

Robust Text Classification: Analyzing Prototype-Based Networks

Anonymous ACL submission

Abstract

Downstream applications often require text classification models to be accurate and robust. While the accuracy of the state-of-the-art Language Models (LMs) approximates human performance, they often exhibit a drop in performance on noisy data found in the real world. This lack of robustness can be concerning, as even small perturbations in the text, irrelevant to the target task, can cause classifiers to incorrectly change their predictions. 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 has been shown to be robust to noise for computer vision tasks. In this paper, we study whether the robustness properties of PBNs transfer to text classification tasks under both targeted and static adversarial attack settings. Our results show that PBNs, as a mere architectural variation of vanilla LMs, offer more robustness compared to vanilla LMs under both targeted and static settings. We showcase how PBNs' interpretability can help us to understand PBNs' robustness properties. Finally, our ablation studies reveal the sensitivity of PBNs' robustness to how strictly clustering is done in the training phase, as tighter clustering results in less robust PBNs.

1 Introduction

Language models (LMs) are widely used in various NLP tasks and exhibit exceptional performance (Chowdhery et al., 2022; Zoph et al., 2022). In light of the need for real-world applications of these models, the requirements for robustness and interpretability have become urgent for both Large Language Models (LLMs) and fine-tuned LMs. 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, LMs have limited interpretability by design

(Zhao et al., 2023; Gholizadeh and Zhou, 2021), which cannot be fully mitigated by posthoc explainability techniques (Zini and Awad, 2022). Moreover, LMs lack robustness when exposed to text perturbations, noisy data, or distribution shifts (Jin et al., 2020; Moradi and Samwald, 2021). Reportedly, even LLMs 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 this paper, too.

On this ground, NLP research has increasingly focused on benchmarks, methods, and studies that emphasize robustness and interpretability (e.g., Zhou et al., 2020; Jang et al., 2022; Liu et al., 2021). This has also been accompanied by the surge of focus on models that are inherently and architecturally interpretable and robust (e.g., Koh et al., 2020; Papernot and McDaniel, 2018; Keane and Kenny, 2019). An example of such models is the family of Prototype-Based Networks (PBNs) that is designed for robustness and interpretability (Li et al., 2018b). PBNs are based on the theory of categorization in cognitive science (Rosch, 1973), where it is governed by the graded degree of possessing prototypical features of different categories, 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. Aligned with how humans approach classification (Linzen, 2020), classifications in PBNs are expected to have human-like robustness because they classify through distances to prototypical examples found in the data. Leveraging distance between points helps to quantify prototypicality,

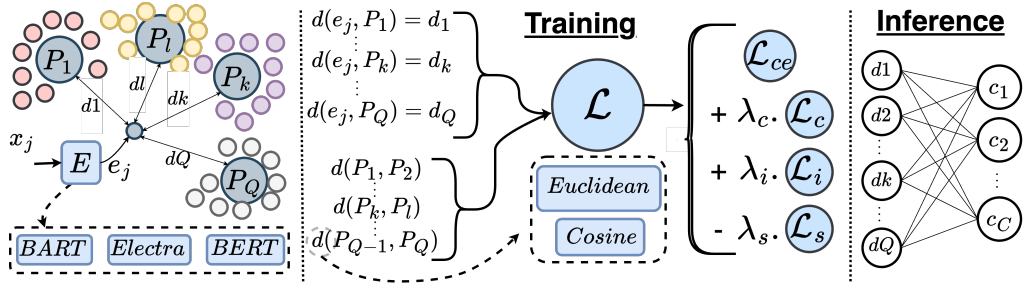


Figure 1: Classification by a PBN. The model computes distances between the new point and prototypes, $d(e_j, P_k)$, and distances within prototypes, $d(P_k, P_l)$, for both inference and training. During training, the model minimizes the loss term, \mathcal{L} , consisting of \mathcal{L}_{ce} , $\lambda_c \mathcal{L}_c$, $\lambda_i \mathcal{L}_i$, $\lambda_s \mathcal{L}_s$, controlling the importance of accuracy, clustering, interpretability, and separation of prototypes, based on all the computed distances; during inference, distances between the new point and prototypes are used for classification by a fully connected layer.

084 which then facilitates identifying noisy or out-of-
085 distribution samples (Yang et al., 2018).

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

103 In this study, we investigate whether the robust-
104 ness properties of PBNs transfer to NLP classifica-
105 tion tasks. In particular, our contributions are: (1)
106 We adopt a modular and comprehensive approach
107 to evaluate PBNs’ robustness properties against
108 various well-known adversarial attacks under both
109 targeted and static adversarial settings; (2) We con-
110 duct a comprehensive analysis of the sensitivity of
111 PBNs’ robustness w.r.t. different hyperparameters.

112 Our experiments show that PBNs’ robustness
113 transfers to realistic perturbations in text classifica-
114 tion tasks under both targeted and static adversarial
115 settings and can, thus, enhance the text classifica-
116 tion robustness of LMs. We note that the robustness
117 boost that adversarial augmented training brings
118 to LMs with access to additional pieces of rele-
119 vant data, is higher than the boost caused by PBNs’
120 architecture. Nevertheless, considering that the

robustness boost in PBNs is only caused by their
architecture without any additional resources (data
or parameters), and this architecture is interpretable
by design, the merits of such models can contribute
to the field. Finally, benefiting from inherent inter-
pretability, we showcase how PBN interpretability
properties help to explain PBNs’ robust behavior.

2 Prototype-Based Networks

PBNs classify data points based on their similarity
to a set of *prototypes* learned during training. These
prototypes summarize prominent semantic patterns
of the dataset through two mechanisms: (1) proto-
types are defined in the same embedding space as
input examples, which makes them interpretable
by leveraging input examples in their proximity;
and (2) prototypes are designed to cluster semanti-
cally similar training examples, which makes them
representative of the prominent patterns embed-
ded in the data and input examples. The PBN’s
decisions, based on quantifiable similarity to proto-
types, are robust as noise and perturbations are bet-
ter reflected in the computed similarity to familiar
prototypical patterns (Hong et al., 2020). Addition-
ally, prototypes can provide insight during infer-
ence by helping users explain the model’s behavior
on input examples through the prototypes utilized
for the model’s prediction (Das et al., 2022).

Inference. Classification in PBNs is done via
a fully connected layer applied on the measured
distances between embedded data points and pro-
totypes. 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 em-
bedding $e_j = E(x_j)$. Next, PBNs compute the

distances between prototypes and e_j using the function d . These distances get fed into a fully connected layer to compute class-wise logits, incorporating the similarities to each prototype. Applying a softmax on top, the final outputs are $\hat{y}_c(x_j)$: 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, clustering, 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:

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

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

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

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.

Original text	Perturbed text
A gentle breeze rustled the leaves.	A gentle wind rustled the leaves.
rescue Engineer Company	Rescue operation Company
embarrassingly foolish	embarrassingly foolish

Table 1: Examples of adversarial perturbations, with the perturbed tokens highlighted.

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 Robustness Evaluation

We assess PBNs’ robustness against adversarial perturbations of original input text that are intended to preserve the text’s original meaning. The perturbations change the classification of the target model upon confronting these perturbed examples from the correct behavior to an incorrect one in an effective and efficient way (Dalvi et al., 2004; Kurakin et al., 2017a,b; Li et al., 2023). Automatic approaches of finding these perturbations vary (Zhang et al., 2020): perturbations can be focused on different **granularities**, i.e., *character-level*, *word-level*, or *sentence-level*; their **generation** can be done in different ways, e.g., *replacing*, *inserting*, *deleting*, *swapping* tokens; they can have different **searching strategies** for their manipulations, such as *context-aware* or *isolated* approaches; and also various **salient token identification strategies** to maximize their adversarial effect.

Orthogonally, these adversarial perturbations are divided into **targeted** and **static**. In the targeted setting, the attacker has access to the target model and can attack it directly (Si et al., 2021). However, in the static setting, the attacker does not have access to the target model. Hence, adversarial perturbations are gathered while attacking external models that the attacker has access to, and the gathered successful perturbations would be used to assess the robustness of the target model (Wang et al., 2022a).

With numerous adversarial perturbation strategies in the literature (Zhang et al., 2020; Wang et al., 2022c), each with unique advantages (e.g., effectiveness vs. efficiency), we use a wide range of existing perturbation strategies in this study. These cover the aforementioned **granularities**, **generation strategies**, **searching strategies**, and **salient token identification strategies**, under both tar-

241 **geted**, and **static** settings. See examples of adver- 286
242 sarial perturbations covered in our study in Table 1. 287

243 4 Experimental Setup 288

244 4.1 Datasets 289

245 PBNs classify instances based on their similarity to 290
246 prototypes learned during training that summarize 291
247 prominent semantic patterns in a dataset. Thus, 292
248 with more classes, we might need more prototypes 293
249 to govern the more complex system between in- 294
250 stances and prototypes (Yang et al., 2018). To study 295
251 the interplay between the number of classes and 296
252 robustness, we employ three datasets: (1) *IMDB* 297
253 *reviews* (Maas et al., 2011): a binary sentiment 298
254 classification dataset; (2) *AG_NEWS* (Gulli): a col- 299
255 lection of news articles that can be associated with 300
256 four categories; (3) *DBPedia*:¹ a dataset with taxo- 301
257 nomic, hierarchical categories for Wikipedia arti- 302
258 cles (Lehmann et al., 2015), with nine classes. We 303
259 use these three datasets to study the robustness of 304
260 PBNs under both targeted and static adversarial set- 305
261 tings. As an additional source of static adversarial 306
262 perturbations, we adopt the SST-2 binary classifi- 307
263 cation split from the existing *Adversarial GLUE* 308
264 (*AdvGLUE*) dataset (Wang et al., 2022a), consist- 309
265 ing of perturbed examples of different granularities, 310
266 filtered both automatically and by human evalua- 311
267 tion for more effectiveness. For statistics of the 312
268 datasets and their perturbations, see Appendix A. 313

