

# Verifiable LLM-Generated Text Detection via Projected Semantic-Structural Distributions

Anonymous ACL submission

## Abstract

The widespread deployment of large language models (LLMs) makes detecting LLM-Generated text a critical security task. Existing methods, primarily relying on output probabilities from proxy models or single semantic features, suffer from distribution misalignment and limited interpretability. We observe that machine-generated text exhibits a directionally consistent systematic translation relative to human-written text within the joint semantic-structural space. Accordingly, we propose ProSSD, a statistical framework utilizing supervised subspace learning to extract compact features and construct conditional semantic distributions based on syntactic structures. By employing a likelihood ratio test, we derive a modified Mahalanobis distance, weighted by the Wasserstein distance, as the discriminative metric. Experiments demonstrate ProSSD’s superior robustness and computational efficiency across cross-domain, cross-model, and adversarial scenarios. Furthermore, we reveal the phenomena of systematic semantic translation and semantic collapse in machine-generated text, offering interpretable statistical insights into LLM generation behaviors.

## 1 Introduction

The rapid advancement of large language models (LLMs) has significantly enhanced the efficiency of various text processing tasks (Demszky et al., 2023; Doshi and Hauser, 2024). However, this progress introduces severe challenges, notably the mass generation of fake news (Ahmed et al., 2021; Hu et al., 2025), academic fabrication (Koike et al., 2024), copyright infringement (Liu et al., 2024), and the contamination of web corpora. These issues not only precipitate widespread trust crises but also threaten the integrity of information ecosystems and human creativity (Lee et al., 2024). Consequently, the detection and governance of machine-generated text (MGT) have become urgent prior-

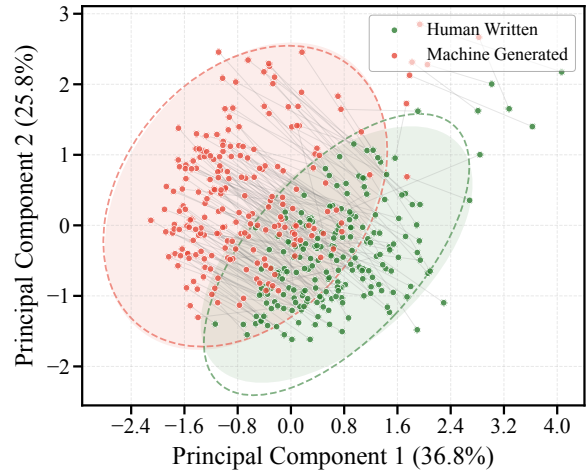


Figure 1: Average semantic positions of HWT and MGT in 2D space. MGT exhibits a systematic directional shift relative to HWT across different syntactic structures.

ities. There is a pressing need for efficient, accurate, and interpretable detection methods to provide robust technical (Wu et al., 2025a) support for forensic analysis, academic integrity, and content moderation.

To address these challenges, extensive research has been conducted, as detailed in Section 2. As a prevailing paradigm, statistical zero-shot detection (Crothers et al., 2023) distinguishes texts by exploiting the tendency of LLMs to select high-probability tokens, utilizing metrics such as perplexity (Solaiman et al., 2019) and logits Curvature (Mitchell et al., 2023). However, these methods face significant practical limitations. First, distributional discrepancies in text generation are not statistically significant across all contexts (Jiang et al., 2025). Second, the closed-source nature of commercial models (e.g., GPT (OpenAI, 2025), Gemini (Google DeepMind, 2025)) forces reliance on open-source proxy models for distribution approximation (Zhou et al., 2025). This inevitably introduces distribution misalignment and incurs

high training and inference costs, hindering large-scale real-time deployment. Furthermore, these methods lack intrinsic interpretability. Relying on model output probabilities, they fail to offer transparent computational processes or a step-by-step verifiable chain of evidence.

A more fundamental limitation is that MGT achieves semantic fluency comparable to human-written text (HWT), logit-based metrics relying on token-level confidence fail to define clear decision boundaries using such single semantic features (Tang et al., 2025). Recent studies further suggest that despite simulating human-like surface semantics, LLMs exhibit relatively constrained syntax-contextualized semantic expressions (Durward and Thomson, 2024). Inspired by this, we extend our perspective from a single semantic dimension to the joint semantic-structural distribution, aiming to quantify the intrinsic differences between HWT and MGT via conditional semantic statistics in a low-dimensional projected subspace.

This shift is driven by a key geometric regularity observed in the projected space. As shown in Figure 1, we calculate the semantic centroids for various local syntactic configurations (e.g., "noun-verb" pairs). The gray mapping lines connecting the centroids of HWT and MGT under identical syntactic structures reveal a striking phenomenon: while distinct syntactic configurations scatter across the feature space, **the semantic centroids of MGT consistently exhibit a Systematic Translation relative to those of HWT**. This implies a systematic bias across syntactic structures during the LLM decoding process, not a local perturbation specific to certain parts of speech, but a global characteristic inherent to the generation mechanism. For a more detailed analysis of systematic semantic translation and the phenomenon of semantic collapse, please refer to Appendix B.

Building upon these geometric insights, we propose the **Projected Semantic-Structural Distributions (ProSSD)** framework. Diverging from methods dependent on internal model states or raw output probabilities, ProSSD establishes a transparent, step-by-step verifiable detection paradigm. The framework proceeds in three distinct steps: first, utilizing supervised subspace learning to extract dense, low-dimensional aggregated semantic features. Second, modeling semantic distributions conditioned on syntactic structures to quantify the intrinsic divergence between HWT and MGT. And finally, deriving the

modified Mahalanobis distance via the likelihood ratio test as the core metric, dynamically weighted by the Wasserstein distance to optimize discrimination. Our main contributions are summarized as follows:

(1) **We propose ProSSD, a novel statistical detection paradigm.** Utilizing supervised subspace projection and joint semantic-structural distribution modeling, we construct a discriminant statistic based on the Wasserstein distance-weighted likelihood ratio test. This design suppresses high-dimensional noise, transforming subtle artifacts into distinguishable statistical features.

(2) **We achieve SOTA detection performance with high efficiency.** On the DetectRL benchmark covering advanced LLMs such as GPT-5.1, ProSSD excels in cross-model, cross-domain, and adversarial settings. Our approach yields robust results even with minimal samples and reduces inference costs by orders of magnitude, demonstrating significant deployment value.

(3) **We establish a step-by-step verifiable interpretability framework.** Supported by rigorous derivations, we provide theoretical guarantees and enhance transparency via statistical evidence. Furthermore, we uncover systematic semantic transition and semantic collapse in MGT under syntactic constraints, offering fresh insights into generative mechanisms.

## 2 Related work

Current machine-generated text detection methods are primarily categorized into supervised classifiers, statistical zero-shot detection, and auxiliary retrieval or watermarking techniques. As retrieval and watermarking approaches rely on external reference corpora or active injection rather than intrinsic distributional features, we discuss them in Appendix A.

**Supervised training-based methods** formulate detection as a binary classification task, heavily relying on labeled data (Zellers et al., 2019; Fagni et al., 2020). Early research predominantly utilized shallow linguistic features (e.g., TF-IDF) combined with traditional classifiers (Solaiman et al., 2019), subsequently transitioning to fine-tuning pre-trained models to enhance performance (Solaiman et al., 2019; Ippolito et al., 2020). To improve robustness against domain shifts and adversarial attacks, recent works introduce advanced optimization objectives: RADAR (Hu et al., 2023)

