# Robust Text Classification: Analyzing Prototype-Based Networks

**Zhivar Sourati**[1,2]**, Darshan Deshpande**[1,2]**, Filip Ilievski**[3]**,**
**Kiril Gashteovski**[4,5] **and Sascha Saralajew**[4]

[1]Information Sciences Institute, University of Southern California, Marina del Rey, CA, USA
[2]Department of Computer Science, University of Southern California, Los Angeles, CA, USA
[3]Department of Computer Science, Vrije Universiteit Amsterdam, The Netherlands
[4]NEC Laboratories Europe, Heidelberg, Germany
[5]CAIR, Ss. Cyril and Methodius University, Skopje, North Macedonia

## Abstract

Downstream applications often require text classification models to be accurate, robust, and interpretable. While the accuracy of the state-of-the-art language models approximates human performance, they are not designed to be interpretable and often exhibit a drop in performance on noisy data. The family of Prototype-Based Networks (PBNs) that classify examples based on their similarity to prototypical examples of a class (prototypes) is natively interpretable and shown to be robust to noise, which enabled its wide usage for computer vision tasks. In this paper, we study whether the robustness properties of PBNs transfer to text classification tasks. We design a modular and comprehensive framework for studying PBNs, which includes different backbone architectures, backbone sizes, and objective functions. Our evaluation protocol assesses the robustness of models against character-, word-, and sentence-level perturbations. Our experiments on three benchmarks show that the robustness of PBNs transfers to NLP classification tasks facing realistic perturbations. Moreover, the robustness of PBNs is supported mostly by the objective function that keeps prototypes interpretable, while the robustness superiority of PBNs over vanilla models becomes more salient as datasets get more complex.

## 1 Introduction

In light of the needs of real-world applications, Natural Language Processing (NLP) research has increasingly focused on benchmarks, methods, and studies that emphasize robustness and interpretability (Zhou et al., 2020; Jang et al., 2022; Liu et al., 2021). The requirements for robustness and interpretability are intuitive for high-stake domains that affect health, racial bias, and safety (Rudin et al., 2022), where model errors can have serious consequences on human lives. More fundamentally, robustness and interpretability are essential components of developing trustworthy technology that can be adopted by experts in any domain (Wagstaff, 2012; Slack et al., 2022). However, the widely adopted pre-trained language models (PLMs) that report exceptional accuracy on NLP classification benchmarks (Chowdhery et al., 2022; Zoph et al., 2022) have limited interpretability by design, which cannot be fully mitigated by posthoc explainability techniques (Zini and Awad, 2022). Moreover, PLMs lack robustness when they are exposed to text perturbations, noisy data, or distribution shifts (Moradi and Samwald, 2021).

Meanwhile, the family of Prototype-Based Networks (PBNs; Li et al., 2018b) is designed for robustness and interpretability. PBNs are based on the theory of categorization in cognitive science (Rosch, 1973): categorization is governed by the graded degree of possessing a prototypical feature of that category, with some members being more central (*prototypical*) than others. Consider, for example, classifying different types of birds. Then, pelican classification can be done through their prototypical tall necks and similarity to a prototypical pelican (Nauta et al., 2021a). This classification approach is both interpretable and robust, because it classifies through distances to prototypical examples found in the data. Simple associations between sample points and central prototypical examples bring interpretability, while leveraging distance between points helps to quantify prototypicality, which then facilitates identifying noisy or out-of-distribution samples (Yang et al., 2018).

PBNs have been popular in computer vision (CV) tasks, including image classification (Angelov and Soares, 2020) and novel class detection (Hase et al., 2019). Inspired by PBNs in CV, NLP researchers have also developed PBN models for text classification, in particular, for sentiment classification (Pluciński et al., 2021; Ming et al., 2019; Hong et al., 2021), few-shot relation extraction (Han et al., 2021; Meng et al., 2023),
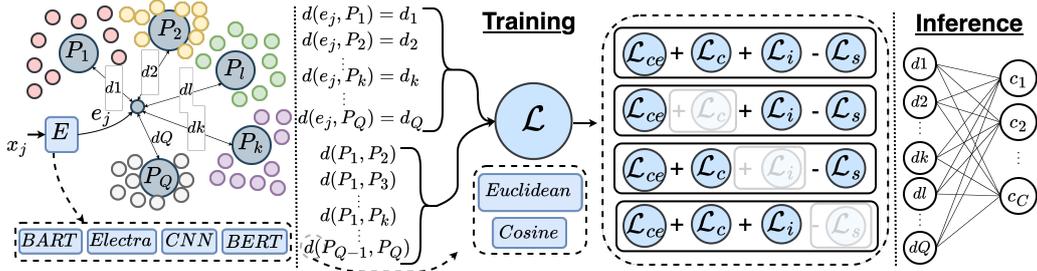
Figure 1: Classification by a PBN. The model computes distances between the new point and prototypes and distances within prototypes to do both inference and training. In training, the model minimizes the loss terms based on the computed distances between the new point and prototypes and within prototypes, and in inference, distances between the new point and prototypes are used for the classification by a fully connected layer. We test variations of the loss terms, different backbones, and distance functions and assess their effect on network's robustness.

and propaganda detection (Das et al., 2022). Yet, while competitive performance and interpretability of PBNs have been studied in both NLP (Das et al., 2022; Hase and Bansal, 2020) and vision (Gu and Ding, 2019; van Aken et al., 2022), their robustness advantages over non-PBNs have only been investigated in CV (Yang et al., 2018; Saralajew et al., 2020; Vorácek and Hein, 2022).

In this study, *we investigate whether the robustness properties of PBNs transfer to NLP classification tasks*. In particular, the contributions of this work are as follows: (1) We propose a **modular and comprehensive framework** to evaluate the robustness properties of PBNs that take into account state-of-the-art backbone architectures, different backbone sizes, and loss terms; (2) We devise an evaluation protocol that employs three **perturbation strategies** to evaluate the robustness of models against character-, word-, and sentence-level perturbations; (3) We perform **extensive experiments** on three benchmarks to compare the robustness of PBNs to non-PBN models, evaluate the effect of key PBN design choices, and assess their sensitivity to varying task complexity. Our experiments show that the robustness of PBNs transfers to realistic perturbations in text classification tasks, and it is impacted by the backbone architecture, the objective function, and the number of prototypes. Our experiments suggest that PBNs can enhance the text classification robustness of PLMs.

