

# EAGLE: A Domain Generalization Framework for AI-generated Text Detection

Amrita Bhattacharjee<sup>1</sup>, Raha Moraffah<sup>1</sup>, Joshua Garland<sup>1</sup>, and Huan Liu<sup>1</sup>

Arizona State University, Tempe AZ 85281, USA  
{abhatt43,rmoraffa,joshua.garland,huanliu}@asu.edu

**Abstract.** With the advancement in capabilities of Large Language Models (LLMs), one major step in the responsible and safe use of such LLMs is to be able to detect text generated by these models. While supervised AI-generated text detectors perform well on text generated by older LLMs, with the frequent release of new LLMs, building supervised detectors for identifying text from such new models would require new labeled training data, which is infeasible in practice. In this work, we tackle this problem and propose a domain generalization framework for the detection of AI-generated text from unseen target generators. Our proposed framework, EAGLE, leverages the labeled data that is available so far from older language models and learns features invariant across these generators, in order to detect text generated by an unknown target generator. EAGLE learns such domain-invariant features by combining the representational power of self-supervised contrastive learning with domain adversarial training. Through our experiments we demonstrate how EAGLE effectively achieves impressive performance in detecting text generated by unseen target generators, including recent state-of-the-art ones such as GPT-4 and Claude, reaching detection scores of within 4.7% of a fully supervised detector.

**Keywords:** Large Language Models · AI-generated Text Detection · Domain Generalization.

## 1 Introduction

Large language models (LLMs) are becoming ubiquitous and an increasing number of people are using these models, often equipped with easy-to-use public facing APIs, for a variety of use cases. Being brilliant productivity aids, these systems are being used for creative writing, homework help, education, and general writing assistants. However, given the human-like quality of text generated by these models, these are also susceptible to being used maliciously to generate misinformation, disinformation and misleading content at scale [39,43,16]. Such misuse is even more harmful during major social or political events such as presidential elections [38,24]. Furthermore, with improvement in model fluency, recent work has shown that human readers often struggle to differentiate be-

tween actual human written text and LLM generated text [10]. This necessitates the development of automated systems to detect such LLM-generated<sup>1</sup> content.

Most existing work on automated detectors for AI-generated text are supervised classifiers trained using labeled data from some text generator(s), but these detectors do not generalize well to newly released, possibly much larger and more capable LLMs [33]. However, given the sheer variety of LLMs available alongside the fast development and release of new LLMs, it is challenging to develop, train and maintain a general purpose AI text detector that would also work for newly emerging LLMs, since this requires continuous collection and/or generation of new labeled training datasets. While recently there have been some zero-shot detection methods proposed [30,17], most of these methods rely on a proxy language model, and are therefore highly sensitive to the choice of this model. Furthermore, these methods do not make use of the available labeled data from older language models. In this work, we aim to leverage the labeled data that is already available, and we propose a novel domain generalization framework to learn domain-invariant (i.e., LLM-invariant) features to perform detection on text generated from a completely unseen domain (i.e., a new LLM). We assume that existing data from older generators must contain crucial discriminative features that we can leverage in a detection framework. Our proposed framework aims to capture the (1) cross-domain invariance: that is, learning the invariant features across different generators, and (2) in-domain invariance: that is, learning better latent representations in order to be robust against minor perturbations in text. For the purposes of this paper, we focus on the critical issue of AI-generated news articles and we demonstrate the effectiveness of our proposed framework on established benchmark datasets as well as data from newer LLMs, including our own GPT-4 generated data. Overall, in this work, our contributions are:

- We propose EAGLE<sup>2</sup>: a novel domain generalization for AI-generated text detection framework to detect text from new, unseen target generators, by leveraging labeled data from pre-existing, possibly older generators.
- Through comprehensive experiments on text from a variety of language models, we demonstrate the effectiveness of EAGLE in learning domain-invariant features.
- Alongside generating our own GPT-4 data, we evaluate the efficacy of EAGLE on detecting text from new state-of-the-art LLMs, by only leveraging data from older, much smaller language models.

## 2 Related Work

**AI-generated Text Detection.** With the rapid progress of language models over the last several years, there have also been approaches proposed for detection of text generated by such models. In the case where plenty of labeled data is

<sup>1</sup> or ‘AI-generated’, used interchangeably throughout this paper.

<sup>2</sup> *inspired by the sharp eyesight of this specific bird of prey*

available, fine-tuned pre-trained language models are often the best performing detectors [20,18]. An example of this is the OpenAI detector that is simply a RoBERTa [27] model fine-tuned on GPT-2 data [36]. A recent supervised method called Ghostbuster uses a series of weaker models followed by a search over combination functions and then a linear classifier [42]. Some recent work also explore the use of LLMs as the detector for detecting AI-generated text [6]. Authors in [5] propose an unsupervised domain adaptation framework to detect AI-generated text by leveraging labeled source and unlabeled target data.

**Zero-shot AI-generated Text Detection.** In addition to supervised methods for detection, recently there has been a lot of effort in the area of zero-shot detection of text generated by AI text generators, i.e., LLMs. While some works [33,29] analyze the zero-shot transfer capabilities of AI-text detectors to text from new generators, some also propose novel zero-shot or unsupervised detection methods. Some of these methods [15,3] leverage statistical measures to identify generation artifacts across common sampling schemes, while some rely on assumptions surrounding the log probabilities of the generated text under a proxy model [30,4,37]. Under a full black-box setting, with no access to the token probabilities, a recent approach [45] uses n-gram analysis to detect such AI-generated text. Another interesting recent zero-shot detection method compares the perplexity of the input text under two related language models as a signal towards detection [17]. Authors in [5] proposes a domain adaptation framework for unsupervised detection. Zero-shot detection approaches have also been proposed for code generated by LLMs [46].