269 4.2 Perturbations 314

270 **Attacking strategies.** We selected five well- 315
271 established adversarial attack methods: BAE (Garg 316
272 and Ramakrishnan, 2020), TextFooler (Jin et al., 317
273 2020), TextBugger (Li et al., 2018a), DeepWord- 318
274 Bug (Gao et al., 2018), and PWWS (Ren et al., 319
275 2019).² As mentioned in Section 3, these at- 320
276 tacks cover a wide range of **granularities** (e.g., 321
277 character-based in DeepWordBug and word-based 322
278 in PWWS), **generation strategies** (e.g., word 323
279 substitution in PWWS and TextFooler and dele- 324
280 tion in TextBugger), **searching strategies** (e.g., 325
281 context-aware in BAE and isolated synonym-based 326
282 in TextFooler), and **salient token identification** 327
283 **strategies** (e.g., finding the important sentences 328
284 first and then words in TextBugger and finding the 329
285 important words to change in BAE). 330

¹<https://bit.ly/3RgX41H>

²We also employed paraphrased-based perturbations (Lei et al., 2019), generated by GPT3.5 (OpenAI, 2022). However, both our baselines and PBNs were robust to these perturbations, and we include them in the Appendix in Table 6. 331

Targeted perturbations. In this setting, the ad- 286
287 versarial attacks are directly conducted against 288
289 PBNs and vanilla LMs trained on original datasets. 290
291 For each attack strategy, we aim for 800 successful 292
293 perturbations and report the robustness of PBNs 294
295 against adversarial attacks by Attack Success Rate 296
297 (ASR; Wu et al., 2021) and Average Percentage of 298
299 Words Perturbed (APWP; Yoo et al., 2020) to reach 300
301 the observed ASR. Successful perturbations are 302
303 those that change the prediction of a target model 304
305 already fine-tuned on that dataset from the correct 306
307 prediction to the wrong prediction. 308

Static perturbations. In this setting, the adver- 298
299 sarial attacks are conducted on external models: 300
301 BERT (Devlin et al., 2018), RoBERTa (Liu et al., 302
303 2019), and DistilBERT (Sanh et al., 2019), which 304
305 are trained on the original datasets, and a compila- 306
307 tion of the successful perturbations on those models 308
309 is used to assess the robustness of PBNs against the 309
310 studied adversarial attacks by their accuracy on the 310
311 perturbations, similar to the study by Wang et al. 311
312 (2022a). To obtain the perturbations, each model 312
313 is fine-tuned on each dataset, and 800 successful 313
314 perturbations for each attack strategy are obtained. 314
315 We focus on examples whose perturbations are pre- 315
316 dicted incorrectly by all three models to maximize 316
317 the generalizability of this static set of perturbations 317
318 to a wider range of unseen target models. In princi- 318
319 ple, the perturbations for each model are different, 319
320 yielding three variations per original example for 320
321 a dataset-perturbation pair. For instance, focusing 321
322 on DBPedia and BAE attack strategy, after 800 322
323 successful perturbations for each of the three target 323
324 models, the perturbations of 347 original examples 324
325 could change all models’ predictions, resulting in a 325
326 total of 1401 (3×347) perturbations compiled for 326
327 BAE attack strategy and DBPedia dataset. 327

328 4.3 PBNs’ Hyperparameters 323

Backbone (E). Prototype alignment and training 324
325 are highly dependent on the quality of the latent 325
326 space created by the backbone encoder E , which 326
327 in turn affects the performance, robustness, and 327
328 interpretability of PBNs. We consolidate previous 328
329 methods for text classification using PBNs (Plu- 329
330 ciński et al., 2021; Das et al., 2022; Ming et al., 330
331 2019; Hong et al., 2020) and consider three back- 331
332 bone architectures: BERT (Devlin et al., 2018), 332
333 BART encoder (Lewis et al., 2019), and Electra 333
334 (Clark et al., 2020). Based on our empirical evi- 334
335 dence, fine-tuning all the layers of the backbone 335

was causing the PBNs’ training not to converge. Hence, we freeze all the layers of the backbones except for the last layer when training.

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 incompatible 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, we consider both Cosine and Euclidean distance functions when training PBNs.

Number of prototypes (Q). Number of prototypes in PBNs is a key factor for mapping difficult data distributions (Yang et al., 2018; Sourati et al., 2023). Hence, to cover a wide range, we consider five values for $Q = \{2, 4, 8, 16, 64\}$.

Objective functions (\mathcal{L}). Given the partly complementary goals of loss terms, we investigate the effect of interpretability, clustering, and separation loss on PBNs’ robustness, keeping the accuracy constraint (\mathcal{L}_{ce}) intact. To do so, we consider three values, $\{0, 0.9, 10\}$ for λ_i , λ_c , and λ_s . 0 value represents the condition where the corresponding loss function is not being utilized in the training process. 0.9 value was empirically found to offer good accuracy, clustering, and interpretability, across datasets and was also motivated by prior works (Das et al., 2022). 10 value was chosen as an upper bound dominating the corresponding loss objective (e.g., interpretability) in the training process.

4.4 Baselines

Since PBNs are architectural enhancements of vanilla LMs using learned prototypes for classification instead of a traditional softmax layer used in vanilla LMs, **vanilla LMs** employed as PBNs’ backbones serve as a baseline for comparing the robustness of PBNs. We also employ **adversarial augmented training** (Goyal et al., 2023) on top of the vanilla LMs as another baseline. Note that

the same layers frozen for PBNs’ training are also frozen for the baselines. As we need additional data for such extra training, we use this baseline under static perturbations, where the set of perturbations has already been compiled beforehand. Finally, although we note that LLMs are more appropriate choices for generic chat and text generation due to their decoder-only architecture, and fine-tuned LMs are still superior to LLMs when it comes to task-oriented performance (Chang et al., 2024), we compare PBNs with two LLMs, namely, GPT4o (AI, 2024) and Llama3 (AI@Meta, 2024).

5 Results

5.1 Robustness of PBNs

The robustness report of PBNs under both targeted adversarial attacks and static attacks under different experimental setups (i.e., datasets, backbones, and attack strategies), using the best hyperparameters is presented in Table 2.^{3 4} Best hyperparameters were chosen among the permutation of all hyperparameters presented in Section 4.3 to yield the highest robustness (lowest ASR or highest accuracy). Under the targeted adversarial attack setting, our results showed that PBNs are more robust than vanilla LMs (having lower ASR) regardless of the utilized backbone, dataset, or attacking strategy. We also saw similar trends analyzing the robustness of PBNs compared to vanilla LMs, averaging over all PBN hyperparameters (find the details in Table 8). Focusing on the APWP metric, we observed that in 71.0% of the conditions, the PBNs’ robustness was greater than vanilla LMs (having higher APWP), and this superiority dropped to 31.0% of the conditions when averaging over all the hyperparameters (find the details in Table 7), which suggested that PBNs’ robustness is sensitive to hyperparameters involved in training.

We observed similar trends under static adversarial attacks, where the PBNs’ robustness was higher than vanilla LMs (having higher accuracy under attack) in the majority of the conditions (93.7% of all variations of experimental setups and hyperparameters). We observed that in every experimental condition (dataset and attack strategy), a PBN exists with a robustness outperforming LLMs like GPT4o (AI, 2024) and Llama3 (AI@Meta,

³The semantic similarity between original and perturbed texts using OpenAI text-embedding-ada-002 across all datasets and attack types was 0.97 ($SD = 0.01$).

⁴Our results showed that adversarial perturbations from TextFooler and PWWS were more effective than others.

Targeted Attacks; Attack Success Rate (ASR %) reported

	AG_News					DBPedia					IMDB				
	BAE	DWB	PWWS	TB	TF	BAE	DWB	PWWS	TB	TF	BAE	DWB	PWWS	TB	TF
BART	14.8	53.2	53.6	31.8	76.5	18.9	28.3	43.1	21.1	71.9	74.1	74.7	99.3	78.5	100.0
+ PBN	11.1	32.3	41.3	23.1	62.2	15.2	14.7	28.7	12.6	45.5	36.1	41.0	75.9	41.3	73.1
BERT	17.0	78.0	69.8	45.7	88.8	13.9	24.8	31.6	22.0	61.3	82.5	79.7	99.9	83.9	99.9
+ PBN	7.7	42.6	47.0	30.4	70.5	9.8	17.3	21.6	13.0	41.0	42.8	41.0	79.7	57.7	79.8
ELEC.	24.8	89.5	69.1	87.8	87.9	14.5	42.8	45.6	42.3	75.3	52.5	49.2	95.3	67.8	99.3
+ PBN	14.0	34.9	42.9	51.8	70.2	7.8	11.5	17.8	19.1	35.6	28.9	27.4	66.6	36.8	78.0