## 2   Prototype-Based Networks

PBNs classify data points based on their similarity to a set of *prototypes* learned during training. These prototypes summarize the prominent semantic patterns of the dataset through two mechanisms: (1) prototypes are defined in the same embedding

space as input examples, which makes them interpretable by leveraging input examples close to them; and (2) prototypes are designed to cluster semantically similar training examples, which makes them representative of the prominent patterns embedded in the data and input examples. The PBN's decisions are inherently interpretable, because prototypes are trained to be aligned with previous observations (Hong et al., 2020). This enables insights into the behavior of the model during inference by looking at the closest activated prototypes (Das et al., 2022). Prototypes being in the same embedding space as input examples allows them to be represented as either the training examples (Das et al., 2022) or parts of training examples, such as key phrases (Pluciński et al., 2021) or key sequences (Ming et al., 2019; Hong et al., 2021) extracted from training examples. These prototypes can be associated with semantic patterns of particular classes from their initialization or be trained freely and subsequently associated with the prominent semantic patterns of the whole dataset.

**Inference.**  As shown in Figure 1, given a set of data points $x_j, j \in \{1, \ldots, N\}$ with labels $y_j \in \{1, \ldots, C\}$, and $Q$ prototypes, PBNs first encode examples with a backbone $E$, resulting in the embedding $e_j = E(x_j)$. Next, PBNs compute the distances between prototypes and $e_j$ using the function $d$. These distances get fed into a linear layer to compute class-wise logits by incorporating the similarities to each prototype. Applying a softmax on top of logits, the final outputs are $\hat{y}_c(x_j)$: a probability that $x_j$ belongs to class $c \in \{1, \ldots, C\}$.

**Training.**   To compute a total loss term $\mathcal{L}$, PBNs use the computed distances within prototypes $d(P_k, P_l)_{k \neq l}$, distances between all $Q$ prototypes and $N$ training examples given by

$d(e_j, P_k)_{j\in\{1,...,N\};k\in\{1,...,Q\}}$, and the computed probabilities $\hat{y}_c$. The prototypes and the weights in the backbone are adjusted according to $\mathcal{L}$. The total loss $\mathcal{L}$ consists of different inner loss terms that ensure high accuracy, high interpretability, and low redundancy among prototypes; i.e., the classification loss $\mathcal{L}_{ce}$, the clustering loss $\mathcal{L}_c$ (Li et al., 2018b), the interpretability loss $\mathcal{L}_i$ (Li et al., 2018b), and separation loss $\mathcal{L}_s$ (Hong et al., 2020):

$$\mathcal{L} = \mathcal{L}_{ce} + \lambda_c\mathcal{L}_c + \lambda_i\mathcal{L}_i - \lambda_s\mathcal{L}_s, \qquad (1)$$

where $\lambda_c, \lambda_i, \lambda_s \geq 0$ are regularization factors to adjust the contribution of the auxiliary loss terms.

*Classification loss* $\mathcal{L}_{ce}$ is defined as the cross-entropy loss between predicted and true labels:

$$\mathcal{L}_{ce} = -\sum_{j=1}^{N} \log(\hat{y}_{y_j}(x_j)). \qquad (2)$$

*Clustering loss* $\mathcal{L}_c$ ensures that the training examples close to each prototype form a cluster of similar examples. In practice, $\mathcal{L}_c$ keeps all the training examples as close as possible to at least one prototype and minimizes the distance between training examples and their closest prototypes:

$$\mathcal{L}_c = \frac{1}{N}\sum_{j=1}^{N} \min_{k\in\{1,...,Q\}} d(P_k, x_j). \qquad (3)$$

*Interpretability loss* $\mathcal{L}_i$ ensures that the prototypes are interpretable by minimizing the distance to their closest training sample:

$$\mathcal{L}_i = \frac{1}{Q}\sum_{k=1}^{Q} \min_{j\in\{1,...,N\}} d(P_k, x_j). \qquad (4)$$

Keeping the prototypes close to training samples allows PBNs to represent a prototype by its closest training samples that are domain-independent and make prototypes understandable by task experts.

*Separation loss* $\mathcal{L}_s$ maximizes the inter-prototype distance to reduce the probability of redundant prototypes:

$$\mathcal{L}_s = \frac{2}{Q(Q-1)} \sum_{k,l\in\{1,...,Q\};k\neq l;} d(P_k, P_l). \qquad (5)$$

## 3 PBN Evaluation Framework

The robustness of PBNs may be affected by architectural choices in the design and implementation of PBNs, including the backbone encoder $E$, the distance function $d$, the number of prototypes $Q$, and the regularization factors of the objective functions in $\mathcal{L}$. Inspired by prior work that studied the impact of some of these choices in computer vision (Yang et al., 2018), we design a framework

to systematically investigate their impact on the robustness of PBNs in text classification tasks.

**Backbone** ($E$). Prototype alignment and training are highly dependent on the quality of the latent space created by the backbone encoder $E$, which in turn affects the performance, robustness, and interpretability of PBNs. We consolidate previous methods for text classification using PBNs (Pluciński et al., 2021; Das et al., 2022; Ming et al., 2019; Hong et al., 2020) and consider four backbone architectures: CNN, BERT (Devlin et al., 2018), BART encoder (Lewis et al., 2019), and Electra (Clark et al., 2020). Besides including PLMs due to their general language abilities (Min et al., 2021), we experiment with CNN as a backbone for two reasons. First, the network is more interpretable since embeddings extracted from shallow CNN-based networks can be mapped with n-grams. Second, they are faster and smaller compared to Transformers (Vaswani et al., 2017), which makes them a better option for models deployed on production systems (Pluciński et al., 2021). We also include models with a different number of parameters to analyze the effect of scaling up the backbones.

**Distance function** ($d$). The pairwise distance calculation quantifies how closely the prototypes are aligned to the training examples (Figure 1). In recent work, Euclidean distance has been shown to be better than Cosine distance for similarity calculation (van Aken et al., 2022; Snell et al., 2017) as it helps to align prototypes closer to the training examples in the latent space of the encoder's output. However, given uncomparable design choices, alongside some utilizing Cosine distance (Chen et al., 2019) while others prioritizing Euclidean distance (Mettes et al., 2019), the comparison would be difficult and the choice of distance function is usually treated as a hyperparameter. Accordingly, we hypothesize that the impact of $d$ will be significant in our study of robustness, and hence, our framework considers both Cosine and Euclidean distance functions when training PBNs.

**Number of prototypes** ($Q$). The number of prototypes $Q$ in PBNs is a key factor for mapping difficult data distributions (Yang et al., 2018; Sourati et al., 2023). Hence, we compare the effect with three values for $Q$: number of classes in the dataset, 16, and 64, covering a small, medium, and large number of prototype parameters.

**Objective functions** ($\mathcal{L}$). Given the partly complementary goals of the four loss terms, we in-

vestigate whether all are necessary for training an accurate and robust model. However, keeping the accuracy constraint ($\mathcal{L}_{ce}$) intact, to assess the effect of clustering and interpretability on the robustness of PBNs, we train the following model variations: a model employing all loss terms, a model dropping the interpretability loss while keeping other loss terms intact, and a model that omits clustering loss from training (see Figure 1). Moreover, to assess how the segregation within prototypes affects PBN's robustness, we train PBNs using five values for $\lambda_s \in \{0, 2, 4, 10, 20\}$, with lower values representing higher tolerance for prototypes being close to each other. We chose these values to cover small, medium, and large tolerance.

