. Specifically, for a dataset that contains K ID classes, these modules allocate an additional OOD class for all the OOD samples and essentially perform a $K + 1$ class classification task. The following metrics are used to evaluate the overall performance of these modules:

F1: The macro F1 score is used to evaluate classification performance, which keeps the balance between precision and recall. Usually, F1 scores are calculated over all samples to estimate the overall performance. Some studies also compute F1 scores over ID and OOD samples, respectively, to evaluate fine-grained performances (Zhang et al., 2021).

Acc: The accuracy score is also used to evaluate classification performance (Zhan et al., 2021). See Appendix B for more details of various metrics.

6 Discussion

6.1 Pros and Cons for Different Settings

Labeled OOD data provide valuable information for OOD distributions, and thus models trained using these OOD samples usually achieve high performance in different applications. However, the collection of labeled OOD samples requires additional efforts that are extremely time-consuming and labor-extensive. Moreover, due to the infinite compositions of language, it is generally impractical to collect OOD samples for all unseen cases. Using only a small subset of OOD samples may

lead to serious selection bias issues and thus hurt the generalization of the learned model. Therefore, it is important to develop OOD detection methods that do not rely on labeled OOD samples.

OOD detection using only labeled ID data fits the above requirements. The representation learning and detecting approaches decompose the OOD detection process in this setting into two stages so that we can separately optimize each stage. Specifically, the representation learning stage attempts to learn distinct feature spaces for ID/OOD samples. Results show that this stage benefit from recent advances in pre-training and semi-supervised learning schemes on unlabeled data. Recent research also shows that a good ID classifier benefits the OOD detection (Vaze et al., 2021). OOD scoring functions aim to produce reliable scores for OOD detection. Various approaches generate the OOD score with different distance measurements and distributions. Another way to tackle the problem of lacking annotated OOD data is to generate pseudo OOD samples. Approaches in this category benefit from the strong language modeling prior and the generation ability of pre-trained models. Promising results are reported by applying the mixup strategy.

In some applications, we can only obtain a set of ID data without any labels. This situation is commonly encountered in non-classification tasks where we also need to detect OOD inputs. Compared to NLP, this setting is more widely investigated in other fields like machine learning and computer vision (CV). Popular approaches involve using estimated distribution densities or reconstruction losses as the OOD scores.

6.2 Comparison between NLP and CV in OOD Detection

OOD detection is an active research field in CV communities (Yang et al., 2021). Compared to CV, models in NLP need to tackle discrete input spaces and handle complex output structures. Therefore, additional efforts should be paid to develop algorithms for OOD detection in NLP. We summarize the differences in OOD detection between NLP and CV in the following three aspects:

Discrete Input NLP handles token sequences that lie in discrete spaces. Therefore distorting ID samples in their surface space (Ouyang et al., 2021; Choi et al., 2021; Shu et al., 2021) produces high-quality OOD samples if a careful filtering process is designed. On the contrary, CV tackles

669 inputs from continuous spaces, where it is hard to
670 navigate on the manifold of the data distribution.
671 [Du et al. \(2022b,a\)](#) showed OOD synthesizing in
672 the pixel space with a noise-additive manner led to
673 limited performance.

674 **Complex Output** Most OOD detection methods
675 in CV are proposed for K -way classification tasks.
676 However, in NLP, conditional language generation
677 tasks need to predict token sequences that lie in se-
678 quentially structured distributions, such as seman-
679 tic parsing ([Lukovnikov et al., 2021](#)), abstractive
680 summarization, and machine translation ([Ren et al.,](#)
681 [2022](#)). Hence, the perils of OOD are arguably more
682 severe as (a) errors may propagate and magnify in
683 sequentially structured output, and (b) the space of
684 low-quality outputs is greatly increased as arbitrary
685 text sequences can be generated. OOD detection
686 methods for these conditional language generation
687 tasks should consider the internal dependency of
688 input-output samples.

689 **Contextual Information** Some datasets in NLP
690 contain contextual information. It is important to
691 properly model this extra information for OOD de-
692 tection in these tasks. For example, STAR ([Mosig](#)
693 [et al., 2020](#)) is a multi-turn dialogue dataset, and ef-
694 fective OOD detectors should consider multi-turn
695 contextual knowledge in their modeling process
696 ([Chen and Yu, 2021](#)). However, most CV models
697 only consider single images as their inputs.