Static Attacks; Accuracy (%) reported

	AG_News					DBPedia					IMDB					SST2
	BAE	DWB	PWWS	TB	TF	BAE	DWB	PWWS	TB	TF	BAE	DWB	PWWS	TB	TF	GLUE
BART	53.2	76.7	83.2	77.5	85.8	55.5	68.6	58.4	72.5	71.3	74.1	80.5	83.6	85.8	87.6	29.8
+ PBN	<u>57.6</u>	80.6	<u>84.8</u>	79.2	<u>88.8</u>	<u>65.0</u>	<u>71.6</u>	<u>65.7</u>	<u>78.4</u>	<u>74.8</u>	<u>80.4</u>	<u>81.3</u>	<u>86.3</u>	<u>89.3</u>	<u>90.4</u>	50.4
+ Aug.	71.7	<u>78.4</u>	85.5	<u>77.6</u>	90.1	84.0	79.6	89.7	88.8	94.0	85.7	86.7	92.9	89.9	96.5	-
BERT	47.8	64.0	75.9	69.4	80.7	62.3	61.4	75.4	78.4	82.0	75.1	77.1	85.0	83.4	85.9	42.0
+ PBN	<u>52.9</u>	<u>70.4</u>	78.5	73.8	<u>84.3</u>	<u>66.9</u>	<u>66.6</u>	<u>80.3</u>	<u>82.0</u>	<u>85.8</u>	<u>77.6</u>	79.1	<u>85.3</u>	<u>85.0</u>	<u>86.5</u>	51.1
+ Aug.	58.3	71.6	<u>78.3</u>	<u>71.2</u>	85.4	75.5	70.9	84.1	90.5	91.0	83.2	<u>77.6</u>	91.7	90.8	89.2	-
ELEC.	50.4	<u>65.0</u>	<u>73.5</u>	<u>63.9</u>	<u>77.8</u>	<u>79.7</u>	66.9	<u>80.9</u>	81.4	84.4	<u>89.7</u>	90.3	<u>94.6</u>	94.5	95.6	44.3
+ PBN	64.6	74.1	85.1	77.2	89.0	<u>78.7</u>	<u>69.8</u>	<u>79.3</u>	<u>82.5</u>	<u>85.8</u>	90.0	<u>90.8</u>	<u>94.6</u>	95.5	96.3	65.6
+ Aug.	<u>55.0</u>	59.5	71.7	61.6	<u>79.5</u>	86.2	73.8	88.1	84.5	92.8	89.4	93.7	95.3	<u>94.9</u>	<u>95.8</u>	-
GPT4o	57.1	73.3	73.0	76.5	79.9	66.0	63.4	61.0	69.0	44.0	87.0	89.5	91.2	93.7	94.2	59.8
Llama3	57.6	56.4	55.0	65.9	62.8	44.0	53.7	37.8	45.0	44.4	82.0	86.0	93.2	89.0	91.5	56.0

Table 2: Comparison of PBNs and vanilla LMs (+ vanilla LMs with adversarial augmented training under static attack setting) under both targeted and static adversarial attack perturbations, using the best hyperparameters for PBNs, on IMDB, AG_News, DBPedia (+ SST-2 from AdvGLUE under static attack setting) datasets, under BAE, DeepWordBug (DWB), PWWS, TextBugger (TB), TextFooler (TF). The highest accuracy and lowest ASR showing the superior model for each architecture is **boldfaced**, and the second best model is underlined for static attacks.

2024) that have orders of magnitude more parameters and are not interpretable by design as opposed to PBNs. Vanilla LMs with adversarial augmented training demonstrated greater robustness than PBNs in 71.2% of the conditions. This highlighted the more effective role of additional data in adversarial augmented training compared to PBNs’ robust architecture and makes PBNs a preferable choice when efficiency is prioritized (Goodfellow et al., 2014). Analyzing PBNs’ robustness under the static adversarial setting averaging over all PBNs’ hyperparameters, our results showed that in only 31.2% of the conditions, PBNs have greater robustness compared to vanilla LMs (find the details in Table 8), which similar to observations on APWP, suggested that PBNs’ robustness is sensitive to hyperparameters involved in the training.

To sum up, we observed that PBNs consistently and over different metrics were more robust compared to vanilla LMs and LLMs, using the best hyperparameters without sacrificing performance on the original unperturbed samples (find performance on original datasets in Table 6). We believe that the observed 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), espe-

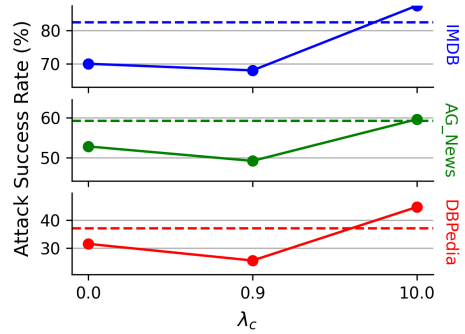


Figure 2: Attack Success Rate (ASR %) of PBNs with different λ_c values adjusting the importance of clustering in the trained PBNs, with other hyperparameters set to their best values, and averaged across other possible variables (e.g., backbone and attack type). The dotted line represents the ASR for the non-PBN model.

cially 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. Finally, we observed that the robustness superiority of PBNs compared to vanilla LMs dropped when averaging over all the possible hyperparameters, which is what we investigate further in Section 5.2.

5.2 Sensitivity to Hyperparameters

We studied the sensitivity of PBNs’ robustness to the hyperparameters involved in training, covering

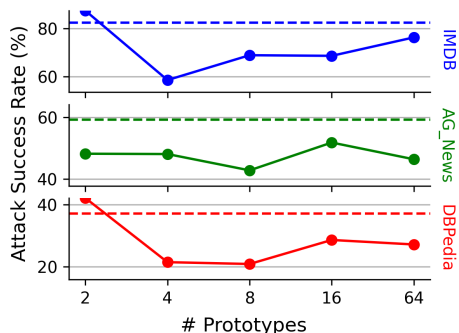


Figure 3: Attack Success Rate (ASR %) of PBNs with different numbers of prototypes, with other hyperparameters set to their best values, and averaged across other possible variables (e.g., backbone and attack type). Dotted line represents the ASR for the non-PBN model.

values discussed in Section 4.3. Focusing on each hyperparameter, the value for the other ones was selected to yield the best performance so that, overall, we could better depict the sensitivity and limiting effect of the hyperparameter of interest. We did not observe any sensitivity from PBNs with respect to the backbone, interpretability term (λ_i ; see Section C.5), separation term (λ_s ; see Section C.7), and the distance function (d ; see Section C.4).

However, as presented in Figure 2, we observed that higher values of λ_c , promoting tighter clustering of input examples around prototypes, hinder PBNs’ robustness. Clustering loss is a regularization term that encourages samples to be close to 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 $(-0.24 \pm 1.7) \times 10^{-7}$ with $\lambda_c = 0.9$, and $(-0.18 \pm 1.5) \times 10^{-6}$ with $\lambda_c = 10$, showing less diverse prototypes indicated by smaller measured distances caused by stronger clustering.

Additionally, as depicted in Figure 3, we observed poor robustness from PBNs when the number of prototypes is as low as two, which is intuitive as a low number of prototypes also means a lower number of semantic patterns learned, which constraints the PBNs’ abilities to distinguish between different classes. Noting that more prototypes add to the complexity and size of the network as a whole, the observed stable trend of the robustness with the higher number of prototypes (> 2)

Proto.	Representative Training Examples	Label
P_0	Handly’s Lessee v. Anthony (1820) : Determined Indiana-Kentucky boundary.	UnitWork
	Rasul v. Bush (2004) : Decided jurisdiction over Guantanamo detainees.	UnitWork
P_1	Marine Corps Air Station Futenma : U.S. Marine Corps base, Ginowan, Okinawa; regional military hub.	Place
	Özdere : Turkish coastal resort town in İzmir Province, popular among tourists.	Place
P_2	Yevgeni Viktorovich Balyaikin : Russian footballer for FC Tom Tomsk.	Agent
	Gigi Morasco : Fictional character on ABC’s One Life to Live.	Agent

Table 3: Examples of prototypes, their closest training examples, alongside their label derived from their closest training examples, extracted from a PBN with 16 prototypes and a BART backbone on DBpedia. Note that the presented training examples are the summarization of their longer version for easier interpretation.

suggests that as long as the number of prototypes is not too low, PBNs with lower number of prototypes can be preferred. This corroborates with the studies performed by Yang et al. (2018). Finally, note that the same analysis using other metrics (e.g., APWP) and under static adversarial setting (using accuracy as the studied metric) depicted the same trend and can be found in Section C.6 and Section C.8.

5.3 PBNs’ Interpretability w.r.t. Robustness

PBNs are interpretable by design, and we can understand their behavior through the distance of input examples to prototypes and the importance of these distances, extracted by the last fully connected layer of PBNs transforming vector of distances to log probabilities for classes. Examples of learned prototypes that can be represented by their closest training input examples are shown in Table 3. These input examples help the user identify the semantic features that the prototypes are associated with, which by our observations in our case, were mostly driven by the class label of the closest training examples.

We can also benefit from interpretable properties of PBNs to better understand their robustness properties, regardless of the success of perturbations. Table 4 illustrates predictions of a PBN on three original and perturbed examples from the DBpedia dataset, alongside the top-2 prototypes that were utilized by the PBN’s fully connected layer for prediction and prototypes’ associated label (by their closest training examples). In the first two examples, PBN correctly classifies both the original and perturbed examples, and from the top-2 prototypes, we observe that this is due to unchanged prototypes

Text	Activ. Proto.s	Proto.s Labels	Pred.	Label
Roman Catholic Diocese of Barra: Diocese in Barra, Feira de Santana province, Brazil.	P_1, P_{14}	Place, Place	Place	Place
Roman Catholic Bishop of Barra: Episcopal seat in Barra, Feira de Santana province, Brazil.	P_1, P_{14}	Place, Place	Place	Place
Inta Ezergailis: Latvian American professor emerita at Cornell University.	P_2, P_8	Agent, Agent	Agent	Agent
Inta Ezergailis: Latvian American poet and scholar at Cornell University.	P_2, P_7	Agent, Work	Agent	Agent
Saint Eigrad: 6th-century Precongregational North Wales saint and Patron Saint of Llaneigrad.	P_2, P_8	Agent, Agent	Agent	Agent
<i>St Eigrad:</i> 6th-century Precongregational street of North Wales and Patron Saint of Llaneigrad.	P_1, P_{14}	Place, Place	Place	Agent

Table 4: Examples of the original test (top) and adversarially perturbed examples (bottom) of DBpedia using TextFooler, classified by a PBN, alongside the top-2 activated prototypes by the PBN’s fully connected layer and their associated labels. Incorrectly predicted examples are in *italic*.

utilized in prediction. However, in the last example, the model incorrectly classifies an example that is associated with an Agent as a Place. Interestingly, this incorrect behavior can be explained by the change in the top-2 activated prototypes, where they are changing from Agent-associated to Place-associated prototypes because of the misspelling of "saint" with "street." Thus, the use of prototypes not only enhances our understanding of the model’s decision-making process but also unveils how minor perturbations influence the model’s predictions.