## 4 Robustness Evaluation Protocol

Text classification models are often misled by perturbations that differ from their original version in an imperceptible way to the human eye (Dalvi et al., 2004; Kurakin et al., 2017a,b). We analyze the robustness of PBNs in our framework by designing perturbations that keep the label unchanged, preserve the meaning of the original example, and maintain fluency as formalized by Jin et al. (2020). We consider perturbations that explore vulnerabilities of models in three levels: character-, word-, and sentence-level (see examples in Appendix A).

**Character-level perturbations.** In many text classification applications, specific keywords play an important role in the model's prediction. However, they might be unintentionally disguised and mislead the model, e. g., a hateful content detector might be misled by the misspelling of the word *Women* as *Wmen*. To assess the robustness of the models against such typos, we use TextBugger (Li et al., 2018a) as an effective representative for typo-based perturbations (Wang et al., 2022a). This character-level perturbation identifies the most important words in a text and then performs substitution, insertion, or deletion of characters in words (e. g., Citrix → Citri×).

**Word-level perturbations.** Machine learning models may learn spurious word correlations rather than how to solve a task (Wang and Culotta, 2020). To analyze whether models are sensitive to word choices, we use a strategy that modifies words in a text, while keeping the meaning unchanged. We use TextFooler (Jin et al., 2020), which strongly affects PLMs (Wang et al., 2022a). It perturbs text by replacing words with their synonyms according to

their distance in the embedding space (e. g., film → cinematographers). Similar to TextBugger, all important words are identified, and embedding-based transformations are applied on them.

**Sentence-level perturbations.** To assess if the models learn the underlying semantics behind the sentence instead of using certain phrases as clues, we employ paraphrasing as a sentence-level perturbation strategy. We obtain sentence-level perturbations by prompting GPT3.5 (OpenAI, 2022).

## 5 Experimental Setup

To assess the robustness of PBNs to real-world perturbations, our experimental framework tracks the F1 scores of the models on three text classification datasets. We train each model on the original training set without any perturbation or adversarial training (Goodfellow et al., 2014) and test it on both the original test examples and their perturbed versions. Please see Appendix B for further details.

**Datasets.** PBNs classify instances based on their similarity to prototypes learned during training that summarize prominent semantic patterns in a dataset. Thus, with more classes, we might need more prototypes to govern the more complex system between instances and prototypes (Yang et al., 2018). To study the interplay between the number of classes and robustness, we employ three datasets: (1) *IMDB reviews*:[1] a binary sentiment classification dataset; (2) *AG_NEWS* (Zhang et al., 2015): a collection of news articles that can be associated with four categories: world, sports, business, and science; (3) *DBPedia*:[2] a dataset with taxonomic, hierarchical categories for Wikipedia articles, with 9 classes: Agent, Work, Place, Species, UnitOfWork, Event, SportsSeason, Device, and TopicalConcept. Moreover, we adopt the SST-2 binary classification split from the existing *Adversarial GLUE (AdvGLUE)* dataset (Wang et al., 2022a). The AdvGLUE SST-2 benchmark consists of 131 examples perturbed using various word- and sentence-level perturbations. For statistics of the datasets and their perturbations, see Appendix A.

**Perturbations.** Perturbations are designed to simulate real noisy data that is generalizable to any *victim model* (model facing perturbations). Generalizable noisy data can affect any model regardless of the architecture and also without having

---

[1] https://ai.stanford.edu/~amaas/data/sentiment
[2] https://www.huggingface.co/datasets/DeveloperOats/DBPedia_Classes

access to the model to directly attack it. We generate character- and word-level perturbations for IMDB, AG_News, and DBPedia, following Wang et al. (2022a). For a given dataset, we randomly perturb examples from its original test split until we obtain at most 800 successful perturbations. *Successful* perturbations of a dataset are those that change the prediction of a victim model already fine-tuned on that dataset from the correct prediction to the wrong prediction. We consider three fine-tuned victim models: BERT (Devlin et al., 2018), RoBERTa (Liu et al., 2019), and DistilBERT (Sanh et al., 2019). We preserve examples whose perturbations are predicted wrongly by three models to ensure that the perturbations are not model-specific and are difficult for multiple models. In principle, the perturbations for each model are different, yielding three variations per original example for a dataset-perturbation pair. For instance, this procedure successfully creates three character-level perturbations of 461 original examples from the AG_News dataset, resulting in a total of 1,383 ($3 \times 461$) data points in its character-level test set. To avoid introducing additional noise to the data for the sentence-level perturbations, we randomly select 1,000 samples from the testing set of each dataset and paraphrase them with GPT3.5 using the prompt: *Paraphrase the paragraph below concisely and accurately: {input text}*. In our main experiments, we design generalizable perturbations that simulate noisy data and test robustness without assuming access to any specific model to directly attack it. In an auxiliary experiment, we assume to have access to the victim model and we attack PBNs directly.

**Baselines.** We compare the performance of our PBN framework variants with their vanilla counterparts, i.e., CNN, BERT, BART encoder, and Electra. The only difference between PBNs and their non-PBN counterparts is the model wrapped around their backbones: while PBNs compute distances to prototypes and classify data points based on these distances, non-PBNs leverage a fully connected layer for classification instead.

**Implementation details.** In our experiments, we vary the PBN framework dimensions (see Section 3) and fix other implementation decisions. All prototypes are initialized randomly for a fair comparison, and PLM backbones are also trainable. The prototypes are trained without being constrained to a certain class from the beginning,

and their corresponding class can be identified after training. The transformation from distances to class logits is done through a simple linear layer without intercept to avoid introducing additional complexity and keep the prediction interpretable through prototype distances. Apart from CNN, which was trained from scratch, both the backbone of PBNs and their non-PBN counterparts leveraged the same PLM and were fine-tuned separately to show the difference that is only attributed to the models' architecture. Focusing on the BERT-based baselines, since BERT-base is one of the models from which we extract perturbations by directly attacking it, to ensure generalization of the experiments on different backbones, we use BERT-Medium (Turc et al., 2019), which being smaller than other Transformer baselines (BART and Electra), allows us to assess the sensitivity of PBNs to the backbone size too.

## 6 Results

This section compares the robustness of PBNs to baselines and in relation to task complexity. Moreover, it assesses the impact of the design choices within our framework from Section 3.