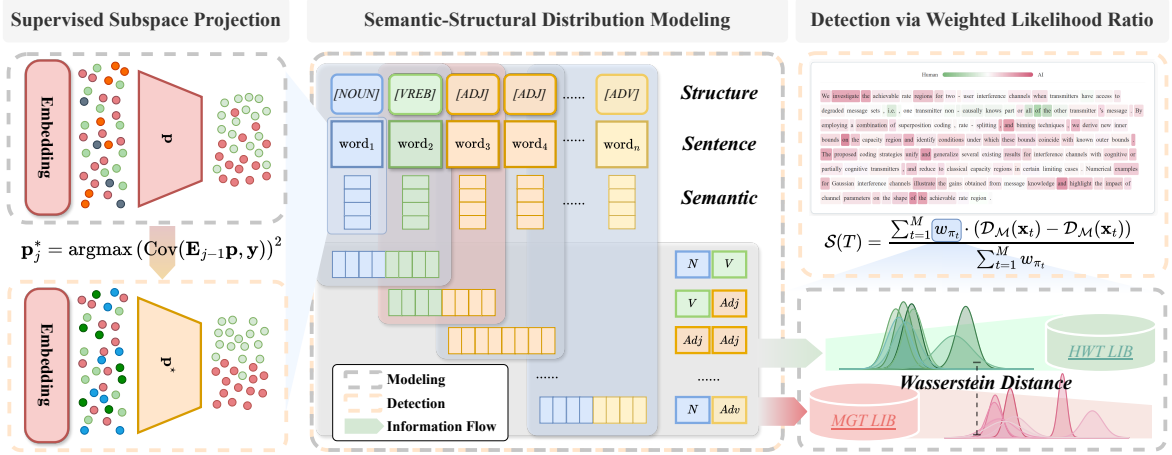


Figure 2: **The architecture of ProSSD.** It extracts features via Supervised Subspace Projection (Left), models Semantic-Structural Distributions (Middle), and calculates the final score through Detection via Weighted Likelihood Ratio (Right).

leverages proximal policy optimization (PPO) (Schulman et al., 2017) to strengthen defense capabilities; ImBD (Chen et al., 2025a) employs direct preference optimization (DPO) (Rafailov et al., 2023) to align stylistic discrepancies between humans and machines; and DetectAnyLLM (Fu et al., 2025) proposes direct difference learning (DDL) to maximize distributional distances. Despite these advancements, supervised methods face two inherent limitations: severe domain dependency, which leads to significant performance degradation across models or topics (Sarvazyan et al., 2023), and high maintenance costs, as the rapid iteration of LLMs necessitates frequent classifier retraining that consumes substantial computational and data resources.

**Statistical Detection methods** aim to exploit distributional discrepancies between MGT and HWT using specific metrics. In terms of scalar metrics, early work utilized basic statistics such as entropy and log-probability (Gehrmann et al., 2019; Solaiman et al., 2019), later expanding to more comprehensive measures like log-lank ratio (LRR) (Su et al., 2023), n-gram distributions (Yang et al., 2024), and intrinsic text dimension (Tulchinskii et al., 2023). Perturbation-based methods leverage probability curvature for discrimination; DetectGPT (Mitchell et al., 2023) hypothesizes that MGT tends to occupy negative curvature regions of the log-probability surface. To address high sampling costs, Fast-DetectGPT (Bao et al., 2024) employs conditional probability curvature for efficient approximation, while AdaDetectGPT (Zhou et al., 2025) introduces an adaptive witness function to

provide theoretical guarantees. Regarding high-order features and representations, recent research explores deeper patterns: Binoculars (Hans et al., 2024) measures prediction surprise via perplexity ratios; ReprGuard (Chen et al., 2025b) analyzes hidden representations to capture activation differences; and GECScore (Wu et al., 2025b) performs detection from a syntactic perspective using grammatical error correction distance.

### 3 Methodology

#### 3.1 Problem Definition

While LLMs generate text that is semantically indistinguishable from human writing on the surface, the generation process remains fundamentally constrained by the probabilistic decoding algorithms (Fröhling and Zubiaga, 2021), leading to statistical discrepancies in the underlying joint syntactic-semantic distribution. Drawing on the observations in Figure 1, we posit a core hypothesis: for a given local syntactic structure  $\pi \in \Pi$ , the conditional distributions of human-written text and machine-generated text within the semantic embedding space  $\mathbf{e} \in \mathbb{R}^D$  diverge significantly. Formally, let  $y \in \{H, M\}$  denote the class labels for HWT and MGT respectively, we hypothesize:

$$P(\mathbf{e}|\pi, y = H) \neq P(\mathbf{e}|\pi, y = M). \quad (1)$$

The objective of this paper is to achieve zero-shot detection by modeling and quantifying the divergence in this joint semantic-structural Distribution.

### 3.2 Supervised Subspace Projection

Directly estimating the aforementioned joint semantic-structural distribution using raw embeddings  $\mathbf{e}$  from pre-trained models (e.g. RoBERTa (Liu et al., 2019)) presents significant challenges. First, the semantic masking effect obscures the subtle signals distinguishing HWT from MGT, as raw embeddings are overwhelmingly dominated by variance from general semantic topics irrelevant to the generation source (Nagata et al., 2023). Furthermore, high-dimensional sparsity poses a computational barrier; since the embedding dimension  $D$  far exceeds the local sample size, direct parametric estimation of the covariance matrix becomes numerically highly unstable (Ledoit and Wolf, 2004).

To address these challenges, we propose the supervised subspace learning algorithm. The goal is to learn a projection matrix  $\mathbf{P} \in \mathbb{R}^{D \times k}$  ( $k \ll D$ ) to map the raw vectors into low-dimensional compact semantic features  $\mathbf{v} = \mathbf{P}^T \mathbf{e}$ . Unlike principal component analysis (PCA), which aims to maximize total variance, our objective is to maximize the statistical correlation between features and class labels  $\mathbf{y}$ , thereby suppressing noise and preserving discriminative information. We model this process as a recursive optimization problem. Let  $\mathbf{E}_0$  be the centered raw embedding matrix, and  $\mathbf{E}_{j-1}$  be the residual feature matrix at step  $j$ . We solve for the  $j$ -th projection basis vector  $\mathbf{p}_j$  via the following objective function:

$$\mathbf{p}_j^* = \operatorname{argmax}_{\mathbf{p}: \|\mathbf{p}\|_2=1} (\operatorname{Cov}(\mathbf{E}_{j-1} \mathbf{p}, \mathbf{y}))^2. \quad (2)$$

After obtaining the optimal direction  $\mathbf{p}_j^*$ , we perform feature deflation:  $\mathbf{E}_j = \mathbf{E}_{j-1} - \mathbf{s}_j (\mathbf{p}_j^*)^T$ , where  $\mathbf{s}_j = \mathbf{E}_{j-1} \mathbf{p}_j^*$  is the projection score of the current dimension. This step ensures the informational orthogonality of projection directions across different dimensions; detailed derivation is provided in Appendix C.1. Finally, for the  $i$ -th word of the input text, we obtain its compact semantic vector  $\mathbf{v}_i \in \mathbb{R}^k$ .

### 3.3 Semantic-Structural Distribution Modeling

After obtaining the low-dimensional compact semantic features  $\mathbf{v}$ , we extend our perspective from single semantic point estimation to Joint semantic-structural distribution modeling.

**Local Feature Construction.** To capture semantic transition patterns within local syntactic environ-

ments, we define the meta-semantics vector  $\mathbf{x}_t$  and meta-structure  $\pi_t$ . At time step  $t$ , we concatenate projected features from adjacent positions and pair them with corresponding pos tags to form observation pairs:

$$\mathbf{x}_t = \begin{bmatrix} \mathbf{v}_t \\ \mathbf{v}_{t+1} \end{bmatrix}, \quad \pi_t = (\text{pos}_t, \text{pos}_{t+1}). \quad (3)$$

**Distribution Modeling.** Consequently, any text sequence can be parsed into a set of local observations  $\mathcal{D} = \{(\mathbf{x}_t, \pi_t)\}_{t=1}^{T-1}$ . Based on the central limit theorem and empirical observations detailed in Appendix C.2, we assume that conditioned on the meta-structure  $\pi$ , the vector  $\mathbf{x}$  follows a multivariate Gaussian distribution. That is, for class  $y \in \{H, M\}$ :

$$P(\mathbf{x}|\pi, y) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_y^\pi, \boldsymbol{\Sigma}_y^\pi). \quad (4)$$

Here,  $\boldsymbol{\mu}$  characterizes the average semantic information under this syntactic structure, while the block covariance matrix  $\boldsymbol{\Sigma}$  encodes the semantic dependencies between adjacent tokens.