While these various approaches have shown promising results in the detection of AI-generated text both in fully supervised and zero-shot scenarios, there is currently no work that leverages labeled data from older generators to perform unsupervised detection of text from newer generators. To the best of our knowledge, this is the first work to propose a domain generalization framework for AI-generated text detection, whereby we investigate if we can tackle the real-world scenario of detecting text from an unseen generator while leveraging labeled data from older, potentially much smaller generators.

### 3 Background & Preliminaries

Operationalizing machine learning frameworks for real-world use requires such methods to have the ability to perform well under domain shift. Examples of such cases in text classification are product review sentiment classifiers that are trained for a specific category of products (e.g., Books) and need to perform well to the unseen category of Electronics. There have been substantial work done to tackle such scenarios by domain adaptation techniques. Typically such domain adaptation methods assume access to one or multiple source domains as well as *unlabeled data* from the target domain. The model is then trained to optimize an objective function that retains performance on the source domain while learning the invariance across the source and target domains.

In our task of AI-generated text detection, we treat each generator or LLM as a ‘domain’; that is, texts generated by different LLMs follow a different distribution. In this work, we consider the case where we have access to labeled data from  $k$  such generators, and we aim to detect text generator by an unseen  $(k + 1)$ th generator. This is a typical domain generalization setting. While numerous approaches have been proposed for this setting in the field of computer vision, text is relatively under explored [21,25,40]. A widely used method is DANN - Domain Adversarial Neural Networks [13,14] - that uses domain adversarial training to make the model learn domain invariant features. Such a method has also been adapted to text applications for sentiment classification, albeit with unlabeled data from the target. In this work, we build on top of this framework by: (1) modifying it to only use source domain data from multiple sources and *no* target data, and (2) adding an in-domain contrastive regularization function via which the model aims to learn more robust, better representations of the text. As supported by prior work [44], this additional regularization component should provide more generalizability to the trained model.

## 4 Problem Definition

In this work, we address the task of AI-generated text detection. Specifically, we assume we have access to labeled data from  $k$  different AI-text generators, or LLMs, denoted by  $\{D^1, D^2, \dots, D^k\}$ . Each domain is defined as the combination of the input (or feature) space and the label space:  $D^i = \{X^i, Y^i\}, \forall i \in [1, k]$ . Typically, for most domain generalization use-cases, including our problem setting, the label space is constant across all the different domains, that is,  $Y^1 = Y^2 = \dots = Y^k$ . The goal of domain generalization in a classification setting is to learn a mapping  $f : X \rightarrow Y$  that captures the domain-invariant features across the  $k$  domains in order to minimize the predictive loss on data from a new unseen domain  $D^{k+1}$ .

In the task of AI-generated text detection, with text from different generators, we formulate the problem as follows: Each generator  $i$  is a domain  $D^i$ ,  $X^i$  consists of text samples either generated by the generator  $i$  or written by a human, the label space  $Y = \{0, 1\}$  denoting human-written and AI-generated text respectively. Hence the function  $f(\cdot)$  that we aim to learn should ideally capture the features invariant across different text generators, while being discriminative enough to be able to distinguish between generated vs. human-written text. In the following section, we describe our proposed framework to learn this function to perform the task of AI-generated text detection on unknown target domains.

## 5 Proposed Framework

In this section, we introduce our proposed EAGLE framework, as showed in Figure 1, and describe each component in detail.

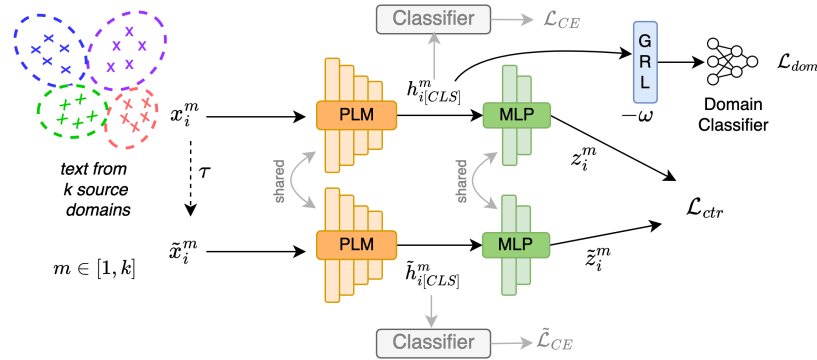


Fig. 1. Our proposed EAGLE framework.

**Classification Backbone.** Given the superior performance of pre-trained language models on a variety of tasks, and the ease of fine-tuning such models on task-specific data, we opt to use a pre-trained language model (PLM) RoBERTa [27] (roberta-base) from Huggingface transformers<sup>3</sup>, along with a classifier head on top of it. Our input text consists of labeled text data from  $k$  different domains:  $x_i^m \in \{D^m\}_{m=1}^k$ . Each text  $x_i^m$  is fed into the PLM and we obtain the final hidden layer embedding  $h_{i[CLS]}^m$ . Since we have the labeled data from each of the  $k$  sources, we feed this embedding into the linear classifier layer and compute the binary-cross entropy loss  $\mathcal{L}_{CE}$ . Since we also have a perturbed version of the input (as explained later in the contrastive loss section), we compute a similar loss for this and denote it as  $\tilde{\mathcal{L}}_{CE}$ .

$$\mathcal{L}_{CE} = -\frac{1}{b} \sum_{i=1}^b [y_i \log p(y_i | h_{i[CLS]}^m) + (1 - y_i) \log(1 - p(y_i | h_{i[CLS]}^m))] \quad (1)$$

where  $y_i$  is the ground truth label for input  $x_i^m$  and  $b$  is the batch size. We have a similar loss for the perturbed version of the text as well.

