

A CRITICAL ANALYSIS OF OUT-OF-DISTRIBUTION DETECTION FOR DOCUMENT UNDERSTANDING

Anonymous authors

Paper under double-blind review

ABSTRACT

Large-scale pretraining is widely used in recent document understanding models. During deployment, one may expect that large-scale pretrained models should trigger a conservative fallback policy when encountering out-of-distribution (OOD) samples, which suggests the importance of OOD detection. However, most existing OOD detection methods focus on single-modal inputs such as images or texts. While documents are multi-modal in nature, it is underexplored if and how multi-modal information in documents can be exploited for OOD detection. In this work, we first provide a systematic and in-depth analysis on OOD detection for document understanding models. We study the effects of model modality, pretraining, and finetuning across various types of OOD inputs. In particular, we find that spatial information is critical for document OOD detection. To better exploit spatial information, we propose a simple yet effective spatial-aware adapter, which serves as an add-on module to adapt transformer-based language models to document domain. Extensive experiments show that our method consistently improves ID accuracy and OOD detection performance compared to baselines. We hope our findings can help inspire future works on understanding OOD robustness for documents.

1 INTRODUCTION

The recent success of large-scale pretrained models has led to the widespread deployment of deep models in various applications. In the document domain, model predictions are increasingly used to help humans make decisions in important applications ranging from tax form processing, machine learning assistant medical reports analysis, deep analyses from financial forms, *etc.* However, in most cases, models are pretrained on collected data but are then deployed in an environment with a different distribution over the observed data (Cui et al., 2021). For example, with the outbreak of COVID-19 (Velavan & Meyer, 2020), machine-assisted medical document analysis systems have to face continually changing data distributions. This motivates the need for reliable methods in the document domain to detect out-of-distribution (OOD) inputs.

The goal of OOD detection is to categorize in-distribution (ID) test samples into one of the known categories and detect instances that do not belong to any known classes (Huang & Li, 2021; Bendale & Boulton, 2016). Generally, a model is optimized on a particular task (*e.g.*, image classification (Deng et al., 2009)), and a companion OOD detector is built as a safeguard for the classifier. Recently, large-scale pretrained models have demonstrated promising results in multiple domains (Dosovitskiy et al., 2021; Hendrycks et al., 2020) as pretraining enables models to learn powerful and transferable feature representations (Radford et al., 2021). In particular, the models obtained by finetuning large-scale pretrained models are significantly better at OOD detection even with a simple distance metric (Lee et al., 2018; Radford et al., 2021).

It is underexplored whether existing OOD detection methods that demonstrate success for images or text can be naturally extended to documents. The main challenges posed in document OOD detection stem from the fact that document understanding is inherently multi-modal, thus, it is suboptimal to rely on a single

modality. The majority of recent OOD detection approaches focus on single-modal learning (Hsu et al., 2020; Zhou et al., 2021; Xu et al., 2021a; Jin et al., 2022), and they are not compatible with document understanding tasks which require multi-modal learning. The spatial relationship of text blocks in documents further differentiated the document multimodal learning from the multimodal learning in vision-language domain (Lu et al., 2019; Li et al., 2020). In addition, recent document pretraining methods have demonstrated remarkable performance on various downstream document understanding tasks (Xu et al., 2020; 2021b; Huang et al., 2022; Li et al., 2021; Cui et al., 2021; Hong et al., 2022; Gu et al., 2022; Wang et al., 2022). However, existing pretraining datasets for documents are limited and lack diversity, in sharp contrast to common pre-training datasets for natural images. Therefore, it is not obvious which OOD detection methods are reliable in the document domain and how pretraining impacts OOD robustness.

This paper investigates the OOD robustness in the document domain through the following questions: (1) Are pretrained models robust to OOD examples? Is further pretraining beneficial? How do the pretraining data and tasks affect the performance? (2) How does multimodality (textual, visual, and spatial) affect OOD robustness? (3) Are existing OOD detection methods developed for natural images and texts transferrable to documents? We present a large-scale evaluation of recent approaches. We focus on models pretrained on different data types and evaluate them on a diverse range of document understanding benchmarks across visual, textual, and spatial modalities. Our key contributions are summarized as follows:

- We show that pretraining datasets and tasks significantly impact OOD detection performance. Through extensive pretraining and finetuning experiments, we find that higher finetuning performance on ID data does not usually translate to better performance on OOD data. This observation emphasizes the importance of considering metrics beyond ID performance for measuring model reliability.
- We propose a spatial-aware adapter, which can serve as an add-on module to transformer-based models and learn the spatial-aware representation. Our method can easily transfer the pretrained language models to the document domain. Extensive experiments show that our method can consistently improve ID accuracy and OOD detection performance across a broad spectrum of datasets.
- We show that recent conclusions drawn from OOD detection methods are valid for images and texts but do not always transfer to documents. For a wide range of document models, we observe that OOD samples are easier to identify in the feature space than in the logit space.

The rest of the paper is organized as follows. Sec. 2 provides the preliminaries and related works. In Sec. 3, we provide a comprehensive analysis of OOD robustness for document models and conclude in Sec. 4.

2 PRELIMINARIES AND RELATED WORKS

2.1 DOCUMENT MODELS AND PRETRAINING

Large-scale pretrained models have attracted a lot of attention in the document domain. In vision or natural language processing (NLP) tasks, pretraining has shown great success in producing generic representations that learn from large-scale unlabeled corpora (Devlin et al., 2018; Lu et al., 2019; Su et al., 2019; He et al., 2020). Document pretraining also seeks to find universal representations suitable for any downstream task. However, the unique characteristics of document images distinguish document pretraining works from previous ones in vision or language domains. For documents, the contents are spatially distributed, and visual and textual information co-occurs within the semantic regions. In contrast, inputs in the language domain are pure texts, and inputs in the vision-language domain are image-text pairs.

Recent document pretraining models differ in architecture and objectives, as depicted in Fig. 2. LayoutLM (Xu et al., 2020) extends BERT to learn contextualized word representations for document images through multi-task learning. It takes a sequence of Optical Character Recognition (OCR) (Smith, 2007)

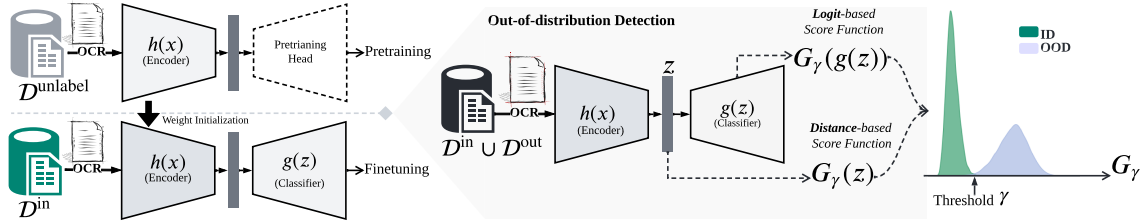


Figure 1: Schematic description of OOD detection for document classification. The left part shows the pretraining and finetuning pipelines. During inference time, for a given input document image, we calculate the OOD detection score G_γ according to different methods (logit-base or distance-base). The OOD detector will identify the input document as OOD if the OOD score is smaller than the threshold value γ .

words and word bounding boxes as inputs during pretraining and finetuning. LayoutLMv2 (Xu et al., 2021b) improves on the prior work by including an image encoder in pretraining and training them jointly. Like LayoutLMv2, DocFormer (Appalaraju et al., 2021) also adopts a CNN model to extract image grid features. It fuses the spatial information as an inductive bias for the self-attention module. The latest version, LayoutLMv3 (Huang et al., 2022), shares similar ideas as LayoutLMv2 and further enhances the visual and spatial characteristics by introducing two other tasks: masked image modeling and word-patch alignment. Another line of works for document pretraining focuses on different granularities of document images and takes region-level text blocks as the basic input elements, such as SelfDoc (Li et al., 2021) and UDoc (Gu et al., 2021). The pretraining tasks of SelfDoc and UDoc are based on feature space. They adopt a cross-modal encoder to model the relationship between visual and textual features. Instead of using the spatial information at the input layer, SelfDoc and UDoc encode the 2D spatial information with a linear mapping and fuse the position embeddings at the output layer of the image encoder and sentence encoder. Despite the promising performance of those pretrained models on downstream applications, it remains largely underexplored whether recent document pretraining models are robust to various types of OOD data, the role of pretraining and finetuning, and the key factors for document OOD detection.

2.2 OUT-OF-DISTRIBUTION DETECTION

Many OOD detection methods have been proposed for deep models, including generative model-based methods (Ge et al., 2017; Oza & Patel, 2019; Nalisnick et al., 2019; Ren et al., 2019; Xiao et al., 2020; Morteza & Li, 2022), and discriminative-model based methods. For the latter category, an OOD scoring function can be derived based on the softmax output or logit space (Liu et al., 2020; Hsu et al., 2020; Huang & Li, 2021; Liang et al., 2018; Sun et al., 2021), gradient information (Huang et al., 2021), or the feature space (Sastry & Oore, 2020; Sehwan et al., 2021; Winkens et al., 2020; Sun et al., 2022) of a classifier. Despite their impressive performance, most of the scores are developed for natural images and text inputs. A recent work (Larson et al., 2022) studies OOD detection performance for documents, but only explores a limited number of models and OOD detection methods. Furthermore, the relationship between pretraining, finetuning, and spatial information is underexplored. In this work, we provide a finer-grained and comprehensive analysis and hope to shed light on the key factors of OOD robustness for documents.

Notations We denote the input and label space \mathcal{X}^{in} and $\mathcal{Y}^{\text{in}} = \{1, \dots, K\}$, respectively. Let $\mathcal{D}^{\text{in}} = \{(\mathbf{x}_i^{\text{in}}, y_i^{\text{in}})\}_{i=1}^N$ denote an ID dataset, where $\mathbf{x} \in \mathcal{X}^{\text{in}}$ is the input feature vector, and $y^{\text{in}} \in \mathcal{Y}^{\text{in}}$ denotes the semantic label for K -way classification. Let $\mathcal{D}^{\text{out}} = \{(\mathbf{x}_i^{\text{out}}, y_i^{\text{out}})\}_{i=1}^M$ denote an OOD test set where $y^{\text{out}} \in \mathcal{Y}^{\text{out}}$, and $\mathcal{Y}^{\text{out}} \cap \mathcal{Y}^{\text{in}} = \emptyset$. OOD detection can be formulated as a binary classification problem, which aims to distinguish between ID and OOD data. We express the neural network model $f := g \circ h$ as a composition of a feature extractor $h : \mathcal{X}^{\text{in}} \rightarrow \mathbb{R}^d$ and a classifier $g : \mathbb{R}^d \rightarrow \mathbb{R}^K$, which maps the feature

embedding of an input to K real-valued numbers known as logits. During inference time, OOD detection can be performed by exercising a thresholding mechanism $G_\gamma(\mathbf{x}) = \mathbb{1}\{S(\mathbf{x}) \geq \gamma\}$ where by convention samples with higher scores $S(\mathbf{x})$ are classified as ID and vice versa. The threshold γ is typically chosen so that a high fraction of ID data (e.g., 95%) is correctly classified.

We group OOD detection methods into two major categories: logit-based scores are derived from the logit layer of the model, while distance-based methods are directly based on the feature embedding layer, as shown in Fig. 1. We describe a few popular OOD detection methods for each category as follows.

