

# Robust Text Classification: Analyzing Prototype-Based Networks

Anonymous ACL submission

## Abstract

Downstream applications often require text classification models to be accurate, robust, and interpretable. While the accuracy of the state-of-the-art (large) language models approximates human performance, they are not designed to be interpretable and often exhibit a drop in performance on noisy data in the real world. This lack of robustness is particularly concerning in critical domains; e.g., toxicity detection in social media, where it may harm readers by, for example, misidentifying harmful content. A potential solution can be the family of Prototype-Based Networks (PBNs) that classifies examples based on their similarity to prototypical examples of a class (prototypes) and 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 the robustness of PBNs, which includes different backbone architectures, distance functions, and objective functions. The proposed evaluation protocol assesses the robustness of models against character-, word-, and sentence-level perturbations. Our experiments show that PBNs consistently enhance the robustness of vanilla language models, supported by the objective function that keeps prototypes interpretable.

## 1 Introduction

The requirements for robustness and interpretability have become urgent for high-stake tasks, such as toxicity detection (Davidson et al., 2017) and fake news detection (Rubin et al., 2016), where model errors can lead to serious consequences such as increased bias and misinformation (Rudin et al., 2022). 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). 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). However, the widely adopted pre-trained language models (PLMs) and large language models (LLMs), which report exceptional accuracy on NLP classification benchmarks (Chowdhery et al., 2022; Zoph et al., 2022; Zhao et al., 2023), 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). Reportedly, LLMs also lack robustness when faced with out-of-distribution data and noisy inputs (Wang et al., 2023), a finding that is supported by the empirical findings of the present paper.

Meanwhile, the family of Prototype-Based Networks (PBNs) is designed for robustness and interpretability (Li et al., 2018b). 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). Computationally, this idea is implemented by finding prototypical points or examples in the shared embedding space of data points and using the distance between prototypes and data points to accomplish the classification task. This classification approach is both interpretable and robust because it classifies through distances to prototypical examples found in the data. Simple associations between data points and central prototypical examples bring

043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083

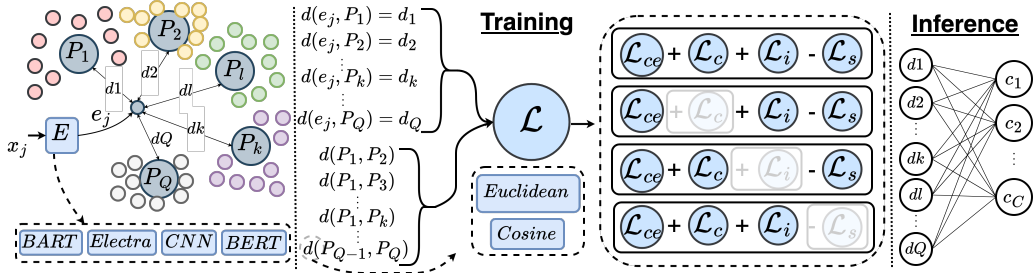


Figure 1: Classification by a PBN. The model computes distances between the new point and prototypes,  $d(e_j, P_k)$ , and distances within prototypes,  $d(P_k, P_l)$ , for both inference and training. During training, the model minimizes the loss term,  $\mathcal{L}$ , based on the distances between the new point and prototypes as well as within prototypes; during inference, distances between the new point and prototypes are used for classification by a fully connected layer. We test variations of the loss terms ( $\mathcal{L}$ ), different encoder backbones ( $E$ ), and distance functions ( $d$ ) and assess their effect on the PBN’s robustness.

084 interpretability while leveraging distance between  
 085 points helps to quantify prototypicality, which then  
 086 facilitates identifying noisy or out-of-distribution  
 087 samples (Yang et al., 2018).

088 PBNs have been popular in Computer Vision  
 089 (CV) tasks, including image classification (An-  
 090 gelov and Soares, 2020) and novel class detec-  
 091 tion (Hase et al., 2019). Inspired by PBNs in CV,  
 092 NLP researchers have also developed PBN mod-  
 093 els for text classification, in particular, for senti-  
 094 ment classification (Pluciński et al., 2021; Ming  
 095 et al., 2019; Hong et al., 2021), few-shot relation  
 096 extraction (Han et al., 2021; Meng et al., 2023),  
 097 and propaganda detection (Das et al., 2022). Yet,  
 098 while competitive performance and interpretability  
 099 of PBNs have been studied in both NLP (Das et al.,  
 100 2022; Hase and Bansal, 2020) and CV (Gu and  
 101 Ding, 2019; van Aken et al., 2022), their robust-  
 102 ness advantages over vanilla models have only been  
 103 investigated in CV (Yang et al., 2018; Saralajew  
 104 et al., 2020; Voráček and Hein, 2022).

105 In this study, we investigate whether the robust-  
 106 ness properties of PBNs transfer to NLP classifica-  
 107 tion tasks. In particular, the contributions of this  
 108 work are as follows: (1) We propose a **modular**  
 109 **and comprehensive framework** to evaluate the ro-  
 110 bustness properties of PBNs, which combines dif-  
 111 ferent backbone architectures, distance functions,  
 112 and loss terms; (2) We devise an evaluation proto-  
 113 col that employs **three perturbation strategies** to  
 114 evaluate the robustness of models against charac-  
 115 ter-, word-, and sentence-level noise; (3) We per-  
 116 form **extensive experiments** on four benchmarks  
 117 to compare the robustness of PBNs to vanilla mod-  
 118 els, evaluate the effect of key PBN design choices,  
 119 and assess their sensitivity to varying task com-

120 plexity. Our experiments show that the robustness  
 121 of PBNs transfers to realistic perturbations in text  
 122 classification tasks, and can, thus, enhance the text  
 123 classification robustness of PLMs and outperform  
 124 state-of-the-art LLMs.