**Domain Adversarial Training.** For the task of domain generalization, we would require our framework to learn features that are invariant across the different domains, i.e., text generators. To facilitate this, we employ domain adversarial training [13,14]. To do this, we use the final hidden layer embedding  $h_{i[CLS]}^m$  for each input text, and feed this into a domain classifier that computes a domain loss. The objective of the domain classifier is to minimize the classification error of distinguishing the source and target domains, while the objective of the entire classification model is to learn domain-invariant i.e., transferable features. Through this adversarial training process, the model learns better, more discriminative representations of the input text specific to the downstream task that are also domain-invariant. To facilitate gradient-based optimization of

<sup>3</sup> <https://huggingface.co/FacebookAI/roberta-base>

these conflicting objectives, a gradient reversal layer (GRL) is used [14]. More precisely, the hidden layer embedding  $h_{i[CLS]}^m$  is fed into the gradient reversal layer, and then into the domain classifier. The gradient reversal layer does not have any trainable parameters. In the forward pass, it essentially functions as an identity transform, and in the backpropagation, the GRL layer reverses the gradient by multiplying the gradient from the subsequent layer by a negative scalar value  $-\omega$ . Therefore, in the forward pass, the output of this layer is

$$GRL(h_{i[CLS]}^m) = h_{i[CLS]}^m \quad (2)$$

This is now input into the domain classifier which is simply a dropout layer followed by a linear layer with dimension  $768 \times num\_domains$ , where  $num\_domains$  is the number of source domains we are using. Say the predicted domain label for this input is  $d_i^m$ . We compute the domain loss as the cross-entropy loss as

$$\mathcal{L}_{dom} = -\frac{1}{k} \sum_{m=1}^k \sum_{i=1}^b d_i^m \log p(d_i^m) \quad (3)$$

where  $k$  is the number of source domains,  $b$  is the batch size. During backpropagation of this loss  $L_{dom}$  through the GRL layer, the gradient becomes:

$$GRL\left(\frac{\partial L_{dom}}{\partial \theta_f}\right) = -\omega \frac{\partial L_{dom}}{\partial \theta_f} \quad (4)$$

where  $\theta_f$  are the parameters of the feature extractor network, which is the pre-trained RoBERTa language model in our case, and  $-\omega$  is the negative scalar with which the gradient is reversed. Training the model using the GRL ensures that the domain classifier is optimized to predict domain labels correctly, while the RoBERTa parameters are optimized to deceive the domain classifier, thereby learning domain invariant features.

**Contrastive Learning for Better Representations.** Inspired by previous work, we add a contrastive loss component to learn better representations of the input text. A contrastive loss component would act as a regularizer to learn more robust representations of the input. Ideally our framework should be robust to small perturbations in the input, to facilitate in-domain invariance to such noisy perturbations. To facilitate that, we use a loss such as in [5]. Following previous work, for each input, we apply a perturbation  $\tau$  which is synonym replacement. We use the hidden layer embeddings of the original and perturbed texts:  $h_{i[CLS]}^m$  and  $\tilde{h}_{i[CLS]}^m$ , and pass these through a projection layer in order to compute the contrastive loss in a lower dimensional space, as done in [5], and compute a SimCLR-style [9] contrastive loss between the original and the perturbed views of the input text. Here the original and perturbed views are the positives, and all other samples in the mini-batch are considered as negatives. Therefore, the contrastive loss is

$$\mathcal{L}_{ctr} = - \sum_{i \in b} \log \frac{\exp(\text{sim}(z_i^m, \tilde{z}_i^m)/t)}{\sum_{j=1}^{2|b|} \mathbb{1}_{[j \neq i]} \exp(\text{sim}(z_i^m, z_j)/t)} \quad (5)$$

where  $z_i^m$  and  $\tilde{z}_i^m$  refer to the projected embeddings in the lower dimension space, corresponding to the original and the perturbed views of the input text respectively,  $t$  is the temperature.

The final training objective is

$$\mathcal{L} = \frac{\lambda_1}{2}(\mathcal{L}_{CE} + \tilde{\mathcal{L}}_{CE}) + \lambda_2\mathcal{L}_{ctr} + \lambda_3\mathcal{L}_{dom} \quad (6)$$

where  $\lambda_1, \lambda_2$  and  $\lambda_3$  are hyper-parameters,  $\tilde{\mathcal{L}}_{CE}$  is the cross-entropy loss from the perturbed version of the text.

## 6 Experimental Settings

In this section, we describe the datasets we trained and evaluated our framework on, along with the baselines and experimental settings used.