6 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 various benchmarks (Wang et al., 2019; Sarlin et al., 2020), concerns have shifted towards models’ performance facing noisy data (Wang et al., 2022a,b). Studies have designed novel and effective adversarial attacks (Jin et al., 2020; Zhang et al., 2020), defense mechanisms (Goyal et al., 2023; Liu et al., 2020), and evaluations to better understand the robustness properties of NLP models (Wang et al., 2022a; Morris et al., 2020a). These evaluations are also being extended to LLMs, as they similarly lack robustness (Wang et al., 2023; Shi et al., 2023). While prior work has studied LMs’ robustness, to our knowledge, PBNs’ robustness properties have not been explored yet. Our study bridges this gap.

Prototype-based networks. PBNs are widely used in CV (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 detec-

tion (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 interpretability. ProtoryNet (Hong et al., 2020) leverages RNN-extracted prototype trajectories and deploys a pruning procedure for prototypes, 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 in NLP. Our paper consolidates PBN components used in prior studies and studies their robustness in different adversarial settings.

7 Conclusions

Inspired by the lack of robustness to noisy data of state-of-the-art LMs and LLMs, we study the robustness of PBNs, as an architecturally robust variation of LMs, against both targeted and static adversarial attacks. We find that PBNs are more robust than vanilla LMs and even LLMs such as Llama3, both under targeted and static adversarial attack settings. Our results suggest that this robustness can be sensitive to hyperparameters involved in PBNs’ training. More particularly, we note that a low number of prototypes and tight clustering conditions limit the robustness capacities of PBNs. Additionally, benefiting from the inherently interpretable architecture of PBNs, we showcase how learned prototypes can be utilized for robustness and also for gaining insights about their behavior facing adversarial perturbations, even when PBNs are wrong. In summary, our work provides encouraging results for the potential of PBNs to enhance the robustness of LMs across a variety of text classification tasks and quantifies the impact of architectural components on PBN robustness.

621 Limitations

622 Although we cover a wide range of adversarial per-
623 turbations and strategies for their generation, we
624 acknowledge that more complicated perturbations
625 can also be created that are more effective and
626 help the community have a more complete under-
627 standing of the models’ robustness. Hence, we do
628 not comment on the generalizability of our study
629 to all possible textual perturbations besides our
630 evaluation on AdvGLUE. Moreover, although it is
631 customary in the field to utilize prototype-based
632 networks for classification tasks, their application
633 and robustness on other tasks remain to be explored.
634 Furthermore, while we attempt to use a wide vari-
635 ety of backbones for our study, we do not ascertain
636 similar patterns for all possible PBN backbones
637 and leave this study for future work. Finally, we
638 encourage more exploration of the interpretability
639 of these models under different attacks to better
640 understand the interpretability benefits of models
641 when analyzing robustness.

642 Ethical Considerations

643 Although the datasets and domains we focus on
644 do not pose any societal harm, the potential harm
645 that is associated with using the publicly available
646 tools we used in this study to manipulate models in
647 other critical domains should be considered. Issues
648 surrounding anonymization and offensive content
649 hold importance in data-driven studies, particularly
650 in fields like natural language processing. Since we
651 utilize datasets like IMDB, AG_News, DBPedia,
652 and AdvGLUE that are already content-moderated,
653 there is no need for anonymization of data before
654 testing for robustness in this study.

655 References

656 Open AI. 2024. Hello GPT-4o. <https://openai.com/index/hello-gpt-4o/>. [Accessed 15-06-2024].
657
658 AI@Meta. 2024. [Llama 3 model card](#).
659 Plamen Angelov and Eduardo Soares. 2020. Towards
660 explainable deep neural networks (xdnn). *Neural*
661 *Networks*, 130:185–194.
662 Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu,
663 Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi,
664 Cunxiang Wang, Yidong Wang, et al. 2024. A sur-
665 vey on evaluation of large language models. *ACM*
666 *Transactions on Intelligent Systems and Technology*,
667 15(3):1–45.

Chaofan Chen, Oscar Li, Chaofan Tao, Alina Jade Bar-
668 nett, Jonathan Su, and Cynthia Rudin. 2019. *This*
669 *Looks like That: Deep Learning for Interpretable Im-*
670 *age Recognition*. Curran Associates Inc., Red Hook,
671 NY, USA. 672
Aakanksha Chowdhery, Sharan Narang, Jacob Devlin,
673 Maarten Bosma, Gaurav Mishra, Adam Roberts,
674 Paul Barham, Hyung Won Chung, Charles Sutton,
675 Sebastian Gehrmann, et al. 2022. Palm: Scaling
676 language modeling with pathways. *arXiv preprint*
677 *arXiv:2204.02311*. 678
Kevin Clark, Minh-Thang Luong, Quoc V Le, and
679 Christopher D Manning. 2020. Electra: Pre-training
680 text encoders as discriminators rather than generators.
681 *arXiv preprint arXiv:2003.10555*. 682
Nilesh Dalvi, Pedro Domingos, Mausam, Sumit Sang-
683 hai, and Deepak Verma. 2004. *Adversarial classifica-*
684 *tion*. KDD ’04, page 99–108, New York, NY, USA.
685 Association for Computing Machinery. 686
Anubrata Das, Chitrant Gupta, Venelin Kovatchev,
687 Matthew Lease, and Junyi Jessy Li. 2022. [ProtoTEx:](#)
688 [Explaining model decisions with prototype tensors](#).
689 In *Proceedings of the 60th Annual Meeting of the*
690 *Association for Computational Linguistics (Volume*
691 *1: Long Papers)*, pages 2986–2997, Dublin, Ireland.
692 Association for Computational Linguistics. 693
Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
694 Kristina Toutanova. 2018. Bert: Pre-training of deep
695 bidirectional transformers for language understand-
696 ing. *arXiv preprint arXiv:1810.04805*. 697
Javid Ebrahimi, Anyi Rao, Daniel Lowd, and Dejing
698 Dou. 2018. [HotFlip: White-box adversarial exam-](#)
699 [ples for text classification](#). In *Proceedings of the 56th*
700 *Annual Meeting of the Association for Computational*
701 *Linguistics (Volume 2: Short Papers)*, pages 31–36,
702 Melbourne, Australia. Association for Computational
703 Linguistics. 704
Ji Gao, Jack Lanchantin, Mary Lou Soffa, and Yanjun
705 Qi. 2018. Black-box generation of adversarial text
706 sequences to evade deep learning classifiers. In *2018*
707 *IEEE Security and Privacy Workshops (SPW)*, pages
708 50–56. IEEE. 709
Siddhant Garg and Goutham Ramakrishnan. 2020.
710 [BAE: BERT-based adversarial examples for text clas-](#)
711 [sification](#). In *Proceedings of the 2020 Conference on*
712 *Empirical Methods in Natural Language Processing*
713 *(EMNLP)*, pages 6174–6181, Online. Association for
714 Computational Linguistics. 715
Shafie Gholizadeh and Nengfeng Zhou. 2021. [Model ex-](#)
716 [plainability in deep learning based natural language](#)
717 [processing](#). 718
Ian J Goodfellow, Jonathon Shlens, and Christian
719 Szegedy. 2014. Explaining and harnessing adver-
720 sarial examples. *arXiv preprint arXiv:1412.6572*. 721

722	Shreya Goyal, Sumanth Doddapaneni, Mitesh M Khapra, and Balaraman Ravindran. 2023. A survey of adversarial defenses and robustness in nlp. <i>ACM Computing Surveys</i> , 55(14s):1–39.	777
723		778
724		779
725		780
726	Xiaowei Gu and Weiping Ding. 2019. A hierarchical prototype-based approach for classification. <i>Information Sciences</i> , 505:325–351.	781
727		782
728		783
729	Antonio Gulli. AG’s Corpus of News Articles. groups.di.unipi.it/~gulli/AG_corpus_of_news_articles.html . Accessed 15 June 2024.	784
730		785
731		786
732	Jiale Han, Bo Cheng, and Wei Lu. 2021. Exploring task difficulty for few-shot relation extraction. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 2605–2616, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	787
733		788
734		789
735		790
736		791
737		792
738	Peter Hase and Mohit Bansal. 2020. Evaluating explainable AI: Which algorithmic explanations help users predict model behavior? In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics – ACL2020</i> , pages 5540–5552, Online. Association for Computational Linguistics.	793
739		794
740		795
741		796
742		797
743		798
744	Peter Hase, Chaofan Chen, Oscar Li, and Cynthia Rudin. 2019. Interpretable image recognition with hierarchical prototypes. In <i>Proceedings of the AAAI Conference on Human Computation and Crowdsourcing</i> , volume 7, pages 32–40.	799
745		800
746		801
747		802
748		803
749	Junxian He, Taylor Berg-Kirkpatrick, and Graham Neubig. 2020. Learning sparse prototypes for text generation. <i>Advances in Neural Information Processing Systems</i> , 33:14724–14735.	804
750		805
751		806
752		807
753	Dan Hendrycks, Xiaoyuan Liu, Eric Wallace, Adam Dziedzic, Rishabh Krishnan, and Dawn Song. 2020. Pretrained transformers improve out-of-distribution robustness. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 2744–2751, Online. Association for Computational Linguistics.	808
754		809
755		810
756		811
757		812
758		813
759		814
760	Dat Hong, Stephen S Baek, and Tong Wang. 2020. Interpretable sequence classification via prototype trajectory. <i>arXiv preprint arXiv:2007.01777</i> .	815
761		816
762		817
763	Dat Hong, Stephen S. Baek, and Tong Wang. 2021. Interpretable sequence classification via prototype trajectory.	818
764		819
765		820
766	Myeongjun Jang, Deuk Sin Kwon, and Thomas Lukasiewicz. 2022. Becel: Benchmark for consistency evaluation of language models. In <i>Proceedings of the 29th International Conference on Computational Linguistics</i> , pages 3680–3696.	821
767		822
768		823
769		824
770		825
771	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 <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 34, pages 8018–8025.	826
772		827
773		828
774		829
775		830
776		831
		832
	Mark T Keane and Eoin M Kenny. 2019. How case-based reasoning explains neural networks: A theoretical analysis of xai using post-hoc explanation-by-example from a survey of ann-cbr twin-systems. In <i>Case-Based Reasoning Research and Development: 27th International Conference, ICCBR 2019, Otzenhausen, Germany, September 8–12, 2019, Proceedings 27</i> , pages 155–171. Springer.	
	Eunji Kim, Siwon Kim, Minji Seo, and Sungroh Yoon. 2021. Xprotonet: Diagnosis in chest radiography with global and local explanations. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 15719–15728.	
	Pang Wei Koh, Thao Nguyen, Yew Siang Tang, Stephen Mussmann, Emma Pierson, Been Kim, and Percy Liang. 2020. Concept bottleneck models. In <i>International conference on machine learning</i> , pages 5338–5348. PMLR.	
	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.	
	Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch, Dimitris Kontokostas, Pablo N Mendes, Sebastian Hellmann, Mohamed Morsey, Patrick Van Kleef, Sören Auer, et al. 2015. Dbpedia—a large-scale, multilingual knowledge base extracted from wikipedia. <i>Semantic web</i> , 6(2):167–195.	
	Qi Lei, Lingfei Wu, Pin-Yu Chen, Alex Dimakis, Inderjit S Dhillon, and Michael J Witbrock. 2019. Discrete adversarial attacks and submodular optimization with applications to text classification. <i>Proceedings of Machine Learning and Systems</i> , 1:146–165.	
	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.	
	Ang Li, Fangyuan Zhang, Shuangjiao Li, Tianhua Chen, Pan Su, and Hongtao Wang. 2023. Efficiently generating sentence-level textual adversarial examples with seq2seq stacked auto-encoder. <i>Expert Systems with Applications</i> , 213:119170.	
	Jinfeng Li, Shouling Ji, Tianyu Du, Bo Li, and Ting Wang. 2018a. Textbugger: Generating adversarial text against real-world applications. <i>arXiv preprint 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 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</i> . AAAI Press.	