- *Logit-based*: Maximum Softmax Probability (MSP) score (Hendrycks & Gimpel, 2017) $S_{\text{MSP}} = \max_{i \in [K]} e^{f_i(\mathbf{x})} / \sum_{j=1}^K e^{f_j(\mathbf{x})}$ naturally arises as a classic baseline since logits can be converted to a categorical distribution $p(y|\mathbf{x})$; Energy score (Liu et al., 2020): $S_{\text{Energy}} = \log \sum_{i \in [K]} e^{f_i(\mathbf{x})}$ utilizes the Helmholtz free energy of the data and theoretically aligns with the logarithm of the ID density; MaxLogit score (Hendrycks et al., 2022): $S_{\text{Maxlogit}} = \max_{i \in [K]} f_i(\mathbf{x})$ removes the softmax function in MSP and demonstrates promising performance on large-scale natural image datasets recently.
- *Distance-based*: Distance-based methods directly leverage feature embeddings h based on the idea that OOD inputs are relatively far away from ID centroids or prototypes. Depending on the distributional assumption of feature embeddings, methods can be characterized as 1) parametric methods such as Mahalanobis score (Lee et al., 2018; Sehwag et al., 2021) which assumes ID embeddings follow class-conditional Gaussian distributions and use Mahalanobis distance from the ID centroid as the distance metric; 2) non-parametric methods such as KNN+ (Sun et al., 2022) which uses cosine similarity as the distance metric.

Evaluation Metrics To evaluate OOD detection performance, we adopt two commonly used metrics (Hendrycks & Gimpel, 2017): Area Under the Receiver Operating Characteristic (AUROC) and False Positive Rate at 95% Recall (FPR95). For ID test sets, we report Accuracy (Acc), F1 score, and Mean Average Precision (mAP).

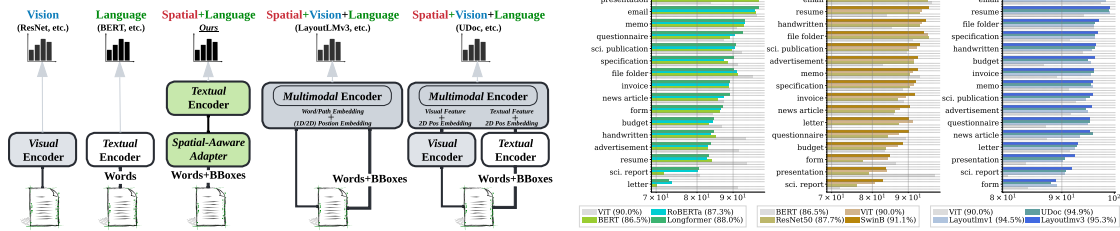
3 ANALYZING OOD ROBUSTNESS FOR DOCUMENT MODELS

In this section, we consider the task of document classification, where models are expected to classify documents into categories such as *scientific papers*, *resumes*, etc. However, it is underexplored whether models are robust to OOD samples at test time. Most document classification datasets exist in the form of images (Harley et al., 2015). Usually, the first step is to pass the image through an OCR system to obtain a set of text blocks along with their coordinates in the image. Given the input image, extracted words, and coordinates, models can utilize single-modal or multi-modal information to classify the document.

Models Fig. 2 (a) shows common structures for document image pretraining and classification models¹. According to the input modalities, we categorize them into the following groups:

(1) *Vision-based*: Since current document datasets exist as images, we can treat document classification as the standard image classification problem. In our experiments, we consider ResNet-50 (He et al., 2016) and ViT (Fort et al., 2021) as exemplar document image classification models. As for pretrained weights, we consider two settings: pretrained on ImageNet (Deng et al., 2009) and further pretrained on IIT-CDIP (Lewis et al., 2006). We adopt masked image modeling (MIM) for image pretraining with a mask ratio of 0.6. Note that the document classification dataset we used in this paper, RVL-CDIP, is a subset of IIT-CDIP. Hence, unless otherwise specified, the IIT-CDIP pretraining data used in this paper excludes RVL-CDIP.

¹See Appendix A.1.2 for further details about the models and hyperparameters.



(a) Illustration of common structures for document pre-training/classification. (b) A detailed comparison of per-category accuracy on the RVL-CDIP test set.

Figure 2: (Left) Illustration of models for document pretraining and classification. Our proposed spatial-aware pretraining and finetuning models are the network architectures in green blocks. We also show the modality information on top of each architecture. (Right): Evaluating finetuning performance for document classification of pretrained models. We group models into three groups (from left to right): language-only, vision-only, and multimodal. For each group, we also present the performance of a model in another group (shown in grey) for better reference. The average accuracy for each model is shown in the parenthesis.

(2) *Text-based*: Alternatively, we can define the classification as a text classification problem since documents typically contain words. In our experiments, we consider RoBERTa (Liu et al., 2019) as the backbone and append a classifier for finetuning. Since some documents such as scientific papers consist of sentences with more than 512 tokens, we also consider Longformer (Beltagy et al., 2020), which can handle a maximum of 4,096 input tokens. Similar to the vision-based models, we further pretrain the language models with masked language modeling (MLM) on IIT-CDIP extracted text corpus.

(3) *Text+Spatial*: Layout information plays a crucial role in the document domain. As shown in Fig. 3, a document is composed of words or images with some specific layouts. To investigate the effect of layout information, we adopt LayoutLM as a baseline. It is specifically designed for documents and trained on the full IIT-CDIP data. Inspired by the promising OOD detection performance of spatial-aware models (Sec 3.3) and the recent advances in adapter-based transformers (Pfeiffer et al., 2020), we propose a new spatial-aware adaptor, a small learned module that can be inserted within a pretrained language model. Besides the simplicity, our adaptor is competitive for both ID classification and OOD detection (Sec 3.4).

(4) *Visual+Textual+Spatial*: Current state-of-the-art methods tailored to documents consider various input granularity and modality and utilize textual, visual, and spatial information for document tasks. Despite the promising performance, such models are large in size and computationally heavy. We select two representative models to evaluate upon: LayoutLMv3 and UDoc. For a fair comparison, both models are pretrained on full IIT-CDIP.

Constructing ID and OOD Datasets We construct ID datasets from RVL-CDIP (Harley et al., 2015). Specifically, we specify 12 out of 16 classes as ID classes. For OOD datasets, we consider two scenarios:

(1) *In-domain OOD*: To determine the OOD categories, we extensively analyze the performance of recent document classification models. Fig. 2(b) shows a detailed comparison of per-category accuracy on the RVL-CDIP test set. Naturally, for the classes the model performs poorly on, we may expect models to detect such inputs as OOD instead of assigning a specific ID class with low confidence. We observe that the 4 categories *letter*, *form*, *scientific report*, *presentation* result in the worst performance across most of models with different modality, which we use as OOD categories and construct the OOD datasets accordingly. The ID dataset is constructed from the remaining 12 categories. We refer to these OOD datasets as *in-domain*, as they are also constructed from RVL-CDIP.

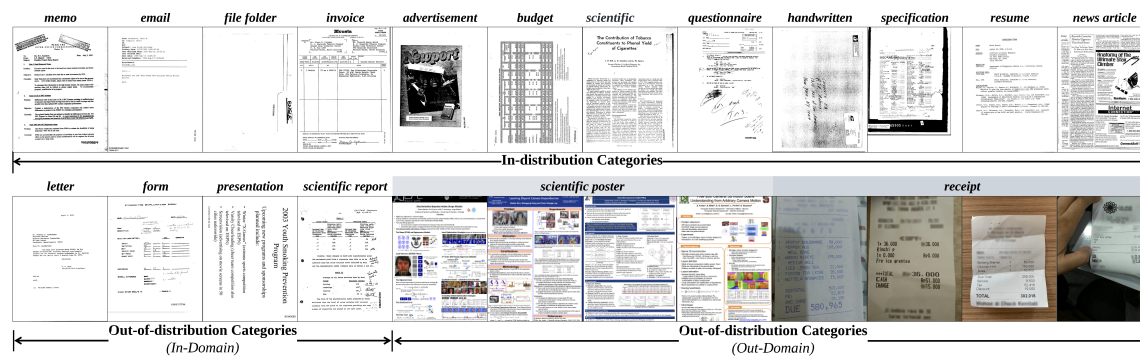


Figure 3: **(Top)** Examples of ID inputs sampled from RVL-CDIP (top). **(Bottom)** In-domain OOD from RVL-CDIP, and out-domain OOD from *Scientific Poster* and *Receipts*.

(2) *Out-domain OOD*: In the open-world setting, test inputs can have significantly different color schemes and layouts compared to ID samples. To mimic such scenarios, we use two public datasets as the *out-domain* OOD test sets. Specifically, NJU-Fudan Paper-Poster Dataset (Qiang et al., 2019) contains scientific posters in the digital PDF format, and we extract the document contents with ². CORD (Park et al., 2019) is a receipt understanding dataset that contains significantly different inputs from that in RVL-CDIP. As shown in Fig. 3, document images in CORD are receipt images without creases or warping, which requires the model to be capable of handling text information but also visual and spatial information.

In the following sections, we provide detailed analysis and share insights on various aspects of OOD detection performance for document understanding models under different OOD detection methods. Further details on the setup are provided in Appendix A.

3.1 ARE PRETRAINED MODELS SUFFICIENT FOR OOD DETECTION?

As shown in Sec. 2.1, most domain processing models deployed in the real world are pretrained on a large-scale dataset. Naturally, one may expect pre-trained models to be robust to OOD data when equipped with competitive OOD detection methods. To better understand the role of pretraining, we first provide more nuanced discussions on the following questions: 1) Are models equally robust to in-domain and out-domain OOD inputs? 2) How does model modality impact OOD detection performance?

We consider a wide range of models pretrained on pure-text/image data (*e.g.*, ImageNet, Wikipedia, *etc*). A detailed description of these models can be found in Appendix A.1.2. During finetuning, we combine the pretrained model with a classifier and finetune on RVL-CDIP (ID). For models before and after finetuning, we extract the final feature outputs as the feature embeddings and use the same KNN+ score (Sun et al., 2022) for OOD detection. The results are shown in Figure 4. We observe the following trends. First, finetuning largely improves OOD detection performance for both in-domain and out-domain OOD data. Pretrained models, despite the fact that they have “seen” a diverse collection of data during pretraining, do not yield sufficient OOD robustness. The same trend holds broadly across models with different modalities. Second, the improvement of finetuning is less significant for out-domain OOD data. For example, the AUROC on Receipt (out-domain OOD) for pretrained ViT model is 97.13, whereas finetuning only improves by 0.79%. This suggests that pretrained models do have the potential to separate data from different domains due to the diversity of data used for pretraining, while it remains hard for pretrained models to perform finer-grained separation for in-domain OOD inputs. Therefore, finetuning is beneficial for improving both types of OOD detection performance as a consequence of improved feature representation.

²<https://github.com/pymupdf/PyMuPDF>

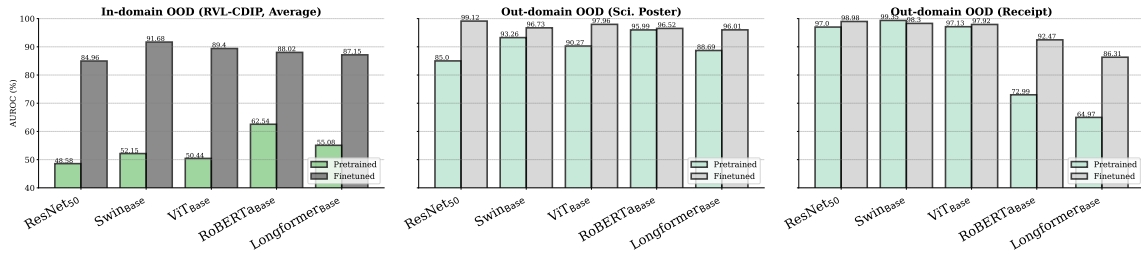


Figure 4: OOD detection performance for pretrained models w. and w.o. finetuning based on distance-based score KNN+(k=10). Finetuning significantly improves performance for both in and out-domain OOD.