698 6.3 Future Research Challenges

699 **OOD Detection and Domain Generalization** In
700 most practical applications, we are not only inter-
701 ested in detecting OOD inputs that are semanti-
702 cally shifted, but also required to build more robust
703 ID classifiers that can tackle covariate shifted data
704 ([Yang et al., 2021](#)). We believe there are oppor-
705 tunities to tackle problems of OOD detection and
706 domain generalization in a unified framework. Fu-
707 ture research opportunities can be explored to equip
708 OOD detectors with better text representation ex-
709 tractors since recent results demonstrate that a good
710 ID classifier improves the OOD detection perfor-
711 mance ([Vaze et al., 2021](#)). Both new task design
712 and algorithm development can be investigated.

713 **OOD Detection with Extra Information Sources**
714 Humans usually consider OOD inputs easily dis-
715 tinguishable because they can access external in-
716 formation besides plain texts (e.g., images, audio,
717 and videos). OOD detectors are expected to per-

718 form better if we can equip them with inputs from
719 different sources. Although various works are pro-
720 posed to model each single information source,
721 such as text or image, few works are dedicated
722 to combining different sources, and no studies try
723 to equip OOD detectors with external knowledge,
724 such as structured knowledge graphs. We envision
725 great performance improvements if we can prop-
726 erly model external knowledge in OOD detectors.

727 Moreover, Internet search engines are common
728 approaches for humans to obtain external knowl-
729 edge ([Komeili et al., 2021](#)). More research opportu-
730 nities can be explored to build Internet-augmented
731 OOD detectors that can utilize rich and updated
732 knowledge yielded by search engines to enhance
733 the OOD detection performance.

734 **OOD Detection and Lifelong Learning** All pre-
735 vious approaches focus on detecting OOD inputs
736 so that we can safely ignore them. However, OOD
737 inputs usually represent new tasks that the current
738 system does not support. Systems deployed in an
739 ever-evolving environment are usually expected to
740 continuously learn from these OOD inputs rather
741 than ignoring them ([Liu and Mazumder, 2021](#)).
742 However, humans exhibit outstanding abilities in
743 learning new tasks from OOD inputs. We believe
744 OOD detectors are essential components in a life-
745 long learning system, and it is helpful to combine
746 OOD detection with a downstream lifelong learn-
747 ing process to build stronger systems.

748 **Theoretical Analysis of OOD Detection** De-
749 spite impressive empirical results that OOD studies
750 have achieved, theoretical investigation of OOD de-
751 tection is far behind the empirical success ([Morteza](#)
752 [and Li, 2022](#); [Fang et al., 2022](#)). We hope more
753 attention can be paid to theoretical analysis for
754 OOD detection and provide insights to guide the
755 development of better algorithms and applications.

756 7 Conclusion

757 In this survey, we provide a comprehensive review
758 of OOD detection methods in NLP. We formalize
759 the OOD detection tasks and identify the major
760 challenges of OOD detection in NLP. A taxonomy
761 of existing OOD detection methods is also pro-
762 vided. We hope this survey helps researchers lo-
763 cate their target problems and find the most suitable
764 datasets, metrics, and baselines. Moreover, we also
765 provide some promising directions that can inspire
766 future research and exploration.

767 Limitations

768 There are several limitations of this work. First,
769 this survey mainly focuses on OOD detection ap-
770 proaches for NLP domains. Despite the restrictive
771 scope, our work well complements the existing
772 survey on OOD detection in CV tasks, and hence
773 will benefit a well-targeted research community
774 in NLP. Second, some OOD detection methods
775 mentioned in this paper are not extended in this
776 survey due to space limitations. We include details
777 that are necessary to outline the development of
778 OOD detection methods so that readers can get
779 a comprehensive overview of this field. Our sur-
780 vey provides an elaborate starting point for readers
781 who want to dive deep into OOD detection for NLP.
782 Moreover, The term “OOD detection” has vari-
783 ous alias, such as “Anomaly Detection”, “Outlier
784 Detection”, “One-class Classification”, “Novelty
785 Detection”, and “Open Set Recognition”. These
786 notations represent similar tasks with subtle differ-
787 ences in detailed experiment settings. We do not
788 extensively discuss these differences due to space
789 limitations. Readers can refer to other papers for
790 more detailed discussions (Yang et al., 2021). Fi-
791 nally, we do not present any new empirical results.
792 It would be helpful to perform comparative experi-
793 ments over different OOD detection methods (Yang
794 et al., 2022). We leave this to future work.