833	Tal Linzen. 2020. How can we accelerate progress towards human-like linguistic generalization? <i>arXiv preprint arXiv:2005.00955</i> .	
834		
835		
836	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 Linguistics (ACL), System Demonstrations</i> .	
837		
838		
839		
840		
841		
842	Xiaodong Liu, Hao Cheng, Pengcheng He, Weizhu Chen, Yu Wang, Hoifung Poon, and Jianfeng Gao. 2020. Adversarial training for large neural language models. <i>arXiv preprint arXiv:2004.08994</i> .	
843		
844		
845		
846	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. <i>arXiv preprint arXiv:1907.11692</i> .	
847		
848		
849		
850		
851	Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. 2011. Learning word vectors for sentiment analysis . In <i>Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies</i> , pages 142–150, Portland, Oregon, USA. Association for Computational Linguistics.	
852		
853		
854		
855		
856		
857		
858		
859	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. <i>arXiv preprint arXiv:2310.15743</i> .	
860		
861		
862		
863		
864	Pascal Mettes, Elise Van der Pol, and Cees Snoek. 2019. Hyperspherical prototype networks. <i>Advances in neural information processing systems</i> , 32.	
865		
866		
867	Yao Ming, Panpan Xu, Huamin Qu, and Liu Ren. 2019. Interpretable and steerable sequence learning via prototypes . In <i>Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining</i> . ACM.	
868		
869		
870		
871		
872	Milad Moradi and Matthias Samwald. 2021. Evaluating the robustness of neural language models to input perturbations .	
873		
874		
875	John X Morris, Eli Lifland, Jack Lanchantin, Yangfeng Ji, and Yanjun Qi. 2020a. Reevaluating adversarial examples in natural language. <i>arXiv preprint arXiv:2004.14174</i> .	
876		
877		
878		
879	John X. Morris, Eli Lifland, Jin Yong Yoo, Jake Grigsby, Di Jin, and Yanjun Qi. 2020b. Textattack: A framework for adversarial attacks, data augmentation, and adversarial training in nlp .	
880		
881		
882		
883	Meike Nauta, Annemarie Jutte, Jesper Provoost, and Christin Seifert. 2021a. This looks like that, because ... explaining prototypes for interpretable image recognition . In <i>Communications in Computer and Information Science</i> , pages 441–456. Springer International Publishing.	
884		
885		
886		
887		
888		
	Meike Nauta, Ron van Bree, and Christin Seifert. 2021b. Neural prototype trees for interpretable fine-grained image recognition . 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.	889
		890
		891
		892
		893
		894
	OpenAI. 2022. Chatgpt. https://openai.com/blog/chatgpt . Accessed: April 30, 2023.	895
		896
	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.	897
		898
		899
		900
		901
	Nicolas Papernot and Patrick McDaniel. 2018. Deep k-nearest neighbors: Towards confident, interpretable and robust deep learning. <i>arXiv preprint arXiv:1803.04765</i> .	902
		903
		904
		905
	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.	906
		907
		908
		909
		910
		911
		912
	Shuhuai Ren, Yihe Deng, Kun He, and Wanxiang Che. 2019. Generating natural language adversarial examples through probability weighted word saliency . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 1085–1097, Florence, Italy. Association for Computational Linguistics.	913
		914
		915
		916
		917
		918
		919
	Eleanor H. Rosch. 1973. Natural categories . <i>Cognitive Psychology</i> , 4(3):328–350.	920
		921
	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> .	922
		923
		924
		925
	Sascha Saralajew, Lars Holdijk, and Thomas Villmann. 2020. Fast adversarial robustness certification of nearest prototype classifiers for arbitrary seminorms . In <i>Advances in Neural Information Processing Systems</i> , pages 13635–13650. Curran Associates, Inc.	926
		927
		928
		929
		930
	Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. 2020. Superglue: Learning feature matching with graph neural networks .	931
		932
		933
		934
	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 .	935
		936
		937
		938
	Chenglei Si, Zhengyan Zhang, Fanchao Qi, Zhiyuan Liu, Yasheng Wang, Qun Liu, and Maosong Sun. 2021. Better robustness by more coverage: Adversarial and mixup data augmentation for robust	939
		940
		941
		942

943	finetuning. In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 1569–1576.	995
944		996
945		997
946	Dylan Slack, Satyapriya Krishna, Himabindu Lakkaraju, and Sameer Singh. 2022. Talktomodel: Explaining machine learning models with interactive natural language conversations.	998
947		999
948		1000
949		
950	Jake Snell, Kevin Swersky, and Richard Zemel. 2017. Prototypical networks for few-shot learning . In <i>Advances in Neural Information Processing Systems</i> , volume 30. Curran Associates, Inc.	1001
951		1002
952		1003
953		
954	Eduardo Soares, Plamen Angelov, and Neeraj Suri. 2022. Similarity-based deep neural network to detect imperceptible adversarial attacks . In <i>2022 IEEE Symposium Series on Computational Intelligence (SSCI)</i> , pages 1028–1035.	1008
955		1009
956		1010
957		1011
958		1012
959	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. <i>Knowledge-Based Systems</i> , 266:110418.	1013
960		1014
961		1015
962		1016
963		
964		
965	Iulia Turc, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Well-read students learn better: On the importance of pre-training compact models .	1017
966		1018
967		1019
968		1020
969		1021
970	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 .	1022
971		1023
972		1024
973	Václav Voráček and Matthias Hein. 2022. Provably adversarially robust nearest prototype classifiers . 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.	1025
974		1026
975		
976		
977		
978		
979	Kiri Wagstaff. 2012. Machine learning that matters. <i>arXiv preprint arXiv:1206.4656</i> .	1027
980		1028
981		1029
982	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 .	1030
983		1031
984		
985	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 .	1032
986		1033
987		1034
988		
989		
990	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 .	1035
991		1036
992		1037
993		1038
994		
	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 .	995
		996
		997
		998
		999
		1000
	Xuezhi Wang, Haohan Wang, and Diyi Yang. 2022c. Measure and improve robustness in nlp models: A survey .	1001
		1002
		1003
	Jing Wu, Mingyi Zhou, Ce Zhu, Yipeng Liu, Mehrtash Harandi, and Li Li. 2021. Performance evaluation of adversarial attacks: Discrepancies and solutions. <i>arXiv preprint arXiv:2104.11103</i> .	1004
		1005
		1006
		1007
	Hong-Ming Yang, Xu-Yao Zhang, Fei Yin, and Cheng-Lin Liu. 2018. Robust classification with convolutional prototype learning. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 3474–3482.	1008
		1009
		1010
		1011
		1012
	Jin Yong Yoo, John X. Morris, Eli Lifland, and Yanjun Qi. 2020. Searching for a search method: Benchmarking search algorithms for generating nlp adversarial examples .	1013
		1014
		1015
		1016
	Wei Emma Zhang, Quan Z Sheng, Ahoud Alhazmi, and Chenliang Li. 2020. Adversarial attacks on deep-learning models in natural language processing: A survey. <i>ACM Transactions on Intelligent Systems and Technology (TIST)</i> , 11(3):1–41.	1017
		1018
		1019
		1020
		1021
	Haiyan Zhao, Hanjie Chen, Fan Yang, Ninghao Liu, Huiqi Deng, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, and Mengnan Du. 2023. Explainability for large language models: A survey. <i>arXiv preprint arXiv:2309.01029</i> .	1022
		1023
		1024
		1025
		1026
	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. <i>arXiv preprint arXiv:2005.00782</i> .	1027
		1028
		1029
		1030
		1031
	Julia El Zini and Mariette Awad. 2022. On the explainability of natural language processing deep models . <i>ACM Comput. Surv.</i> , 55(5).	1032
		1033
		1034
	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 .	1035
		1036
		1037
		1038
	A Dataset Details	1039
	The statistics of the datasets we used in this study to test the robustness of PBNs against perturbations are demonstrated in Table 5 . We present both statistics about the original dataset and statistics and details about the number of perturbations that we have gathered on each dataset with different attack strategies. All the original datasets we use in this study are gathered by other researchers and	1040
		1041
		1042
		1043
		1044
		1045
		1046
		1047

1048 have been made public by them, mentioning non-
1049 commercial use, which aligns with how we use
1050 these datasets. We have included information on
1051 their descriptions and how they were gathered:

1052 **IMDB.** This dataset is compiled from a set of
1053 50000 reviews sourced from IMDB in English, lim-
1054 iting each movie to a maximum of 30 reviews. It
1055 has maintained an equal count of positive and neg-
1056 ative reviews, ensuring a 50% accuracy through
1057 random guessing. To align with prior research
1058 on polarity classification, the authors specifically
1059 focus on highly polarized reviews. A review is
1060 considered negative if it scores ≤ 4 out of 10 and
1061 positive if it scores ≥ 7 out of 10. Neutral reviews
1062 are excluded from this dataset. Authors have made
1063 the dataset publicly available, and you can find
1064 more information about this dataset at [https://
1065 ai.stanford.edu/~amaas/data/sentiment/](https://ai.stanford.edu/~amaas/data/sentiment/).