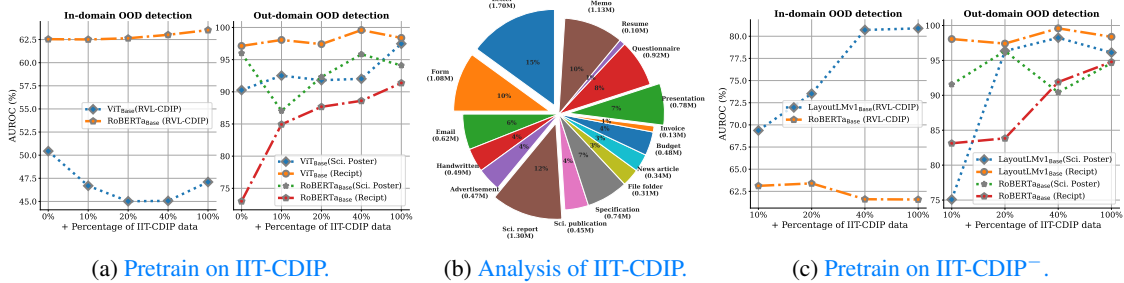


Figure 5: The impact of pretraining data on zero-shot OOD detection performance. IIT-CDIP⁻ denotes the filtered pretraining data after removing the “OOD” categories.

To make the analysis more thorough, we have two additional in-domain OOD settings: (1) select the classes the model performs well on, as in-domain OOD categories; (2) randomly select classes as OOD categories. As shown in Appendix (Table 10 and Table 11), we can see that finetuning also improves both types of OOD detection, which further reaffirm our conclusion. We also visualize the optimal transport dataset distance (OTDD) (Alvarez-Melis & Fusi, 2020) between in-domain and out-domain OOD datasets in Appendix (Fig. 10(b) and Fig. 10(c)). Please refer to the Appendix for more details.

3.2 THE IMPACT OF PRETRAINING DATA ON ZERO-SHOT OOD DETECTION

In the previous section, we analyze the impacts of finetuning for OOD detection where the pretraining dataset is fixed and unrelated to documents. Next, we dive deeper and study the impacts of pretraining dataset on zero-shot OOD detection. For each model, we adopt the same pretraining objective while adjusting the amount of pretraining data. Specifically, we increase the data diversity by appending 10, 20, 40, and 100% of randomly sampled data from IIT-CDIP dataset (around 11M) and pretrain each model. After pretraining, we measure OOD detection performance with KNN+ score based on feature embeddings.

For out-domain OOD data (Fig. 5, right), increasing the amount of pretraining data can significantly improve the zero-shot OOD detection performance (w.o. finetuning) for models across different modalities. This further verifies our previous hypothesis that pretraining with diverse data is beneficial for coarse-grained OOD detection, such as inputs from different domains (e.g., color schemes). On the other hand, for in-domain OOD inputs, even increasing the amount of pretraining data by over 40% provides negligible improvements (Fig. 5, left). This also suggests the necessity of finetuning for improving in-domain OOD detection.

We further explore zero-shot OOD detection by removing the potential OOD categories from IIT-CDIP. In practice, we first adopt the LayoutLMv1 finetuned on RVL-CDIP as the classifier for predicting labels for all IIT-CDIP document images. Fig. 5(b) shows the distribution of the predicted classes on IIT-CDIP. Next, we

remove the “*OOD*” categories from the IIT-CDIP data and pretrain two models (RoBERTa and LayoutLMv1) with 10, 20, 40, and 100% of randomly sampled data from the filtered IIT-CDIP. Fig. 5(c) shows our zero-shot OOD performance. Note that we do not show 0% in Fig. 5(c) since we pretrain LayoutLMv1 from scratch. For RoBERTa, we start from the public pretrained model and see a similar trend in Fig. 5(c) – the influence of pretraining for those well-pretrained language models is minor for in-domain OOD detection since there is a considerable gap between OCR words and pure-text data. *E.g.*, words in a document are spatially arranged, while words in text corpus are arranged sequentially. In contrast, pretraining data has a bigger impact on those models trained from scratch – the zero-shot performance of LayoutLMv1 increases when more pretraining data is added. We provide more details in the Appendix (Table 4 and Table 5).

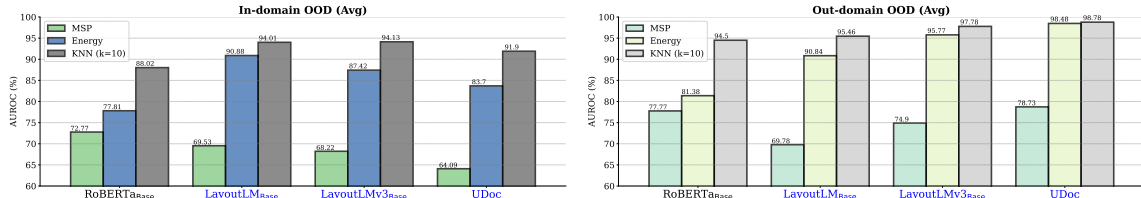


Figure 6: Comparison between representative feature-based scores and logit-based scores for spatial-aware and non-spatial-aware models. Spatial-aware models are colored in blue.

3.3 INVESTIGATING SPATIAL-AWARE MODELS FOR OOD DETECTION

In previous sections, we mainly focus on mainstream text-based and vision-based models to analyze the effects of finetuning and pretraining on in- and out-domain OOD detection. Next, we largely expand the scope of our study by incorporating models tailored to document processing, which we refer to as spatial-aware models, such as layoutLM, LayoutLMv3, and UDoc. Moreover, given finetuned models, we are able to compare the performance of logit-based scores and distance-based scores. Some key comparisons are shown in Figure 6. Please refer to the Appendix for full results.

Distance-based vs. logit-based score

We can see that simple KNN-based score outperforms logit-based scores for both in-domain and out-domain OOD data. The trend holds consistently across models with different modalities. Moreover, spatial-aware models demonstrate both stronger OOD detection performance for in- and out-domain OOD. For example, with the best scoring function (KNN+), compared to RoBERTa, LayoutLMv3 improves the average AUROC by 7.09% for out-domain OOD and 7.54%. The significant improvement suggests the value of spatial and visual information in improving OOD robustness for document data. **Note**

that despite this paper mainly comparing the logit-based and distance-based scores, we need to be aware that Gradient-based score has also been proposed for OOD detection. We also report the GradNorm (Huang et al., 2021) OOD detection score in Appendix A.3 and achieves similar performance as logit-based scores.

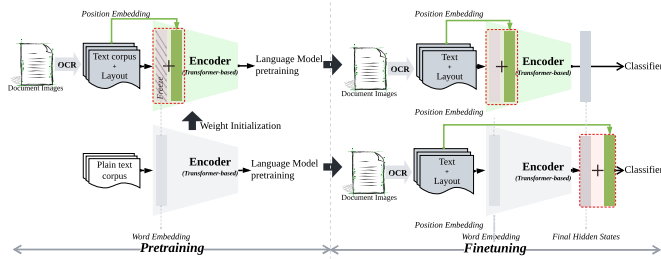


Figure 7: Illustration of our spatial-aware adapter design for language models. We have two adapter designs (marked in red box): (1) insert the adapter into the word embedding layer during pretraining and finetuning; (2) insert the adapter into the output layer for finetuning only. For the first design, we freeze the word embedding layer and learn the adapter and transformer layers.

3.4 TOWARDS SIMPLE AND EFFECTIVE SPATIAL-AWARE ADAPTORS

Spatial-aware models tailored for documents such as LayoutLM rely on spatial information and demonstrates superior OOD detection performance. This brings us a question: given the abundance of well-pretrained large-scale language models on text data such as RoBERTa, is there a simple and effective method that allows us to exploit the pretrained language model to document inputs for effective OOD detection? Next, we show that by enhancing transformer-based pretrained models with a spatial-aware adapter module, we can achieve good performance with minimal code edits.

Spatial-aware adapter Given a public pretrained RoBERTa model, depending on the position of the adaptor, we consider two architectures: 1) appending the adaptor to the word embedding layer, denoted as Spatial-RoBERTa (pre). It requires pretraining and finetuning, as illustrated in the top row in Fig. 7; 2) appending the adaptor to the final layer, denoted as Spatial-BoBERTa (post). As the model can utilize pretrained textual encoder, it only requires finetuning (illustrated in the bottom row). In the following, we only discuss Spatial-RoBERTa (pre). Full results for both Spatial-RoBERTa variants are in the Appendix.

We freeze the word embedding layer during pretraining for the following considerations: 1) word embeddings learned from large-scale corpus already cover most of those words from documents; 2) pretraining on documents without strong language dependency may not help improve word embeddings. For example, in semi-structured documents (*e.g.*, forms, receipts), language dependencies are not be as strong as rich-text documents (*e.g.*, letters, resumes), which may degenerate the learned word representations. In practice, each word has a normalized bounding box (x_0, y_0, x_1, y_1), where $(x_0, y_0) / (x_1, y_1)$ corresponds to the position of the upper left / lower right in the bounding box. Each coordinate is fed into an embedding layer and outputs a position embedding. All position embeddings are added to the initial word embedding to form a new spatial-aware embedding.

Spatial-RoBERTa significantly outperforms RoBERTa To verify the effectiveness of Spatial-RoBERTa, We compare OOD detection performance for pre-trained and fine-tuned models. The results are shown in Fig. 8. Spatial-RoBERTa significantly improves the OOD detection performance, especially after finetuning. Specifically, compared to RoBERTa, Spatial-RoBERTa improves AUROC by 9.20% for in-domain OOD and 4.95% for out-domain OOD data. This further verifies the importance of spatial-awareness for OOD detection in the document domain.

Spatial-RoBERTa is competitive for both ID classification and OOD detection Beyond OOD detection performance, we also examine ID classification accuracy and plot the two metrics for all the models with different modalities in Fig. 9. We find that there exists a positive correlation between ID accuracy and OOD detection performance (measured by AUROC) for both in-domain and out-domain OOD. Moreover, spatial-aware models display superior ID accuracy and OOD robustness compared to text-only and vision-only models. Finally, our Spatial-RoBERTa provides a simple and effective solution that greatly improves upon RoBERTa and matches the performance of models with specific architecture such as LayoutLM. Specifically, Spatial-RoBERTa_{Large} achieves 97.37 ID accuracy, which is even higher than LayoutLM (97.28) and UDoc (97.36). Furthermore, since our Spatial-RoBERTa_{Base} and Spatial-RoBERTa_{Large} freeze the word embed-

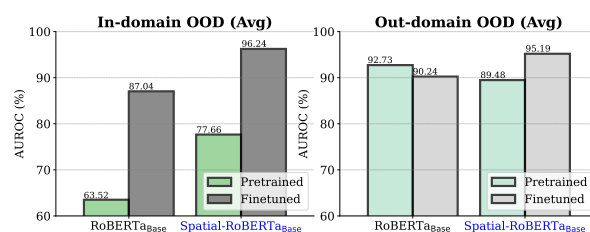


Figure 8: OOD detection performance between Spatial-RoBERTa and RoBERTa. All models are initialized with public pretrained checkpoints on pure-text data and further pretrained on IIT-CDIP with the same pretraining tasks. The only difference here is that Spatial-RoBERTa has an additional spatial-aware adapter and takes the word bounding boxes as the additional inputs.

Figure 8: OOD detection performance between Spatial-RoBERTa and RoBERTa. All models are initialized with public pretrained checkpoints on pure-text data and further pretrained on IIT-CDIP with the same pretraining tasks. The only difference here is that Spatial-RoBERTa has an additional spatial-aware adapter and takes the word bounding boxes as the additional inputs.

ding during pretraining, it can learn spatial-aware feature target document data while keeping the word embedding fixed, thus reduce the trainable model size and reduce the training cost.