## 2 Prototype-Based Networks

125 PBNs classify data points based on their similar-  
 126 ity to a set of *prototypes* learned during training.  
 127 These prototypes summarize the prominent seman-  
 128 tic patterns of the dataset through two mechanisms:  
 129 (1) prototypes are defined in the same embedding  
 130 space as input examples, which makes them in-  
 131 terpretable by leveraging input examples close to  
 132 them; and (2) prototypes are designed to cluster se-  
 133 mantically similar training examples, which makes  
 134 them representative of the prominent patterns em-  
 135 bedded in the data and input examples. The PBN’s  
 136 decisions are inherently interpretable because pro-  
 137 totypes are trained to be aligned with previous ob-  
 138 servations (Hong et al., 2020). This enables in-  
 139 sights into the behavior of the model during infer-  
 140 ence by looking at the closest activated prototypes  
 141 (Das et al., 2022). Prototypes being in the same  
 142 embedding space as input examples allows them  
 143 to be represented as either the training examples  
 144 (Das et al., 2022) or parts of training examples,  
 145 such as key phrases (Pluciński et al., 2021) or key  
 146 sequences (Ming et al., 2019; Hong et al., 2021) ex-  
 147 tracted from training examples. These prototypes  
 148 can be associated with semantic patterns of partic-  
 149 ular classes from their initialization or be trained  
 150 freely and subsequently associated with the promi-  
 151 nent semantic patterns of the whole dataset.

**Inference.** Classification in PBNs is done via  
 152 a fully connected layer applied on the measured  
 153  
 154

distances between embedded data points and prototypes. As shown in Figure 1, given a set of data points  $x_j, j \in \{1, \dots, N\}$  with labels  $y_j \in \{1, \dots, C\}$ , 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  $c \in \{1, \dots, C\}$ .

**Training.** The model is trained using objectives that simultaneously tweak the backbone parameters and the (randomly initialized) prototypes, thus promoting high performance and meaningful prototypes. 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  $d(e_j, P_k)_{j \in \{1, \dots, N\}; k \in \{1, \dots, Q\}}$ , 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):

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

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:

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

*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:

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

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

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

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

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

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

### 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 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 Transformer-based PLMs (Vaswani et al., 2017) due to their general language abilities (Min et al., 2021), we experiment with a CNN as a backbone since their architecture, being simpler than Transformers, may make PBNs more interpretable by mapping the extracted embeddings from shallow CNN-based networks to n-grams (Pluciński et al., 2021). This can bring more insight into the interplay between robustness and interpretability. 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 with 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 encoder’s latent space. However, with some utilizing Cosine distance (Chen et al., 2019) while others prioritizing Euclidean distance (Mettes et al., 2019), and the two having incomparable experimental setups, conclusive arguments about the superiority of one over the other cannot be justified, 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 investigate 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 ( $\mathcal{L}$  w/o  $\mathcal{L}_i$ ), and a model that omits clustering loss from training ( $\mathcal{L}$  w/o  $\mathcal{L}_c$ ). 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 to character-, word-, and sentence-level noise (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 may consequently mislead the model. For instance, 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  $\rightarrow$  Citrix).

**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  $\rightarrow$  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 four 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 about our experimental setup and training choices.

**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* (Maas et al., 2011): a binary sen-

349 timent classification dataset; (2) *AG\_NEWS* (Zhang  
350 et al., 2015):<sup>1</sup> a collection of news articles that  
351 can be associated with four categories; (3) *DB-*  
352 *Pedia* (Zhang et al., 2015):<sup>2</sup> a dataset with taxo-  
353 nomic, hierarchical categories for Wikipedia arti-  
354 cles, with nine classes. Moreover, we adopt the  
355 SST-2 binary classification split from the exist-  
356 ing *Adversarial GLUE (AdvGLUE)* dataset (Wang  
357 et al., 2022a). The AdvGLUE SST-2 benchmark  
358 consists of 131 examples perturbed using various  
359 word- and sentence-level perturbations that are fil-  
360 tered both automatically and by human evaluation  
361 for more effectiveness. For statistics of the datasets  
362 and their perturbations, see [Appendix A](#).

363 **Perturbations.** The perturbations are designed  
364 to simulate real noisy data. To generate character-  
365 and word-level perturbations for IMDB, AG\_News,  
366 and DBPedia, we follow Wang et al. (2022a). For  
367 each dataset, we randomly perturb examples from  
368 its original test split until we obtain 800 *successful*  
369 perturbations. *Successful* perturbations of a dataset  
370 are those that change the prediction of a victim  
371 model (model facing perturbations) already fine-  
372 tuned on that dataset from the correct prediction  
373 to the wrong prediction. We consider three fine-  
374 tuned victim models: BERT (Devlin et al., 2018),  
375 RoBERTa (Liu et al., 2019), and DistilBERT (Sanh  
376 et al., 2019). We preserve examples whose pertur-  
377 bations are predicted wrongly by the three models.  
378 Based on prior work (Wang et al., 2022a), we ex-  
379 pect these perturbations to be *generalizable*, i.e., to  
380 represent noisy data that can affect a wide range of  
381 victim models without any assumption or informa-  
382 tion about them. In principle, the perturbations for  
383 each model are different, yielding three variations  
384 per original example for a dataset-perturbation pair.  
385 For instance, this procedure successfully creates  
386 three character-level perturbations of 461 original  
387 examples from the AG\_News dataset, resulting in a  
388 total of 1,383 ( $3 \times 461$ ) data points in its character-  
389 level test set. To generate sentence-level perturba-  
390 tions, we use GPT3.5 (OpenAI, 2022). To avoid in-  
391 troducing additional noise to data, we randomly se-  
392 lect 1,000 samples from the test set of each dataset  
393 and paraphrase them using the prompt: *Paraphrase*  
394 *the paragraph below concisely and accurately: {in-*  
395 *put text}*. While our main experiments evaluate  
396 PBNs on the simulated real noisy data, in an auxil-  
397 iary experiment, we consider PBNs as victim mod-

els and adversarially attack them.

398 **Baselines.** We compare the performance of our  
399 PBN framework variants with their vanilla counter-  
400 parts, i.e., CNN, BERT, BART encoder, and Elec-  
401 tra. The only difference between PBNs and their  
402 vanilla counterparts is the model wrapped around  
403 their backbones: while PBNs compute distances to  
404 prototypes and classify data points based on these  
405 distances, the vanilla models leverage a fully con-  
406 nected layer for classification instead. We also test  
407 the performance of ChatGPT (OpenAI, 2022) on  
408 both original and perturbed instances as a repre-  
409 sentative LLM baseline in a zero-shot setting. We  
410 do not combine ChatGPT with PBNs because that  
411 would require access to the model, and the Chat-  
412 GPT model is not released to date.

## 414 6 Results

415 We study the effect of perturbations that simulate  
416 noisy data, on the robustness of PBNs and their  
417 vanilla counterparts. We investigate the impact  
418 of the design choices within our framework from  
419 [Section 3](#) overall and in relation to task complexity.