795 Ethics Statement

796 This work does not present any direct ethical issues.
797 In this survey, we provide a comprehensive review
798 of OOD detection methods in NLP, and we believe
799 this study leads to intellectual merits that benefit
800 from a reliable application of NLU models.

801 References

802 Shikha Agrawal and Jitendra Agrawal. 2015. Survey
803 on anomaly detection using data mining techniques.
804 *Procedia Computer Science*, 60:708–713.

805 Dario Amodei, Chris Olah, Jacob Steinhardt, Paul
806 Christiano, John Schulman, and Dan Mané. 2016.
807 Concrete problems in ai safety. *arXiv preprint*
808 *arXiv:1606.06565*.

809 Udit Arora, William Huang, and He He. 2021. Types
810 of out-of-distribution texts and how to detect them.
811 In *Proceedings of the 2021 Conference on Empiri-
812 cal Methods in Natural Language Processing*, pages
813 10687–10701.

814 Abhijit Bendale and Terrance Boulton. 2015. Towards
815 open world recognition. In *Proceedings of the IEEE*

*conference on computer vision and pattern recogni-
tion*, pages 1893–1902. 816
817

Abhijit Bendale and Terrance E Boulton. 2016. Towards 818
open set deep networks. In *Proceedings of the IEEE
conference on computer vision and pattern recogni-
tion*, pages 1563–1572. 819
820
821

Christophe Bertero, Matthieu Roy, Carla Sauvanaud, 822
and Gilles Trédan. 2017. Experience report: Log
mining using natural language processing and ap-
plication to anomaly detection. In *2017 IEEE 28th
International Symposium on Software Reliability En-
gineering (ISSRE)*, pages 351–360. IEEE. 823
824
825
826
827

John Blitzer, Ryan McDonald, and Fernando Pereira. 828
2006. Domain adaptation with structural correspon-
dence learning. In *Proceedings of the 2006 con-
ference on empirical methods in natural language
processing*, pages 120–128. 829
830
831
832

Terrance E Boulton, Steve Cruz, Akshay Raj Dhamija, 833
Manuel Gunther, James Henrydoss, and Walter J
Scheirer. 2019. Learning and the unknown: Survey-
ing steps toward open world recognition. In *Proceeed-
ings of the AAAI conference on artificial intelligence*,
volume 33, pages 9801–9807. 834
835
836
837
838

Markus M Breunig, Hans-Peter Kriegel, Raymond T 839
Ng, and Jörg Sander. 2000. Lof: identifying density-
based local outliers. In *Proceedings of the 2000 ACM
SIGMOD international conference on Management
of data*, pages 93–104. 840
841
842
843

Nicholas Carlini and David Wagner. 2017. Adver- 844
sarial examples are not easily detected: Bypassing
ten detection methods. In *Proceedings of the 10th
ACM workshop on artificial intelligence and security*,
pages 3–14. 845
846
847
848

Raghavendra Chalapathy and Sanjay Chawla. 2019. 849
Deep learning for anomaly detection: A survey.
arXiv preprint arXiv:1901.03407. 850
851

Derek Chen and Zhou Yu. 2021. Gold: Improving 852
out-of-scope detection in dialogues using data aug-
mentation. In *Proceedings of the 2021 Conference on
Empirical Methods in Natural Language Processing*,
pages 429–442. 853
854
855
856

Qian Chen, Xiaodan Zhu, Zhen-Hua Ling, Si Wei, Hui 857
Jiang, and Diana Inkpen. 2017. Enhanced lstm for
natural language inference. In *Proceedings of the
55th Annual Meeting of the Association for Computa-
tional Linguistics (Volume 1: Long Papers)*, pages
1657–1668. 858
859
860
861
862

Ting Chen, Simon Kornblith, Mohammad Norouzi, and 863
Geoffrey Hinton. 2020. A simple framework for
contrastive learning of visual representations. In *In-
ternational conference on machine learning*, pages
1597–1607. PMLR. 864
865
866
867

DongHyun Choi, Myeong Cheol Shin, EungGyun Kim, 868
and Dong Ryeol Shin. 2021. Outflip: Generating
869