### 6.1 Datasets

In this work, we specifically focus on AI-generated news articles. The nature of our domain generalization task also requires having AI-generated text from separate AI-text generators (i.e., language models or ‘domains’). Therefore, in our experiments, we use the benchmark dataset TuringBench. TuringBench [41] is an extensive dataset consisting of new-style text from a total of 19 different generators along with human-written text. These 19 generators comprise various sizes of 10 different model architectures: {GPT-1 [34], GPT-2 [35], GPT-3 [7], GROVER [48], CTRL [22], XLNET [47], XLM [26], TRANSFORMER-XL [11], FAIR [31,8], and PPLM [12]}. Following previous work [5], we use a subset of 6 of these generators as our domains. This choice ensures that we have a good coverage of different architectures and model families of generators. The generators are:

**CTRL:** This is a 1.5B parameter transformer-based language model that is capable of controlled generation of text. Generation can be conditioned on control codes such as style, sentiment, etc. Authors in [41] generate text from CTRL using the *News* control code.

**FAIR\_wmt19:** This is a 656M parameter transformer-based model developed by FAIR as part of their submission to the WMT19 news translation task. Texts are generated using the English version of the model using the FAIRSEQ [32] toolkit.

**GPT2\_xl:** This is the 1.5B version of the GPT-2 model, the precursor to GPT-3 and the more recent GPT-3.5 and GPT-4 models.

**GPT-3:** This is the largest model in the TuringBench dataset. It is a 175B model which is a successor of the previous GPT-2 model.

**GROVER\_mega:** This is a 1.5B parameter transformer-based model, which is trained to generate news style text.

**XLM:** This 550M parameter model is the smallest model we use in our evaluation. This is a transformer-based model that is designed for cross-lingual tasks.

Given the rapid development and release of new text generators or LLMs, it is necessary to evaluate AI-text detection systems on new language models, that are not covered by the TuringBench dataset. Therefore, we also use news-style data generated by new language models such as **GPT-3.5**, **GPT-4** from OpenAI [1] and **Claude** from Anthropic [2]. GPT-3.5 (or often referred to as ChatGPT) and its newer variant GPT-4 have garnered immense attention from academics, industry, educators and practitioners for its impressive instruction-following and text generation capabilities, as well as superior performance on complex tasks. Given the human-like quality of text generated by GPT-3.5 and GPT-4, there have also been concerns surrounding whether such text can be detected by automated AI-text detection tools. To investigate this, we use such data in our evaluation. For GPT-3.5, we use the data from [5]. For Claude, we use data from [42]. For GPT-4, we generate our own news-style data as described below.

Following a similar data generation pipeline as in [5], we first collect a set of 2,000 human written articles from CNN and Washington Post (as done in previous work [41]). For each article, we use the **headline** of the article to generate an article using GPT-4 by simply prompting the model with the following prompt: "Generate a news article with the headline '<headline>'." Due to resource constraints, we only generate 2,000 articles using GPT-4. Therefore, for GPT-4, we have a balanced set of 2,000 human-written articles and 2,000 GPT-4 generated articles.

## 6.2 Baseline Detection Methods

Similar to previous work [5,19], we use a variety of unsupervised baselines:

**GLTR** [15]: This is a set of four statistical measures to identify whether a text is AI generated or not. These are computed based on token-wise log probabilities. These measures are: (1) log probability ( $\log p(x)$ ): that assumes that higher log probability indicates that the text is AI-generated, (2) average token rank and (3) token log-rank: that are based on the assumption that tokens with lower rank are possibly AI-generated, and (4) predictive entropy: that is based on the assumption that AI-generated text often has less diversity than human-written, and therefore less entropy.

**DetectGPT** [30]: This method relies on a proxy language model to compute log probabilities of generated tokens. The main assumption in this method is the minor perturbations of AI-generated text would fall in the negative curvature region of the log-likelihood curve, while such perturbations to human-written text does not follow this trend. This assumption on the curvature of the log probabilities is used as the discriminative feature for prediction.

Apart from these unsupervised baselines, in order to test the transferability across different domains (i.e., generators) we also use standard baselines as in [33]:

**Data Mix**: Here we combine the data from the training splits of the  $k$  known generators, and train one detector on this unified dataset. During evaluation, we evaluate this trained model on the test split of the unseen  $(k + 1)$ th generator.



**Ensemble:** Here for each of the  $k$  known generators, we train a separate detector, and then during inference on the unknown  $(k + 1)$ th generator, we evaluate all  $k$  detectors, and use an aggregation of the  $k$  prediction scores to get the final label. For the aggregation, we use max and average.

For both Data Mix and Ensemble baselines, we use a RoBERTa model (roberta-base) with a classification head on top as the detector. We train this for 3 epochs with Adam optimizer [23] and a learning rate of  $2 \times 10^{-5}$ .

### 6.3 Implementation Details

All our experiments were performed using PyTorch, on an NVIDIA A100 GPU with 40 GB memory. For hyperparameter values in the detection experiments, we use grid search and use the best performing value in our final experiments. For generating the GPT-4 data, we use temperature of 0.5 and top\_p of 1. To facilitate reproducibility, all code and data will be made available at <link-to-be-inserted-after-blind-review>.

Target Generator	Fully Supervised	GLTR				Det.GPT	DataMix	Ensemble		Ours (k=5)
		log p(x)	rank	log rank	entropy			Avg	Max	
CTRL	1	0.951	0.849	<u>0.956</u>	0.379	0.793	0.902	0.801	0.761	<b>0.984</b>
FAIR_wmt19	0.999	0.558	0.618	0.546	0.656	0.5045	0.797	<u>0.836</u>	0.728	<b>0.896</b>
GPT2_xl	0.998	0.485	0.508	0.48	0.631	0.529	<b>0.995</b>	<u>0.993</u>	0.982	0.94
GPT3	0.988	0.362	0.356	0.341	0.756	0.5485	<u>0.976</u>	0.951	0.947	<b>0.983</b>
GROVER_mega	0.996	0.434	0.469	0.434	0.592	0.5415	<u>0.895</u>	0.702	0.705	<b>0.912</b>
XLM	1	0.473	0.762	0.442	0.696	0.7355	0.864	<u>0.978</u>	0.952	<b>0.98</b>

**Table 1.** Performance of our framework on TuringBench dataset. Scores are AUROC values averaged over three different seeds. Best values are in **bold** and second best values are underlined. Det.GPT refers to the DetectGPT baseline [30].