1066 **AG_News.** This dataset comprises over 1 million
1067 English news articles sourced from 2000+ news
1068 outlets over a span of more than a year by Come-
1069 ToMyHead, an academic news search engine op-
1070 erational since July 2004. Provided by the aca-
1071 demic community, this dataset aids research in
1072 data mining, information retrieval, data compres-
1073 sion, data streaming, and non-commercial activi-
1074 ties. This news topic classification dataset features
1075 four classes: world, sports, business, and science.
1076 The details about the intended use and access condi-
1077 tions are provided at [http://www.di.unipi.it/
1078 ~gulli/AG_corpus_of_news_articles.html](http://www.di.unipi.it/~gulli/AG_corpus_of_news_articles.html).

1079 **DBPedia.** DBPedia⁵ seeks to extract organized
1080 information from Wikipedia’s vast content. The
1081 gathered subset of data we used offers hierar-
1082 chical categories for 342782 Wikipedia articles.
1083 These classes are distributed across three lev-
1084 els, comprising 9, 70, and 219 classes, respec-
1085 tively. We used the version that has nine classes:
1086 Agent, Work, Place, Species, UnitOfWork, Event,
1087 SportsSeason, Device, and TopicalConcept. Al-
1088 though the articles are in English, specific names
1089 (e.g., the name of a place or person) can be
1090 non-English. Find more information about this
1091 dataset at [https://huggingface.co/datasets/
1092 DeveloperOats/DBPedia_Classes](https://huggingface.co/datasets/DeveloperOats/DBPedia_Classes).

1093 **AdvGLUE.** Adversarial GLUE (AdvGLUE)
1094 (Wang et al., 2022a) introduces a multi-task En-
1095 glish benchmark designed to investigate and assess

⁵<https://www.dbpedia.org/>

1096 the vulnerabilities of modern large-scale language
1097 models against various adversarial attacks. It en-
1098 compasses five corpora, including SST-2 sentiment
1099 classification, QQP paraphrase test dataset, and
1100 QNLI, RTE, and MNLI, all of which are natural lan-
1101 guage inference datasets. To assess robustness, per-
1102 turbations are applied to these datasets through both
1103 automated and human-evaluated methods, span-
1104 ning word-level, sentence-level, and human-crafted
1105 examples. Our focus primarily centers on SST-2
1106 due to its alignment with the other covered datasets
1107 in our study and its classification nature. This
1108 dataset has been made public by the authors and is
1109 released with CC BY-SA 4.0 license.

1110 **B Implementation Details**

1111 **B.1 Experimental Environment**

1112 For all the experiments that involved training or
1113 evaluating PBNs or vanilla LMs, we used three
1114 GPU NVIDIA RTX A5000 devices with Python
1115 v3.9.16 and CUDA v11.6, and each experiment
1116 took between 10 minutes to 2 hours, depending
1117 on the dataset and model used. All Transformer
1118 models were trained using the Transformers pack-
1119 age v4.30.2 and Torch package v2.0.1+cu117. We
1120 used TextAttack v0.3.10 (Morris et al., 2020b) for
1121 implementing the employed attack strategies and
1122 perturbations.

1123 **B.2 Training Details**

1124 All prototypes are initialized randomly for a fair
1125 comparison, and only the last layer of LM back-
1126 bones are trainable. The prototypes are trained
1127 without being constrained to a certain class from
1128 the beginning, and their corresponding class can
1129 be identified after training. The transformation
1130 from distances to class logits is done through a
1131 simple fully connected layer without intercept to
1132 avoid introducing additional complexity and keep
1133 the prediction interpretable through prototype dis-
1134 tances. Both the backbone of PBNs and their
1135 vanilla counterparts leveraged the same LM and
1136 were fine-tuned separately to show the difference
1137 that is only attributed to the models’ architecture.
1138 Focusing on the BERT-based PBN for evaluation,
1139 since BERT-base is one of the models from which
1140 we extract static perturbations by directly attacking
1141 it, to ensure generalization of the experiments on
1142 different backbones in the evaluation step, we use
1143 BERT-Medium (Turc et al., 2019) as the backbone
1144 for BERT-based PBN and its vanilla counterpart.

Dataset	#Classes	#Tokens	#Train	#Val	#Test	BAE	DWB	PWWS	TB	TF	Other
IMDB	2	234	22,500	2,500	25,000	1784	1584	2816	2408	2880	-
AG_News	4	103	112,400	7,600	7,600	663	1287	1533	1383	1893	-
DBPedia	9	38	240,942	36,003	60,794	1041	1143	1401	1281	1836	-
SST-2	2	14	67,349	872	1,821	-	-	-	-	-	148

Table 5: Dataset statistics: number of classes, the average number of tokens, and size of the perturbed datasets under BAE, DeepWordBug (DWB), PWWS, TextBugger (TB), TextFooler (TF), obtained. 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.

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 were trained with the number of epochs adjusted according to an early stopping module with patience of 4 and a threshold value of 0.01 for change in accuracy.

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 static perturbations were also pre-trained Transformer models gathered from <https://huggingface.co/>. 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., 2020b) 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 target model.

B.3 GPT4o and Llama3 Baseline

We used GPT4o and Llama3 (AI@Meta, 2024) as baselines 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 these models to generate the baseline responses and the reported performance in Table 2. 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

C.1 Robustness of PBNs Against Paraphrased-Based Perturbations

Comparison between PBNs and vanilla LMs on the original and paraphrased version of texts from AG_News, DBPedia, and IMDB datasets that GPT3.5 generated are shown in Table 6, which

	AG_News		DBPedia		IMDB	
	Orig	Adv	Orig	Adv	Orig	Adv
BART	93.7	92.6	91.2	91.3	97.5	96.0
+ PBN	93.2	93.8	92.0	91.6	97.2	97.0
BERT	92.5	91.0	90.8	90.5	95.5	94.2
+ PBN	92.8	91.2	90.3	90.8	95.2	95.0
ELEC.	93.0	92.1	90.5	90.0	96.0	94.5
+ PBN	93.5	91.8	90.8	89.7	95.8	95.0

Table 6: Comparison between PBNs and vanilla LMs on the original and paraphrased version of texts from AG_News, DBPedia, and IMDB datasets that GPT3.5 generated.

illustrated that both PBNs and vanilla LMs are robust to such perturbations.

C.2 Robustness of PBNs’ w.r.t. Average Percentage of Words Perturbed

The Comparison of PBNs and vanilla LMs’ robustness with respect to the Average Percentage of Words Perturbed (APWP) under different adversarial settings, different datasets, and perturbation strategies is shown in Table 7. We observed that while using the best hyperparameters, PBNs are more robust than vanilla LMs in the majority of the cases, this superiority is less salient when averaging over all hyperparameters involved in PBNs’ training, which entails how PBNs’ robustness is sensitive to hyperparameters.

C.3 Robustness of PBNs’ Averaged over Hyperparameters

The comparison of PBNs and vanilla LMs under different adversarial settings, on different datasets, and different attacking strategies, averaged over all hyperparameters of PBNs, is shown in Table 8. Comparing the observed trends with the same trends when using the best hyperparameters for PBNs, our results suggested that PBNs’ robustness is sensitive to hyperparameters that are involved in their training.

C.4 Effect of Distance Function on Robustness

Figure 6, Figure 4, and Figure 5 illustrate the robustness of PBNs compared to vanilla LMs, using different distance functions, showing that PBNs’ robustness is not sensitive to this hyperparameter.

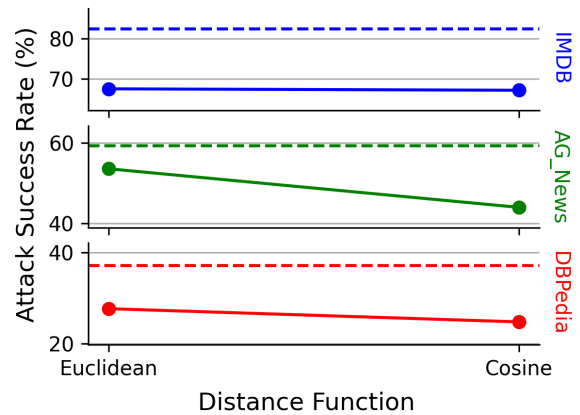


Figure 4: Attack Success Rate (ASR %) of PBNs with different distance functions and other hyperparameters set to their best values and averaged across other possible variables (e.g., backbone and attack type). The dotted line represents the ASR for the vanilla LMs.

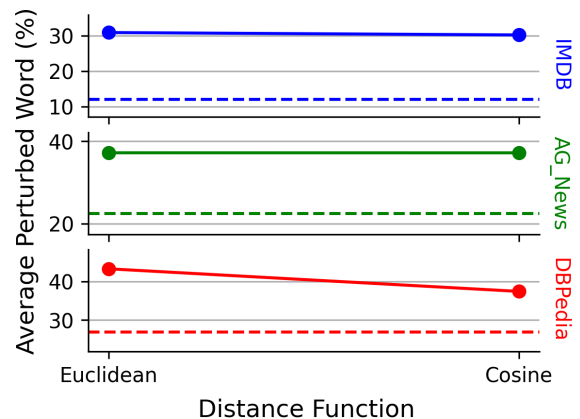


Figure 5: Average Percentage of Words Perturbed (APWP) of PBNs with different distance functions and other hyperparameters set to their best values and averaged across other possible variables (e.g., backbone and attack type). The dotted line represents the APWP for the vanilla LMs.