**Robustness of PBNs.** We found that PBNs are consistently more robust to perturbations compared to their non-PBN counterparts (see Table 1). We observed this trend both for our custom perturbations and for the existing benchmark, AdvGLUE's SST-2 (Wang et al., 2022a). Character-level and word-level perturbations lead to a larger robustness gap between PBNs and non-PBNs, while both PBN and non-PBNs are robust against sentence-level perturbations.[3] PBNs' performance on perturbations improve by up to 24%, 33%, and 30%, with BERT, BART, and Electra, respectively, relative to their non-PBN variants. We note that the perturbations are general enough to be more challenging for all models, including PBNs, as both character- and word-level perturbations decrease the average model performance by 20% compared to the scores on the original test sets. Finally, we note that PBNs perform slightly worse than non-PBNs on original datasets, and we attribute this to PBNs being optimized to satisfy multiple objectives, including interpretability and clustering, whereas non-PBNs are only optimized to have high accuracy.

**Effect of task complexity.** Both PBNs and non-

---

[3]As we observe a low effect from sentence-level perturbations across all models, we focus on character- and word-level perturbations for brevity.

| | IMDB | | | | AG_News | | | | DBPedia | | | | SST-2 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Model** | **Orig** | **Char** | **Word** | **Sent** | **Orig** | **Char** | **Word** | **Sent** | **Orig** | **Char** | **Word** | **Sent** | **Orig** | **Adv** |
| CNN | 85.7 | **75.5** | **76.9** | **86.2** | 90.4 | **58.9** | **67.7** | 75.9 | 97.9 | **64.0** | **63.9** | **94.4** | 75.0 | 48.0 |
| +PBN | 82.6 | 73.4 | 72.5 | 83.2 | 80.1 | 48.5 | 57.2 | **76.3** | 85.0 | 50.6 | 49.7 | 83.1 | 75.0 | **49.0** |
| Δ | -3.1 | -2.1 | -4.4 | -3.0 | -10.3 | -10.4 | -10.5 | +0.4 | -12.9 | -13.3 | -14.2 | -11.3 | 00.0 | +1.0 |
| BERT | 94.7 | **84.2** | 85.5 | **92.2** | 93.9 | 71.1 | 78.8 | 88.3 | 98.4 | 66.2 | 60.5 | **98.0** | 83.9 | 40.9 |
| +PBN | 90.5 | 83.5 | **85.6** | 85.9 | 92.1 | **78.4** | **80.2** | **88.5** | 96.5 | **69.0** | **75.5** | 97.4 | 77.8 | **46.3** |
| Δ | -4.2 | -0.7 | +0.1 | -6.3 | -1.8 | +7.3 | +1.4 | +0.2 | -1.9 | +2.8 | +15.0 | -0.6 | -6.1 | +5.4 |
| BART | 98.1 | 89.3 | 92.1 | 94.9 | 97.4 | 74.2 | 79.2 | **90.1** | 96.3 | 69.5 | 68.8 | 97.0 | 93.1 | 29.7 |
| +PBN | 96.9 | **89.5** | **93.4** | **95.0** | 93.7 | **75.6** | **81.2** | 88.3 | 93.6 | **72.9** | **71.5** | **97.4** | 90.0 | **39.6** |
| Δ | -1.2 | +0.2 | +1.3 | +0.1 | -3.7 | +1.4 | +2.0 | -1.8 | -2.7 | +3.4 | +2.7 | +0.4 | -3.1 | +9.9 |
| ELEC. | 98.4 | **92.8** | 94.0 | **94.5** | 95.0 | 71.8 | 78.7 | **88.7** | 93.7 | 61.8 | 65.0 | 94.5 | 87.6 | 43.5 |
| +PBN | 94.9 | 91.9 | **94.2** | 91.8 | 90.7 | **73.9** | **80.4** | 82.8 | 99.1 | **76.7** | **70.5** | **98.4** | 98.5 | **56.6** |
| Δ | -3.5 | -0.9 | +0.2 | -2.7 | -4.3 | +2.1 | +1.7 | -5.9 | +5.4 | +14.9 | +5.5 | +3.9 | +10.9 | +13.1 |

Table 1: F1 scores of PBNs and their non-PBN counterparts on three datasets (also, on the SST-2 dataset and its perturbed version) and three perturbation types: character-level (Char), word-level (Word), and sentence-level (Sent). The **boldfaced** numbers indicate top performance within a particular test-split of PBN vs. non-PBN model; Δ's indicate the difference between PBN vs. non-PBN counterpart. Performance on the original unperturbed examples (Orig) is shown as the original scores on main test splits because of the low variance across different attack splits.

PBN methods perform worse on more complicated datasets with more classes. Meanwhile, the superior robustness of PBNs compared to their non-PBN counterparts gets more pronounced as datasets get more complicated. While the performance of PBNs and non-PBNs is on par for the IMDB dataset, their gap widens on the datasets with more classes: AG_News and DBPedia. We believe that this robust behavior is due to the design of the PBN architecture. Standard neural networks for text classification distinguish classes by drawing hyperplanes between samples of different classes that are prone to noise (Yang et al., 2018), especially when dealing with several classes. Instead, PBNs are inherently more robust since they perform classification based on the similarity of data points to prototypes, acting as class centroids.

**Effect of backbone design** ($E$). The performance of PBNs is sensitive to the choice of the backbone architecture (Table 1). Transformer-based architectures (like BERT, BART, and Electra) yield higher absolute F1 scores compared to their vanilla backbones with differences ($\Delta$) of up to 15%. Yet, this trend cannot be seen when using CNN as a backbone, as it performs worse than vanilla CNN. Comparing the results gathered from models with different sizes and also different embedding properties (Transformers capturing context better), we attribute the disparity to the embedding space properties of the backbone and consider a strong backbone more favorable to the robustness of PBNs. In a further experiment, we measure the impact of the size of the backbone on the robustness properties of PBNs (see Figure 2). The results show a boost
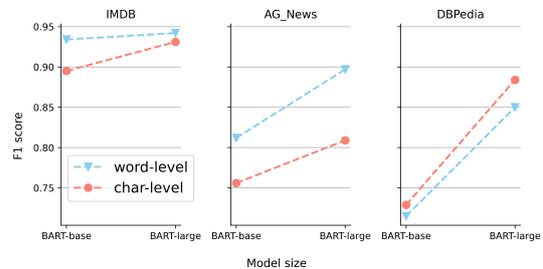


Figure 2: F1 of PBN with BART-base and BART-large.

in performance when scaling up the backbone used in PBNs regardless of the training dataset, which is in line with theoretical work (Tripuraneni et al., 2021) proving that overparameterized models exhibit enhanced robustness to perturbations. The same observation holds across different architectures, too, with larger models being more robust. We conclude that larger Transformer architectures are optimal to ensure the robust behavior of PBNs. In the remaining experiments, we use BART, as its performance is similar to Electra and its number of parameters is 2x lower, while outperforming BERT.