## 7 Experimental Results

To investigate the effectiveness of our EAGLE framework, we perform a comprehensive set of experiments to answer the following research questions:

- **RQ1:** Can our framework learn domain-invariant features and perform well on unseen generators?
- **RQ2:** Can our framework learn to detect text from new generators after being trained only on text from older generators?

Apart from these main research questions, we also explore the effectiveness of the different components in the framework, as well as the effect of number of source domains.

### 7.1 RQ1: Performance on Unseen Targets

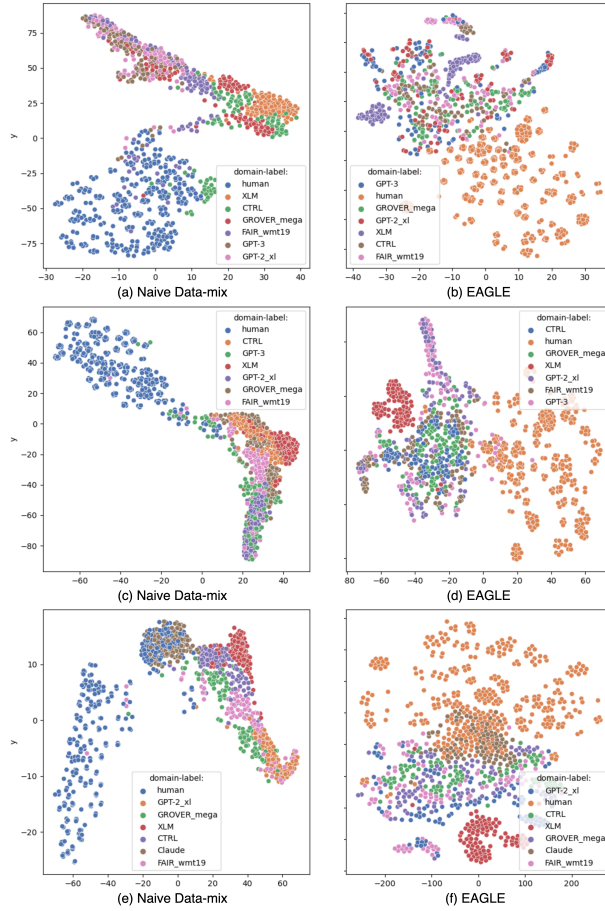
We evaluate our framework on the TuringBench dataset and report these results in Table 1. Here we compare the performance of our framework with unsupervised baselines as described in Section 6.2. We also compare performance with a fully supervised detection model that serves as an upper bound. For the fully supervised model we use a RoBERTa [27] (roberta-base<sup>4</sup>) model with a classification head on top and the model is fine-tuned on labeled data from the target generator. This essentially serves as an upper bound for the performance. We report the AUROC scores from this experiment in Table 1. We see that for all of the generators in our experiment, except GPT2\_xl, our proposed framework performs better than the unsupervised baselines, with no labeled target data.

To understand if our framework is effectively learning the domain invariant features, we visualize the t-SNE plots [28] of the learned representations from our EAGLE framework vs. a naive data mix framework with the same set of source domains, for a fair comparison. We produce these visualizations for three target generators (top row to bottom row): CTRL, GPT3 and Claude, and show these in Figure 2. For each target generator, we see that the data mix plot (left) shows somewhat homogeneous clusters for each of the generators, implying texts from these different source generators and the target generators lie in disjoint spaces in the latent representation space. However, for the EAGLE plots (right), we see more overlap between text from the different generators. This is encouraging since this implies our framework has successfully learned domain-invariant features, thereby making the data points inseparable based on the domain label. For Claude, we see some overlap between human written articles and Claude generated ones, implying that Claude-generated text is still somewhat challenging to distinguish from human-written text.

Target	Data Mix		Ensemble		EAGLE (Ours)	
	Avg	Max	Avg	Max	Avg	Max
Claude	36.67 ± 22.18	69	0.4 ± 0.43	0.8	<b>40.67 ± 29.85</b>	<b>97</b>
GPT-3.5	75 ± 19.79	98	54.72 ± 45.22	96	<b>86.17 ± 14.66</b>	<b>99</b>
GPT-4	32.5 ± 20.98	64	1.62 ± 1.63	3.1	<b>58 ± 37.18</b>	<b>99</b>

**Table 2.** Performance of our model on data from newer generators. Values are TPR. We report both the average across all possible sets of 5 sources (average ± standard deviation) and the maximum across the choices. Best performance is in **bold**.

<sup>4</sup> <https://huggingface.co/FacebookAI/roberta-base>



**Fig. 2.** t-SNE visualizations for representations learned by our EAGLE model, vs the naive data-mix baseline with the same set of sources. (a)-(b) denote plots for data mix and EAGLE representations for target generator CTRL, respectively, (c)-(d) for target generator GPT-3 and (e)-(f) for target generator Claude.