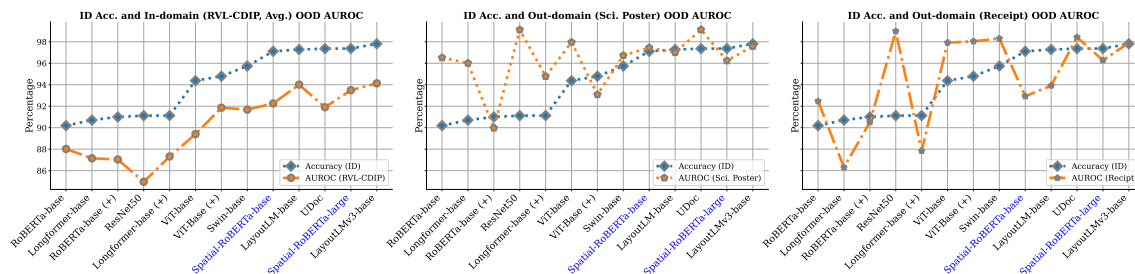


Figure 9: Correlation between ID accuracy and OOD detection performance. For most models, ID accuracy is positively correlated with OOD detection performance. Spatial-aware models display both higher ID accuracy and stronger OOD robustness (in AUROC).

4 CONCLUSION AND OUTLOOK

This paper presents a large-scale study of various methods for quantifying OOD robustness across different data modalities and models for document domains. Our key novelties include a large-scale investigation of OOD robustness in the document domain and a simple yet powerful spatial-aware adapter for transformer-based language models. We start from document classification and explore the pretrained models for document OOD robustness. A variety of substantial experiments in different settings demonstrates that pretraining datasets and tasks greatly impact OOD detection performance. Notably, OOD samples in the document domain are more accessible to identify in the feature space than in the logit space. Investigations from various perspectives explain certain intriguing phenomena and inspire more research on evaluating OOD robustness towards more reliable document understanding models.

REFERENCES

- David Alvarez-Melis and Nicolo Fusi. Geometric dataset distances via optimal transport. In *NeurIPS*, 2020.
- Srikar Appalaraju, Bhavan Jasani, Bhargava Urala Kota, Yusheng Xie, and R Manmatha. Docformer: End-to-end transformer for document understanding. In *ICCV*, 2021.
- Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*, 2020.
- Abhijit Bendale and Terrance E Boult. Towards open set deep networks. In *CVPR*, 2016.
- Lei Cui, Yiheng Xu, Tengchao Lv, and Furu Wei. Document ai: Benchmarks, models and applications. *arXiv preprint arXiv:2111.08609*, 2021.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*, 2018.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and

- Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021.
- Xuefeng Du, Zhaoning Wang, Mu Cai, and Yixuan Li. Vos: Learning what you don't know by virtual outlier synthesis. In *ICLR*, 2022.
- Stanislav Fort, Jie Ren, and Balaji Lakshminarayanan. Exploring the limits of out-of-distribution detection. In *NeurIPS*, 2021.
- ZongYuan Ge, Sergey Demyanov, Zetao Chen, and Rahil Garnavi. Generative openmax for multi-class open set classification. *arXiv preprint arXiv:1707.07418*, 2017.
- Jiuxiang Gu, Jason Kuen, Vlad I Morariu, Handong Zhao, Rajiv Jain, Nikolaos Barmpalios, Ani Nenkova, and Tong Sun. Unified pretraining framework for document understanding. In *NeurIPS*, 2021.
- Zhangxuan Gu, Changhua Meng, Ke Wang, Jun Lan, Weiqiang Wang, Ming Gu, and Liqing Zhang. Xylay-outlm: Towards layout-aware multimodal networks for visually-rich document understanding. In *CVPR*, 2022.
- Adam W Harley, Alex Ufkes, and Konstantinos G Derpanis. Evaluation of deep convolutional nets for document image classification and retrieval. In *ICDAR*, 2015.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *cvpr*, 2020.
- Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. In *ICLR*, 2017.
- Dan Hendrycks, Xiaoyuan Liu, Eric Wallace, Adam Dziedzic, Rishabh Krishnan, and Dawn Song. Pre-trained transformers improve out-of-distribution robustness. In *ACL*, 2020.
- Dan Hendrycks, Steven Basart, Mantas Mazeika, Mohammadreza Mostajabi, Jacob Steinhardt, and Dawn Song. Scaling out-of-distribution detection for real-world settings. In *ICML*, 2022.
- Teakgyu Hong, Donghyun Kim, Mingi Ji, Wonseok Hwang, Daehyun Nam, and Sungrae Park. Bros: A pre-trained language model focusing on text and layout for better key information extraction from documents. In *AAAI*, 2022.
- Yen-Chang Hsu, Yilin Shen, Hongxia Jin, and Zsolt Kira. Generalized odin: Detecting out-of-distribution image without learning from out-of-distribution data. In *CVPR*, 2020.
- Rui Huang and Yixuan Li. Mos: Towards scaling out-of-distribution detection for large semantic space. In *CVPR*, 2021.
- Rui Huang, Andrew Geng, and Yixuan Li. On the importance of gradients for detecting distributional shifts in the wild. In *NeurIPS*, 2021.
- Yupan Huang, Tengchao Lv, Lei Cui, Yutong Lu, and Furu Wei. Layoutlmv3: Pre-training for document ai with unified text and image masking. *arXiv preprint arXiv:2204.08387*, 2022.
- Guillaume Jaume, Hazim Kemal Ekenel, and Jean-Philippe Thiran. Funsd: A dataset for form understanding in noisy scanned documents. In *ICDAR Workshop*, 2019.

- Di Jin, Shuyang Gao, Seokhwan Kim, Yang Liu, and Dilek Hakkani-Tur. Towards textual out-of-domain detection without in-domain labels. *TASLP*, 2022.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2014.
- Stefan Larson, Gordon Lim, Yutong Ai, David Kuang, and Kevin Leach. Evaluating out-of-distribution performance on document image classifiers. In *NeurIPS*, 2022.
- Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. In *NeurIPS*, 2018.
- D. Lewis, G. Agam, S. Argamon, O. Frieder, D. Grossman, and J. Heard. Building a test collection for complex document information processing. In *SIGIR*, 2006.
- Gen Li, Nan Duan, Yuejian Fang, Daxin Jiang, and Ming Zhou. Unicoder-vl: A universal encoder for vision and language by cross-modal pre-training. In *AAAI*, 2020.
- Peizhao Li, Jiuxiang Gu, Jason Kuen, Vlad Morariu, Handong Zhao, Rajiv Jain, Varun Manjunatha, and Hongfu Liu. Selfdoc: Self-supervised document representation learning. In *CVPR*, 2021.
- Shiyu Liang, Yixuan Li, and Rayadurgam Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. In *ICLR*, 2018.
- Weitang Liu, Xiaoyun Wang, John Owens, and Yixuan Li. Energy-based out-of-distribution detection. In *NeurIPS*, 2020.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *NeurIPS*, 2019.
- Ajoy Mondal, Peter Lipps, and CV Jawahar. Iiit-ar-13k: a new dataset for graphical object detection in documents. In *International Workshop on Document Analysis Systems*, 2020.
- Peyman Morteza and Yixuan Li. Provable guarantees for understanding out-of-distribution detection. In *AAAI*, 2022.
- Eric Nalisnick, Akihiro Matsukawa, Yee Whye Teh, Dilan Gorur, and Balaji Lakshminarayanan. Do deep generative models know what they don’t know? In *ICLR*, 2019.
- Poojan Oza and Vishal M Patel. C2ae: Class conditioned auto-encoder for open-set recognition. In *CVPR*, 2019.
- Seunghyun Park, Seung Shin, Bado Lee, Junyeop Lee, Jaeheung Surh, Minjoon Seo, and Hwalsuk Lee. Cord: A consolidated receipt dataset for post-ocr parsing. In *NeurIPS Workshop*, 2019.
- Jonas Pfeiffer, Andreas Rücklé, Clifton Poth, Aishwarya Kamath, Ivan Vulić, Sebastian Ruder, Kyunghyun Cho, and Iryna Gurevych. Adapterhub: A framework for adapting transformers. *arXiv preprint arXiv:2007.07779*, 2020.
- Yu-Ting Qiang, Yan-Wei Fu, Xiao Yu, Yan-Wen Guo, Zhi-Hua Zhou, and Leonid Sigal. Learning to generate posters of scientific papers by probabilistic graphical models. *JCST*, 2019.

- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021.
- Jie Ren, Peter J Liu, Emily Fertig, Jasper Snoek, Ryan Poplin, Mark Depristo, Joshua Dillon, and Balaji Lakshminarayanan. Likelihood ratios for out-of-distribution detection. In *NeurIPS*, 2019.
- Chandramouli Shama Sastry and Sageev Oore. Detecting out-of-distribution examples with Gram matrices. In *International Conference on Machine Learning (ICML)*, 2020.
- Vikash Sehwal, Mung Chiang, and Prateek Mittal. Ssd: A unified framework for self-supervised outlier detection. In *ICLR*, 2021.
- Ray Smith. An overview of the tesseract ocr engine. In *ICDAR*, 2007.
- Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. Vi-bert: Pre-training of generic visual-linguistic representations. In *ICLR*, 2019.
- Yiyu Sun, Chuan Guo, and Yixuan Li. React: Out-of-distribution detection with rectified activations. In *NeurIPS*, 2021.
- Yiyu Sun, Yifei Ming, Xiaojin Zhu, and Yixuan Li. Out-of-distribution detection with deep nearest neighbors. In *ICML*, 2022.
- Thirumalaisamy P Velavan and Christian G Meyer. The covid-19 epidemic. *Tropical medicine & international health*, 25(3):278, 2020.
- Wenjin Wang, Zhengjie Huang, Bin Luo, Qianglong Chen, Qiming Peng, Yinxiu Pan, Weichong Yin, Shikun Feng, Yu Sun, Dianhai Yu, et al. mmlayout: Multi-grained multimodal transformer for document understanding. In *ACMMM*, 2022.
- Jim Winkens, Rudy Bunel, Abhijit Guha Roy, Robert Stanforth, Vivek Natarajan, Joseph R Ledsam, Patricia MacWilliams, Pushmeet Kohli, Alan Karthikesalingam, Simon Kohl, et al. Contrastive training for improved out-of-distribution detection. *arXiv preprint arXiv:2007.05566*, 2020.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Huggingface’s transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*, 2019.
- Yuxin Wu, Alexander Kirillov, Francisco Massa, Wan-Yen Lo, and Ross Girshick. Detectron2. <https://github.com/facebookresearch/detectron2>, 2019.
- Zhisheng Xiao, Qing Yan, and Yali Amit. Likelihood regret: An out-of-distribution detection score for variational auto-encoder. In *NeurIPS*, 2020.
- Keyang Xu, Tongzheng Ren, Shikun Zhang, Yihao Feng, and Caiming Xiong. Unsupervised out-of-domain detection via pre-trained transformers. In *ACL*, 2021a.
- Yang Xu, Yiheng Xu, Tengchao Lv, Lei Cui, Furu Wei, Guoxin Wang, Yijuan Lu, Dinei Florencio, Cha Zhang, Wanxiang Che, et al. Layoutlmv2: Multi-modal pre-training for visually-rich document understanding. In *ACL*, 2021b.
- Yiheng Xu, Minghao Li, Lei Cui, Shaohan Huang, Furu Wei, and Ming Zhou. Layoutlm: Pre-training of text and layout for document image understanding. In *SIGKDD*, 2020.

Xu Zhong, Jianbin Tang, and Antonio Jimeno Yepes. Publaynet: largest dataset ever for document layout analysis. In *ICDAR*, 2019.

Wenxuan Zhou, Fangyu Liu, and Muhao Chen. Contrastive out-of-distribution detection for pretrained transformers. *EMNLP*, 2021.

A APPENDIX

A.1 DOCUMENT CLASSIFICATION

A.1.1 DATASETS

RVL-CDIP consists of 320K/40K/40K training/validation/testing images under 16 categories. We select 12 categories and treat them as ID (In-Domain) data. We extract the text and layout information with Google OCR engine³ which provides both tokens and text blocks along with their corresponding bounding boxes. Most recent models take the full IIT-CDIP as pretraining data and finetuning on RVL-CDIP. However, it is not reasonable for the OOD setting since RVL-CDIP is a subset of IIT-CDIP. To make OOD results more reliable, in our experiments, we exclude the RVL-CDIP from the IIT-CDIP during pretraining.