**Effect of loss functions** ($\mathcal{L}$). Removing the interpretability loss consistently leads to a drop in performance on perturbed examples, especially on word-level perturbations (Figure 3). This means that the model training process enables higher robustness to perturbations if the prototypes are kept close to at least one training example. Meanwhile, while the separation loss often enhances the robustness against perturbations and out-of-distribution data in computer vision tasks (Yang et al., 2018),
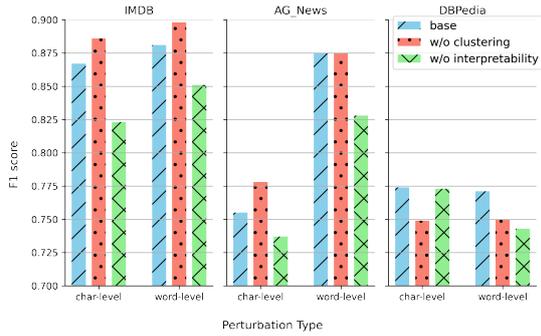
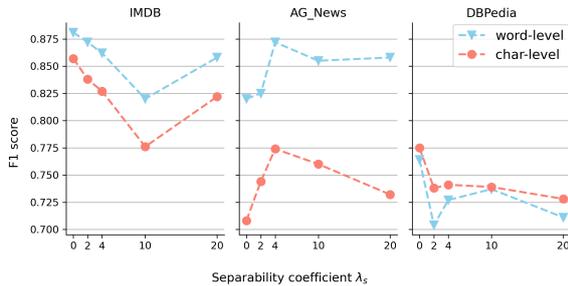Figure 3: F1 score of PBN (BART) with ablated loss functions across datasets and perturbations.



Figure 4: Performance of PBN (BART) under different perturbations with different separability ($\lambda_s$).
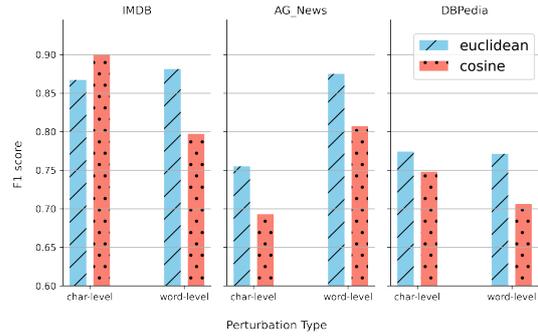


Figure 5: Performance of PBN (BART) using Cosine and Euclidean distance transformations.



Figure 6: Performance of the PBN (BART) model with varying numbers of prototypes ($Q$).

we show that the separation between prototypes does not have a clear impact on the robustness of the PBNs for text classification (Figure 4). However, our results suggest that keeping the prototypes close to each other (lower $\lambda_s$ coefficients in Figure 4) helps robustness for 2 out of 3 datasets. We conclude that the interpretability loss should be included to enhance the robustness of PBNs on text classification tasks, whereas the separation loss parameter should be tuned carefully according to the task at hand. We also observe that clustering loss is hindering the robustness. The clustering loss is a regularization term that encourages samples to be all close to certain prototypes in the embedding space, enhancing interpretability, but potentially reducing accuracy by narrowing the diversity in the embedding space; which is a common phenomenon in loss terms of competing goals. The mean and standard deviation over the (transformed) distances between prototypes and each sample can be used to describe the spread of the embedded data points around the prototype. These values are $(-0.18 \pm 1.5) \times 10^{-6}$ with the clustering loss and $(-0.24 \pm 1.7) \times 10^{-7}$ without, showing less diverse prototypes indicated by smaller values of measured distances caused by stronger clustering.

**Effect of distance functions** ($d$). In line with prior observation by Snell et al. (2017), we found that computing similarities using Euclidean distance makes PBNs more robust than using Cosine distance. In other words, forcing the samples to be close to a point (the prototype) in the embedding space is more robust than forcing the samples to be in the same direction as the prototypes (Figure 5).

**Effect of the number of prototypes** ($Q$). We found that $Q$ highly affects PBNs' robustness (Figure 6): the F1-scores of PBN (BART) change as a function of $Q$ for all datasets and perturbations. We observe that PBNs perform poorly when the number of prototypes is as low as the number of classes in the dataset and when there is an excessive number of prototypes (we consider 64 as an upper bound). Similar to prior findings in vision (Yang et al., 2018) and in logical fallacy identification (Sourati et al., 2023), we find that the optimal number of prototypes is higher than the number of classes, while too many prototypes are detrimental to the model performance since it is forced to learn a more complex embedding space. It is known that the optimal number of prototypes is non-trivial

| Model | IMDB | | AG_News | | DBPedia | |
|---|---|---|---|---|---|---|
| | Char | Word | Char | Word | Char | Word |
| BART | 94.01 (23.8) | 99.80 (05.7) | 56.34 (**33.8**) | 90.19 (**24.6**) | 39.62 (**45.0**) | 68.00 (24.0) |
| + PBN | **43.96** (**28.4**) | **88.30** (**12.3**) | **24.44** (30.6) | **62.60** (24.1) | **12.62** (43.0) | **53.33** (**26.0**) |

Table 2: Character and word perturbation success rates (lower=better) and average perturbed word percentages (higher=better; values in parentheses) for BART and its PBN counterpart.

and not necessarily the number of classes or training data points (Crammer et al., 2003). Therefore, $Q = 16$ is in accordance with prior work, and we found it empirically to yield optimal performance.

**Effect of direct attacks on PBNs.** As opposed to the previous experiments that do not assume having access to the victim model (PBNs), in an additional experiment, we compare the adversarial perturbations directly on both PLMs and PBNs, following standard robustness evaluation procedure (Wang et al., 2022a, 2023; Xiao et al., 2018). Using character- and word-level perturbations, we attack each model until we have 800 successful adversarial perturbations, reporting the attack success rate and the average number of perturbed tokens for each victim model (see Table 2). We observe that attacks on PBN (BART) are successful only 27% and 68% of the times on average across all datasets as compared to 63% and 86% for the vanilla BART on character- and word-level perturbations, respectively. We observed the opposite pattern in terms of the average number of perturbed words, where more tokens need to be perturbed in PBNs compared to vanilla PLMs. These results demonstrate the superiority of PBNs in a dynamic adversarial setting, too, where models are directly attacked.

## 7 Related Work

**Robustness evaluation.** Robustness in NLP is defined as models' ability to perform well under noisy (Ebrahimi et al., 2018) and out-of-distribution data (Hendrycks et al., 2020). With the wide adoption of NLP models in different domains and their near-human performance on natural language benchmarks (Wang et al., 2019; Sarlin et al., 2020), concerns have shifted towards the NLP models' performance facing noisy data (Wang et al., 2022a). Wang et al. (2023) evaluated the adversarial robustness of ChatGPT and found that, although it is more robust than other LLMs, it is far from perfect. Shi et al. (2023) studied the effect of irrelevant context on LLMs and found its dramatic effect on the model's performance. Wang et al. (2022b) presented ReCode to evaluate the robust-

ness properties of code generation models. While prior work has studied PLMs' robustness, to our knowledge, PBNs' robustness properties have not been explored for text classification. Our study bridges this gap.