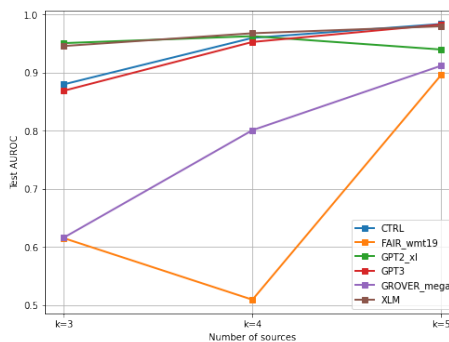
Framework variant	CTRL	GRO_m	XLM
EAGLE	<b>0.984</b>	<b>0.912</b>	<b>0.980</b>
EAGLE $-\mathcal{L}_{ctr}$	0.763	0.883	0.932
EAGLE $-\mathcal{L}_{dom}$	0.948	0.814	0.968
EAGLE $-\mathcal{L}_{ctr} - \mathcal{L}_{dom}$	0.902	0.895	0.970

**Table 3.** Performance comparison across different variants of EAGLE. GRO\_m refers to GROVER\_mega. Scores are AUROC and best performance is in **bold**.

## 7.2 RQ2: Detection of Text Generated by Newer LLMs

Here we are interested in evaluating the performance of our framework on new state-of-the-art LLMs: Claude, GPT-3.5 and GPT-4, *without* using data generated by any of these in the training set. For this experiment, we use the same framework and  $k = 5$  number of source domains. For each new LLM (Claude, GPT-3.5 and GPT-4), we only consider  $k = 5$  *older* generators as the source domains. This source domain data is the same TuringBench data as used in the previous experiment. Essentially, we are interested in evaluating how well our proposed framework can transfer discriminative features from older generators with possibly easily available data, to newer generators from which text is more difficult to identify. We report the true positive rate of detection in Table 2 from the best performing set of sources, along with the average across different 5-set choices. We also compare results from our framework with two baseline methods: (1) Data Mix - where the training data comes from  $k = 5$  source domains (from generators in TuringBench), and (2) Ensemble - where  $k = 5$  models, each trained on a single source are used. We keep the same set of sources across all 3 settings: our framework, data mix and ensemble. This is because we are interested in understanding how our framework can learn domain invariant features from the same data, above and beyond what is already do-able by using all of the same sources in a naive way.

We report the true positive rate of detection of LLM-generated text for both ‘avg’ and ‘max’ settings: where we take the average and maximum across all combinations of 5 out of the 6 sources, respectively (as described in Section 6.1). We see that our framework performs better than both the Data Mix and Ensemble baselines by a significant margin. Interesting, we see that the Ensemble approach performs the worst, and has negligible performance for detecting Claude and GPT-4 generated text. However, performance varies significantly across different choices of sources as depicted by the high values of standard deviation in most cases.



**Fig. 3.** Variation of test performance with respect to number  $k$  of source domains or generators used in training.

## 7.3 Ablation: Effectiveness of Framework Components

In this experiment, we evaluate the effectiveness of the different components in our framework. We remove one component at a time, train and evaluate each variant and report these results for

a randomly chosen set of three target domains in Table 3. EAGLE  $-\mathcal{L}_{ctr}$  removes the contrastive loss component (along with the cross-entropy loss for the perturbed version of the input text), EAGLE  $-\mathcal{L}_{dom}$  removes the domain loss, and EAGLE  $-\mathcal{L}_{ctr} - \mathcal{L}_{dom}$  removes both. This last variant is essentially similar to the Data Mix baseline. We see performance drops for each of those variants, while  $\mathcal{L}_{ctr}$  and  $\mathcal{L}_{dom}$  having different degrees of impact on the performance for different target generators. Overall, we see our full framework EAGLE performs the best across all the target generators.

#### 7.4 Hyperparameter Analysis: Effect of Number of Sources

We are also interested in evaluating how sensitive our model is to the number of source generators we include in our training. For this analysis, we vary the number of source domains/generators  $k = \{3, 4, 5\}$ . We show the variation of performance across the different choices of  $k$  in Figure 3. We see that in general, the performance increases with increase in the value of  $k$ . This is possibly due to the increase in variability and diversity of the data as we increase  $k$ , since the sources are all different model families with different architectural backbones, with the exception of GPT2\_xl and GPT3, albeit differing in the type of attention used in the 2 models [7]. Therefore, we use  $k = 5$  as the number of source domains in all our experiments.

## 8 Conclusion and Future Work

In this work, we propose a novel framework called EAGLE, to perform AI-generated text detection from unseen target generators, by effectively leveraging data from older, possibly much smaller generators. Our framework learns domain invariant features via a gradient reversal layer and adversarial training, thereby retaining task-specific discriminative features, while learning to ignore domain specific features. We demonstrate the effectiveness of the proposed framework via experiments on a vast variety of language models, ranging from smaller language models such as XLM, to recent state of the art models such as GPT-4 and Claude. We also generate our own GPT-4 data that we will make available for research purposes upon request. Our experiments show that EAGLE can effectively leverage data from older generators, learn transferable features and perform detection on unseen target LLMs, in an unsupervised manner. Our framework and findings pave the way for building more detectors for newly emerging LLMs, simply by leveraging data from older, smaller generators.