420 **Robustness of PBNs.** We found that PBNs are con-  
421 sistentlly more robust to perturbations compared to  
422 their vanilla counterparts (see [Table 1](#)). We ob-  
423 served this trend both for our custom perturbations  
424 and for the more difficult benchmark, AdvGLUE’s  
425 SST-2 (Wang et al., 2022a). Character-level and  
426 word-level perturbations lead to a larger robust-  
427 ness gap between PBNs and vanilla models, while  
428 both PBN and vanilla models are robust against our  
429 sentence-level perturbations.<sup>3</sup> PBNs’ performance  
430 on perturbations improve by up to 24%, 33%, and  
431 30%, with BERT, BART, and Electra, respectively,  
432 relative to their vanilla variants. We note that the  
433 perturbations are general enough to be more chal-  
434 lenging for all models, including PBNs, as both  
435 character- and word-level perturbations decrease  
436 the average model performance by 20% compared  
437 to the scores on the original test sets. Finally, we  
438 note that PBNs perform slightly worse than the  
439 vanilla models on the original datasets, and we  
440 attribute this to PBNs being optimized to satisfy  
441 multiple objectives, including interpretability and  
442 clustering, whereas the vanilla models are only  
443 tuned for accuracy. Interestingly, we saw a low  
444 performance of ChatGPT on both original and per-

<sup>1</sup><https://bit.ly/47Uzha5>

<sup>2</sup><https://bit.ly/3RgX41H>

<sup>3</sup>As we observe a low effect from sentence-level perturba-  
tions across all models, we focus on character- and word-level  
perturbations and AdvGLUE for brevity.

Model	IMDB				AG_News				DBPedia				SST-2	
	Orig	Char	Word	Sent	Orig	Char	Word	Sent	Orig	Char	Word	Sent	Orig	Adv
CNN	85.7	<b>75.5</b>	<b>76.9</b>	<b>86.2</b>	90.4	<b>58.9</b>	<b>67.7</b>	75.9	97.9	<b>64.0</b>	<b>63.9</b>	<b>94.4</b>	75.0	48.0
+PBN	82.6	73.4	72.5	83.2	80.1	48.5	57.2	<b>76.3</b>	85.0	50.6	49.7	83.1	75.0	<b>49.0</b>
$\Delta$	<b>-3.1</b>	<b>-2.1</b>	<b>-4.4</b>	<b>-3.0</b>	<b>-10.3</b>	<b>-10.4</b>	<b>-10.5</b>	<b>+0.4</b>	<b>-12.9</b>	<b>-13.4</b>	<b>-14.2</b>	<b>-11.3</b>	0.0	<b>+1.0</b>
BERT	94.7	<b>84.2</b>	85.5	<b>92.2</b>	93.9	71.1	78.8	88.3	98.4	66.2	60.5	<b>98.0</b>	83.9	40.9
+PBN	90.5	83.5	<b>85.6</b>	85.9	92.1	<b>78.4</b>	<b>80.2</b>	<b>88.5</b>	96.5	<b>69.0</b>	<b>75.5</b>	97.4	77.8	<b>46.3</b>
$\Delta$	<b>-4.2</b>	<b>-0.7</b>	<b>+0.1</b>	<b>-6.3</b>	<b>-1.8</b>	<b>+7.3</b>	<b>+1.4</b>	<b>+0.2</b>	<b>-1.9</b>	<b>+2.8</b>	<b>+15.0</b>	<b>-0.6</b>	<b>-6.1</b>	<b>+5.4</b>
BART	98.1	89.3	92.1	94.9	<b>97.4</b>	74.2	79.2	<b>90.1</b>	96.3	69.5	68.8	97.0	93.1	29.7
+PBN	96.9	<b>89.5</b>	<b>93.4</b>	<b>95.0</b>	93.7	<b>75.6</b>	<b>81.2</b>	88.3	93.6	<b>72.9</b>	<b>71.5</b>	<b>97.4</b>	90.0	<b>39.6</b>
$\Delta$	<b>-1.2</b>	<b>+0.2</b>	<b>+1.3</b>	<b>+0.1</b>	<b>-3.7</b>	<b>+1.4</b>	<b>+2.0</b>	<b>-1.8</b>	<b>-2.7</b>	<b>+3.4</b>	<b>+2.7</b>	<b>+0.4</b>	<b>-3.1</b>	<b>+9.9</b>
ELEC.	<b>98.4</b>	<b>92.8</b>	94.0	<b>94.5</b>	95.0	71.8	78.7	<b>88.7</b>	93.7	61.8	65.0	94.5	87.6	43.5
+PBN	94.9	91.9	<b>94.2</b>	91.8	90.7	<b>73.9</b>	<b>80.4</b>	82.8	<b>99.1</b>	<b>76.7</b>	<b>70.5</b>	<b>98.4</b>	<b>98.5</b>	<b>56.6</b>
$\Delta$	<b>-3.5</b>	<b>-0.9</b>	<b>+0.2</b>	<b>-2.7</b>	<b>-4.3</b>	<b>+2.1</b>	<b>+1.7</b>	<b>-5.9</b>	<b>+5.4</b>	<b>+14.9</b>	<b>+5.5</b>	<b>+3.9</b>	<b>+10.9</b>	<b>+13.1</b>
ChatGPT	84.0	78.4	80.3	90.1	71.0	63.9	71.3	72.5	41.8	42.7	41.8	46.0	94.6	<b>62.5</b>
$\Delta$ (best)	<b>-14.4</b>	<b>-14.4</b>	<b>-13.9</b>	<b>-4.9</b>	<b>-26.4</b>	<b>-14.5</b>	<b>-9.9</b>	<b>-17.6</b>	<b>-57.3</b>	<b>-34.0</b>	<b>-33.7</b>	<b>-52.4</b>	<b>-3.9</b>	<b>+5.9</b>

Table 1: F1 scores of PBNs, their vanilla counterparts, and ChatGPT 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. vanilla model;  $\Delta$ 's indicate the difference between PBN vs. vanilla counterpart;  $\Delta$  (best) indicates the difference between ChatGPT and best performance among other models; gray highlighted cells indicate the highest performance among all models on a particular dataset and test-split. Performance on the original examples (Orig) is shown as the original scores on main test splits because of low variance across different attack splits.

turbed examples compared to other models, despite the fact that ChatGPT is a much larger model, has seen much more data (probably including perturbed instances), and is reported to be more robust than PLMs (Wang et al., 2023). Furthermore, the perturbations, although not targeted at ChatGPT, decrease its performance, which, alongside its overall lower performance, further accentuates the need for robust text classification models going beyond vanilla PLMs and LLMs. Find the complete results of ChatGPT in Table 6.