Using the best hyperparameters

	AG_News					DBPedia					IMDB				
	BAE	DWB	PWWS	TB	TF	BAE	DWB	PWWS	TB	TF	BAE	DWB	PWWS	TB	TF
BART	8.7	26.9	20.8	35.7	25.0	9.1	27.3	16.9	50.1	26.2	4.1	6.4	4.2	33.3	5.9
+ PBN	9.0	24.8	22.2	37.7	27.6	10.1	17.1	15.9	43.3	26.0	4.7	6.6	8.1	33.4	13.6
BERT	7.4	26.8	21.6	37.4	24.1	9.7	27.9	19.4	53.8	28.8	4.0	5.7	4.4	30.1	5.0
+ PBN	7.7	26.6	24.1	37.7	28.8	10.9	27.9	22.4	50.0	30.6	4.6	6.7	9.3	35.9	15.4
ELEC.	8.2	23.7	17.5	32.7	20.8	10.9	24.6	17.7	58.0	22.9	5.4	8.1	8.8	44.7	11.2
+ PBN	8.1	21.2	18.9	31.8	24.0	11.9	25.1	19.4	48.5	26.8	5.6	8.4	13.3	38.6	18.5

Averaged over all hyperparameters

	AG_News					DBPedia					IMDB				
	BAE	DWB	PWWS	TB	TF	BAE	DWB	PWWS	TB	TF	BAE	DWB	PWWS	TB	TF
BART	8.7	26.9	20.8	35.7	25.0	9.1	27.3	16.9	50.1	26.2	4.1	6.4	4.2	33.3	5.9
+ PBN	8.3	19.3	20.5	32.6	25.2	9.7	17.1	15.9	40.4	24.7	4.4	6.1	6.5	29.5	10.1
BERT	7.4	26.8	21.6	37.4	24.1	9.7	27.9	19.4	53.8	28.8	4.0	5.7	4.4	30.1	5.0
+ PBN	7.2	24.1	21.9	35.0	25.9	9.5	24.1	19.3	43.1	27.6	4.1	5.5	5.0	27.3	7.1
ELEC.	8.2	23.7	17.5	32.7	20.8	10.9	24.6	17.7	58.0	22.9	5.4	8.1	8.8	44.7	11.2
+ PBN	7.7	15.3	16.1	26.1	20.1	10.2	18.2	16.6	40.2	23.7	5.4	6.7	10.1	31.3	13.6

Table 7: Comparison of PBNs and vanilla LMs’ robustness with respect to Average Percentage of Words Perturbed (APWP) under targeted adversarial attack perturbations, both using the best hyperparameters and averaged over all hyperparameters for PBNs, on IMBD, AG_News, and DBPeida datasets, under BAE, DeepWordBug (DWB), PWWS, TextBugger (TB), TextFooler (TF). The highest APWP showing the superior model for each architecture is **boldfaced**.

Targeted Attacks; Attack Success Rate (ASR %) reported

	AG_News					DBPedia					IMDB				
	BAE	DWB	PWWS	TB	TF	BAE	DWB	PWWS	TB	TF	BAE	DWB	PWWS	TB	TF
BART	14.8	53.2	53.6	31.8	76.5	18.9	28.3	43.1	21.1	71.9	74.1	74.7	99.3	78.5	100.0
+ PBN	14.8	40.4	50.7	29.8	76.2	17.0	14.7	28.7	12.7	49.4	55.5	49.2	86.2	49.7	88.5
BERT	17.0	78.0	69.8	45.7	88.8	13.9	24.8	31.6	22.0	61.3	82.5	79.7	99.9	83.9	99.9
+ PBN	14.0	64.7	57.0	39.3	82.1	13.5	23.4	27.6	19.6	51.3	68.4	61.8	91.3	74.0	92.4
ELEC.	24.8	89.5	69.1	87.8	87.9	14.5	42.8	45.6	42.3	75.3	52.5	49.2	95.3	67.8	99.3
+ PBN	18.5	50.4	55.7	63.6	80.0	12.6	19.4	26.1	27.1	46.5	41.0	35.9	77.7	55.6	86.2

Static Attacks; Accuracy (%) reported

	AG_News					DBPedia					IMDB					SST2 GLUE
	BAE	DWB	PWWS	TB	TF	BAE	DWB	PWWS	TB	TF	BAE	DWB	PWWS	TB	TF	
BART	<u>53.2</u>	<u>76.7</u>	<u>83.2</u>	<u>77.5</u>	<u>85.8</u>	55.5	68.6	58.4	<u>72.5</u>	<u>71.3</u>	74.1	80.5	83.6	85.8	87.6	29.8
+ PBN	50.4	68.3	75.7	70.5	79.6	<u>56.4</u>	65.8	<u>58.7</u>	70.9	69.5	69.2	78.7	79.7	81.9	78.3	37.6
+ Aug.	71.7	78.4	85.5	77.6	90.1	84.0	79.6	89.7	88.8	94.0	85.7	86.7	92.9	89.9	96.5	-
BERT	47.8	64.0	75.9	69.4	80.7	62.3	61.4	75.4	<u>78.4</u>	<u>82.0</u>	<u>75.1</u>	<u>77.1</u>	85.0	83.4	<u>85.9</u>	42.0
+ PBN	49.5	66.2	76.4	71.3	<u>82.3</u>	<u>63.5</u>	61.1	73.9	76.9	79.4	71.0	73.9	81.1	80.2	<u>79.2</u>	47.1
+ Aug.	58.3	71.6	78.3	<u>71.2</u>	85.4	75.5	70.9	84.1	90.5	91.0	83.2	77.6	91.7	90.8	89.2	-
ELECTRA	50.4	65.0	<u>73.5</u>	<u>63.9</u>	77.8	<u>79.7</u>	<u>66.9</u>	<u>80.9</u>	<u>81.4</u>	<u>84.4</u>	89.7	<u>90.3</u>	94.6	<u>94.5</u>	<u>95.6</u>	44.3
+ PBN	<u>52.7</u>	<u>63.9</u>	73.7	67.1	77.8	73.4	64.1	73.0	76.4	80.6	80.6	79.4	79.9	80.2	86.8	56.4
+ Aug.	55.0	59.5	71.7	61.6	79.5	86.2	73.8	88.1	84.5	92.8	<u>89.4</u>	93.7	95.3	94.9	95.8	-

Table 8: Comparison of PBNs and vanilla LMs (+ vanilla LMs with adversarial augmented training under static attack setting) under both targeted and static adversarial attack perturbations, averaged over all hyperparameters for PBNs, on IMBD, AG_News, DBPeida (+ SST-2 AdvGLUE under static attack setting) datasets, under BAE, DeepWordBug (DWB), PWWS, TextBugger (TB), TextFooler (TF). The highest accuracy and lowest ASR showing the superior model for each architecture is **boldfaced**, and the second best model is underlined for static attacks.

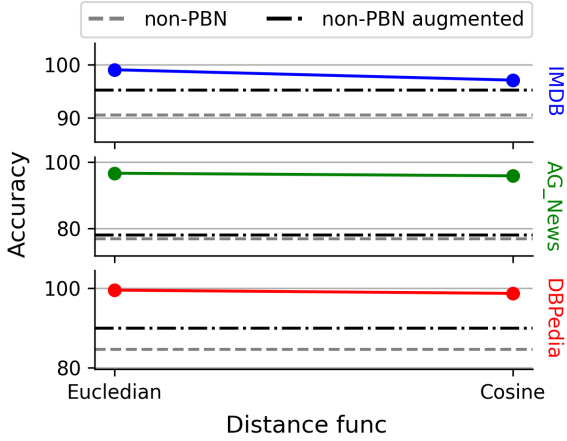


Figure 6: Accuracy of PBNs under static adversarial settings, with different distance functions, with other hyperparameters set to their best values and averaged across other possible variables (e.g., backbone and attack type). The dotted line represents the ASR for the vanilla LMs.

C.5 Effect of Interpretability on Robustness

Figure 9, Figure 7, and Figure 8 illustrate the robustness of PBNs compared to vanilla LMs, using different values of λ_i adjusting the importance of interpretability, showing that overall, PBNs' robustness is not sensitive to this hyperparameter.

C.6 Effect of Clustering on Robustness

Figure 10, Figure 11 illustrate the robustness of PBNs compared to vanilla LMs, using different values of λ_c adjusting the importance of clustering, that alongside the trends observed using ASR (see Figure 2), show that overall, PBNs' robustness degrades with tighter clustering in PBNs' training.

C.7 Effect of Separation on Robustness

Figure 14, Figure 12, and Figure 13 illustrate the robustness of PBNs compared to vanilla LMs, using different values of λ_s adjusting the importance of separability between prototypes, showing that overall, PBNs' robustness is not sensitive to this hyperparameter.

C.8 Effect of Number of Prototypes on Robustness

Figure 15, Figure 16 illustrate the robustness of PBNs compared to vanilla LMs, using different numbers of prototypes, that alongside the trends observed using ASR (see Figure 3), show that overall, PBNs' robustness degrades with low number of prototypes as PBNs can capture lower number of semantic patterns in such conditions.

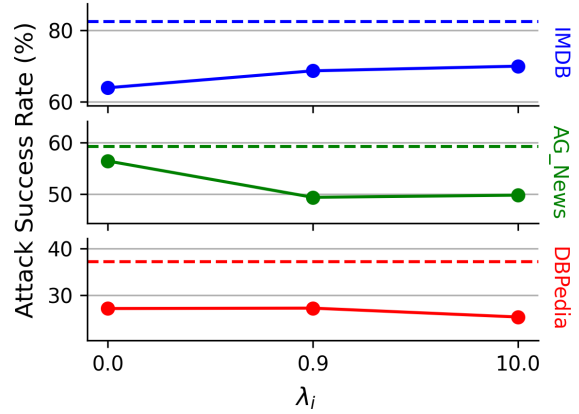


Figure 7: Attack Success Rate (ASR %) of PBNs with different λ_i values adjusting the importance of interpretability of the prototypes in training, with other hyperparameters set to their best values, and averaged across other possible variables (e.g., backbone and attack type). The dotted line represents the ASR for the non-PBN model.