**Prototype-based networks.** PBNs are widely used in computer vision (Chen et al., 2019; Hase et al., 2019; Kim et al., 2021; Nauta et al., 2021b; Pahde et al., 2021) because of their interpretability and robustness properties (Soares et al., 2022; Yang et al., 2018). While limited work has been done in the NLP domain, PBNs have recently found application in text classification tasks such as propaganda detection (Das et al., 2022), logical fallacy detection (Sourati et al., 2023), sentiment analysis (Pluciński et al., 2021), and few-shot relation extraction (Meng et al., 2023). ProseNet (Ming et al., 2019), a prototype-based text classifier, uses several criteria for constructing prototypes (He et al., 2020), and a special optimization procedure for better explainability. ProtoryNet (Hong et al., 2020) leverages RNN-extracted prototype trajectories and deploys a pruning procedure for prototypes. ProtoCNN (Pluciński et al., 2021) uses phrase-based prototypes to provide explanations using n-grams, and ProtoTex (Das et al., 2022) employs the concept of negative prototypes for handling the absence of features for classification. While PBNs are expected to be robust to perturbations, this property has not been systematically studied for text classification tasks. Our paper consolidates prior PBN methods like ProtoTex and ProtoCNN into a comprehensive study framework.

## 8 Conclusions

We study the robustness of PBNs with the help of character-, word-, and sentence-level perturbations. We find that PBNs are typically more robust than baseline models both in static and dynamic perturbation settings. Our experiments show that the choice of the encoder backbone is critical, with Tranformer-based backbones being relatively robust compared to CNN-based PBNs that are not. Here, scaling the size of the backbone improves

robustness further. We study the impact of PBN components like individual loss functions, the number of prototypes, and the distance functions on robustness. We find that the interpretability loss contributes the most to robustness, robust PBNs require the number of prototypes to be higher than the number of classes, and adopting a Euclidean distance calculation instead of Cosine can be more effective in terms of robustness. For future work, we plan to study the interplay between robustness and interpretability for text classification via user studies.

## Limitations

While we perform a systematic study of the robustness of PBNs, we do not analyze how the explanations offered by the model change when faced with perturbations. Additionally, we limit our study to one of each kind of character-, word-, and sentence-level perturbations. While the specific attack implementations are popular text perturbation strategies within these categories and have been shown to affect language models in the past, we acknowledge that more complicated perturbations can also be created that are more effective and help the community have a more complete understanding of models' robustness, hence we do not comment on the generalizability of our study to all possible textual perturbations besides our evaluation on AdvGLUE. We leave the interpretability analysis, evaluation of tasks beyond text classification, and generalizability study to future work.

## References

Plamen Angelov and Eduardo Soares. 2020. Towards explainable deep neural networks (xdnn). *Neural Networks*, 130:185–194.

Chaofan Chen, Oscar Li, Chaofan Tao, Alina Jade Barnett, Jonathan Su, and Cynthia Rudin. 2019. *This Looks like That: Deep Learning for Interpretable Image Recognition*. Curran Associates Inc., Red Hook, NY, USA.

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*.

Kevin Clark, Minh-Thang Luong, Quoc V Le, and Christopher D Manning. 2020. Electra: Pre-training text encoders as discriminators rather than generators. *arXiv preprint arXiv:2003.10555*.

Koby Crammer, Ran Gilad-Bachrach, Amir Navot, and Naftali Tishby. 2003. Margin analysis of the LVQ algorithm. In *Advances in Neural Information Processing Systems 15: Proceedings of the Neural Information Processing Systems Conference – NIPS 2002*, pages 479–486, Vancouver, BC, Canada. MIT Press.

Nilesh Dalvi, Pedro Domingos, Mausam, Sumit Sanghai, and Deepak Verma. 2004. Adversarial classification. KDD '04, page 99–108, New York, NY, USA. Association for Computing Machinery.

Anubrata Das, Chitrank Gupta, Venelin Kovatchev, Matthew Lease, and Junyi Jessy Li. 2022. ProtoTEx: Explaining model decisions with prototype tensors. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2986–2997, Dublin, Ireland. Association for Computational Linguistics.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

Javid Ebrahimi, Anyi Rao, Daniel Lowd, and Dejing Dou. 2018. HotFlip: White-box adversarial examples for text classification. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 31–36, Melbourne, Australia. Association for Computational Linguistics.

Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. 2014. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*.

Xiaowei Gu and Weiping Ding. 2019. A hierarchical prototype-based approach for classification. *Information Sciences*, 505:325–351.

Jiale Han, Bo Cheng, and Wei Lu. 2021. Exploring task difficulty for few-shot relation extraction. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 2605–2616, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Peter Hase and Mohit Bansal. 2020. Evaluating explainable AI: Which algorithmic explanations help users predict model behavior? In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics – ACL2020*, pages 5540–5552, Online. Association for Computational Linguistics.

Peter Hase, Chaofan Chen, Oscar Li, and Cynthia Rudin. 2019. Interpretable image recognition with hierarchical prototypes. In *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, volume 7, pages 32–40.

Junxian He, Taylor Berg-Kirkpatrick, and Graham Neubig. 2020. Learning sparse prototypes for text generation. *Advances in Neural Information Processing Systems*, 33:14724–14735.

Dan Hendrycks, Xiaoyuan Liu, Eric Wallace, Adam Dziedzic, Rishabh Krishnan, and Dawn Song. 2020. Pretrained transformers improve out-of-distribution robustness. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2744–2751, Online. Association for Computational Linguistics.

Dat Hong, Stephen S Baek, and Tong Wang. 2020. Interpretable sequence classification via prototype trajectory. *arXiv preprint arXiv:2007.01777*.

Dat Hong, Stephen S. Baek, and Tong Wang. 2021. Interpretable sequence classification via prototype trajectory.

Myeongjun Jang, Deuk Sin Kwon, and Thomas Lukasiewicz. 2022. Becel: Benchmark for consistency evaluation of language models. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 3680–3696.

Di Jin, Zhijing Jin, Joey Tianyi Zhou, and Peter Szolovits. 2020. Is bert really robust? a strong baseline for natural language attack on text classification and entailment. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 8018–8025.

Eunji Kim, Siwon Kim, Minji Seo, and Sungroh Yoon. 2021. Xprotonet: Diagnosis in chest radiography with global and local explanations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15719–15728.