870	examples for unknown intent detection with natural language attack. In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 504–512.	925
871		926
872		927
873		928
874	Soham Dan and Dan Roth. 2021. On the effects of transformer size on in-and out-of-domain calibration. In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 2096–2101.	929
875		
876		930
877		931
878	Jesse Davis and Mark Goadrich. 2006. The relationship between precision-recall and roc curves. In <i>Proceedings of the 23rd international conference on Machine learning</i> , pages 233–240.	932
879		933
880		934
881		935
882	Shrey Desai and Greg Durrett. 2020. Calibration of pre-trained transformers. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 295–302.	936
883		937
884		938
885		939
886	Xuefeng Du, Xin Wang, Gabriel Gozum, and Yixuan Li. 2022a. Unknown-aware object detection: Learning what you don’t know from videos in the wild. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> .	940
887		941
888		942
889		943
890		944
891	Xuefeng Du, Zhaoning Wang, Mu Cai, and Yixuan Li. 2022b. Vos: Learning what you don’t know by virtual outlier synthesis. In <i>Proceedings of the International Conference on Learning Representations</i> .	945
892		946
893		947
894		948
895	Zhen Fang, Yixuan Li, Jie Lu, Jiahua Dong, Bo Han, and Feng Liu. 2022. Is out-of-distribution detection learnable? In <i>Advances in Neural Information Processing Systems</i> .	949
896		950
897		951
898		952
899	Geli Fei and Bing Liu. 2016. Breaking the closed world assumption in text classification. In <i>Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 506–514.	953
900		954
901		955
902		956
903		957
904	Yarin Gal and Zoubin Ghahramani. 2016. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In <i>international conference on machine learning</i> , pages 1050–1059. PMLR.	958
905		959
906		960
907		961
908	Varun Gangal, Abhinav Arora, Arash Einolghozati, and Sonal Gupta. 2020. Likelihood ratios and generative classifiers for unsupervised out-of-domain detection in task oriented dialog. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 34, pages 7764–7771.	962
909		963
910		964
911		965
912		966
913		967
914	Yonatan Geifman and Ran El-Yaniv. 2017. Selective classification for deep neural networks. <i>Advances in neural information processing systems</i> , 30.	968
915		969
916		970
917	Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2020. Generative adversarial networks. <i>Communications of the ACM</i> , 63(11):139–144.	971
918		972
919		973
920		974
921		975
922	Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. 2014. Explaining and harnessing adversarial examples. <i>arXiv preprint arXiv:1412.6572</i> .	976
923		977
924		
	Jiuxiang Gu, Zhenhua Wang, Jason Kuen, Lianyang Ma, Amir Shahroudy, Bing Shuai, Ting Liu, Xingxing Wang, Gang Wang, Jianfei Cai, et al. 2018. Recent advances in convolutional neural networks. <i>Pattern recognition</i> , 77:354–377.	
	Xiaoyi Gu, Leman Akoglu, and Alessandro Rinaldo. 2019. Statistical analysis of nearest neighbor methods for anomaly detection. <i>Advances in Neural Information Processing Systems</i> , 32.	
	Beliz Gunel, Jingfei Du, Alexis Conneau, and Ves Stoyanov. 2020. Supervised contrastive learning for pre-trained language model fine-tuning. <i>arXiv preprint arXiv:2011.01403</i> .	
	Dan Hendrycks and Kevin Gimpel. 2016. A baseline for detecting misclassified and out-of-distribution examples in neural networks. <i>arXiv preprint arXiv:1610.02136</i> .	
	Dan Hendrycks, Xiaoyuan Liu, Eric Wallace, Adam Dziedzic, Rishabh Krishnan, and Dawn Song. 2020. Pretrained transformers improve out-of-distribution robustness. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 2744–2751.	
	Dan Hendrycks, Mantas Mazeika, and Thomas Dietterich. 2018. Deep anomaly detection with outlier exposure. <i>arXiv preprint arXiv:1812.04606</i> .	
	Dan Hendrycks, Mantas Mazeika, Saurav Kadavath, and Dawn Song. 2019. Using self-supervised learning can improve model robustness and uncertainty. <i>Advances in neural information processing systems</i> , 32.	
	Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. <i>Neural computation</i> , 9(8):1735–1780.	
	Amita Kamath, Robin Jia, and Percy Liang. 2020. Selective question answering under domain shift. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 5684–5696.	
	Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In <i>Proceedings of NAACL-HLT</i> , pages 4171–4186.	
	Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. 2020. Supervised contrastive learning. <i>Advances in Neural Information Processing Systems</i> , 33:18661–18673.	
	Joo-Kyung Kim and Young-Bum Kim. 2018. Joint learning of domain classification and out-of-domain detection with dynamic class weighting for sacrificing false acceptance rates. <i>arXiv preprint arXiv:1807.00072</i> .	