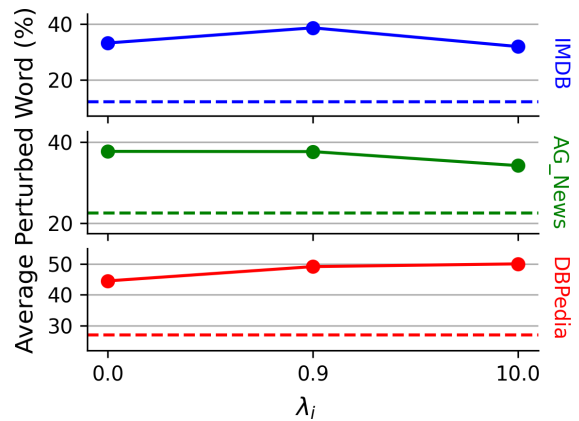


Figure 8: Average Percentage of Words Perturbed (APWP) of PBNs with different λ_i values adjusting the importance of interpretability of the prototypes in training, with other hyperparameters set to their best values, and averaged across other possible variables (e.g., backbone and attack type). The dotted line represents the ASR for the non-PBN model.

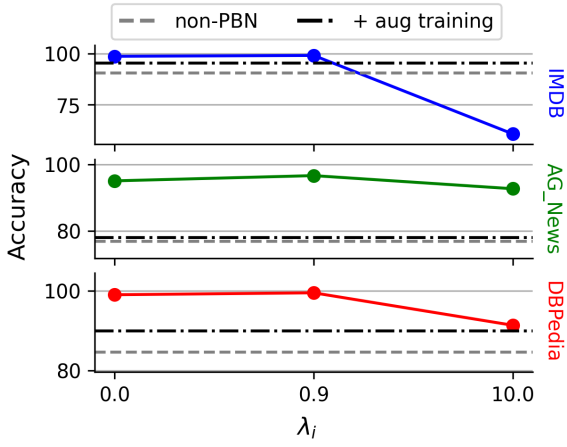


Figure 9: Accuracy of PBNs under static adversarial settings, with different λ_i values adjusting the level of interpretability in PBNs, with other hyperparameters set to their best values and averaged across other possible variables (e.g., backbone and attack type). The dotted line represents the ASR for the vanilla LMs.

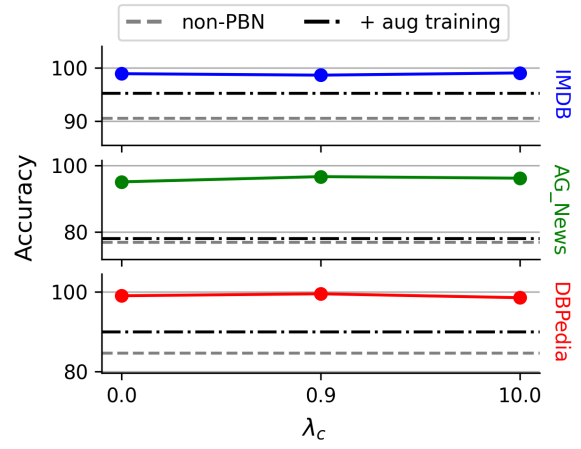


Figure 11: Accuracy of PBNs under static adversarial settings, with different λ_c values adjusting the level of clustering in PBNs, with other hyperparameters set to their best values and averaged across other possible variables (e.g., backbone and attack type). The dotted line represents the ASR for the vanilla LMs.

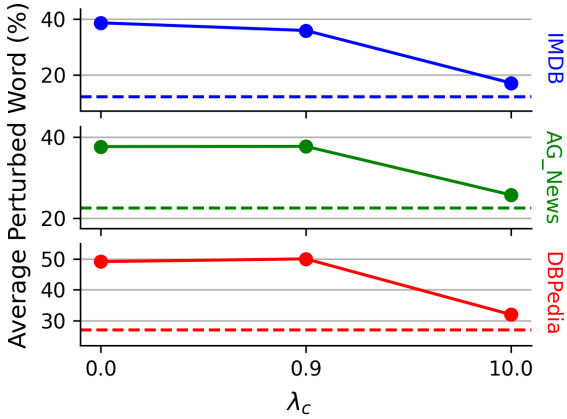


Figure 10: Average Percentage of Words Perturbed (APWP) of PBNs with different λ_c values adjusting the importance of clustering of examples in PBNs, with other hyperparameters set to their best values, and averaged across other possible variables (e.g., backbone and attack type). The dotted line represents the ASR for the non-PBN model.

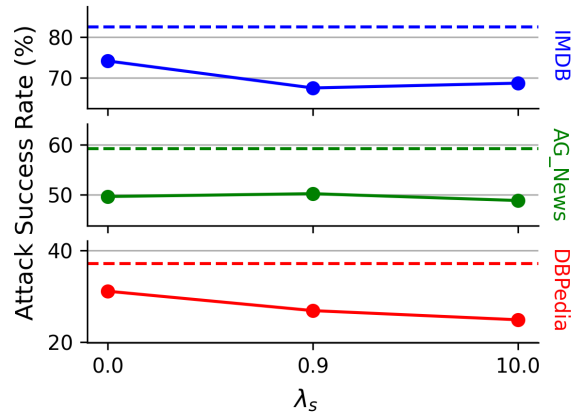


Figure 12: Attack Success Rate (ASR %) of PBNs with different λ_s values adjusting the level of separation between the prototypes, with other hyperparameters set to their best values and averaged across other possible variables (e.g., backbone and attack type). The dotted line represents the ASR for the vanilla LMs.

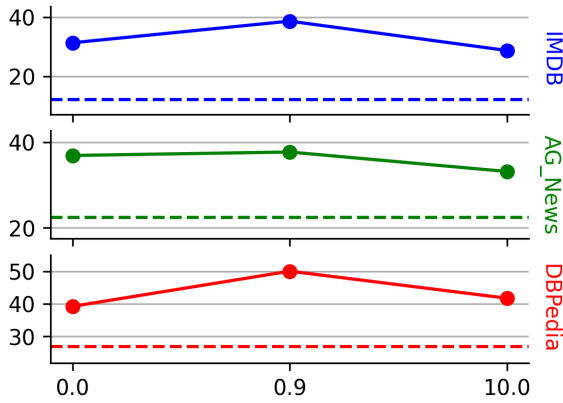


Figure 13: Average Percentage of Words Perturbed (APWP) of PBNs with different λ_s values adjusting the level of separation between the prototypes, with other hyperparameters set to their best values and averaged across other possible variables (e.g., backbone and attack type). The dotted line represents the ASR for the vanilla LMs.

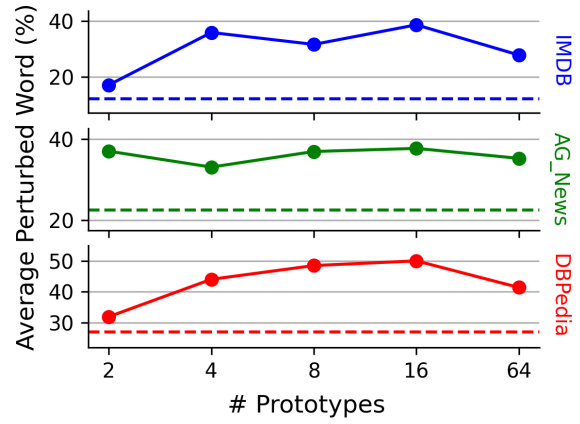


Figure 15: Average Percentage of Words Perturbed (APWP) of PBNs with different numbers of prototypes, with other hyperparameters set to their best values, and averaged across other possible variables (e.g., backbone and attack type). The dotted line represents the ASR for the non-PBN model.

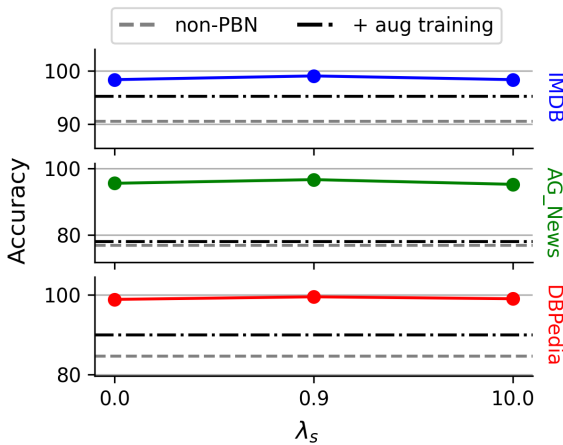


Figure 14: Accuracy of PBNs under static adversarial settings, with different λ_s values adjusting the level of separation between the prototypes, with other hyperparameters set to their best values and averaged across other possible variables (e.g., backbone and attack type). The dotted line represents the ASR for the vanilla LMs.

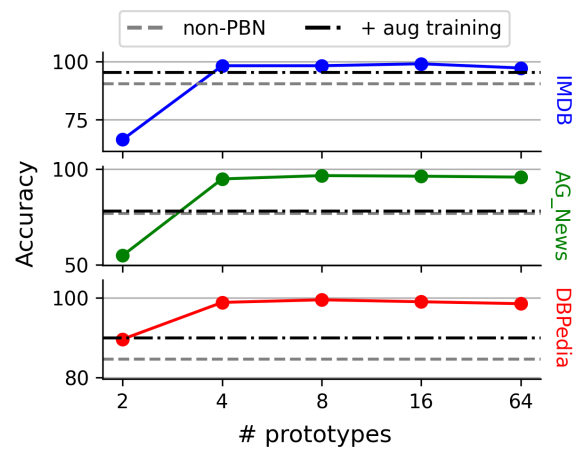


Figure 16: Accuracy of PBNs under static adversarial settings, with different numbers of prototypes, with other hyperparameters set to their best values and averaged across other possible variables (e.g., backbone and attack type). The dotted line represents the ASR for the vanilla LMs.