Here, we assume the test data only comes from one generator. Future work could explore the more challenging setting of multiple unseen test generators, perhaps assuming some mixture of distributions over the texts. Furthermore, we only use new-style text since this is the only type of text that is widely available from a variety of different generators. Future work could also explore how our detector performs on other types of text. Another direction could be generalization in two dimensions: first over the text generators, and second, over the type of text, such as news, scientific article, finance, etc.

## 9 Ethical Statement

With the prevalence of LLMs everywhere, we are increasingly exposed to AI-generated text. While automatic detection of such AI-generated text is a first step to evaluate content authenticity, care must be taken while using such automatic systems in practice. Most of these detectors are black-box models, do not provide explanations for their predictions, and are therefore challenging to be used in practice. For high-stakes applications where users may be penalized heavily for using AI-generated content, false positives can be highly undesirable. Similarly, for applications such as evaluating content authenticity on a creative writing website, false negatives would be highly undesirable, since real human authors may potentially miss out on the credit they deserve. Furthermore, alongside flagging generated content, systems also need to identify the intent behind the content generation - whether it is being used maliciously or not. However, such determination is highly subjective and therefore a vastly unsolved issue. For the purposes of this work, we do not condone misuse of our proposed detection framework especially for high-stakes applications and urge users of such applications to instead resort to a human-in-the-loop, hybrid solution.

## Acknowledgements

This work is supported by the DARPA SemaFor project (HR001120C0123), Army Research Lab (W911NF2020124) and Army Research Office (W911NF2110030). The views, opinions and/or findings expressed are those of the authors.

## References

1. Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F.L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al.: Gpt-4 technical report. arXiv preprint arXiv:2303.08774 (2023)
2. Bai, Y., Kadavath, S., Kundu, S., Askell, A., Kernion, J., Jones, A., Chen, A., Goldie, A., Mirhoseini, A., McKinnon, C., et al.: Constitutional ai: Harmlessness from ai feedback. arXiv preprint arXiv:2212.08073 (2022)
3. Bakhtin, A., Gross, S., Ott, M., Deng, Y., Ranzato, M., Szlam, A.: Real or fake? learning to discriminate machine from human generated text. arXiv preprint arXiv:1906.03351 (2019)
4. Bao, G., Zhao, Y., Teng, Z., Yang, L., Zhang, Y.: Fast-detectgpt: Efficient zero-shot detection of machine-generated text via conditional probability curvature. arXiv preprint arXiv:2310.05130 (2023)
5. Bhattacharjee, A., Kumarage, T., Moraffah, R., Liu, H.: Conda: Contrastive domain adaptation for ai-generated text detection. arXiv preprint arXiv:2309.03992 (2023)
6. Bhattacharjee, A., Liu, H.: Fighting fire with fire: Can chatgpt detect ai-generated text? arXiv preprint arXiv:2308.01284 (2023)
7. Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al.: Language models are few-shot learners. *Advances in neural information processing systems* **33**, 1877–1901 (2020)

8. Chen, P.J., Lee, A., Wang, C., Goyal, N., Fan, A., Williamson, M., Gu, J.: Facebook ai’s wmt20 news translation task submission. arXiv preprint arXiv:2011.08298 (2020)
9. Chen, T., Kornblith, S., Norouzi, M., Hinton, G.: A simple framework for contrastive learning of visual representations. In: International conference on machine learning. pp. 1597–1607. PMLR (2020)
10. Clark, E., August, T., Serrano, S., Haduong, N., Gururangan, S., Smith, N.A.: All that’s human is not gold: Evaluating human evaluation of generated text. arXiv preprint arXiv:2107.00061 (2021)
11. Dai, Z., Yang, Z., Yang, Y., Carbonell, J., Le, Q.V., Salakhutdinov, R.: Transformer-xl: Attentive language models beyond a fixed-length context. arXiv preprint arXiv:1901.02860 (2019)
12. Dathathri, S., Madotto, A., Lan, J., Hung, J., Frank, E., Molino, P., Yosinski, J., Liu, R.: Plug and play language models: A simple approach to controlled text generation. arXiv preprint arXiv:1912.02164 (2019)
13. Ganin, Y., Lempitsky, V.: Unsupervised domain adaptation by backpropagation. In: International conference on machine learning. pp. 1180–1189. PMLR (2015)
14. Ganin, Y., Ustinova, E., Ajakan, H., Germain, P., Larochelle, H., Laviolette, F., March, M., Lempitsky, V.: Domain-adversarial training of neural networks. *Journal of machine learning research* **17**(59), 1–35 (2016)
15. Gehrmann, S., Strobelt, H., Rush, A.M.: Gltr: Statistical detection and visualization of generated text. arXiv preprint arXiv:1906.04043 (2019)
16. Goldstein, J.A., Chao, J., Grossman, S., Stamos, A., Tomz, M.: Can ai write persuasive propaganda? *SocArXiv*. April **8** (2023)
17. Hans, A., Schwarzschild, A., Cherepanova, V., Kazemi, H., Saha, A., Goldblum, M., Geiping, J., Goldstein, T.: Spotting llms with binoculars: Zero-shot detection of machine-generated text. arXiv preprint arXiv:2401.12070 (2024)
18. He, X., Shen, X., Chen, Z., Backes, M., Zhang, Y.: Mgtbench: Benchmarking machine-generated text detection. arXiv preprint arXiv:2303.14822 (2023)
19. Hu, X., Chen, P.Y., Ho, T.Y.: Radar: Robust ai-text detection via adversarial learning. arXiv preprint arXiv:2307.03838 (2023)
20. Ippolito, D., Duckworth, D., Callison-Burch, C., Eck, D.: Automatic detection of generated text is easiest when humans are fooled. arXiv preprint arXiv:1911.00650 (2019)
21. Jia, C., Zhang, Y.: Prompt-based distribution alignment for domain generalization in text classification. In: Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing. pp. 10147–10157 (2022)
22. Keskar, N.S., McCann, B., Varshney, L.R., Xiong, C., Socher, R.: Ctrl: A conditional transformer language model for controllable generation. arXiv preprint arXiv:1909.05858 (2019)
23. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
24. Klepper, D., Swenson, A.: Ai-generated disinformation poses threat of misleading voters in 2024 election (May 2023), <https://www.pbs.org/newshour/politics/ai-generated-disinformation-poses-threat-of-misleading-voters-in-2024-election>
25. Krishnan, J., Purohit, H., Rangwala, H.: Diversity-based generalization for neural unsupervised text classification under domain shift. In: ECML-PKDD (2020)
26. Lample, G., Conneau, A.: Cross-lingual language model pretraining. arXiv preprint arXiv:1901.07291 (2019)

27. Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., Stoyanov, V.: Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692 (2019)
28. Van der Maaten, L., Hinton, G.: Visualizing data using t-sne. *Journal of machine learning research* **9**(11) (2008)
29. Mireshghallah, F., Mattern, J., Gao, S., Shokri, R., Berg-Kirkpatrick, T.: Smaller language models are better black-box machine-generated text detectors. arXiv preprint arXiv:2305.09859 (2023)
30. Mitchell, E., Lee, Y., Khazatsky, A., Manning, C.D., Finn, C.: Detectgpt: Zero-shot machine-generated text detection using probability curvature. arXiv preprint arXiv:2301.11305 (2023)
31. Ng, N., Yee, K., Baevski, A., Ott, M., Auli, M., Edunov, S.: Facebook fair’s wmt19 news translation task submission. arXiv preprint arXiv:1907.06616 (2019)
32. Ott, M., Edunov, S., Baevski, A., Fan, A., Gross, S., Ng, N., Grangier, D., Auli, M.: fairseq: A fast, extensible toolkit for sequence modeling. arXiv preprint arXiv:1904.01038 (2019)
33. Pu, X., Zhang, J., Han, X., Tsvetkov, Y., He, T.: On the zero-shot generalization of machine-generated text detectors. arXiv preprint arXiv:2310.05165 (2023)
34. Radford, A., Narasimhan, K., Salimans, T., Sutskever, I., et al.: Improving language understanding by generative pre-training (2018)
35. Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., et al.: Language models are unsupervised multitask learners. *OpenAI blog* **1**(8), 9 (2019)
36. Solaiman, I., Brundage, M., Clark, J., Askill, A., Herbert-Voss, A., Wu, J., Radford, A., Krueger, G., Kim, J.W., Kreps, S., et al.: Release strategies and the social impacts of language models. arXiv preprint arXiv:1908.09203 (2019)
37. Su, J., Zhuo, T.Y., Wang, D., Nakov, P.: Detectllm: Leveraging log rank information for zero-shot detection of machine-generated text. arXiv preprint arXiv:2306.05540 (2023)
38. Swenson, A., Chan, K.: Election disinformation takes a big leap with ai being used to deceive worldwide (Mar 2024), <https://apnews.com/article/artificial-intelligence-elections-disinformation-chatgpt-bc283e7426402f0b4baa7df280a4c3fd>
39. Tamkin, A., Brundage, M., Clark, J., Ganguli, D.: Understanding the capabilities, limitations, and societal impact of large language models. arXiv preprint arXiv:2102.02503 (2021)
40. Tan, Q., He, R., Bing, L., Ng, H.T.: Domain generalization for text classification with memory-based supervised contrastive learning. In: *Proceedings of the 29th International Conference on Computational Linguistics*. pp. 6916–6926 (2022)
41. Uchendu, A., Ma, Z., Le, T., Zhang, R., Lee, D.: Turingbench: A benchmark environment for turing test in the age of neural text generation. arXiv preprint arXiv:2109.13296 (2021)
42. Verma, V., Fleisig, E., Tomlin, N., Klein, D.: Ghostbuster: Detecting text ghost-written by large language models. arXiv preprint arXiv:2305.15047 (2023)
43. Vykopal, I., Pikuliak, M., Srba, I., Moro, R., Macko, D., Bielikova, M.: Disinformation capabilities of large language models. arXiv preprint arXiv:2311.08838 (2023)
44. Xu, H., Mannor, S.: Robustness and generalization. *Machine learning* **86**, 391–423 (2012)
45. Yang, X., Cheng, W., Petzold, L., Wang, W.Y., Chen, H.: Dna-gpt: Divergent n-gram analysis for training-free detection of gpt-generated text. arXiv preprint arXiv:2305.17359 (2023)
46. Yang, X., Zhang, K., Chen, H., Petzold, L., Wang, W.Y., Cheng, W.: Zero-shot detection of machine-generated codes. arXiv preprint arXiv:2310.05103 (2023)



47. Yang, Z., Dai, Z., Yang, Y., Carbonell, J., Salakhutdinov, R.R., Le, Q.V.: Xlnet: Generalized autoregressive pretraining for language understanding. *Advances in neural information processing systems* **32** (2019)
48. Zellers, R., Holtzman, A., Rashkin, H., Bisk, Y., Farhadi, A., Roesner, F., Choi, Y.: Defending against neural fake news. *Advances in neural information processing systems* **32** (2019)