**Parameter Estimation.** We estimate the distribution parameters for HWT and MGT on a reference corpus using maximum likelihood estimation (MLE) for each appearing meta-structure  $\pi \in \Pi$ :

$$\begin{aligned} \hat{\boldsymbol{\mu}}_y^\pi &= \frac{1}{N_{\pi,y}} \sum_{j=1}^{N_{\pi,y}} \mathbf{x}_j, \\ \hat{\boldsymbol{\Sigma}}_y^\pi &= \frac{1}{N_{\pi,y} - 1} \sum_{j=1}^{N_{\pi,y}} (\mathbf{x}_j - \hat{\boldsymbol{\mu}}_y^\pi)(\mathbf{x}_j - \hat{\boldsymbol{\mu}}_y^\pi)^T. \end{aligned} \quad (5)$$

Thus, the detection model is parameterized as a set of structured distributions  $\mathcal{M} = \{(\mathcal{N}_H^\pi, \mathcal{N}_M^\pi) \mid \pi \in \Pi\}$ .

**Adaptive Weighting via Wasserstein Distance.** Different syntactic structures carry significantly different amounts of discriminative information. To quantify this structural non-uniformity, we introduce the Wasserstein distance (Villani et al., 2008) as the discriminative weight  $w_\pi$  for the meta-structure  $\pi$ . For two Gaussian distributions  $\mathcal{N}_H^\pi$  and  $\mathcal{N}_M^\pi$ , the Wasserstein distance has an analytical solution. As a core metric in optimal transport theory, the Wasserstein distance defines a strict geometric metric in the probability distribution space (Ge et al., 2021; Just et al., 2023). It quantifies intrinsic shifts between continuous distributions more robustly than traditional divergence metrics.

The specific form is as follows:

$$w_\pi = \left( \|\boldsymbol{\mu}_H^\pi - \boldsymbol{\mu}_M^\pi\|_2^2 + \text{Tr} \left( \boldsymbol{\Sigma}_H^\pi + \boldsymbol{\Sigma}_M^\pi - 2(\boldsymbol{\Sigma}_H^{\pi \frac{1}{2}} \boldsymbol{\Sigma}_M^\pi \boldsymbol{\Sigma}_H^{\pi \frac{1}{2}})^{\frac{1}{2}} \right) \right)^{1/2}. \quad (6)$$

The derivation of Equation 6 is provided in the Appendix C.3. We select the Wasserstein distance not only because it captures changes in both the first moment and second moment, but also because it theoretically bounds the performance of the discriminative function.

**Theorem 1. (Wasserstein Discrimination Bound)**

For any discriminative function  $f$  satisfying the  $K$ -Lipschitz continuity condition, the difference in expected scores on the HWT and MGT distributions is strictly bounded by the Wasserstein distance between these two distributions:

$$|\mathbb{E}_{P_H}[f(\mathbf{x})] - \mathbb{E}_{P_M}[f(\mathbf{x})]| \leq K \cdot W_2(P_H, P_M). \quad (7)$$

*Proof detailed in Appendix C.4*

This theorem indicates that the value of  $w_\pi$  directly reflects the theoretical possibility of distinguishing the two types of text under that structure. A larger  $W_2$  distance implies a clearer decision boundary.

**3.4 Detection via Weighted Likelihood Ratio**

Grounded in the Neyman-Pearson lemma (Larsen and Marx, 2005), we formulate the detection problem as a statistical test based on the log-likelihood Ratio (LLR) for each local feature, followed by a weighted aggregation strategy to derive the final decision.

**Local Statistic Derivation.** Given a text  $T$  to be detected, we obtain the sequence  $\{(\mathbf{x}_t, \pi_t)\}_{t=1}^M$  after projection, slicing, and feature construction. For each local slice  $t$ , we decide between the null hypothesis  $H_0 : \mathbf{x}_t \sim P_H^{\pi_t}$  and the alternative hypothesis  $H_1 : \mathbf{x}_t \sim P_A^{\pi_t}$ . The optimal test statistic is the log-likelihood ratio:

$$s_t = \ln \frac{P(\mathbf{x}_t | \pi_t, \theta_M)}{P(\mathbf{x}_t | \pi_t, \theta_H)}. \quad (8)$$

Based on the multivariate Gaussian assumption mentioned above, substituting the density function reveals that the local anomaly score  $s_t$  is equivalent to the difference in the modified Mahalanobis distance:

$$s_t = \frac{1}{2} [\mathcal{D}_M(\mathbf{x}_t; \boldsymbol{\mu}_H^{\pi_t}, \boldsymbol{\Sigma}_H^{\pi_t}) - \mathcal{D}_M(\mathbf{x}_t; \boldsymbol{\mu}_M^{\pi_t}, \boldsymbol{\Sigma}_M^{\pi_t})], \quad (9)$$

where  $\mathcal{D}_M$  includes a penalty term for the covariance determinant:

$$\mathcal{D}_M(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) + \ln |\boldsymbol{\Sigma}|. \quad (10)$$

The detailed derivation of this metric is provided in Appendix C.5. This equation indicates that if the sample deviates from the human distribution significantly more than from the machine distribution,  $s_t$  becomes positive, indicating MGT characteristics. **Global Aggregation.** To obtain a document-level decision score and ensure interpretability, we map window-level scores  $s_t$  back to the word level. We define the anomaly score of the  $i$ -th word as the mean of adjacent windows:  $o_i = \frac{1}{2}(s_{i-1} + s_i)$ , and assuming boundary terms are 0. Finally, we use Wasserstein distance weights to aggregate all scores, obtaining the final detection statistic  $\mathcal{S}(T)$  for text  $T$ :

$$\mathcal{S}(T) = \frac{\sum_{t=1}^M w_{\pi_t} \cdot s_t}{\sum_{t=1}^M w_{\pi_t}}. \quad (11)$$

If  $\mathcal{S}(T)$  exceeds the preset threshold  $\tau$ , the text is classified as MGT.

**4 Experiments**

**4.1 Experimental Settings**

**Benchmark Datasets.** To rigorously evaluate model performance in realistic and challenging detection scenarios, we adopt DetectRL (Wu et al., 2024) as the foundational evaluation benchmark. This benchmark encompasses four typical high risk misuse domains: academic writing (ArXiv<sup>1</sup>), news summarization (XSum, (Narayan et al., 2018)), creative writing (Writing Prompts, (Fan et al., 2018)), and social media reviews (Yelp Review (Zhang et al., 2015)). Given that the generation models in the original benchmark (e.g., GPT-3.5 (OpenAI, 2023), PaLM-2 (Anil et al., 2023)) fail to capture the frontier capabilities of current LLMs in semantic alignment and reasoning, we upgrade the generation sources to verify detector effectiveness against SOTA capabilities. Strictly following the DetectRL data construction protocol, we employ GPT-5.1 (OpenAI, 2025), Claude-Sonnet-4 (Anthropic, 2025), and Gemini-3-Flash (Google DeepMind, 2025), Grok-4.1 (xAI, 2025). Using this upgraded dataset, we execute Model Generalization, Domain Generalization, and OOD detection

<sup>1</sup><https://www.kaggle.com/datasets/spsayakpaul/axiv-paper-abstracts/data>

Table 1: Performance comparison across different target models and text domains. Results are reported in terms of AUROC (%) and F1 (%). The best and second-best results in each column are highlighted in **bold** and underlined.