We measure the distance between in-domain and out-domain datasets via OTDD⁴. We first visualize the OTDD distance between ID and the OOD data (in-domain and out-domain datasets in our main paper) in Fig. 10(a). During analysis, we sample the maximum number of 1000 images from each data and calculate the distance between datasets. It can be seen that there is a clear gap between in-domain data and out-domain data. To make the analysis more thorough, we have two additional in-domain OOD settings: (1) select the classes the model performs well as OOD data; (2) randomly select classes as OOD data. Fig. 10(a) and Fig. 10(c) show the dataset distance. As for the other two selection strategies, we can see that the domain gap is not as clear as the subset we selected for the main experiments. Interestingly, for those rare-word documents, such as file-folder and advertisements, the dataset distances are larger than documents with rich words. The background colors and layouts may yield a big distinction for those documents.

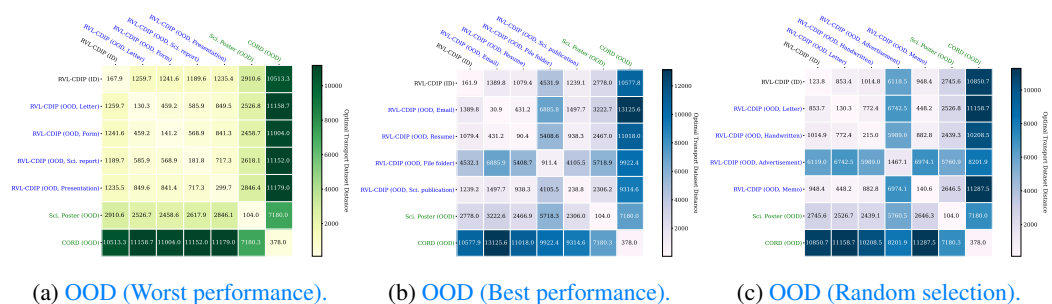


Figure 10: Visualization of optimal transport dataset distance for ID and OOD (in-domain and out-domain) datasets. We highlight the in-domain data in blue and the out-domain in green.

A.1.2 MODELS

All models reported in Fig. 2(b), except for UDoc, are initialized with model pretrained weights from Huggingface⁵ and finetune on the full RVL-CDIP training set. During finetuning, we train those models on RVL-CDIP with the cross-entropy loss. Models are optimized with Adam optimizer (Kingma & Ba, 2014) for 30 epochs with a batch size of 50 and a learning rate of 2×10^{-5} on 8 A100 GPUs. We list the hyperparameters of models used in our paper as follows:

³<https://cloud.google.com/vision/docs/ocr>

⁴<https://github.com/microsoft/otdd>

⁵<https://huggingface.co/models>

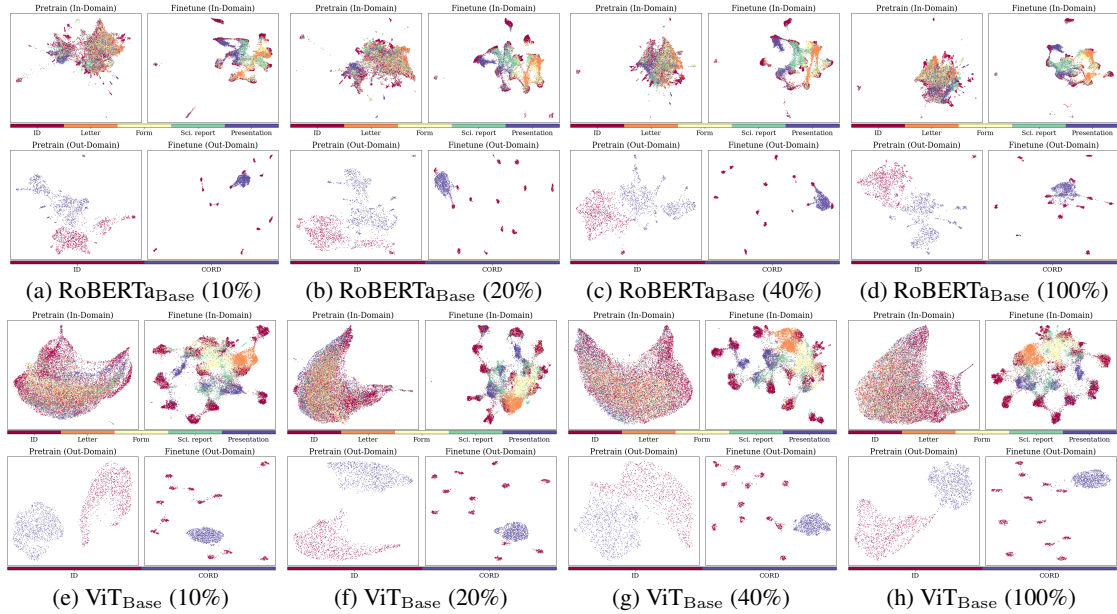


Figure 11: Feature visualization for pretrained (with different numbers of pretraining data) and finetuned models. We show both In-Domain (RVL-CDIP) and Out-Domain (CORD) OOD datasets.

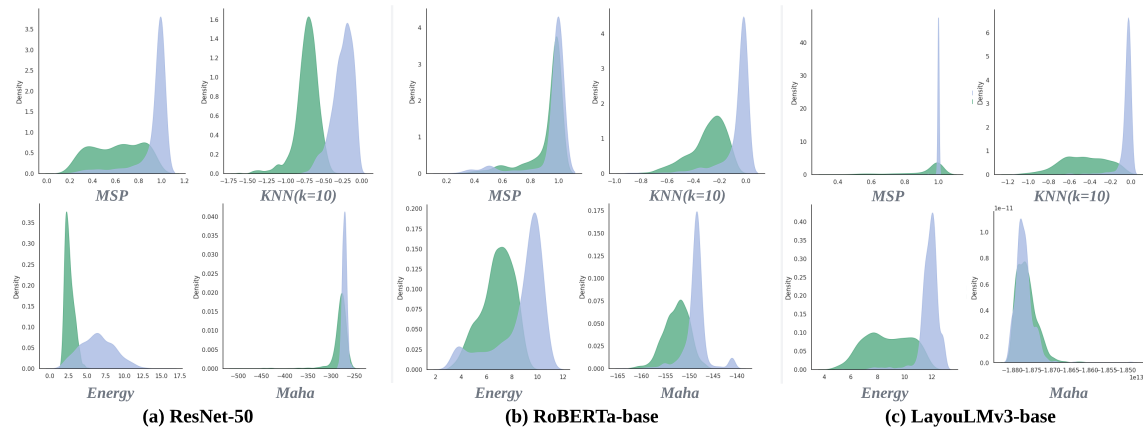


Figure 12: MSP, Energy, KNN, and Maha score histogram distributions of ID (blue) and OOD (green) inputs derived from finetuned ResNet-50, RoBERTa, and LayoutLMv3. The KNN scores calculated from both vision and language models naturally form smooth distributions. In contrast, MSP and Maha scores for both in- and out-of-distribution data concentrate on high values. Overall our experiments show that using feature space makes the scores more distinguishable between and out-of-distributions and, as a result, enables more effective OOD detection.

Language-only: (1) BERT and RoBERTa. We adopt the RoBERTa_{Base} (12 layer, 768 hidden size) and BERT_{Base} (12 layer, 768 hidden size) as the backbone and set the maximum sequence length to 512. For RoBERTa, the classifier is composed of two linear layers followed by a tanh activation function. (2) Longformer. We also adopts the Longformer_{Base} (12 layer, 768 hidden size) as the backbone and sets the maximum sequence length to 4,096.

Vision-only: (1) ResNet50: This model adopts the ResNet50 (pretrained on ImageNet-1k) as the backbone. We finetune it at resolution 224×224 . (2) ViT: This model adopts ViT_{Base} (vit-base-patch16-224, pretrained on ImageNet-21k) as the visual backbone. We finetune it at resolution 224×224 . (3) SwinB: This model adopt swin transformer (swin-base-patch4-window7-224-in22k, pretrained on ImageNet-21k) as backbone. We finetune it at resolution 224×224 .

Layout+Language: (1) LayoutLMv1: This model adopts the LayoutLM (layoutlm-base-uncased, 12 layer, 768 hidden size, pretrained on IIT-CDIP) as the backbone. We set the maximum sequence length to 512. (2) Spatial-RoBERTa_{Base} (Pre): This model is combines our spatial-aware adapter to the pretrained RoBERTa_{Base} model. The adapter is applied to the word embedding layer. During pretraining, we freeze the pretrained word embedding and optimize the spatial-aware adapter and transformers. During finetuning, we optimize all the parameters. (3) Spatial-RoBERTa_{Base} (Post): Instead of inserting the spatial-aware adapter in the input layer, this model combines the spatial-aware adapter at the transformers’ output layer.

Vision+Language+Layout: (1) LaytoulMv3: This model adopt the LayoutLMv3 (layoutlmv3-base, 12 layer, 768 hidden size, pretrained on IIT-CDIP) as the backbone. 2) UDoc: This model follows the design of UDoc. The only difference is the sentence encoder, we adopt a smaller version of the pretrained sentence encoder (all-MiniLM-L6-v2, 6 layer, 384 hidden size) instead of the sentence encoder (bert-base-nli-mean-tokens, 12 layer, 768 hidden size) used in their paper.

A.2 BEYOND DOCUMENT CLASSIFICATION

Going beyond document classification, we explore OOD detection for two entity-level tasks: document entity recognition and document object detection. Basic units such as text, tables, and figures in the document are the objects that need to be detected and recognized. Document entity recognition aims to predict the label for each semantic entity with given bounding boxes. Document object detection is an object detection task for document images. Specifically, we denote the input as x , the bounding box coordinates associated with object instances in the image as $\mathbf{b} \in \mathbb{R}^4$, and use the model with parameters θ to model the bounding box regression $p_{\theta}(\mathbf{b}|x)$ and the label classification $p_{\theta}(y|x, \mathbf{b})$. Given a test input \hat{x} , the OOD detection scoring function for entity detection and recognition can be unified as $S(\hat{x}, \hat{\mathbf{b}})$, where $\hat{\mathbf{b}}$ denotes the object instance predicted by the object detector. In particular, for document entity recognition, since the bounding boxes are provided, the OOD score can be simplified as $S(\hat{x}, \hat{\mathbf{b}})$, where $\hat{\mathbf{b}}$ is the given object instance.

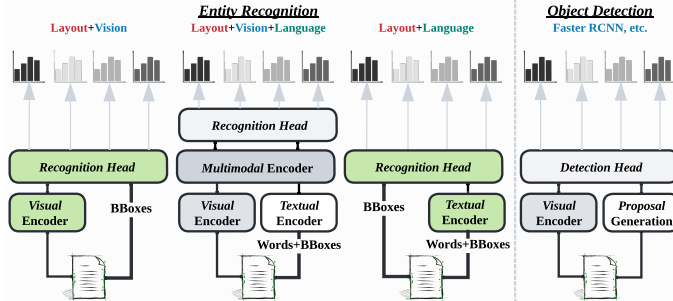


Figure 13: The network architectures in green blocks are our proposed models. We also show the modality information on top of each architecture.

A.2.1 DATASETS

Document Entity Recognition The original FUNSD (Jaume et al., 2019) dataset contains 149/50 training/testing images. we treat entities with category *other/header* as the OOD entities. After doing the split, if we treat *other* as OOD, we have a total number of 8,330/1,019 ID/OOD entities in total. Otherwise, if we treat *header* as OOD, we have 8,981/368 ID/OOD entities in total.

Document Object Detection PubLayNet (Zhong et al., 2019) contains 336K/11K training/validation images with 6 category labels (*text, title, list, figure, and table*). The original IIIT-AR-13K (Mondal et al., 2020) contains (*table, figure, natural image, logo, and signature*). In our paper, considering the overlap between IIIT-AR-13K and PubLayNet, we select those images that contain *Natural Image* as the OOD test set. After filtering, we have 2,880 OOD entities across 1,837 document images.