**Effect of task complexity.** Both the PBNs and their vanilla models perform worse on perturbed datasets that have more classes or more complex perturbations, as apparent from DBPedia (with 9 classes) and SST-2 from AdvGLUE (with 9 types of word- and sentence-level perturbations that have been filtered both automatically and by human evaluation). Meanwhile, the superior robustness of PBNs compared to their vanilla counterparts gets more pronounced as datasets get more complicated. While the performance of PBNs and vanilla models is on par for the IMDB dataset, their gap widens on the datasets with more classes: AG\_News and DBPedia. The same holds for more complex perturbations. Although both IMDB and SST-2 are binary datasets, the superiority of PBNs is more pronounced on SST-2, which contains more complex perturbations. We believe that this robust behavior is due to the design of the PBN architecture. Standard neural networks for text classification dis-

tinguish 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 PLM architectures (like BERT, BART, and Electra) yield higher absolute F1 scores compared to their vanilla backbones with differences ( $\Delta$ ) of up to 33%. 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 the remaining experiments, we use BART, as its performance is higher than BERT and similar to Electra while having half as many parameters.

**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 2). This means that the model training process enables higher robustness to perturbations if prototypes are kept close to at least one training example. Meanwhile,





Figure 2: F1 score of PBN (BART) with ablated loss functions across datasets and perturbations.

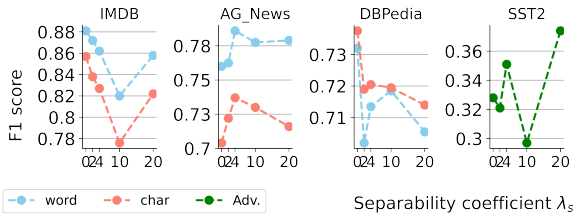


Figure 3: Performance of PBN (BART) under different perturbations with different separability ( $\lambda_s$ ).

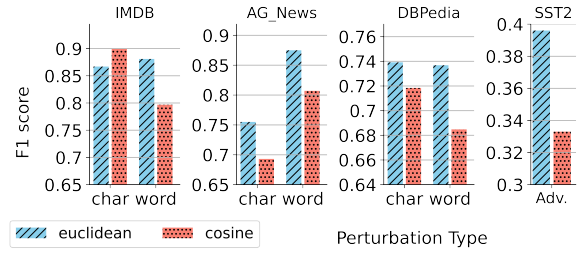


Figure 4: Performance of PBN (BART) using Cosine and Euclidean distance transformations.

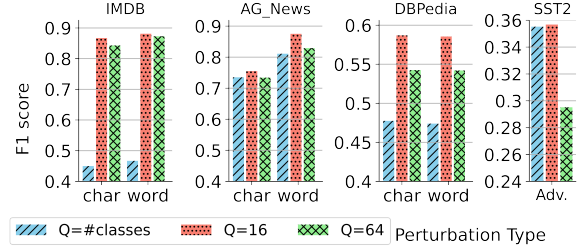


Figure 5: Performance of the PBN (BART) model with varying numbers of prototypes ( $Q$ ).

while the separation loss often enhances robustness against perturbations and out-of-distribution data in CV tasks (Yang et al., 2018), we show that the separation between prototypes does not have a clear impact on the robustness of the PBNs for text classification (Figure 3). 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 robustness. The clustering loss is a regularization term that encourages samples to be close to certain prototypes in the embedding space, further enhancing interpretability, but potentially reducing accuracy by narrowing the diversity in embedding space; which is a common phenomenon in loss terms of competing goals. The mean and standard deviation over (transformed) distances between prototypes and samples can be used to describe the spread of embedded data points around prototypes. These values are  $(-0.24 \pm 1.7) \times 10^{-7}$  without, and  $(-0.18 \pm 1.5) \times 10^{-6}$  with clustering loss, showing less diverse prototypes indicated by smaller 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 4).

**Effect of the number of prototypes ( $Q$ ).** We found that  $Q$  highly affects PBNs' robustness (Figure 5): 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 CV (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 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 previous experiments that do not assume having access to the victim model (PBNs), in an additional experiment common in robustness analysis (Wang et al., 2022a, 2023; Xiao et al., 2018), we adversarially attack both PLMs and PBNs and compare the results. Using character- and word-level perturbations, we attack each model until we have 800

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	<b>43.96 (28.4)</b>	<b>88.30 (12.3)</b>	<b>24.44</b> (30.6)	<b>62.60</b> (24.1)	<b>12.62</b> (43.0)	<b>53.33 (26.0)</b>

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.

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 vanilla BART on character- and word-level perturbations, respectively. We observed the opposite pattern in terms of 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 an 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,b). Wang et al. (2023) evaluated adversarial robustness of ChatGPT and found that, although it is more robust than PLMs, it is far from perfect. Similarly, Shi et al. (2023) studied the effect of irrelevant context on LLMs and found its dramatic effect on the model’s performance. 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 extrac-

tion (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) uses 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 components of prior PBNs like ProtoTex and ProtoCNN into a comprehensive framework to study their robustness.

## 8 Conclusions

Inspired by the lack of robustness to noisy data of state-of-the-art PLMs and LLMs, we study the sensitivity of PBNs to character-, word-, and sentence-level perturbations. We find that PBNs are typically more robust than vanilla models, both under simulated real noisy data and adversarial perturbations. Our experiments show that the choice of the encoder backbone is critical, with Transformer-based backbones being relatively robust compared to CNN-based PBNs that are not. 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. In summary, our work provides encouraging results for the potential of PBNs to enhance the robustness of PLMs across a variety of text classification tasks and quantifies the impact of architectural components on PBN robustness.



## 9 Limitations & Ethical Considerations

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. The specific attack implementations are popular text perturbation strategies within these categories and have been shown to affect language models in the past. However, we have also included AdvGLUE as a complementary perturbation resource in our study that provides more effective perturbations. Nevertheless, we acknowledge that more complicated perturbations can also be created that are more effective and help the community have a more complete understanding of the 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. Moreover, focusing on the ethical considerations of this study, although the datasets and domains we focus on do not pose any societal harm, the potential harm that is associated with using the publicly available tools we used in this study to manipulate models in other critical domains should be considered. Issues surrounding anonymization and offensive content hold immense importance in data-driven studies, particularly in fields like natural language processing. Since we utilize datasets like IMDB, AG\_News, DBpedia, and AdvGLUE, due to the inherent nature of the datasets that do not involve human participation or potentially offensive or harmful data, and subsequent manual inspections, applying anonymization techniques or addressing offensive content is unwarranted.