978	Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, et al. 2021. Wilds: A benchmark of in-the-wild distribution shifts. In <i>International Conference on Machine Learning</i> , pages 5637–5664. PMLR.	1033
979		1034
980		1035
981		1036
982		
983		1037
984		1038
		1039
985	Mojtaba Komeili, Kurt Shuster, and Jason Weston. 2021. Internet-augmented dialogue generation. <i>arXiv preprint arXiv:2107.07566</i> .	1040
986		
987		
988	Lingkai Kong, Haoming Jiang, Yuchen Zhuang, Jie Lyu, Tuo Zhao, and Chao Zhang. 2020. Calibrated language model fine-tuning for in-and out-of-distribution data. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 1326–1340.	1041
989		1042
990		1043
991		1044
992		1045
993		
994	Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. 2017. Simple and scalable predictive uncertainty estimation using deep ensembles. <i>Advances in neural information processing systems</i> , 30.	1046
995		1047
996		1048
997		1049
998		
999	Hao Lang, Yinhe Zheng, Jian Sun, Fei Huang, Luo Si, and Yongbin Li. 2022. Estimating soft labels for out-of-domain intent detection. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 261–276.	1050
1000		1051
1001		1052
1002		1053
1003		1054
1004	Stefan Larson, Anish Mahendran, Joseph J Peper, Christopher Clarke, Andrew Lee, Parker Hill, Jonathan K Kummerfeld, Kevin Leach, Michael A Laurenzano, Lingjia Tang, et al. 2019. An evaluation dataset for intent classification and out-of-scope prediction. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 1311–1316.	1055
1005		1056
1006		1057
1007		1058
1008		1059
1009		
1010		1060
1011		1061
1012		1062
1013		
1014	Kimin Lee, Honglak Lee, Kibok Lee, and Jinwoo Shin. 2018a. Training confidence-calibrated classifiers for detecting out-of-distribution samples. In <i>International Conference on Learning Representations</i> .	1063
1015		1064
1016		1065
1017		1066
1018		1067
1019		1068
1020		1069
1021		
1022		1070
1023		1071
1024		1072
1025		1073
1026		
1027		1074
1028		1075
1029		1076
1030		
1031		1077
1032		1078
		1079
		1080
		1081
		1082
		1083
		1084
		1085