We consider three ID datasets in this experiment. (1) PubLayNet: This is the original PubLayNet dataset. We treat all the entities in training/validation images as ID entities. (2) Considering the domain shift between ID data (PubLayNet) and OOD data (IIIT-AR-13K). We combine the PubLayNet training data with the images from IIIT-AR-13K with overlapping annotations (*table and figure*) and train the object detection model.

A.2.2 MODELS

Document Entity Recognition Fig. 13 illustrates the entity recognition models used in this paper. We consider the entity on regions instead of tokens since regions contain richer semantic information. As for the pretrained model, we adopt UDoc (pretrained on IIT-CDIP) since it models the inputs at the region level. Based on UDoc framework, we develop the following models.

Vision/Vision+Layout: (1) ResNet-50: This model is composed of the ResNet-50 from pretrained UDoc. It adopts the RoI pooling followed by a classifier to extract the entity features. (2) ResNet-50+Position: This model also adapts UDoc pretrained ResNet-50. It further improves the RoI features to be spatial-aware by adding position embedding, where the position embeddings are mapped from bounding boxes via a linear mapping layer.

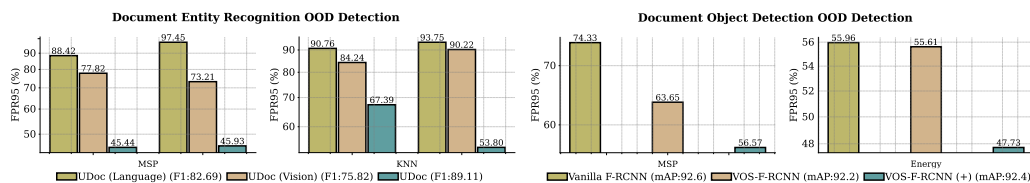
Language/Language+Layout: (1) Sentence BERT: This model adopts the language branch of UDoc and appends the classifier to the output of the sentence encoder. (2) Sentence BERT+Position: This model is close to the above model but adds position embedding to the sentence embeddings.

Vision+Language+Layout: (1) ResNet-50+sentence BERT: This model follows the same framework as UDoc, but replaces the sentence encoder in their original design with a more miniature sentence encoder (all-MiniLM-L6-v2). (2) SwinT+Sentence BERT: This model replaces the ResNet-50 visual backbone with a pretrained Swin tiny model (swin-tiny-patch4-window7-224) adopted from the Huggingface.

All the models are finetuned with cross-entropy loss for 100 epochs with a learning rate of 10^{-5} and batch size of 8 on one A100 GPU.

Document Object Detection Two object detection models are considered in this paper: (1) Vanilla Faster-RCNN: This model is the Faster-RCNN with ResNet-50 visual backbone. (2) Faster-RCNN with VOS⁶: This model enhances above model with VOS⁶. Following their paper, we use 1,000 samples for each ID class to estimate the class-conditional Gaussians. We train detection models with the Detectron2 framework (Wu et al., 2019). Models are trained for 180k iterations with a base learning rate of 0.01 and a batch size of 8. Mean average precision (MAP) @ intersection over union (IOU) [0.50:0.95] of bounding boxes is used to measure the performance.

⁶<https://github.com/deeplearning-wisc/vos>



(a) Comparison of OOD detection methods on different models on two OOD classes: *other* and *header*. (b) OOD detection results from different object detection methods.

Figure 14: Ablation on document entity recognition and object detection. Numbers are reported in FPR95.

A.2.3 EXPERIMENTAL SETUP

For document entity recognition, we construct ID and OOD datasets from FUNSD. Each semantic entity includes a list of words, a label, and a bounding box. The standard label set for this dataset contains 4 categories: *question*, *answer*, *header*, and *other*. In this paper, we select the entity labeled as *other* or *header* category as OOD. And the entities belonging to the other three categories are ID. Instead of treating entity recognition as a named-entity recognition problem, we follow UDoc and solve this problem at the semantic region level. We replace the sentence encoder in UDoc with a smaller sentence encoder (all-MiniLM-L6-v2⁷) from Huggingface (Wolf et al., 2019). We also have the following model variants to verify the effectiveness of each modality combination: textual-only, visual-only, textual+spatial, visual+spatial, and visual+textual+spatial.

For document object detection, we use PubLayNet as the ID dataset. We construct the OOD dataset from IIIT-AR-13K. Unlike PubLayNet, where the documents are scientific articles, IIIT-AR-13K is a dataset for graphical object detection in business documents (*e.g.*, annual reports). Hence, there exists an obvious domain gap between these two datasets. We select *natural images* as the OOD entity and filter images that contain the OOD entity. We first adopt the vanilla Faster RCNN with ResNet-50 backbone for document object detection as the baseline model. We also enhance Faster RCNN with VOS (Du et al., 2022), a recent unknown-aware learning framework to improve OOD detection performance for natural images.

A.2.4 OBSERVATIONS

To identify the entity type, models should not only understand the words but also require spatial and visual reasoning ability. We summarize our findings on document entity recognition in Fig. 14 (a) and describe them in more detail in Table 1. We can see that models can better predict the entity type with the help of the spatial position and also achieves better OOD robustness. Considering the weak language dependency between entities, it is not surprising that vision-based models achieve better performance than text-based models. We can see that UDoc with ResNet-50 achieves the best performance on two OOD test sets, illustrating that visual information plays a major role in increasing the discrimination of entities with similar semantics. We summarize our findings on document object detection in Fig. 14 (b) and describe them in more detail in Table 2. We can see that the OOD detection performance is further improved by introducing document images from IIIT-AR-13K with the same ID annotations as training data.

In Fig. 15, we visualize some document entity recognition OOD detection results. In Fig. 16, we visualize the prediction on sample OOD images, using object detection models trained without VOS (top) and with VOS (bottom), respectively. There is a clear difference between PubLayNet and IIIT-AR-13K – *natural image* annotations and entities rarely exist in PubLayNet. We can see that vanilla Faster RCNN trained on PubLayNet produces false positives when applied to the OOD document image from IIIT-AR-13K. After

⁷<https://huggingface.co/sentence-transformers>

introducing the unknown-aware learning method optimized for both ID and OOD, as shown in Table 2, the FPR95 reduces while preserving the mAP on the ID data. This experiment indicates that bringing uncertainty estimation into the entity detection training procedure can improve the reliability of the document object detection system.

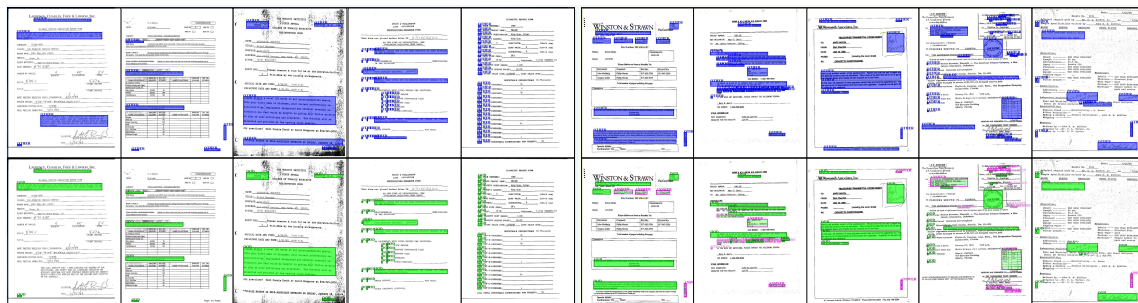


Figure 15: Visualization of detected OOD entities on the form images. The top part shows the entities in blue are entities annotated as *other*. The bottom part shows the detected OOD entities (green). We also show failure cases on the right part.

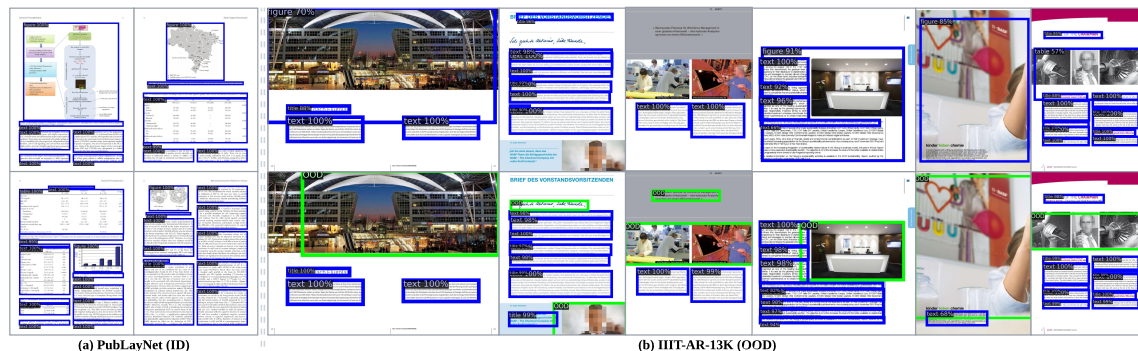


Figure 16: Visualization of detected objects on the OOD images (from IIIT-AR-13K) by a vanilla Faster-RCNN (top) and Faster-RCNN with VOS (bottom). Objects in blue boxes are detected and classified as one of the ID classes. The detected OOD objects (green) reduce false positives among detected objects. We also visualization of detected objects on the ID images. There is a clear difference between PubLayNet and IIIT-AR-13K – *natural image* annotations and entities rarely exist in PubLayNet.

A.3 ADDITIONAL EXPERIMENTAL RESULTS

- Table 1 corresponds to the results shown in Fig. 15 and Fig. 14(a).
- Table 2 corresponds to the results shown in Fig. 16 and Fig. 14(b).
- Table 3 and Table 7 correspond to the results shown in Fig. 5(a).
- Table 4 and Table 5 correspond to the results shown in Fig. 5(c).
- Table 6 corresponds to the results shown in Fig. 8 and Fig. 9.
- Table 9 and Table 8 correspond to the results shown in Fig. 4 and Fig. 9.
- Table 10 and Table 11 correspond to the analysis in Sec. 3.1.
- Table 12 corresponds to the results shown in Fig. 9.

Table 1: Comparison with different models on FUNSD OOD setting. All models are initialized with UDoc pretrained on IIT-CDIP and finetuned on FUNSD data with ID entities. All values are percentages. A lower FPR95 or higher AUROC value indicates better performance.