## 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*. 702  
703  
704

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*. 705  
706  
707  
708

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. 709  
710  
711  
712  
713  
714

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. 715  
716  
717  
718

Anubrata Das, Chitrang 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. 719  
720  
721  
722  
723  
724  
725

Thomas Davidson, Dana Warmusley, Michael Macy, and Ingmar Weber. 2017. Automated hate speech detection and the problem of offensive language. In *Proceedings of the international AAAI conference on web and social media*, volume 11, pages 512–515. 726  
727  
728  
729  
730

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*. 731  
732  
733  
734

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. 735  
736  
737  
738  
739  
740  
741

Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. 2014. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*. 742  
743  
744

Xiaowei Gu and Weiping Ding. 2019. A hierarchical prototype-based approach for classification. *Information Sciences*, 505:325–351. 745  
746  
747

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. 748  
749  
750  
751  
752  
753

Peter Hase and Mohit Bansal. 2020. [Evaluating explainable AI: Which algorithmic explanations help users predict model behavior?](#) In *Proceedings of the* 754  
755  
756

757					
758					
759					
760	Peter Hase, Chaofan Chen, Oscar Li, and Cynthia Rudin.				
761	2019. Interpretable image recognition with hierar-				
762	chical prototypes. In <i>Proceedings of the AAAI Con-</i>				
763	<i>ference on Human Computation and Crowdsourcing</i> ,				
764	volume 7, pages 32–40.				
765	Junxian He, Taylor Berg-Kirkpatrick, and Graham Neu-				
766	big. 2020. Learning sparse prototypes for text gen-				
767	eration. <i>Advances in Neural Information Processing</i>				
768	<i>Systems</i> , 33:14724–14735.				
769	Dan Hendrycks, Xiaoyuan Liu, Eric Wallace, Adam				
770	Dziedzic, Rishabh Krishnan, and Dawn Song. 2020.				
771	<a href="#">Pretrained transformers improve out-of-distribution</a>				
772	<a href="#">robustness</a> . In <i>Proceedings of the 58th Annual Meet-</i>				
773	<i>ing of the Association for Computational Linguistics</i> ,				
774	pages 2744–2751, Online. Association for Computa-				
775	tional Linguistics.				
776	Dat Hong, Stephen S Baek, and Tong Wang. 2020. In-				
777	terpretable sequence classification via prototype tra-				
778	jectory. <i>arXiv preprint arXiv:2007.01777</i> .				
779	Dat Hong, Stephen S. Baek, and Tong Wang. 2021.				
780	<a href="#">Interpretable sequence classification via prototype</a>				
781	<a href="#">trajectory</a> .				
782	Myeongjun Jang, Deuk Sin Kwon, and Thomas				
783	Lukasiewicz. 2022. Becel: Benchmark for consis-				
784	tency evaluation of language models. In <i>Proceedings</i>				
785	<i>of the 29th International Conference on Computa-</i>				
786	<i>tional Linguistics</i> , pages 3680–3696.				
787	Di Jin, Zhijing Jin, Joey Tianyi Zhou, and Peter				
788	Szolovits. 2020. Is bert really robust? a strong base-				
789	line for natural language attack on text classification				
790	and entailment. In <i>Proceedings of the AAAI con-</i>				
791	<i>ference on artificial intelligence</i> , volume 34, pages				
792	8018–8025.				
793	Eunji Kim, Siwon Kim, Minji Seo, and Sungroh Yoon.				
794	2021. Xprotonet: Diagnosis in chest radiography				
795	with global and local explanations. In <i>Proceedings of</i>				
796	<i>the IEEE/CVF Conference on Computer Vision and</i>				
797	<i>Pattern Recognition (CVPR)</i> , pages 15719–15728.				
798	Alexey Kurakin, Ian Goodfellow, and Samy Bengio.				
799	2017a. <a href="#">Adversarial examples in the physical world</a> .				
800	Alexey Kurakin, Ian Goodfellow, and Samy Bengio.				
801	2017b. <a href="#">Adversarial machine learning at scale</a> .				
802	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan				
803	Ghazvininejad, Abdelrahman Mohamed, Omer Levy,				
804	Ves Stoyanov, and Luke Zettlemoyer. 2019. <a href="#">Bart: De-</a>				
805	<a href="#">noising sequence-to-sequence pre-training for natural</a>				
806	<a href="#">language generation, translation, and comprehension</a> .				
807	Jinfeng Li, Shouling Ji, Tianyu Du, Bo Li, and Ting				
808	Wang. 2018a. Textbugger: Generating adversarial				
809	text against real-world applications. <i>arXiv preprint</i>				
810	<i>arXiv:1812.05271</i> .				
	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 <i>Proceedings of the Thirty-Second</i>				
	<i>AAAI Conference on Artificial Intelligence and Thirti-</i>				
	<i>eth Innovative Applications of Artificial Intelli-</i>				
	<i>gence Conference and Eighth AAAI Symposium</i>				
	<i>on Educational Advances in Artificial Intelligence</i> ,				
	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. Explain-				
	aBoard: An Explainable Leaderboard for NLP. In				
	<i>Annual Meeting of the Association for Computational</i>				
	<i>Linguistics (ACL), System Demonstrations</i> .				
	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-				
	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,				
	Luke Zettlemoyer, and Veselin Stoyanov. 2019.				
	Roberta: A robustly optimized bert pretraining ap-				
	proach. <i>arXiv preprint arXiv:1907.11692</i> .				
	Andrew L. Maas, Raymond E. Daly, Peter T. Pham,				
	Dan Huang, Andrew Y. Ng, and Christopher Potts.				
	2011. <a href="#">Learning word vectors for sentiment analysis</a> .				
	In <i>Proceedings of the 49th Annual Meeting of the</i>				
	<i>Association for Computational Linguistics: Human</i>				
	<i>Language Technologies</i> , pages 142–150, Portland,				
	Oregon, USA. Association for Computational Lin-				
	guistics.				
	Shiao Meng, Xuming Hu, Aiwei Liu, Shu’ang Li, Fukun				
	Ma, Yawen Yang, and Lijie Wen. 2023. <a href="#">RAPL:</a>				
	<a href="#">A Relation-Aware Prototype Learning Approach</a>				
	<a href="#">for Few-Shot Document-Level Relation Extraction</a> .				
	<i>arXiv preprint arXiv:2310.15743</i> .				
	Pascal Mettes, Elise Van der Pol, and Cees Snoek. 2019.				
	Hyperspherical prototype networks. <i>Advances in</i>				
	<i>neural information processing systems</i> , 32.				
	Bonan Min, Hayley Ross, Elior Sulem, Amir				
	Pouran Ben Veyseh, Thien Huu Nguyen, Oscar Sainz,				
	Eneko Agirre, Ilana Heinz, and Dan Roth. 2021. <a href="#">Re-</a>				
	<a href="#">cent advances in natural language processing via</a>				
	<a href="#">large pre-trained language models: A survey</a> .				
	Yao Ming, Panpan Xu, Huamin Qu, and Liu Ren. 2019.				
	<a href="#">Interpretable and steerable sequence learning via pro-</a>				
	<a href="#">totypes</a> . In <i>Proceedings of the 25th ACM SIGKDD</i>				
	<i>International Conference on Knowledge Discovery</i>				
	<i>&amp; Data Mining</i> . ACM.				
	Milad Moradi and Matthias Samwald. 2021. <a href="#">Evaluating</a>				
	<a href="#">the robustness of neural language models to input</a>				
	<a href="#">perturbations</a> .				
	John X. Morris, Eli Lifland, Jin Yong Yoo, Jake Grigsby,				
	Di Jin, and Yanjun Qi. 2020. <a href="#">Textattack: A frame-</a>				
	<a href="#">work for adversarial attacks, data augmentation, and</a>				
	<a href="#">adversarial training in nlp</a> .				
	Meike Nauta, Annemarie Jutte, Jesper Provoost, and				
	Christin Seifert. 2021a. <a href="#">This looks like that, be-</a>				
	<a href="#">cause ... explaining prototypes for interpretable im-</a>				
	<a href="#">age recognition</a> . In <i>Communications in Computer</i>				