Alexey Kurakin, Ian Goodfellow, and Samy Bengio. 2017a. Adversarial examples in the physical world.

Alexey Kurakin, Ian Goodfellow, and Samy Bengio. 2017b. Adversarial machine learning at scale.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension.

Jinfeng Li, Shouling Ji, Tianyu Du, Bo Li, and Ting Wang. 2018a. Textbugger: Generating adversarial text against real-world applications. *arXiv preprint arXiv:1812.05271*.

Oscar Li, Hao Liu, Chaofan Chen, and Cynthia Rudin. 2018b. Deep learning for case-based reasoning through prototypes: A neural network that explains its predictions. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence and Thirtieth Innovative Applications of Artificial Intelligence Conference and Eighth AAAI Symposium on Educational Advances in Artificial Intelligence*, AAAI'18/IAAI'18/EAAI'18. AAAI Press.

Pengfei Liu, Jinlan Fu, Yanghua Xiao, Weizhe Yuan, Shuaichen Chang, Junqi Dai, Yixin Liu, Zihuiwen Ye, Zi-Yi Dou, and Graham Neubig. 2021. ExplainaBoard: An Explainable Leaderboard for NLP. In *Annual Meeting of the Association for Computational Linguistics (ACL), System Demonstrations*.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.

Shiao Meng, Xuming Hu, Aiwei Liu, Shu'ang Li, Fukun Ma, Yawen Yang, and Lijie Wen. 2023. RAPL: A Relation-Aware Prototype Learning Approach for Few-Shot Document-Level Relation Extraction. *arXiv preprint arXiv:2310.15743*.

Pascal Mettes, Elise Van der Pol, and Cees Snoek. 2019. Hyperspherical prototype networks. *Advances in neural information processing systems*, 32.

Bonan Min, Hayley Ross, Elior Sulem, Amir Pouran Ben Veyseh, Thien Huu Nguyen, Oscar Sainz, Eneko Agirre, Ilana Heinz, and Dan Roth. 2021. Recent advances in natural language processing via large pre-trained language models: A survey.

Yao Ming, Panpan Xu, Huamin Qu, and Liu Ren. 2019. Interpretable and steerable sequence learning via prototypes. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM.

Milad Moradi and Matthias Samwald. 2021. Evaluating the robustness of neural language models to input perturbations.

John X. Morris, Eli Lifland, Jin Yong Yoo, Jake Grigsby, Di Jin, and Yanjun Qi. 2020. Textattack: A framework for adversarial attacks, data augmentation, and adversarial training in nlp.

Meike Nauta, Annemarie Jutte, Jesper Provoost, and Christin Seifert. 2021a. This looks like that, because ... explaining prototypes for interpretable image recognition. In *Communications in Computer and Information Science*, pages 441–456. Springer International Publishing.

Meike Nauta, Ron van Bree, and Christin Seifert. 2021b. Neural prototype trees for interpretable fine-grained image recognition. In *Proceedings of the 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition – CVPR 2021*, pages 14933–14943, Nashville, TN, USA. IEEE.

OpenAI. 2022. Chatgpt. https://openai.com/blog/chatgpt. Accessed: April 30, 2023.

Frederik Pahde, Mihai Puscas, Tassilo Klein, and Moin Nabi. 2021. Multimodal prototypical networks for few-shot learning. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 2644–2653.

Kamil Pluciński, Mateusz Lango, and Jerzy Stefanowski. 2021. Prototypical convolutional neural network for a phrase-based explanation of sentiment

classification. In *Machine Learning and Principles and Practice of Knowledge Discovery in Databases*, pages 457–472, Cham. Springer International Publishing.

Eleanor H. Rosch. 1973. Natural categories. *Cognitive Psychology*, 4(3):328–350.

Cynthia Rudin, Chaofan Chen, Zhi Chen, Haiyang Huang, Lesia Semenova, and Chudi Zhong. 2022. Interpretable machine learning: Fundamental principles and 10 grand challenges. *Statistic Surveys*, 16:1–85.

Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*.

Sascha Saralajew, Lars Holdijk, and Thomas Villmann. 2020. Fast adversarial robustness certification of nearest prototype classifiers for arbitrary seminorms. In *Advances in Neural Information Processing Systems*, pages 13635–13650. Curran Associates, Inc.

Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. 2020. Superglue: Learning feature matching with graph neural networks.

Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed Chi, Nathanael Schärli, and Denny Zhou. 2023. Large language models can be easily distracted by irrelevant context.

Dylan Slack, Satyapriya Krishna, Himabindu Lakkaraju, and Sameer Singh. 2022. Talktomodel: Explaining machine learning models with interactive natural language conversations.

Jake Snell, Kevin Swersky, and Richard Zemel. 2017. Prototypical networks for few-shot learning. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

Eduardo Soares, Plamen Angelov, and Neeraj Suri. 2022. Similarity-based deep neural network to detect imperceptible adversarial attacks. In *2022 IEEE Symposium Series on Computational Intelligence (SSCI)*, pages 1028–1035.

Zhivar Sourati, Vishnu Priya Prasanna Venkatesh, Darshan Deshpande, Himanshu Rawlani, Filip Ilievski, Hông-Ân Sandlin, and Alain Mermoud. 2023. Robust and explainable identification of logical fallacies in natural language arguments. *Knowledge-Based Systems*, 266:110418.

Nilesh Tripuraneni, Ben Adlam, and Jeffrey Pennington. 2021. Overparameterization improves robustness to covariate shift in high dimensions. In *Advances in Neural Information Processing Systems*, volume 34, pages 13883–13897. Curran Associates, Inc.

Iulia Turc, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Well-read students learn better: On the importance of pre-training compact models.

Betty van Aken, Jens-Michalis Papaioannou, Marcel G. Naik, Georgios Eleftheriadis, Wolfgang Nejdl, Felix A. Gers, and Alexander Löser. 2022. This patient looks like that patient: Prototypical networks for interpretable diagnosis prediction from clinical text.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need.

Václav Vorácek and Matthias Hein. 2022. Provably adversarially robust nearest prototype classifiers. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of the Proceedings of Machine Learning Research, pages 22361–22383, Baltimore, MD, USA. PMLR.

Kiri Wagstaff. 2012. Machine learning that matters. *arXiv preprint arXiv:1206.4656*.

Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019. Glue: A multi-task benchmark and analysis platform for natural language understanding.

Boxin Wang, Chejian Xu, Shuohang Wang, Zhe Gan, Yu Cheng, Jianfeng Gao, Ahmed Hassan Awadallah, and Bo Li. 2022a. Adversarial glue: A multi-task benchmark for robustness evaluation of language models.