	Test F1	Method	Other (OOD)		ID F1	Header (OOD)		ID F1	Test F1	Method	Other (OOD)		ID F1	Header (OOD)		ID F1
			FPR95	AUROC		FPR95	AUROC				FPR95	AUROC		FPR95	AUROC	
ResNet-50	75.15	MahaNorm	88.42	59.21		92.12	41.49		75.82	MahaNorm	74.48	69.86		97.28	46.61	
		MahaUnNorm	94.11	29.14		99.46	24.06			MahaUnNorm	89.11	33.51		99.73	34.1	
		KNN10	59.47	79.14		81.79	63.97			KNN10	73.21	73.19		90.22	61.42	
		KNN20	69.97	78.15		81.25	63.66			KNN20	72.91	73.44		88.04	61.54	
		KNN50	84.49	77.40		82.61	62.86			KNN50	75.96	74.43		82.88	60.93	
		KNN100	97.94	77.08	77.65	84.24	61.62	78.04		KNN100	79.69	74.85	77.65	83.70	59.39	77.98
		KNN200	97.84	77.15		94.29	59.74			KNN200	86.06	75.14		91.58	57.42	
		KNN400	97.15	76.09		94.84	57.53			KNN400	87.93	74.92		95.92	55.37	
		MSP	50.54	75.80		75.82	76.55			MSP	77.82	67.60		84.24	66.58	
		MaxLogit	52.40	73.70		73.64	76.72			MaxLogit	76.94	67.05		84.24	65.41	
		Energy	52.50	73.70		75.82	76.55			Energy	76.64	66.93		84.51	64.98	
Sentence BERT	77.15	MahaNorm	93.33	54.99		88.32	67.06		82.69	MahaNorm	95.88	51.73		92.66	64.25	
		MahaUnNorm	93.33	55.05		88.59	67.67			MahaUnNorm	94.90	56.61		96.47	50.46	
		KNN10	93.72	48.44		92.66	60.99			KNN10	97.45	41.24		93.75	62.38	
		KNN20	93.92	47.65		92.93	59.00			KNN20	97.55	39.91		93.48	61.51	
		KNN50	93.62	48.94		93.21	57.90			KNN50	97.15	39.56		92.39	61.76	
		KNN100	93.92	48.79	82.12	93.21	55.07	82.41		KNN100	97.06	41.67	87.08	91.85	60.99	87.01
		KNN200	93.92	47.85		93.48	52.86			KNN200	96.57	41.85		89.67	59.08	
		KNN400	94.11	46.21		95.38	49.86			KNN400	97.25	40.83		90.22	54.03	
		MSP	93.62	54.91		94.29	52.14			MSP	88.42	61.11		90.76	59.58	
		MaxLogit	93.72	54.75		94.57	56.51			MaxLogit	89.70	60.19		88.86	60.92	
		Energy	93.23	54.88		93.21	58.22			Energy	90.48	59.61		89.95	61.12	
ResNet-50+Sentence BERT	89.11	MahaNorm	33.27	95.02		31.25	94.65		86.00	MahaNorm	59.57	88.56		78.80	77.33	
		MahaUnNorm	94.70	27.16		57.07	79.68			MahaUnNorm	75.56	58.18		92.66	56.17	
		KNN10	45.93	87.85		53.80	87.97			KNN10	63.30	83.64		81.52	64.08	
		KNN20	53.58	86.71		55.71	87.06			KNN20	66.73	82.53		81.52	61.50	
		KNN50	73.21	84.36		62.77	85.49			KNN50	70.17	80.21		82.34	57.77	
		KNN100	89.70	83.01	93.13	69.02	83.60	93.18		KNN100	83.91	77.71	90.82	83.15	54.97	90.40
		KNN200	96.66	81.90		75.54	80.85			KNN200	95.39	75.79		95.38	50.57	
		KNN400	98.82	81.00		91.58	77.42			KNN400	96.76	75.49		99.73	47.45	
		MSP	45.44	87.82		67.39	72.85			MSP	69.28	70.70		80.71	52.02	
		MaxLogit	45.53	90.58		63.04	72.39			MaxLogit	67.12	74.41		81.79	52.77	
		Energy	45.53	90.57		63.86	72.37			Energy	67.22	74.41		81.79	52.77	

Table 2: Comparison with different training and detection methods.

Models	ID	Method	IIIT-AR-13K (Natural Image as OOD)			PubLayNet (ID)
			FPR95	AUROC	AUPR	mAP
Vanilla Faster-RCNN	PubLayNet	MSP	74.33	79.12	98.41	92.6
		Energy	55.96	83.55	98.73	
Faster-RCNN with VOS	PubLayNet	MSP	63.65	79.37	98.57	92.2
		Energy	55.61	80.60	98.67	
Faster-RCNN with VOS	PubLayNet+IIIT-AR-13K(ID)	MSP	56.57	82.94	98.59	92.4
		Energy	47.73	84.04	98.67	

Table 8: OOD detection performance for document classification. Longformer₄₀₉₆ denotes the original model adopted from the Huggingface model hub. Longformer₄₀₉₆ (+) denotes the additional pretraining on IIT-CDIP.

ID Acc	Method	OOD Dataset (In-Domain)										OOD Dataset (Out-Domain)				
		Sci. Report		Presentation		Form		Letter		Average		Sci. Poster		Receipt		
		FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	
Finetune on RVL-CDIP (ID)																
Longformer ₄₀₉₆	90.71	MSP	95.00	64.32	95.62	62.17	95.89	60.53	93.95	66.89	95.12	63.48	88.37	77.50	98.60	54.72
		MaxLogit	97.12	72.84	97.07	75.22	98.24	70.39	95.82	77.57	97.06	74.00	90.70	86.62	99.60	68.10
		Energy	97.48	72.82	97.35	75.21	98.36	70.37	96.59	77.56	97.44	73.99	91.86	86.63	99.80	68.08
		Maha _{Norm}	58.21	88.86	64.93	87.16	69.19	84.31	57.55	89.54	62.47	87.47	19.77	97.02	78.50	87.44
		Maha _{inNorm}	61.21	85.11	67.01	81.93	71.99	78.19	63.19	81.94	65.85	81.79	24.42	94.13	76.80	86.47
		KNN ₁₀	58.45	88.21	65.65	86.88	67.80	83.99	56.78	89.53	62.17	87.15	27.91	96.01	82.10	86.31
		KNN ₂₀	58.97	88.04	65.57	86.60	68.12	83.80	57.35	89.34	62.50	86.94	29.07	95.82	82.60	85.93
		KNN ₅₀	60.25	87.64	66.57	86.25	68.91	83.41	58.81	88.96	63.64	86.56	30.23	95.46	82.70	85.27
		KNN ₁₀₀	61.97	87.19	68.14	85.81	70.15	82.95	60.47	88.60	65.18	86.14	34.88	95.04	82.80	84.75
		KNN ₂₀₀	64.29	86.71	69.43	85.40	71.79	82.42	62.74	88.31	67.06	85.71	43.02	94.63	83.60	84.27
		KNN ₄₀₀	66.33	86.11	71.03	84.87	73.82	81.64	64.94	88.07	69.03	85.17	45.35	94.27	84.20	83.50
No finetune																
-	-	KNN ₁₀	98.04	55.45	97.63	59.97	98.76	51.75	98.13	53.16	98.14	55.08	70.93	88.69	100.00	64.97
		KNN ₂₀	98.12	55.19	97.67	59.64	98.80	51.27	98.17	52.71	98.19	54.70	70.93	88.51	100.00	64.08
		KNN ₅₀	98.00	54.82	97.63	59.13	98.80	50.57	98.30	52.07	98.18	54.15	73.26	88.29	100.00	62.82
		KNN ₁₀₀	97.92	54.48	97.67	58.62	98.84	50.00	98.34	51.62	98.19	53.68	74.42	88.14	100.00	61.70
		KNN ₂₀₀	97.96	53.92	97.75	57.85	98.92	49.30	98.38	51.05	98.25	53.03	74.42	87.99	100.00	60.18
KNN ₄₀₀	97.84	53.04	97.71	56.77	98.92	48.33	98.42	50.27	98.22	52.10	70.93	87.86	100.00	57.80		
Pretrain on IIT-CDIP → Finetune on RVL-CDIP (ID)																
Longformer ₄₀₉₆ (+)	91.13	MSP	95.20	64.08	95.62	61.38	96.05	59.47	94.48	63.13	95.34	62.02	90.70	67.26	98.00	55.52
		MaxLogit	96.96	75.41	96.54	76.03	97.89	70.15	96.71	74.56	97.02	74.04	100.00	78.65	99.70	72.88
		Energy	97.28	75.40	96.54	76.03	98.28	70.14	97.16	74.55	97.32	74.03	100.00	78.59	99.70	72.86
		Maha _{Norm}	98.92	42.39	97.83	48.23	99.40	49.47	97.40	52.81	98.39	48.22	100.00	40.91	99.70	65.88
		Maha _{inNorm}	99.20	37.22	97.91	48.12	99.44	43.39	97.93	45.09	98.62	43.46	100.00	40.00	99.70	68.48
		KNN ₁₀	58.73	89.25	66.21	87.57	72.03	83.76	63.68	88.72	65.16	87.32	48.84	94.78	86.40	87.84
		KNN ₂₀	58.61	89.18	65.97	87.45	71.67	83.69	63.39	88.61	64.91	87.23	48.84	94.62	85.30	87.70
		KNN ₅₀	61.17	88.96	66.97	87.29	72.83	83.47	65.83	88.33	66.70	87.01	55.81	94.25	85.20	87.39
		KNN ₁₀₀	61.73	88.79	66.93	87.11	73.30	83.24	66.15	88.15	67.03	86.82	55.81	94.00	84.70	87.21
		KNN ₂₀₀	62.89	88.60	67.34	86.94	73.74	82.93	66.96	87.99	67.73	86.62	59.30	93.77	84.80	87.00
		KNN ₄₀₀	63.41	88.39	67.38	86.83	74.14	82.53	67.09	87.82	68.00	86.39	59.30	93.57	83.90	86.71
Pretrain on IIT-CDIP (no finetune)																
-	-	KNN ₁₀	95.48	61.40	98.07	53.66	97.73	55.55	98.66	48.70	97.49	54.83	81.40	91.12	97.40	46.27
		KNN ₂₀	95.56	60.92	97.95	52.95	97.49	54.97	98.50	48.21	97.38	54.26	84.88	90.62	97.50	45.55
		KNN ₅₀	95.60	59.94	97.95	51.77	97.41	53.97	98.62	47.29	97.40	53.24	87.21	89.95	98.20	44.18
		KNN ₁₀₀	95.60	59.04	97.99	50.74	97.21	52.99	98.58	46.51	97.34	52.32	88.37	89.52	98.50	43.09
		KNN ₂₀₀	95.68	58.13	98.03	49.68	97.45	52.01	98.58	45.79	97.44	51.40	89.53	88.84	98.50	41.98
		KNN ₄₀₀	96.00	57.52	98.15	48.88	97.65	51.31	98.70	45.45	97.62	50.79	91.86	88.46	98.60	41.02

Table 9: OOD detection performance for document classification. All models are pretrained on ImageNet.