868			
869		<i>and Information Science</i> , pages 441–456. Springer International Publishing.	
870	Meike Nauta, Ron van Bree, and Christin Seifert. 2021b.		
871		<a href="#">Neural prototype trees for interpretable fine-grained image recognition</a> . In <i>Proceedings of the 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition – CVPR 2021</i> , pages 14933–14943, Nashville, TN, USA. IEEE.	
872			
873			
874			
875			
876	OpenAI. 2022. Chatgpt. <a href="https://openai.com/blog/chatgpt">https://openai.com/blog/chatgpt</a> . Accessed: April 30, 2023.		
877			
878	Frederik Pahde, Mihai Puscas, Tassilo Klein, and Moin Nabi. 2021. Multimodal prototypical networks for few-shot learning. In <i>Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)</i> , pages 2644–2653.		
879			
880			
881			
882			
883	Kamil Pluciński, Mateusz Lango, and Jerzy Stefanowski. 2021. Prototypical convolutional neural network for a phrase-based explanation of sentiment classification. In <i>Machine Learning and Principles and Practice of Knowledge Discovery in Databases</i> , pages 457–472, Cham. Springer International Publishing.		
884			
885			
886			
887			
888			
889			
890	Eleanor H. Rosch. 1973. <a href="#">Natural categories</a> . <i>Cognitive Psychology</i> , 4(3):328–350.		
891			
892	Victoria Rubin, Niall Conroy, Yimin Chen, and Sarah Cornwell. 2016. <a href="#">Fake news or truth? using satirical cues to detect potentially misleading news</a> . In <i>Proceedings of the Second Workshop on Computational Approaches to Deception Detection</i> , pages 7–17, San Diego, California. Association for Computational Linguistics.		
893			
894			
895			
896			
897			
898			
899	Cynthia Rudin, Chaofan Chen, Zhi Chen, Haiyang Huang, Lesia Semenova, and Chudi Zhong. 2022. Interpretable machine learning: Fundamental principles and 10 grand challenges. <i>Statistic Surveys</i> , 16:1–85.		
900			
901			
902			
903			
904	Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. <i>arXiv preprint arXiv:1910.01108</i> .		
905			
906			
907			
908	Sascha Saralajew, Lars Holdijk, and Thomas Villmann. 2020. <a href="#">Fast adversarial robustness certification of nearest prototype classifiers for arbitrary seminorms</a> . In <i>Advances in Neural Information Processing Systems</i> , pages 13635–13650. Curran Associates, Inc.		
909			
910			
911			
912			
913	Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. 2020. <a href="#">Superglue: Learning feature matching with graph neural networks</a> .		
914			
915			
916			
917	Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed Chi, Nathanael Schärli, and Denny Zhou. 2023. <a href="#">Large language models can be easily distracted by irrelevant context</a> .		
918			
919			
920			
	Dylan Slack, Satyapriya Krishna, Himabindu Lakkaraju, and Sameer Singh. 2022. <a href="#">Talktomodel: Explaining machine learning models with interactive natural language conversations</a> .		921
			922
			923
			924
	Jake Snell, Kevin Swersky, and Richard Zemel. 2017. <a href="#">Prototypical networks for few-shot learning</a> . In <i>Advances in Neural Information Processing Systems</i> , volume 30. Curran Associates, Inc.		925
			926
			927
			928
	Eduardo Soares, Plamen Angelov, and Neeraj Suri. 2022. <a href="#">Similarity-based deep neural network to detect imperceptible adversarial attacks</a> . In <i>2022 IEEE Symposium Series on Computational Intelligence (SSCI)</i> , pages 1028–1035.		929
			930
			931
			932
			933
	Zhivar Sourati, Vishnu Priya Prasanna Venkatesh, Darshan Deshpande, Himanshu Rawlani, Filip Ilievski, Hông-Ân Sandlin, and Alain Mermoud. 2023. <a href="#">Robust and explainable identification of logical fallacies in natural language arguments</a> . <i>Knowledge-Based Systems</i> , 266:110418.		934
			935
			936
			937
			938
			939
	Nilesh Tripuraneni, Ben Adlam, and Jeffrey Pennington. 2021. <a href="#">Overparameterization improves robustness to covariate shift in high dimensions</a> . In <i>Advances in Neural Information Processing Systems</i> , volume 34, pages 13883–13897. Curran Associates, Inc.		940
			941
			942
			943
			944
	Iulia Turc, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. <a href="#">Well-read students learn better: On the importance of pre-training compact models</a> .		945
			946
			947
	Betty van Aken, Jens-Michalis Papaioannou, Marcel G. Naik, Georgios Eleftheriadis, Wolfgang Nejdl, Felix A. Gers, and Alexander Löser. 2022. <a href="#">This patient looks like that patient: Prototypical networks for interpretable diagnosis prediction from clinical text</a> .		948
			949
			950
			951
			952
	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. <a href="#">Attention is all you need</a> .		953
			954
			955
			956
	Václav Voráček and Matthias Hein. 2022. <a href="#">Provably adversarially robust nearest prototype classifiers</a> . In <i>Proceedings of the 39th International Conference on Machine Learning</i> , volume 162 of the Proceedings of Machine Learning Research, pages 22361–22383, Baltimore, MD, USA. PMLR.		957
			958
			959
			960
			961
			962
	Kiri Wagstaff. 2012. <a href="#">Machine learning that matters</a> . <i>arXiv preprint arXiv:1206.4656</i> .		963
			964
	Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019. <a href="#">Glue: A multi-task benchmark and analysis platform for natural language understanding</a> .		965
			966
			967
			968
	Boxin Wang, Chejian Xu, Shuohang Wang, Zhe Gan, Yu Cheng, Jianfeng Gao, Ahmed Hassan Awadallah, and Bo Li. 2022a. <a href="#">Adversarial glue: A multi-task benchmark for robustness evaluation of language models</a> .		969
			970
			971
			972
			973