Method	Models								Domains							
	GPT-5.1		Claude-S-4		Gemini-3-F		Grok-4.1		Arxiv		Writing		XSum		Review	
	AUR.	F1	AUR.	F1	AUR.	F1	AUR.	F1	AUR.	F1	AUR.	F1	AUR.	F1	AUR.	F1
LRR	68.02	58.08	79.62	71.13	63.30	59.10	62.18	49.50	73.90	67.71	48.20	14.90	70.14	63.49	84.01	75.11
NPR	78.37	68.41	81.40	72.97	80.13	70.64	77.02	68.92	92.45	85.08	67.51	64.77	72.78	61.37	85.03	76.18
RoBERTa	53.92	49.60	57.95	58.16	61.29	53.78	55.85	59.55	61.03	55.62	56.96	58.63	55.76	48.76	56.99	47.82
DetectGPT	59.54	47.60	56.77	38.53	50.84	37.05	57.73	43.34	58.89	61.42	77.39	71.00	17.83	66.69	73.87	67.54
Binoculars	75.55	66.59	92.52	84.66	86.27	77.37	71.40	61.43	75.84	65.29	77.51	69.19	74.90	65.37	97.35	91.64
Fast-DetectGPT	67.07	58.59	78.78	72.73	76.13	68.44	52.91	42.29	74.81	64.93	58.88	47.29	57.65	50.42	82.92	75.56
ImBD	86.95	80.23	94.37	86.67	96.83	90.90	88.48	82.39	93.63	85.77	92.51	84.34	99.14	96.06	97.74	92.55
AdaDetectGPT	61.64	56.87	71.97	69.69	71.98	64.50	48.52	17.02	66.41	59.75	55.84	45.47	53.26	55.23	78.28	71.22
DetectAnyLLM	78.73	71.49	91.58	83.38	96.96	91.63	89.45	82.92	96.25	89.64	81.21	73.26	95.51	88.37	89.82	81.08
RepreGuard	97.50	92.77	99.59	97.38	99.24	96.99	98.07	92.98	99.30	96.98	94.42	87.04	99.76	<b>99.10</b>	99.92	98.80
ProSSD	<b>99.72</b>	<b>98.00</b>	<b>99.97</b>	<b>99.30</b>	<b>99.89</b>	<b>99.50</b>	<b>99.95</b>	<b>98.89</b>	<b>99.97</b>	<b>99.65</b>	<b>99.90</b>	<b>98.60</b>	<b>99.80</b>	<u>98.95</u>	<b>99.99</b>	<b>99.55</b>

tasks. For basic adversarial robustness testing, we retain the original data settings to ensure comparability with prior work. We adopt AUROC and F1 Score as the core evaluation metrics. The data construction protocols and descriptive statistics for the new evaluation sets are detailed in Appendix D.1 and Appendix D.2, while supplementary detection results on additional independent benchmarks are presented in Appendix E.5.

**Baselines.** To establish a comprehensive benchmark, we compare ProSSD against a diverse set of state-of-the-art methods, including LRR (Su et al., 2023), NPR (Su et al., 2023), DetectGPT (Mitchell et al., 2023), Fast-DetectGPT (Bao et al., 2024), AdaDetectGPT (Zhou et al., 2025), Binoculars (Hans et al., 2024), RepreGuard (Chen et al., 2025b), RoBERTa-Base (Solaiman et al., 2019), ImBD (Chen et al., 2025a), and DetectAnyLLM (Fu et al., 2025). Specific implementation details and hyperparameter settings for all baselines are detailed in Appendix D.3.

## 4.2 Main Comparative Results

**Overall Detection Performance.** Table 1 presents the main detection results under Multi-LLM and Multi-Domain settings. Overall, ProSSD demonstrates consistent superiority across all test scenarios, achieving SOTA performance. Notably, in detecting text generated by GPT-5.1 and Grok-4.1, ProSSD achieves F1 scores of 98.00% and 98.89%, respectively, representing improvements of 5.64% and 6.36% over the runner-up method, RepreGuard. In contrast, traditional supervised methods (e.g., RoBERTa) and zero-shot statistical methods (e.g., Binoculars) exhibit significant per-

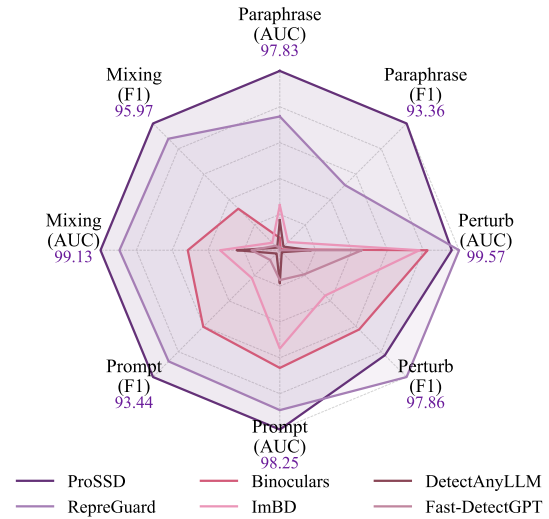


Figure 3: Performance comparison of ProSSD on Paraphrase, Perturbation, and Mixing attacks, as well as the reference Direct Prompt task, in terms of AUROC and F1 scores.

formance volatility when facing frontier LLMs. In the multi-domain evaluation, apart from slightly lower performance on the XSum dataset compared to RepreGuard, ProSSD achieves the best performance across all other text styles (Academic, Creative, Reviews), demonstrating its robust generalizability across diverse semantic styles.

**OOD Performance.** Table 2 reports OOD detection performance, aiming to measure generalization ability when there are significant discrepancies between the generation sources of the training and test sets. The experiment involves using earlier models (e.g., GPT-3.5) for parameter estimation and testing on more advanced models (e.g., GPT-5.1). Results

Table 2: Impact of training data sources on detection generalization. The table compares the performance when training on data generated by GPT-3.5 versus Llama-2-70b and testing on unseen target LLMs. Results are reported as AUROC (%) and F1 (%). The best and second-best scores are highlighted in **bold** and underlined.

Method	Train on GPT-3.5								Train on Llama-2-70b							
	GPT-5.1		Claude-S-4		Gemini-3-F		Grok-4.1		GPT-5.1		Claude-S-4		Gemini-3-F		Grok-4.1	
	AUR.	F1	AUR.	F1	AUR.	F1	AUR.	F1	AUR.	F1	AUR.	F1	AUR.	F1	AUR.	F1
LRR	68.02	58.08	79.62	71.13	63.30	59.10	62.18	49.50	68.02	58.08	79.62	71.13	63.30	59.10	62.18	49.50
Fast-DetectGPT	67.07	58.59	78.78	72.73	76.13	68.44	52.91	42.29	67.07	58.59	78.78	72.73	76.13	68.44	52.91	42.29
ImBD	<u>93.28</u>	83.90	<u>97.29</u>	<u>91.07</u>	<u>97.60</u>	<u>92.19</u>	90.71	83.14	<u>88.59</u>	81.13	92.61	84.38	<u>93.50</u>	86.51	84.51	78.70
AdaDetectGPT	61.32	55.75	74.32	68.67	68.23	55.00	46.90	23.99	58.02	52.83	69.87	65.60	64.35	56.10	43.85	3.12
DetectAnyLLM	89.73	<u>84.27</u>	94.89	89.61	94.35	89.66	<u>92.49</u>	<u>86.75</u>	87.12	<u>84.92</u>	<u>93.68</u>	<u>88.54</u>	91.78	<u>87.38</u>	<u>88.51</u>	<u>86.02</u>
RepreGuard	78.97	73.68	94.49	87.79	89.44	82.38	81.67	74.29	72.83	62.87	88.72	81.90	81.48	74.26	74.85	69.21
ProSSD	<b>94.94</b>	<b>87.42</b>	<b>98.64</b>	<b>93.58</b>	<b>99.31</b>	<b>96.43</b>	<b>97.06</b>	<b>89.63</b>	<b>93.09</b>	<b>86.14</b>	<b>97.95</b>	<b>91.71</b>	<b>97.20</b>	<b>91.82</b>	<b>94.33</b>	<b>86.08</b>

indicate that ProSSD significantly mitigates performance degradation caused by generative distribution shifts. Specifically, When trained only on GPT-3.5 (Table 2 Left), ProSSD maintains an AUROC above 94% on all four target test sets, achieving optimal results. When trained only on Llama-2-70b (Table 2 Right), its advantage remains significant when transferring to closed-source black-box models. For instance, when detecting Grok-4.1, ProSSD achieves an AUROC of 94.33%, significantly surpassing DetectAnyLLM (88.51%) and ImBD (84.51%) under the same setting. More detailed OOD performance results and analyses are available in the Appendix E.1.

**Robustness Against Adversarial Attacks.** Figure 3 displays radar charts illustrating model robustness under three attack scenarios: paraphrase, perturbation, and data mixing. In terms of the performance envelope, ProSSD covers the largest area across all dimensions, indicating superior performance under various attack modes without obvious shortcomings. Conversely, other methods like DetectAnyLLM and Binoculars exhibit significant “star-shaped” contraction, showing drastic performance drops particularly along the most challenging paraphrase axis. This demonstrates that our method maintains exceptional performance and robust reliability when deployed in realistic and complex adversarial environments. Detailed AUROC and F1 results and analyses are provided in the Appendix E.2.