ID Acc	Method	OOD Dataset (In-Domain)								OOD Dataset (Out-Domain)						
		Sci. Report		Presentation		Form		Letter		Average		Sci. Poster		Receipt		
		FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	
ResNet-50	Pretrain on ImageNet → Finetune on RVL-CDIP (ID)															
	91.12	MSP	64.49	87.87	55.89	90.94	66.60	87.31	77.88	80.87	66.22	86.75	51.16	92.76	63.10	90.36
		MaxLogit	64.89	88.59	47.97	92.81	65.40	87.52	77.56	81.87	63.96	87.70	41.86	94.62	54.00	93.29
		Energy	67.09	88.30	47.81	92.86	66.68	87.24	78.53	81.75	65.03	87.54	39.53	94.73	48.50	93.68
		MahaNorm	77.78	83.36	67.66	86.40	84.48	81.43	95.01	72.93	81.23	81.03	1.16	99.75	0.00	99.95
		MahaInNorm	97.16	39.28	93.09	43.28	94.37	50.89	98.99	38.08	95.90	42.88	50.00	82.75	55.50	76.69
		KNN10	73.38	86.82	67.98	87.46	71.31	87.84	92.90	77.74	76.39	84.96	6.98	99.12	5.20	98.98
	KNN20	74.90	86.41	66.29	87.79	73.82	87.21	93.95	76.51	77.24	84.48	6.98	98.96	5.50	98.85	
	KNN50	76.66	86.04	66.41	88.48	78.29	86.39	95.50	74.76	79.22	83.92	5.81	98.68	5.90	98.70	
	KNN100	77.54	85.61	65.41	88.99	82.16	85.43	96.23	73.37	80.33	83.35	6.98	98.34	6.30	98.51	
	KNN200	77.34	84.97	64.89	89.41	84.44	84.40	96.88	71.70	80.89	82.62	9.30	97.89	9.60	98.24	
	KNN400	77.98	84.01	64.48	89.61	87.11	83.04	97.93	69.24	81.88	81.48	15.12	97.29	12.20	97.81	
	Pretrain on ImageNet															
	-	KNN10	96.96	51.14	94.62	51.75	98.76	53.84	99.59	37.60	97.48	48.58	83.56	85.00	20.80	97.00
		KNN20	96.96	50.37	94.34	51.54	98.92	52.98	99.59	36.60	97.45	47.87	83.56	84.49	22.70	96.71
		KNN50	96.92	49.29	94.29	51.30	99.00	51.84	99.59	35.15	97.45	46.90	83.56	84.03	26.70	96.21
		KNN100	97.12	48.60	94.54	51.25	99.16	51.11	99.55	34.36	97.59	46.33	82.19	83.31	29.40	95.67
KNN200		97.48	47.85	94.58	51.12	99.08	50.41	99.68	33.44	97.70	45.70	83.56	82.24	33.30	94.92	
KNN400		97.76	47.00	95.18	50.95	99.04	49.57	99.68	32.27	97.92	44.95	83.56	80.92	40.10	93.86	
SwitchBase	Pretrain on ImageNet → Finetune on RVL-CDIP (ID)															
	95.74	MSP	47.64	88.09	49.90	88.11	58.22	83.14	50.28	88.90	51.51	87.06	49.32	91.31	36.50	93.63
		MaxLogit	42.39	93.11	42.47	93.45	58.62	88.79	45.90	93.18	47.34	92.13	50.68	92.50	32.20	95.65
		Energy	43.15	93.05	42.95	93.40	59.02	88.70	46.71	93.07	47.96	92.06	52.05	92.38	33.60	95.49
		MahaNorm	99.92	28.31	99.88	25.72	99.96	32.66	99.96	22.89	99.93	27.40	100.00	35.29	100.00	29.82
		MahaInNorm	99.96	28.98	100.00	26.25	99.96	33.32	100.00	23.69	99.98	28.06	100.00	35.87	100.00	30.43
		KNN10	49.44	92.82	46.73	92.87	42.90	92.57	72.69	88.45	52.94	91.68	16.44	96.73	6.10	98.30
	KNN20	48.84	92.95	43.27	93.51	44.53	92.32	72.28	88.35	52.23	91.78	17.81	96.52	7.40	98.10	
	KNN50	46.44	93.26	39.25	94.57	47.41	92.09	73.34	87.87	51.61	91.95	26.03	96.15	8.60	97.80	
	KNN100	43.76	93.42	35.03	95.29	50.08	91.72	75.77	87.42	51.16	91.96	28.77	95.94	11.30	97.55	
	KNN200	41.59	93.56	27.84	95.86	52.19	91.45	76.66	87.13	49.57	92.00	27.40	95.75	13.20	97.31	
	KNN400	40.79	93.68	24.71	96.14	53.27	91.17	77.52	86.97	49.07	91.99	26.03	95.72	16.00	97.18	
	Pretrain on ImageNet															
	-	KNN10	98.56	52.75	95.06	55.14	99.36	58.85	99.80	41.86	98.20	52.15	65.75	93.26	2.10	99.35
		KNN20	98.44	51.86	95.18	54.72	99.32	57.88	99.80	40.66	98.18	51.28	68.49	92.52	2.60	99.22
		KNN50	98.52	50.69	95.38	54.13	99.16	56.61	99.76	39.01	98.20	50.11	78.08	91.14	3.40	98.99
		KNN100	98.72	49.96	95.66	53.80	99.16	55.84	99.76	38.16	98.32	49.44	79.45	89.89	4.30	98.77
KNN200		98.68	49.22	95.90	53.39	99.20	55.14	99.76	37.31	98.38	48.76	84.93	88.40	5.40	98.47	
KNN400		98.84	48.36	96.22	52.74	99.40	54.41	99.76	36.26	98.56	47.94	84.93	86.74	8.20	98.05	
ViT-base	Pretrain on ImageNet → Finetune on RVL-CDIP (ID)															
	94.38	MSP	56.81	89.14	52.19	91.80	67.48	84.26	59.90	88.77	59.10	88.49	47.67	92.98	59.50	91.99
		MaxLogit	50.76	91.37	44.60	93.75	68.04	86.94	55.15	91.81	54.64	90.97	40.70	94.20	52.40	93.16
		Energy	51.16	91.31	44.52	93.75	69.43	86.81	56.09	91.77	55.30	90.91	38.37	94.11	53.20	93.11
		MahaNorm	90.63	70.10	91.84	65.75	89.55	70.83	97.81	57.37	92.46	66.01	100.00	58.34	82.20	77.09
		MahaInNorm	90.63	70.10	91.80	65.75	89.55	70.83	97.81	57.37	92.45	66.01	100.00	58.34	82.20	77.09
		KNN10	62.57	90.12	57.73	90.91	53.67	90.36	84.50	86.19	64.62	89.40	12.79	97.96	13.00	97.92
	KNN20	63.01	90.24	56.01	91.51	55.03	90.02	84.38	86.01	64.61	89.44	15.12	97.76	14.90	97.67	
	KNN50	61.97	90.62	53.23	92.62	58.26	89.57	84.25	85.64	64.43	89.61	16.28	97.38	19.80	97.24	
	KNN100	60.29	90.85	49.70	93.53	60.38	89.07	84.01	85.43	63.60	89.72	16.28	97.05	23.60	96.82	
	KNN200	58.45	91.04	45.04	94.36	62.89	88.54	84.17	85.33	62.64	89.82	22.09	96.75	27.50	96.39	
	KNN400	58.01	91.11	40.42	94.94	65.08	87.94	83.93	85.19	61.86	89.80	24.42	96.47	31.60	96.03	
	Pretrain on ImageNet															
	-	KNN10	98.48	52.15	95.02	56.94	99.48	53.77	99.47	38.90	98.11	50.44	93.15	90.27	20.40	97.13
		KNN20	98.48	51.41	95.06	56.61	99.44	52.92	99.55	37.61	98.13	49.64	94.52	89.44	22.60	96.80
		KNN50	98.32	50.43	94.86	56.21	99.40	51.86	99.59	35.82	98.04	48.58	97.26	88.23	26.60	96.25
		KNN100	98.40	49.76	95.06	55.90	99.44	51.15	99.59	34.59	98.12	47.85	98.63	87.24	31.20	95.76
KNN200		98.60	49.01	95.46	55.55	99.48	50.46	99.55	33.24	98.27	47.07	98.63	86.08	36.30	95.15	
KNN400		98.64	48.04	95.50	55.01	99.44	49.60	99.55	31.52	98.28	46.04	100.00	84.82	43.80	94.44	

Table 12: OOD detection performance for document classification. All models are pretrained on IIT-CDIP. For LayoutLM models, we adopt the checkpoints from the Huggingface model hub. For UDoc, we pretrain the model on our side. All models are finetuned on RVL-CDIP ID data.

ID	Method	OOD Dataset (In-Domain)										OOD Dataset (Out-Domain)			
		Sci. Report		Presentation		Form		Letter		Average		Sci. Poster		Receipt	
		FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC
LayoutLMv1Base	MSP	47.48	74.91	59.74	68.72	66.40	65.36	58.89	69.12	58.13	69.53	43.02	77.15	72.40	62.40
	MaxLogit	27.06	92.38	37.97	91.52	45.65	88.36	35.92	91.22	36.65	90.87	24.42	94.96	57.30	86.70
	Energy	27.06	92.40	37.97	91.54	45.65	88.36	35.92	91.23	36.65	90.88	24.42	94.97	57.30	86.70
	MahaNorm	17.73	96.67	33.83	94.27	36.95	92.47	24.55	95.34	28.26	94.69	13.95	97.49	43.60	94.74
	MahaUnNorm	22.86	90.33	37.61	83.70	41.42	82.77	36.16	83.43	34.51	85.06	20.93	90.95	25.80	93.06
	KNN10	20.82	96.09	35.32	93.82	40.06	91.34	28.65	94.80	31.21	94.01	17.44	97.00	49.80	93.92
	KNN20	21.74	95.93	36.20	93.77	41.42	91.12	30.44	94.61	32.45	93.86	17.44	96.82	51.70	93.73
	KNN50	24.34	95.56	38.25	93.41	43.93	90.69	33.64	94.19	35.04	93.46	23.26	96.44	53.80	93.70
	KNN100	25.54	95.30	39.13	93.20	45.17	90.35	34.78	93.99	36.16	93.21	25.58	96.24	54.70	93.45
	KNN200	26.54	95.04	39.53	92.95	46.21	89.97	35.75	93.80	37.01	92.94	25.58	96.04	56.80	93.10
KNN400	27.62	94.76	40.22	92.65	47.29	89.54	36.69	93.61	37.96	92.64	30.23	95.84	57.60	92.62	
LayoutLMv3	MSP	56.16	70.81	63.44	67.17	67.16	65.30	58.60	69.58	61.34	68.22	52.33	72.70	43.60	77.10
	MaxLogit	30.70	89.17	40.42	88.18	42.98	84.09	33.12	88.22	36.80	87.42	19.77	94.50	11.70	97.02
	Energy	30.70	89.18	40.42	88.18	42.98	84.10	33.12	88.23	36.80	87.42	19.77	94.51	11.70	97.03
	MahaNorm	99.16	16.28	98.51	35.13	99.12	20.45	99.19	14.42	99.00	21.57	100.00	7.79	98.50	42.00
	MahaUnNorm	99.76	15.30	99.16	34.20	99.52	19.77	99.92	12.86	99.59	20.53	100.00	7.57	98.30	42.14
	KNN10	21.74	95.03	35.68	93.38	32.88	91.86	18.51	96.26	27.20	94.13	11.63	97.58	8.90	97.97
	KNN20	22.74	94.90	36.56	93.20	33.96	91.66	19.64	96.15	28.22	93.98	12.79	97.44	10.00	97.89
	KNN50	24.62	94.62	38.37	92.71	35.83	91.38	21.63	95.93	30.11	93.66	13.95	97.20	10.70	97.72
	KNN100	25.22	94.38	39.29	92.32	36.55	91.09	22.48	95.79	30.88	93.40	16.28	97.04	11.80	97.59
	KNN200	26.02	94.13	39.82	91.91	37.19	90.78	23.30	95.68	31.58	93.12	18.60	96.91	12.80	97.46
KNN400	26.42	93.86	40.46	91.45	37.95	90.40	23.70	95.58	32.13	92.82	18.60	96.81	14.10	97.33	
UDoc _{Rec+Sci+L50}	MSP	66.13	65.73	69.43	64.09	71.03	63.28	71.06	63.25	69.41	64.09	40.70	78.47	39.80	78.99
	MaxLogit	45.96	82.12	47.21	86.39	49.64	83.16	49.59	83.13	48.10	83.70	2.33	98.57	4.00	98.34
	Energy	45.96	82.12	47.21	86.40	49.64	83.16	49.59	83.13	48.10	83.70	2.33	98.60	4.00	98.36
	MahaNorm	93.31	53.49	94.70	50.13	93.93	53.21	95.37	50.02	94.33	51.71	82.56	71.06	94.90	47.24
	MahaUnNorm	94.16	50.58	94.62	50.19	94.37	49.56	95.29	49.30	94.61	49.91	87.21	69.12	94.90	49.92
	KNN10	30.02	94.47	41.22	88.66	41.90	90.99	36.65	93.48	37.45	91.90	1.16	99.13	5.50	98.42
	KNN20	31.10	94.36	41.98	88.44	42.10	90.90	38.03	93.35	38.30	91.76	1.16	99.04	6.90	98.32
	KNN50	33.95	94.07	43.35	87.89	44.01	90.72	40.71	93.06	40.51	91.43	1.16	98.84	7.40	98.26
	KNN100	34.83	93.84	43.75	87.51	45.01	90.61	41.96	92.90	41.39	91.22	1.16	98.72	8.30	98.16
	KNN200	35.63	93.63	44.11	87.08	45.29	90.57	42.49	92.79	41.88	91.02	1.16	98.65	8.60	98.08
KNN400	36.39	93.29	44.80	86.43	45.65	90.65	42.94	92.60	42.44	90.74	1.16	98.61	9.10	98.04	