1086	Guansong Pang, Chunhua Shen, Longbing Cao, and Anton Van Den Hengel. 2021. Deep learning for anomaly detection: A review. <i>ACM Computing Surveys (CSUR)</i> , 54(2):1–38.	1142
1087		1143
1088		1144
1089		1145
1090	Ellie Pavlick and Joel Tetreault. 2016. An empirical analysis of formality in online communication. <i>Transactions of the Association for Computational Linguistics</i> , 4:61–74.	1146
1091		1147
1092		1148
1093		1149
1094	Alexander Podolskiy, Dmitry Lipin, Andrey Bout, Ekaterina Artemova, and Irina Piontkovskaya. 2021. Revisiting mahalanobis distance for transformer-based out-of-domain detection. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 35, pages 13675–13682.	1150
1095		1151
1096		1152
1097		1153
1098		1154
1099		1155
1100	Alan Ramponi and Barbara Plank. 2020. Neural unsupervised domain adaptation in nlp—a survey. In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 6838–6855.	1156
1101		1157
1102		1158
1103		1159
1104	Mrinal Rawat, Ramya Hebbalaguppe, and Lovekesh Vig. 2021. Pnpood: Out-of-distribution detection for text classification via plug andplay data augmentation. <i>arXiv preprint arXiv:2111.00506</i> .	1160
1105		1161
1106		1162
1107		1163
1108	Jie Ren, Peter J Liu, Emily Fertig, Jasper Snoek, Ryan Poplin, Mark Depristo, Joshua Dillon, and Balaji Lakshminarayanan. 2019. Likelihood ratios for out-of-distribution detection. <i>Advances in neural information processing systems</i> , 32.	1164
1109		1165
1110		1166
1111		1167
1112		1168
1113	Jie Ren, Jiaming Luo, Yao Zhao, Kundan Krishna, Mohammad Saleh, Balaji Lakshminarayanan, and Peter J Liu. 2022. Out-of-distribution detection and selective generation for conditional language models. <i>arXiv preprint arXiv:2209.15558</i> .	1169
1114		1170
1115		1171
1116		1172
1117		1173
1118	Andras Rozsa, Manuel Günther, and Terrance E Boulton. 2017. Adversarial robustness: Softmax versus openmax. <i>arXiv preprint arXiv:1708.01697</i> .	1174
1119		1175
1120		1176
1121	Seonghan Ryu, Seokhwan Kim, Junhwi Choi, Hwanjo Yu, and Gary Geunbae Lee. 2017. Neural sentence embedding using only in-domain sentences for out-of-domain sentence detection in dialog systems. <i>Pattern Recognition Letters</i> , 88:26–32.	1177
1122		1178
1123		1179
1124		1180
1125		1181
1126	Seonghan Ryu, Sangjun Koo, Hwanjo Yu, and Gary Geunbae Lee. 2018. Out-of-domain detection based on generative adversarial network. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 714–718.	1182
1127		1183
1128		1184
1129		1185
1130		1186
1131	Walter J Scheirer, Lalit P Jain, and Terrance E Boulton. 2014. Probability models for open set recognition. <i>IEEE transactions on pattern analysis and machine intelligence</i> , 36(11):2317–2324.	1187
1132		1188
1133		1189
1134		1190
1135	Sebastian Schuster, Sonal Gupta, Rushin Shah, and Mike Lewis. 2019. Cross-lingual transfer learning for multilingual task oriented dialog. In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 3795–3805.	1191
1136		1192
1137		1193
1138		1194
1139		1195
1140		1196
1141		1197
		1198
	Yilin Shen, Yen-Chang Hsu, Avik Ray, and Hongxia Jin. 2021. Enhancing the generalization for intent classification and out-of-domain detection in slu. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 2443–2453.	
	Lei Shu, Yassine Benajiba, Saab Mansour, and Yi Zhang. 2021. Odist: Open world classification via distributionally shifted instances. In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 3751–3756.	
	Lei Shu, Hu Xu, and Bing Liu. 2017. Doc: Deep open classification of text documents. In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , pages 2911–2916.	
	Yiyun Sun, Yifei Ming, Xiaojin Zhu, and Yixuan Li. 2022. Out-of-distribution detection with deep nearest neighbors. In <i>International Conference on Machine Learning</i> .	
	Ming Tan, Yang Yu, Haoyu Wang, Dakuo Wang, Saloni Potdar, Shiyu Chang, and Mo Yu. 2019. Out-of-domain detection for low-resource text classification tasks. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3566–3572.	
	Vladimir Vapnik. 1991. Principles of risk minimization for learning theory. <i>Advances in neural information processing systems</i> , 4.	
	Neeraj Varshney, Swaroop Mishra, and Chitta Baral. 2022. Investigating selective prediction approaches across several tasks in IID, OOD, and adversarial settings. In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 1995–2002, Dublin, Ireland. Association for Computational Linguistics.	
	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>Advances in neural information processing systems</i> , 30.	
	Sagar Vaze, Kai Han, Andrea Vedaldi, and Andrew Zisserman. 2021. Open-set recognition: A good closed-set classifier is all you need. <i>arXiv preprint arXiv:2110.06207</i> .	
	Vikas Verma, Alex Lamb, Christopher Beckham, Amir Najafi, Ioannis Mitliagkas, David Lopez-Paz, and Yoshua Bengio. 2019. Manifold mixup: Better representations by interpolating hidden states. In <i>International Conference on Machine Learning</i> , pages 6438–6447. PMLR.	
	Ricardo Vilalta and Youssef Drissi. 2002. A perspective view and survey of meta-learning. <i>Artificial intelligence review</i> , 18(2):77–95.	