### 4.3 In-depth Analysis

#### 4.3.1 Parameter and Data Sensitivity

**Impact of N-gram Context ( $k$ ).** Figure 4a illustrates the effect of the N-gram context window size  $k$  on detection performance. We observe that

Table 3: Efficiency comparison regarding inference time and GPU memory usage. Lower values are better. The AUROC scores are provided for reference. Best results are highlighted in **bold** and second-best are underlined.

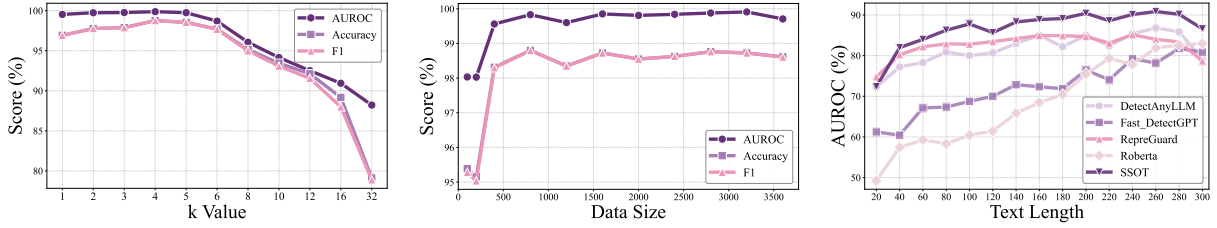
Method	Time (s)	GPU Mem (GB)	AUR.
DetectGPT	11.83	10.87	56.77
Binoculars	<u>0.19</u>	26.24	92.52
Fast-DetectGPT	0.37	33.40	78.78
AdaDetectGPT	3.06	37.98	71.97
DetectAnyLLM	<b>0.08</b>	<u>5.32</u>	91.58
RepreGuard	0.36	30.84	<u>99.59</u>
ProSSD	<b>0.08</b>	<b>1.35</b>	<b>99.97</b>

performance peaks at  $k = 4$ , which achieves an optimal balance between capturing local contextual dependencies and mitigating data sparsity inherent in higher dimensions. Notably, as  $k$  increases further, ProSSD maintains stable AUROC and F1 scores, demonstrating the robustness of the framework to hyperparameter variations.

**Impact of Data Size.** Figure 4b evaluates the impact of the reference set size on model efficacy. The performance curve shows a steep upward trend in the early stage ( $N \leq 400$ ). It then rapidly enters a convergence plateau, where marginal gains from increasing data volume diminish. Unlike deep learning baselines that typically require massive corpora to fit decision boundaries, ProSSD demonstrates extremely high data efficiency. This makes it highly suitable for few-shot application scenarios.

#### 4.3.2 Scalability and Efficiency Analysis

**Impact of Text Length.** Text length determines the information density available for statistical inference. Figure 4c compares the response patterns of different detectors across varying text lengths. Supervised models, represented by RoBERTa, exhibit a clear linear dependency within the 20-300



(a) Impact of N-gram Context Length  $k$  on AUROC, Accuracy and F1 score. (b) Impact of Data Size on AUROC, Accuracy and F1 score. (c) Impact of Text Length on AUROC score.

Figure 4: Impact analysis of key hyperparameters and data variations on detection performance. (a) The AUROC score peaks at  $k = 4$ , identifying the optimal N-gram context length. (b) The model shows high data efficiency, reaching performance saturation with minimal training data size. (c) ProSSD remains robust across varying text lengths, consistently outperforming other detectors from short to long sequences.

Table 4: Ablation studies (AUC) on XSum and Review. See Table 9 for full results. The statistical significance of the performance drop compared to ProSSD is measured by a t-test: \*  $p < 0.05$ , †  $p < 0.01$ , ‡  $p < 0.001$ .

Method	XSum	Review
<b>ProSSD (RoBERTa + SSP)</b>	<b>99.21±0.19</b>	<b>99.67±0.04</b>
<i>Robustness of Semantic Representations</i>		
Qwen-2.5-0.6B-Embed	95.28±0.27‡	99.50±0.06†
Random Projection	93.45±1.04‡	91.49±2.58†
<i>Impact of Subspace Projection</i>		
w/o SSP (PCA)	94.75±0.39‡	97.54±0.08‡
w/o SSP (Rand)	93.39±0.73‡	96.46±0.51‡
w/o SSP (No-Proj)	93.97±0.35‡	96.74±0.37‡
<i>Ablation on Detection Strategies</i>		
w/o Wasserstein (Unif.)	96.48±0.94†	97.57±0.47‡
w/o Contrastive (Human)	93.86±0.34‡	88.70±0.25‡

token range. This indicates they require longer contexts to accumulate sufficient semantic features. In sharp contrast, ProSSD demonstrates a superior information utilization rate, even in the extremely short text range (20-80 tokens), ProSSD achieves a rapid performance climb and quickly converges to a high-performance zone.

**Computational Efficiency.** Computational cost is critical for practical deployment. ProSSD has achieved an inference speed of 0.08s, comparable to lightweight reward models. Notably, it requires only 1.35 GB of VRAM, a 74.62% reduction compared to the runner-up-demonstrating its capability to maintain high-precision detection with minimal resource overhead.

### 4.3.3 Ablation Study

Table 4 verifies the necessity of ProSSD’s core components. First, replacing SSP with unsupervised PCA or random projection causes significant AUC drops (e.g., ~4.5% on XSum), confirming

the critical role of label-aware feature extraction in isolating discriminative signals from noise. Second, removing Wasserstein Distance weights or degenerating to one-sided density estimation consistently weakens performance, establishing the need for adaptive weighting and dual-distribution contrast. Finally, the framework’s robustness across different embedding models suggests it captures intrinsic distributional laws independent of encoder biases. Detailed ablation configurations and results are provided in the Appendix E.3.

### 4.3.4 Interpretability Analysis

Our in-depth analysis reveals that discriminative cues exhibit strong domain dependency: formal corpora rely heavily on logical connectives reflecting deep syntactic coherence, whereas subjective texts hinge on pronouns and adjectives. ProSSD establishes a significantly wider safety margin in the numerical space, compressing the distributional overlap between human and machine text to a negligible level. For comprehensive analysis and visualizations, please refer to the Appendix E.4.

## 5 Conclusion

This paper presents ProSSD, a zero-shot detection framework based on Supervised Subspace Learning and Joint Semantic-Structural Distribution modeling. Leveraging the Wasserstein Distance, ProSSD effectively mitigates high-dimensional noise and captures the statistical characteristics of MGT. Experiments show that ProSSD achieves SOTA performance across various benchmarks while significantly reducing computational costs. Unlike previous methods, our approach formulates detection as an interpretable statistical test, offering a lightweight and transparent solution for LLM governance.

## 570 Limitations

571 Despite the robust performance of ProSSD, sev-  
572 eral limitations remain. First, detection bound-  
573 aries in complex scenarios are limited. While  
574 competitive in handling heavily human-polished  
575 or human-machine hybrid text, our method has not  
576 yet significantly surpassed current SOTA baselines  
577 or achieved near-perfect detection performance.  
578 Second, the theoretical interpretation of internal  
579 mechanisms is insufficient. Our observations of  
580 semantic-structural statistical deviations remain  
581 empirical, lacking a clear theoretical causal link to  
582 internal Transformer mechanisms, such as attention  
583 patterns. Third, the detection paradigm requires  
584 extension. Currently, our framework focuses on  
585 binary classification. As open-source models pro-  
586 liferate, the field is evolving towards fine-grained  
587 source attribution (Park et al., 2025). Our method  
588 has not yet implemented specific attribution for  
589 distinct generative sources.

## 590 References

591 Alim Al Ayub Ahmed, Ayman Aljabouh, Praveen Ku-  
592 mar Donepudi, and Myung Suh Choi. 2021. [Detect-](#)  
593 [ing fake news using machine learning : A systematic](#)  
594 [literature review](#). *CoRR*, abs/2102.04458.

595 Theodore Wilbur Anderson, Theodore Wilbur Ander-  
596 son, Theodore Wilbur Anderson, Theodore Wilbur  
597 Anderson, and Etats-Unis Mathématicien. 1958. *An*  
598 *introduction to multivariate statistical analysis*, vol-  
599 *ume 2*. Wiley New York.