974	Jindong Wang, Xixu Hu, Wenxin Hou, Hao Chen,	are demonstrated in Table 3. We present both statis-	1027
975	Runkai Zheng, Yidong Wang, Linyi Yang, Haojun	tics about the original dataset and statistics and	1028
976	Huang, Wei Ye, Xiubo Geng, Binxin Jiao, Yue Zhang,	details about the number of perturbations that we	1029
977	and Xing Xie. 2023. <a href="#">On the robustness of chatgpt:</a>	have gathered on each dataset in three categories	1030
978	<a href="#">An adversarial and out-of-distribution perspective.</a>	of character-level, word-level, and sentence-level	1031
979	Shiqi Wang, Zheng Li, Haifeng Qian, Chenghao Yang,	perturbations. All the original datasets we use in	1032
980	Zijian Wang, Mingyue Shang, Varun Kumar, Samson	this study are gathered by other researchers and	1033
981	Tan, Baishakhi Ray, Parminder Bhatia, Ramesh Nal-	have been made public by them, mentioning non-	1034
982	lapati, Murali Krishna Ramanathan, Dan Roth, and	commercial use, which aligns with how we use	1035
983	Bing Xiang. 2022b. <a href="#">Recode: Robustness evaluation</a>	these datasets. We have included information on	1036
984	<a href="#">of code generation models.</a>	their descriptions and how they were gathered:	1037
985	Zhao Wang and Aron Culotta. 2020. <a href="#">Identifying spu-</a>	<b>IMDB.</b> This dataset is compiled from a set of	1038
986	<a href="#">rious correlations for robust text classification.</a> In	50000 reviews sourced from IMDB in English, lim-	1039
987	<i>Findings of the Association for Computational Lin-</i>	iting each movie to a maximum of 30 reviews. It	1040
988	<i>guistics: EMNLP 2020</i> , pages 3431–3440, Online.	has maintained an equal count of positive and neg-	1041
989	Association for Computational Linguistics.	ative reviews, ensuring a 50% accuracy through	1042
990	Chaowei Xiao, Bo Li, Jun yan Zhu, Warren He,	random guessing. To align with prior research	1043
991	Mingyan Liu, and Dawn Song. 2018. <a href="#">Generat-</a>	on polarity classification, the authors specifically	1044
992	<a href="#">ing adversarial examples with adversarial networks.</a>	focus on highly polarized reviews. A review is	1045
993	In <i>Proceedings of the Twenty-Seventh International</i>	considered negative if it scores $\leq 4$ out of 10 and	1046
994	<i>Joint Conference on Artificial Intelligence, IJCAI-18,</i>	positive if it scores $\geq 7$ out of 10. Neutral reviews	1047
995	pages 3905–3911. International Joint Conferences on	are excluded from this dataset. Authors have made	1048
996	Artificial Intelligence Organization.	the dataset publicly available, and you can find	1049
997	Hong-Ming Yang, Xu-Yao Zhang, Fei Yin, and Cheng-	more information about this dataset at <a href="https://ai.stanford.edu/~amaas/data/sentiment/">https://</a>	1050
998	Lin Liu. 2018. Robust classification with convolu-	<a href="https://ai.stanford.edu/~amaas/data/sentiment/">ai.stanford.edu/~amaas/data/sentiment/</a> .	1051
999	tional prototype learning. In <i>Proceedings of the IEEE</i>	<b>AG_News.</b> This dataset comprises over 1 million	1052
1000	<i>conference on computer vision and pattern recogni-</i>	English news articles sourced from 2000+ news	1053
1001	<i>tion</i> , pages 3474–3482.	outlets over a span of more than a year by Come-	1054
1002	Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015.	ToMyHead, an academic news search engine op-	1055
1003	Character-level convolutional networks for text classi-	erational since July 2004. Provided by the aca-	1056
1004	fication. <i>Advances in neural information processing</i>	ademic community, this dataset aids research in	1057
1005	<i>systems</i> , 28.	data mining, information retrieval, data compres-	1058
1006	Haiyan Zhao, Hanjie Chen, Fan Yang, Ninghao Liu,	sion, data streaming, and non-commercial activi-	1059
1007	Huiqi Deng, Hengyi Cai, Shuaiqiang Wang, Dawei	ties. This news topic classification dataset features	1060
1008	Yin, and Mengnan Du. 2023. Explainability for	four classes: world, sports, business, and science.	1061
1009	large language models: A survey. <i>arXiv preprint</i>	The details about the intended use and access condi-	1062
1010	<i>arXiv:2309.01029</i> .	tions are provided at <a href="http://www.di.unipi.it/~gulli/AG_corpus_of_news_articles.html">http://www.di.unipi.it/</a>	1063
1011	Pei Zhou, Rahul Khanna, Seyeon Lee, Bill Yuchen	<a href="http://www.di.unipi.it/~gulli/AG_corpus_of_news_articles.html">~gulli/AG_corpus_of_news_articles.html</a> .	1064
1012	Lin, Daniel Ho, Jay Pujara, and Xiang Ren.	<b>DBpedia.</b> DBpedia <sup>4</sup> seeks to extract organized	1065
1013	2020. Rica: Evaluating robust inference capabili-	information from Wikipedia’s vast content. The	1066
1014	ties based on commonsense axioms. <i>arXiv preprint</i>	gathered subset of data we used offers hierar-	1067
1015	<i>arXiv:2005.00782</i> .	chical categories for 342782 Wikipedia articles.	1068
1016	Julia El Zini and Mariette Awad. 2022. <a href="#">On the explain-</a>	These classes are distributed across three lev-	1069
1017	<a href="#">ability of natural language processing deep models.</a>	els, comprising 9, 70, and 219 classes, respec-	1070
1018	<i>ACM Comput. Surv.</i> , 55(5).	tively. We used the version that has nine classes:	1071
1019	Barret Zoph, Irwan Bello, Sameer Kumar, Nan Du, Yan-	Agent, Work, Place, Species, UnitOfWork, Event,	1072
1020	ping Huang, Jeff Dean, Noam Shazeer, and William	SportsSeason, Device, and TopicalConcept. Al-	1073
1021	Fedus. 2022. <a href="#">St-moe: Designing stable and transfer-</a>	though the articles are in English, specific names	1074
1022	<a href="#">able sparse expert models.</a>	(e.g., the name of a place or person) can be	1075
1023	<b>A Dataset Details</b>		
1024	<b>A.1 Dataset Statistics</b>		
1025	The statistics of the datasets we used in this study		
1026	to test the robustness of PBNs against perturbations		