Jindong Wang, Xixu Hu, Wenxin Hou, Hao Chen, Runkai Zheng, Yidong Wang, Linyi Yang, Haojun Huang, Wei Ye, Xiubo Geng, Binxin Jiao, Yue Zhang, and Xing Xie. 2023. On the robustness of chatgpt: An adversarial and out-of-distribution perspective.

Shiqi Wang, Zheng Li, Haifeng Qian, Chenghao Yang, Zijian Wang, Mingyue Shang, Varun Kumar, Samson Tan, Baishakhi Ray, Parminder Bhatia, Ramesh Nallapati, Murali Krishna Ramanathan, Dan Roth, and Bing Xiang. 2022b. Recode: Robustness evaluation of code generation models.

Zhao Wang and Aron Culotta. 2020. Identifying spurious correlations for robust text classification. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3431–3440, Online. Association for Computational Linguistics.

Chaowei Xiao, Bo Li, Jun yan Zhu, Warren He, Mingyan Liu, and Dawn Song. 2018. Generating adversarial examples with adversarial networks. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI-18*, pages 3905–3911. International Joint Conferences on Artificial Intelligence Organization.

Hong-Ming Yang, Xu-Yao Zhang, Fei Yin, and Cheng-Lin Liu. 2018. Robust classification with convolutional prototype learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3474–3482.

| Type | Perturbed Sentence |
|------|--------------------|
| Char. | <u>Disny</u> Rules Out New Deal with Pixar Studios. |
| Word | Disney Rules Out New Deal with <u>Ghibli</u> Studios. |
| Sent. | Disney <u>has decided against pursuing a fresh agreement with Pixar Studios.</u> |

Table 3: Perturbation examples on the following sentence: *Disney Rules Out New Deal with Pixar Studios.* The <u>underlined</u> text indicates the new tokens that replace the original sentence tokens.

Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. *Advances in neural information processing systems*, 28.

Pei Zhou, Rahul Khanna, Seyeon Lee, Bill Yuchen Lin, Daniel Ho, Jay Pujara, and Xiang Ren. 2020. Rica: Evaluating robust inference capabilities based on commonsense axioms. *arXiv preprint arXiv:2005.00782*.

Julia El Zini and Mariette Awad. 2022. On the explainability of natural language processing deep models. *ACM Comput. Surv.*, 55(5).

Barret Zoph, Irwan Bello, Sameer Kumar, Nan Du, Yanping Huang, Jeff Dean, Noam Shazeer, and William Fedus. 2022. St-moe: Designing stable and transferable sparse expert models.

## A  Dataset Details

**Perturbations.** As mentioned in Section 4, we focus on three types of perturbations to assess the robustness of PBNs: character-level, word-level, and sentence-level perturbations. Examples of each type of perturbation from a single sentence are presented in Table 3.

**Dataset statistics.** The statistics of the datasets we used in this study to test the robustness of PBNs against perturbations are demonstrated in Table 4. We present both statistics about the original dataset and statistics and details about the number of perturbations that we have gathered on each dataset in three categories of character-level, word-level, and sentence-level perturbations.

## B  Implementation Details

**Experimental environment.** For all the experiments that involved training or evaluating Transformers or other models like CNN, we used three GPU NVIDIA RTX A5000 devices with Python v3.9.16 and CUDA v11.6. All Transformer models were trained using the Transformers package v4.30.2 and Torch package v2.0.1+cu117. We used TextAttack (Morris et al., 2020) for implementing the character-level and word-level perturbations.

**Training details.** For all the datasets, the training split, validation split, and test splits were used from `https://huggingface.co/`. During training on the IMDB and DBPedia dataset, the batch size was set to 64. This number was 256 on the AG_News dataset. All the models (Table 5) were trained with the number of epochs adjusted according to an early stopping module with the patience of 4 and a threshold value of 0.01 for change in the accuracy. The coefficients controlling the effect of different loss terms in Equation 1 were all set to 0.9 when having all the components in place, and they were set to 0.0 to simulate the situation where they do not contribute to the total loss term.

All the Transformer models were fine-tuned on top of a pre-trained model gathered from `https://huggingface.co/`. Details of the models used in our experiments are presented in the following:

- Electra (Clark et al., 2020): google/electra-base-discriminator;

- BART (Lewis et al., 2019): ModelTC/bart-base-mnli, facebook/bart-base, facebook/bart-large-mnli;

- BERT (Devlin et al., 2018): prajjwal1/bert-medium.

Furthermore, the models that were used in the process of gathering generalizable perturbations were also pre-trained Transformer models gathered from `https://huggingface.co/`. Note that in the process of the mentioned models' pre-training, they were fine-tuned on specific datasets we used in our study before being attacked by the perturbations. Find the details of models used categorized by the dataset below:

- IMDB: textattack/bert-base-uncased-imdb, textattack/distilbert-base-uncased-imdb, textattack/roberta-base-imdb;

- AG_News: textattack/bert-base-uncased-ag-news, andi611/distilbert-base-uncased-ner-agnews, textattack/roberta-base-ag-news;

- DBPedia: dbpedia_bert-base-uncased, dbpedia_distilbert-base-uncased, dbpedia_roberta-base.

| Dataset | #Classes | #Tokens | #Train | #Val | #Test (O) | #Test (C) | #Test (W) | #Test (S) |
|---------|----------|---------|--------|------|-----------|-----------|-----------|-----------|
| IMDB | 2 | 234 | 22,500 | 2,500 | 25,000 | 1,926 | 2,190 | 1,000 |
| AG_News | 4 | 103 | 112,400 | 7,600 | 7,600 | 1,383 | 1,893 | 1,000 |
| DBPedia | 9 | 38 | 240,942 | 36,003 | 60,794 | 1,281 | 1,836 | 1,000 |
| SST-2 | 2 | 14 | 67,349 | 872 | 1,821 | - | 80 | 51 |

Table 4: Dataset statistics: number of classes, the average number of tokens, and partition sizes. #Test (O) shows the size of the original test set, while #Test ({C, W, S}) shows the size of the subsets obtained after successful character-, word-, and sentence-level attacks, repsectively. SST-2 subset comes from the AdvGlue benchmark (Wang et al., 2022a) after removing the human-generated instances that do not belong to either category of perturbation classes.

| Model | Parameters |
|-------|------------|
| BART-base | 53.56 M |
| BART-base (PBN) | 59.86 M |
| BART-large | 205.78 M |
| BART-large (PBN) | 212.08 M |
| Electra-base | 108.89 M |
| Electra-base (PBN) | 115.18 M |
| CNN | 9.3 M |
| CNN (PBN) | 9.3 M |
| BERT-Medium | 41.37 M |
| BERT-Medium (PBN) | 45.56 M |

Table 5: Number of trainable parameters for all experimental models.

Since we could not find models from TextAttack (Morris et al., 2020) library that were fine-tuned on DBPedia, the models that are presented above were fine-tuned by us on that dataset as well and then used as the victim model.