600 Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin John-  
601 son, Dmitry Lepikhin, Alexandre Passos, Siamak  
602 Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng  
603 Chen, Eric Chu, Jonathan H. Clark, Laurent El  
604 Shafey, Yanping Huang, Kathy Meier-Hellstern, Gau-  
605 rav Mishra, Erica Moreira, Mark Omernick, Kevin  
606 Robinson, and 34 others. 2023. [Palm 2 technical](#)  
607 [report](#). *CoRR*, abs/2305.10403.

608 Anthropic. 2025. [Claude 4](#). Anthropic Blog.

609 Guangsheng Bao, Yanbin Zhao, Zhiyang Teng, Linyi  
610 Yang, and Yue Zhang. 2024. [Fast-detectgpt: Effi-](#)  
611 [cient zero-shot detection of machine-generated text](#)  
612 [via conditional probability curvature](#). In *The Twelfth*  
613 *International Conference on Learning Representa-*  
614 *tions, ICLR 2024, Vienna, Austria, May 7-11, 2024*.  
615 OpenReview.net.

616 Yinpeng Cai, Lexin Li, and Linjun Zhang. 2025. [A](#)  
617 [statistical hypothesis testing framework for data mis-](#)  
618 [appropriation detection in large language models](#).  
619 *CoRR*, abs/2501.02441.

Jiaqi Chen, Xiaoye Zhu, Tianyang Liu, Ying Chen, Xin-  
hui Chen, Yiwen Yuan, Chak Tou Leong, Zuchao  
Li, Tang Long, Lei Zhang, Chenyu Yan, Guanghao  
Mei, Jie Zhang, and Lefei Zhang. 2025a. [Imitate be-](#)  
[fore detect: Aligning machine stylistic preference for](#)  
[machine-revised text detection](#). In *AAAI-25, Spon-*  
*sored by the Association for the Advancement of Ar-*  
*tificial Intelligence, February 25 - March 4, 2025,*  
*Philadelphia, PA, USA*, pages 23559–23567. AAAI  
Press. 620 621 622 623 624 625 626 627 628 629

Xin Chen, Junchao Wu, Shu Yang, Runzhe Zhan, Zeyu  
Wu, Ziyang Luo, Di Wang, Min Yang, Lidia S. Chao,  
and Derek F. Wong. 2025b. [Repreguard: Detecting](#)  
[llm-generated text by revealing hidden representation](#)  
[patterns](#). *CoRR*, abs/2508.13152. 630 631 632 633 634

Miranda Christ, Sam Gunn, and Or Zamir. 2024. [Un-](#)  
[detectable watermarks for language models](#). In *The*  
*Thirty Seventh Annual Conference on Learning The-*  
*ory, June 30 - July 3, 2023, Edmonton, Canada*, vol-  
*ume 247 of Proceedings of Machine Learning Re-*  
*search*, pages 1125–1139. PMLR. 635 636 637 638 639 640

Harald Cramér and Herman Wold. 1936. Some theo-  
rems on distribution functions. *Journal of the London*  
*Mathematical Society*, 1(4):290–294. 641 642 643

Evan Crothers, Nathalie Japkowicz, and Herna L. Viktor.  
2023. [Machine-generated text: A comprehensive](#)  
[survey of threat models and detection methods](#). *IEEE*  
*Access*, 11:70977–71002. 644 645 646 647

Dorottya Demszky, Diyi Yang, David S Yeager, Christo-  
pher J Bryan, Margaret Clapper, Susannah Chand-  
hok, Johannes C Eichstaedt, Cameron Hecht, Jeremy  
Jamieson, Meghann Johnson, and 1 others. 2023. Us-  
ing large language models in psychology. *Nature*  
*Reviews Psychology*, 2(11):688–701. 648 649 650 651 652 653

Anil R Doshi and Oliver P Hauser. 2024. Generative  
ai enhances individual creativity but reduces the col-  
lective diversity of novel content. *Science advances*,  
10(28):ead5290. 654 655 656 657

Matthew Durward and Christopher Thomson. 2024.  
[Evaluating vocabulary usage in llms](#). In *Proceed-*  
*ings of the 19th Workshop on Innovative Use of NLP*  
*for Building Educational Applications, BEA 2024,*  
*Mexico City, Mexico, June 20, 2024*, pages 266–282.  
Association for Computational Linguistics. 658 659 660 661 662 663

Tiziano Fagni, Fabrizio Falchi, Margherita Gambini, An-  
tonio Martella, and Maurizio Tesconi. 2020. [Tweep-](#)  
[fake: about detecting deepfake tweets](#). *CoRR*,  
abs/2008.00036. 664 665 666 667

Angela Fan, Mike Lewis, and Yann Dauphin. 2018.  
[Hierarchical neural story generation](#). In *Proceedings*  
*of the 56th Annual Meeting of the Association for*  
*Computational Linguistics (Volume 1: Long Papers)*,  
pages 889–898, Melbourne, Australia. Association  
for Computational Linguistics. 668 669 670 671 672 673

William Feller. 1991. *An introduction to probability the-*  
*ory and its applications, Volume 2*, volume 2. John  
Wiley & Sons. 674 675 676

677	Leon Fröhling and Arkaitz Zubiaga. 2021. <a href="#">Feature-based detection of automated language models: tackling gpt-2, GPT-3 and grover</a> . <i>PeerJ Comput. Sci.</i> , 7:e443.	734
678		735
679		
680		
681	Jiachen Fu, Chun-Le Guo, and Chongyi Li. 2025. <a href="#">Detectanyllm: Towards generalizable and robust detection of machine-generated text across domains and models</a> . <i>CoRR</i> , abs/2509.14268.	736
682		737
683		738
684		739
685	Zheng Ge, Songtao Liu, Zeming Li, Osamu Yoshie, and Jian Sun. 2021. <a href="#">OTA: optimal transport assignment for object detection</a> . In <i>IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021</i> , pages 303–312. Computer Vision Foundation / IEEE.	740
686		741
687		
688		
689		
690		
691	Sebastian Gehrmann, Hendrik Strobelt, and Alexander M. Rush. 2019. <a href="#">GLTR: statistical detection and visualization of generated text</a> . In <i>Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28 - August 2, 2019, Volume 3: System Demonstrations</i> , pages 111–116. Association for Computational Linguistics.	742
692		743
693		744
694		745
695		746
696		747
697		
698		
699	Noah Golowich and Ankur Moitra. 2024. <a href="#">Edit distance robust watermarks via indexing pseudorandom codes</a> . In <i>Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024</i> .	748
700		749
701		750
702		751
703		752
704		753
705	Google DeepMind. 2025. <a href="#">Gemini 3 flash model card</a> . Technical Report.	754
706		755
707	Abhimanyu Hans, Avi Schwarzschild, Valeriia Cherepanova, Hamid Kazemi, Aniruddha Saha, Micah Goldblum, Jonas Geiping, and Tom Goldstein. 2024. <a href="#">Spotting llms with binoculars: Zero-shot detection of machine-generated text</a> . In <i>Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024</i> . OpenReview.net.	756
708		757
709		758
710		
711		
712		
713		
714		
715	Beizhe Hu, Qiang Sheng, Juan Cao, Yang Li, and Danding Wang. 2025. <a href="#">Llm-generated fake news induces truth decay in news ecosystem: A case study on neural news recommendation</a> . In <i>Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2025, Padua, Italy, July 13-18, 2025</i> , pages 435–445. ACM.	759
716		760
717		761
718		762
719		763
720		764
721		765
722		766
723	Xiaomeng Hu, Pin-Yu Chen, and Tsung-Yi Ho. 2023. <a href="#">RADAR: robust ai-text detection via adversarial learning</a> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	767
724		768
725		769
726		770
727		771
728		772
729	Daphne Ippolito, Daniel Duckworth, Chris Callison-Burch, and Douglas Eck. 2020. <a href="#">Automatic detection of generated text is easiest when humans are fooled</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020</i> , pages 1808–1822. Association for Computational Linguistics.	773
730		774
731		775
732		776
733		777
		778
		779
		780
		781
		782
		783
		784
		785
		786
		787
		788
		789
		790