1199	Apoorv Vyas, Nataraj Jammalamadaka, Xia Zhu, Dipankar Das, Bharat Kaul, and Theodore L Willke. 2018. Out-of-distribution detection using an ensemble of self supervised leave-out classifiers. In <i>Proceedings of the European Conference on Computer Vision (ECCV)</i> , pages 550–564.	Jingkang Yang, Pengyun Wang, Dejian Zou, Zitang Zhou, Kunyuan Ding, WENXUAN PENG, Haoqi Wang, Guangyao Chen, Bo Li, Yiyu Sun, et al. 2022. Openood: Benchmarking generalized out-of-distribution detection. In <i>Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track</i> .	1257
1200			1258
1201			1259
1202			1260
1203			1261
1204			1262
1205	Jindong Wang, Cuiling Lan, Chang Liu, Yidong Ouyang, Tao Qin, Wang Lu, Yiqiang Chen, Wenjun Zeng, and Philip Yu. 2022. Generalizing to unseen domains: A survey on domain generalization. <i>IEEE Transactions on Knowledge and Data Engineering</i> .	Jingkang Yang, Kaiyang Zhou, Yixuan Li, and Ziwei Liu. 2021. Generalized out-of-distribution detection: A survey. <i>arXiv preprint arXiv:2110.11334</i> .	1264
1206			1265
1207			1266
1208			
1209			
1210	Wei Wang, Vincent W Zheng, Han Yu, and Chunyan Miao. 2019. A survey of zero-shot learning: Settings, methods, and applications. <i>ACM Transactions on Intelligent Systems and Technology (TIST)</i> , 10(2):1–37.	Lantao Yu, Weinan Zhang, Jun Wang, and Yong Yu. 2017. Seqgan: Sequence generative adversarial nets with policy gradient. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 31.	1267
1211			1268
1212			1269
1213			1270
1214			
1215	Yanan Wu, Keqing He, Yuanmeng Yan, QiXiang Gao, Zhiyuan Zeng, Fujia Zheng, Lulu Zhao, Huixing Jiang, Wei Wu, and Weiran Xu. 2022. Revisit overconfidence for OOD detection: Reassigned contrastive learning with adaptive class-dependent threshold. In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 4165–4179, Seattle, United States. Association for Computational Linguistics.	Zhiyuan Zeng, Keqing He, Yuanmeng Yan, Zijun Liu, Yanan Wu, Hong Xu, Huixing Jiang, and Weiran Xu. 2021a. Modeling discriminative representations for out-of-domain detection with supervised contrastive learning. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)</i> , pages 870–878.	1271
1216			1272
1217			1273
1218			1274
1219			1275
1220			1276
1221			1277
1222			1278
1223			1279
1224			
1225	Ji Xin, Raphael Tang, Yaoliang Yu, and Jimmy Lin. 2021. The art of abstention: Selective prediction and error regularization for natural language processing. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 1040–1051.	Zhiyuan Zeng, Hong Xu, Keqing He, Yuanmeng Yan, Sihong Liu, Zijun Liu, and Weiran Xu. 2021b. Adversarial generative distance-based classifier for robust out-of-domain detection. In <i>ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 7658–7662. IEEE.	1280
1226			1281
1227			1282
1228			1283
1229			1284
1230			1285
1231			1286
1232			
1233	Hu Xu, Bing Liu, Lei Shu, and P Yu. 2019. Open-world learning and application to product classification. In <i>The World Wide Web Conference</i> , pages 3413–3419.	Li-Ming Zhan, Haowen Liang, Bo Liu, Lu Fan, Xiao-Ming Wu, and Albert YS Lam. 2021. Out-of-scope intent detection with self-supervision and discriminative training. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 3521–3532.	1287
1234			1288
1235			1289
1236			1290
1237			1291
1238			1292
1239			1293
1240			1294
1241			
1242			
1243			
1244	Rakesh Bahadur Yadav, P Santosh Kumar, and Sunita Vikrant Dhavale. 2020. A survey on log anomaly detection using deep learning. In <i>2020 8th International Conference on Reliability, Infocom Technologies and Optimization (Trends and Future Directions)(ICRITO)</i> , pages 1215–1220. IEEE.	Bangzuo Zhang and Wanli Zuo. 2008. Learning from positive and unlabeled examples: A survey. In <i>2008 International Symposiums on Information Processing</i> , pages 650–654. IEEE.	1295
1245			1296
1246			1297
1247			1298
1248			
1249			
1250	Guangfeng Yan, Lu Fan, Qimai Li, Han Liu, Xiaotong Zhang, Xiao-Ming Wu, and Albert YS Lam. 2020. Unknown intent detection using gaussian mixture model with an application to zero-shot intent classification. In <i>Proceedings of the 58th annual meeting of the association for computational linguistics</i> , pages 1050–1060.	Hanlei Zhang, Hua Xu, and Ting-En Lin. 2021. Deep open intent classification with adaptive decision boundary. In <i>AAAI</i> , pages 14374–14382.	1299
1251			1300
1252			
1253			
1254			
1255			
1256			
		Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. 2018. mixup: Beyond empirical risk minimization. In <i>International Conference on Learning Representations</i> .	1302
			1303
			1304
			1305
		Yinhe Zheng, Guanyi Chen, and Minlie Huang. 2020. Out-of-domain detection for natural language understanding in dialog systems. <i>IEEE/ACM Transactions on Audio, Speech, and Language Processing</i> , 28:1198–1209.	1306
			1307
			1308
			1309
			1310