<sup>4</sup><https://www.dbpedia.org/>

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 3: 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, respectively. 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.

Type	Perturbed Sentence
Char.	<u>Disney</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 4: Perturbation examples on the following sentence: *Disney Rules Out New Deal with Pixar Studios*. The underlined text indicates the new tokens that replace the original sentence tokens.

non-English. Find more information about this dataset at [https://huggingface.co/datasets/DeveloperOats/DBPedia\\_Classes](https://huggingface.co/datasets/DeveloperOats/DBPedia_Classes).

**AdvGLUE.** Adversarial GLUE (AdvGLUE) (Wang et al., 2022a) introduces a multi-task English benchmark designed to investigate and assess the vulnerabilities of modern large-scale language models against various adversarial attacks. It encompasses five corpora, including SST-2 sentiment classification, QQP paraphrase test dataset, and QNLI, RTE, and MNLI, all of which are natural language inference datasets. To assess robustness, perturbations are applied to these datasets through both automated and human-evaluated methods, spanning word-level, sentence-level, and human-crafted examples. Our focus primarily centers on SST-2 due to its alignment with the other covered datasets in our study and its classification nature. This dataset has been made public by the authors and is released with CC BY-SA 4.0 license.

## A.2 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 4.

## B Implementation Details

### B.1 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, and each experiment took between 10 minutes to 2 hours, depending on the dataset and model used. 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, and GPT3.5 (OpenAI, 2022) for sentence-level perturbations.

### B.2 Training 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 vanilla 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 models for evaluation, 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 in the evaluation step, we use BERT-Medium (Turc et al., 2019) as a model that is different from BERT used in our perturbation gathering step, which being smaller than other Transformer

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.

models included in the evaluation (BART and Electra), also allows us to assess the sensitivity of PBNs to the backbone size.

For all the datasets, the training split, validation split, and test splits were used from <https://huggingface.co/>. During training on the IMDB, SST-2, and DBPedia datasets, 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. The results are shown as the average of three runs.

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 simulating real noisy data 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.

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.

### B.3 GPT-3.5 Baseline

We used GPT3.5 as a baseline in our experiments to compare its performance on original and perturbed examples with PBN and their vanilla models. In this section, we present the prompts that we gave to GPT-3.5 to generate the baseline responses and the reported performance in Table 6. We used the following prompts for the four different datasets:

IMDB: *Identify the binary sentiment of the following text: [text]. Strictly output only "negative" or "positive" according to the sentiment and nothing else. Assistant:*

AG\_News: *Categorize the following news strictly into only one of the following classes: world, sports, business, and science. Ensure that you output only the category name and nothing else. Text: [text]. Assistant:*

DBPedia: *Categorize the following text article strictly into only one taxonomic category from the following list: Agent, Work, Place, Species, UnitOfWork, Event, SportsSeason, Device, and TopicalConcept. Ensure that you output only the category name and nothing else. Text: [text]. Assistant:*

SST-2: *Identify the binary sentiment of the following text: [text]. Strictly output only "negative" or "positive" according to the sentiment and nothing else. Assistant:*

## C Additional Experiments

In a further experiment, we measure the impact of the size of the backbone on the robustness properties of PBNs (see Figure 6). The results show a



	IMDB			AG_News			DBPedia			SST-2
Model	Char	Word	Sent	Char	Word	Sent	Char	Word	Sent	-
GPT-3.5 (Test)	<b>83.5</b>	<b>86.7</b>	82.0	<b>66.0</b>	<b>75.6</b>	71.4	42.5	<b>44.9</b>	38.0	94.6
GPT-3.5 (Adv)	78.4	80.3	<b>90.1</b>	63.9	71.3	<b>72.5</b>	<b>42.7</b>	41.8	<b>46.0</b>	62.5
$\Delta$	<b>-5.1</b>	<b>-6.4</b>	<b>+8.1</b>	<b>-2.1</b>	<b>-4.3</b>	<b>+1.1</b>	<b>+0.2</b>	<b>-3.1</b>	<b>+8.0</b>	<b>-32.1</b>

Table 6: F1 scores of ChatGPT on four datasets and their perturbed versions.  $\Delta$  indicates the difference between the performance on the original data points and perturbed data points.

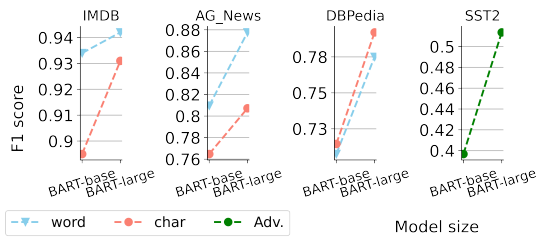


Figure 6: F1 of PBN with BART-base and BART-large.

1225 boost in performance when scaling up the backbone  
1226 used in PBNs regardless of the training dataset,  
1227 which is in line with theoretical work (Tripuraneni  
1228 et al., 2021) proving that overparameterized mod-  
1229 els exhibit enhanced robustness to perturbations.  
1230 The same observation holds across different archi-  
1231 tectures, too, with larger models being more robust.  
1232 We conclude that larger Transformer architectures  
1233 are optimal to ensure the robust behavior of PBNs.