791	Xiaoze Liu, Ting Sun, Tianyang Xu, Feijie Wu, Cunxiang Wang, Xiaoqian Wang, and Jing Gao. 2024. <a href="#">SHIELD: Evaluation and defense strategies for copy-right compliance in LLM text generation</a> . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 1640–1670, Miami, Florida, USA. Association for Computational Linguistics.	848
792		849
793		
794		
795		
796		
797		
798		
799	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. <a href="#">Roberta: A robustly optimized BERT pretraining approach</a> . <i>CoRR</i> , abs/1907.11692.	
800		
801		
802		
803		
804	Eric Mitchell, Yoonho Lee, Alexander Khazatsky, Christopher D. Manning, and Chelsea Finn. 2023. <a href="#">Detectgpt: Zero-shot machine-generated text detection using probability curvature</a> . In <i>International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA</i> , volume 202 of <i>Proceedings of Machine Learning Research</i> , pages 24950–24962. PMLR.	
805		
806		
807		
808		
809		
810		
811		
812	Ryo Nagata, Hiroya Takamura, Naoki Otani, and Yoshifumi Kawasaki. 2023. <a href="#">Variance matters: Detecting semantic differences without corpus/word alignment</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 15609–15622. Association for Computational Linguistics.	
813		
814		
815		
816		
817		
818		
819	Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. <a href="#">Don’t give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization</a> . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 1797–1807, Brussels, Belgium. Association for Computational Linguistics.	
820		
821		
822		
823		
824		
825		
826	OpenAI. 2023. <a href="#">Introducing chatgpt</a> . OpenAI Blog.	
827	OpenAI. 2025. <a href="#">Gpt-5.1</a> . OpenAI Blog.	
828	Hyeonchu Park, Byungjun Kim, and Bugeun Kim. 2025. <a href="#">DART: an AIGT detector using AMR of rephrased text</a> . In <i>Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL 2025 - Volume 2: Short Papers, Albuquerque, New Mexico, April 29 - May 4, 2025</i> , pages 710–721. Association for Computational Linguistics.	
829		
830		
831		
832		
833		
834		
835		
836		
837	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D. Manning, Stefano Ermon, and Chelsea Finn. 2023. <a href="#">Direct preference optimization: Your language model is secretly a reward model</a> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	
838		
839		
840		
841		
842		
843		
844		
845	Vinu Sankar Sadasivan, Aounon Kumar, Sriram Balasubramanian, Wenxiao Wang, and Soheil Feizi. 2025. <a href="#">Can ai-generated text be reliably detected? stress</a>	
846		
847		
	<a href="#">testing AI text detectors under various attacks</a> . <i>Trans. Mach. Learn. Res.</i> , 2025.	
	Areg Mikael Sarvazyan, José Ángel González, Paolo Rosso, and Marc Franco-Salvador. 2023. <a href="#">Supervised machine-generated text detectors: Family and scale matters</a> . In <i>Experimental IR Meets Multilinguality, Multimodality, and Interaction - 14th International Conference of the CLEF Association, CLEF 2023, Thessaloniki, Greece, September 18-21, 2023, Proceedings</i> , volume 14163 of <i>Lecture Notes in Computer Science</i> , pages 121–132. Springer.	850
		851
		852
		853
		854
		855
		856
		857
		858
	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. <a href="#">Proximal policy optimization algorithms</a> . <i>CoRR</i> , abs/1707.06347.	859
		860
		861
	Irene Solaiman, Miles Brundage, Jack Clark, Amanda Askell, Ariel Herbert-Voss, Jeff Wu, Alec Radford, and Jasmine Wang. 2019. <a href="#">Release strategies and the social impacts of language models</a> . <i>CoRR</i> , abs/1908.09203.	862
		863
		864
		865
		866
	Jinyan Su, Terry Yue Zhuo, Di Wang, and Preslav Nakov. 2023. <a href="#">Detectllm: Leveraging log rank information for zero-shot detection of machine-generated text</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 12395–12412. Association for Computational Linguistics.	867
		868
		869
		870
		871
		872
		873
	Chenxia Tang, Jianchun Liu, Hongli Xu, and Liusheng Huang. 2025. <a href="#">Top-<math>\nu\sigma</math>: Eliminating noise in logit space for robust token sampling of LLM</a> . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 10758–10774, Vienna, Austria. Association for Computational Linguistics.	874
		875
		876
		877
		878
		879
		880
	Eduard Tulchinskii, Kristian Kuznetsov, Laida Kushnareva, Daniil Cherniavskii, Sergey I. Nikolenko, Evgeny Burnaev, Serguei Barannikov, and Irina Piontkovskaya. 2023. <a href="#">Intrinsic dimension estimation for robust detection of ai-generated texts</a> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	881
		882
		883
		884
		885
		886
		887
		888
		889
	Cédric Villani and 1 others. 2008. <i>Optimal transport: old and new</i> , volume 338. Springer.	890
		891
	Herman Wold. 1966. Estimation of principal components and related models by iterative least squares. <i>Multivariate analysis</i> , pages 391–420.	892
		893
		894
	Junchao Wu, Shu Yang, Runzhe Zhan, Yulin Yuan, Lidia Sam Chao, and Derek Fai Wong. 2025a. <a href="#">A survey on LLM-generated text detection: Necessity, methods, and future directions</a> . <i>Computational Linguistics</i> , 51(1):275–338.	895
		896
		897
		898
		899
	Junchao Wu, Runzhe Zhan, Derek F. Wong, Shu Yang, Xuebo Liu, Lidia S. Chao, and Min Zhang. 2025b. <a href="#">Who wrote this? the key to zero-shot llm-generated text detection is gecscore</a> . In <i>Proceedings of the</i>	900
		901
		902
		903

904 31st International Conference on Computational Lin-  
 905 guistics, COLING 2025, Abu Dhabi, UAE, January  
 906 19-24, 2025, pages 10275–10292. Association for  
 907 Computational Linguistics.

908 Junchao Wu, Runzhe Zhan, Derek F. Wong, Shu Yang,  
 909 Xinyi Yang, Yulin Yuan, and Lidia S. Chao. 2024.  
 910 Detectrl: Benchmarking llm-generated text detec-  
 911 tion in real-world scenarios. In *Advances in Neural  
 912 Information Processing Systems 38: Annual Confer-  
 913 ence on Neural Information Processing Systems 2024,  
 914 NeurIPS 2024, Vancouver, BC, Canada, December  
 915 10 - 15, 2024*.

916 xAI. 2025. Grok-4.1 model card. Technical Report.

917 Xianjun Yang, Wei Cheng, Yue Wu, Linda Ruth Pet-  
 918 zold, William Yang Wang, and Haifeng Chen. 2024.  
 919 DNA-GPT: divergent n-gram analysis for training-  
 920 free detection of gpt-generated text. In *The Twelfth  
 921 International Conference on Learning Representa-  
 922 tions, ICLR 2024, Vienna, Austria, May 7-11, 2024*.  
 923 OpenReview.net.

924 Rowan Zellers, Ari Holtzman, Hannah Rashkin,  
 925 Yonatan Bisk, Ali Farhadi, Franziska Roesner, and  
 926 Yejin Choi. 2019. Defending against neural fake  
 927 news. In *Advances in Neural Information Processing  
 928 Systems 32: Annual Conference on Neural Informa-  
 929 tion Processing Systems 2019, NeurIPS 2019, De-  
 930 cember 8-14, 2019, Vancouver, BC, Canada*, pages  
 931 9051–9062.

932 Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. 2015.  
 933 Character-level convolutional networks for text clas-  
 934 sification. In *Advances in Neural Information Pro-  
 935 cessing Systems 28: Annual Conference on Neural In-  
 936 formation Processing Systems 2015, December 7-12,  
 937 2015, Montreal, Quebec, Canada*, pages 649–657.

938 Yanzhao Zhang, Mingxin Li, Dingkun Long, Xin Zhang,  
 939 Huan Lin, Baosong Yang, Pengjun Xie, An Yang,  
 940 Dayiheng Liu, Junyang Lin, Fei Huang, and Jingren  
 941 Zhou. 2025. Qwen3 embedding: Advancing text  
 942 embedding and reranking through foundation models.  
 943 *arXiv preprint arXiv:2506.05176*.

944 Hongyi Zhou, Jin Zhu, Pingfan Su, Kai Ye, Ying Yang,  
 945 Shakeel A O. B. Gavioli-Akilagun, and Chengchun  
 946 Shi. 2025. Adadetctgpt: Adaptive detection of llm-  
 947 generated text with statistical guarantees. *CoRR*,  
 948 abs/2510.01268.