1311 Da-Wei Zhou, Han-Jia Ye, and De-Chuan Zhan. 2021a.
1312 Learning placeholders for open-set recognition. In
1313 *Proceedings of the IEEE/CVF Conference on Com-*
1314 *puter Vision and Pattern Recognition*, pages 4401–
1315 4410.

1316 Wenxuan Zhou, Fangyu Liu, and Muhao Chen. 2021b.
1317 Contrastive out-of-distribution detection for pre-
1318 trained transformers. In *Proceedings of the 2021*
1319 *Conference on Empirical Methods in Natural Lan-*
1320 *guage Processing*, pages 1100–1111.

1321 Yunhua Zhou, Peiju Liu, and Xipeng Qiu. 2022. [KNN-](#)
1322 [contrastive learning for out-of-domain intent classifi-](#)
1323 [cation](#). In *Proceedings of the 60th Annual Meeting of*
1324 *the Association for Computational Linguistics (Vol-*
1325 *ume 1: Long Papers)*, pages 5129–5141, Dublin,
1326 Ireland. Association for Computational Linguistics.

1327 Bo Zong, Qi Song, Martin Renqiang Min, Wei Cheng,
1328 Cristian Lumezanu, Daeki Cho, and Haifeng Chen.
1329 2018. Deep autoencoding gaussian mixture model
1330 for unsupervised anomaly detection. In *International*
1331 *conference on learning representations*.

1332 A Surveying Process

1333 In this appendix, we provide more details of how
1334 we select papers for our survey. Specifically, the
1335 selected paper follows at least one criterion listed
1336 below:

- 1337 1. Peer-reviewed papers published in Top-tier
1338 NLP venues, such as ACL, EMNLP, NAACL,
1339 AACL, and IJCAI.
- 1340 2. Peer-reviewed papers that have a significant
1341 impact on the OOD detection area. These pa-
1342 pers are not necessarily limited to NLP tasks.
- 1343 3. Papers that are highly cited in the OOD detec-
1344 tion area.
- 1345 4. Most recently published papers that make
1346 a non-trivial contribution to OOD detection,
1347 such as methods, datasets, metrics, and theo-
1348 retical analysis.
- 1349 5. Papers that initiate each research direction in
1350 the OOD detection area.

1351 B More details of Metrics

1352 Table 1 provides more detailed information of vari-
1353 ous metrics for OOD detection, regarding whether
1354 to consider ID performance, frequency of use, and
1355 applications.

Metric	Definition	Whether to consider ID performance	Frequency of use	Applications	Papers that use this metric (Selected)
AUROC	Area under the Receiver Operating Characteristic curve	No	Very Frequent	NLP, CV, ML	(Hendrycks and Gimpel, 2016; Hendrycks et al., 2018, 2019; Lee et al., 2018a)
AUPR-IN	Area under the Precision-Recall curve (ID samples as positive)	No	Frequent	NLP, CV, ML	(Lee et al., 2018a; Zheng et al., 2020; Shen et al., 2021)
AUPR-OUT	Area under the Precision-Recall curve (OOD samples as positive)	No	Frequent	NLP, CV, ML	(Lee et al., 2018a; Zheng et al., 2020; Shen et al., 2021)
FPR@ N	Value of FPR when TPR is $N\%$	No	Not Frequent	NLP, CV, ML	(Lee et al., 2018a; Zheng et al., 2020; Shen et al., 2021)
F1	Macro F1 score over all testing samples (ID+OOD)	Yes	Very Frequent	NLP	(Xu et al., 2019; Zhan et al., 2021; Shu et al., 2021; Zhou et al., 2022)
Acc	Accuracy score over all testing samples (ID+OOD)	Yes	Very Frequent	NLP	(Zhan et al., 2021; Shu et al., 2017, 2021; Zhou et al., 2022)

Table 1: More detailed information of various metrics for OOD detection.