## 949 A Additional related works

950 Beyond detection methods relying solely on in-  
 951 trinsic features, existing research has extensively  
 952 explored paradigms that incorporate external auxil-  
 953 iary signals.

954 Watermarking employs an active defense strat-  
 955 egy, aiming to implicitly embed verifiable statisti-  
 956 cal signals into text via specific sampling algo-  
 957 rithms during generation (Kirchenbauer et al.,

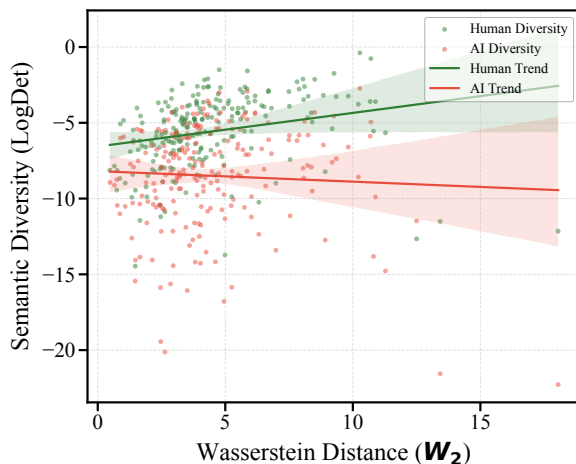


Figure 5: Semantic diversity under different structures. The increasing distance between AI and human centers indicates the rarity of the semantics.

2023). However, this technique faces a fundamen-  
 958 tal trade-off. Balancing the preservation of seman-  
 959 tic integrity with the maximization of robustness  
 960 and statistical detectability has become a central fo-  
 961 cus of recent research (Christ et al., 2024; Golowich  
 962 and Moitra, 2024; Cai et al., 2025).  
 963

964 In contrast, retrieval-based methods incorporate  
 965 external knowledge sources. They utilize sparse or  
 966 dense indexing mechanisms to align and compare  
 967 the input text against large-scale reference corpora.  
 968 By quantifying distributional discrepancies in lin-  
 969 guistic patterns between the input and reference  
 970 samples (human or machine), these methods sig-  
 971 nificantly enhance generalization in cross-domain  
 972 scenarios (Krishna et al., 2023; Sadasivan et al.,  
 973 2025). Furthermore, they offer new perspectives on  
 974 understanding adversarial in-context learning and  
 975 improving detection interpretability (Koike et al.,  
 976 2024, 2025).

977 Notably, although our proposed ProSSD falls  
 978 under the zero-shot detection paradigm, we draw  
 979 inspiration from both aforementioned approaches.

## 980 B Geometric Analysis of HWT and MGT 981 Distributions

982 This section aims to provide an in-depth theoretic-  
 983 al elucidation of the intrinsic differences between  
 984 HWT and MGT revealed in Figure 1 and Figure 5,  
 985 from the perspectives of geometric topology and  
 986 statistical properties.

### 987 B.1 Systematic Translation in Semantic Space

988 Figure 1 visualizes the distribution patterns of  
 989 HWT and MGT within the low-dimensional dis-

crimative subspace. We project high-dimensional embeddings using the supervised subspace learning described in Section 3.2 and calculate the class centroids  $\mu_H^\pi$  and  $\mu_M^\pi$  for each syntactic structure  $\pi$ . For intuitive visualization on a 2D plane, we apply principal component analysis (PCA) to these sets of means, spanning the maximum variance plane defined by the first (PC1) and second (PC2) principal components.

The most significant geometric feature in the plot is the phenomenon of **Systematic Translation**. Although distinct syntactic structures (e.g., conjunction-noun combinations, verb phrases) are scattered across different quadrants of the feature space, the connection lines between each pair of corresponding structural centroids ( $\mu_H^\pi, \mu_M^\pi$ ) exhibit a high degree of directional consistency. This indicates that the stylistic differences of MGT relative to human text are not local, disordered random perturbations, but a systematic bias pervasive throughout the feature space. We formally model this relationship as:

$$\mu_M^\pi = \mu_H^\pi + \delta_{global} + \epsilon_\pi, \quad (12)$$

where  $\delta_{global}$  is the global offset vector independent of specific syntactic structures, and  $\epsilon_\pi$  is a minor residual term for a specific structure (where  $\|\epsilon_\pi\| \ll \|\delta_{global}\|$ ). Notably, since the entire projection transformation is linear, a linear transformation cannot reconstruct anisotropic random noise in high-dimensional space into approximately parallel structured vectors in low-dimensional space. Therefore, this parallelism confirms that the systematic bias  $\delta_{global}$  is an intrinsic essential feature of the generative distribution of LLMs.

## B.2 Duality between Discriminative Distance and Distribution Diversity

Figure 5 further reveals significant differences in second-order statistics between HWT and MGT. While not detailed in the main text, this is crucial for understanding model behavior. The x-axis represents the Wasserstein distance weight  $w_\pi$ , which measuring the deviation of MGT from HWT under that syntactic structure. And the y-axis represents the log-determinant of the covariance matrix  $\ln |\Sigma|$ , which measuring the volume or diversity of the semantic distribution.

The observations reveal two diametrically opposed generation patterns:

- **Human-Written Text:** Exhibits a positive correlation trend. As the syntactic structure

deviates from convention (i.e.,  $w_\pi$  increases), the semantic richness of human text increases rather than decreases. This suggests that humans tend to mobilize a more diverse vocabulary to express precise meanings when navigating complex or rare syntax.

- **Machine-Generated Text:** Exhibits a significant negative correlation trend. When the model faces high-difficulty syntactic structures that deviate from its training priors, the volume of its generated semantic distribution shrinks significantly.

We define this phenomenon in MGT as **Conditional Semantic Collapse**: under the strong constraints of specific syntactic structures, LLMs tend to adopt conservative decoding strategies, converging to high-frequency, generic word combinations. This results in a significant reduction in local variance ( $\ln |\Sigma_M| \ll \ln |\Sigma_H|$ ). This collapse reflects the model’s tendency towards uncertainty avoidance when processing complex syntax.

## B.3 Proof of Discriminative Effectiveness Based on Mahalanobis Distance

Combining the findings of systematic bias  $\delta_{global}$  and semantic collapse, we hereby prove the theoretical validity of the detection statistic (difference in modified Mahalanobis distance) proposed in Section 3.3.

**Proposition 1** (Lower Bound of Expected Detection Statistic). *Assume that in the feature space, the local features  $\mathbf{x}$  of MGT follow the distribution  $\mathcal{N}(\mu_M, \Sigma)$ , and HWT follows  $\mathcal{N}(\mu_H, \Sigma)$ , with a non-zero offset  $\delta = \mu_M - \mu_H$ . Define the single-step detection statistic  $s_t$  as the difference between the distance to the human center and the distance to the machine center:*

$$s_t = \frac{1}{2} \left[ (\mathbf{x} - \mu_H)^T \Sigma^{-1} (\mathbf{x} - \mu_H) - (\mathbf{x} - \mu_M)^T \Sigma^{-1} (\mathbf{x} - \mu_M) \right]. \quad (13)$$

*Then, for machine-generated samples, the expectation of the score is strictly greater than zero.*

**Proof:** Reparameterize the machine-generated sample as  $\mathbf{x} = \mu_M + \mathbf{z}$ , where  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \Sigma)$  is zero-mean noise. Substituting this into the formula for  $s_t$  and expanding:

First term, distance relative to human center:

$$\begin{aligned} & (\mu_M + \mathbf{z} - \mu_H)^T \Sigma^{-1} (\mu_M + \mathbf{z} - \mu_H) \\ &= (\delta + \mathbf{z})^T \Sigma^{-1} (\delta + \mathbf{z}). \end{aligned} \quad (14)$$

Second term, distance relative to machine center:

$$\begin{aligned} & (\boldsymbol{\mu}_M + \mathbf{z} - \boldsymbol{\mu}_M)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_M + \mathbf{z} - \boldsymbol{\mu}_M) \\ & = \mathbf{z}^T \boldsymbol{\Sigma}^{-1} \mathbf{z}. \end{aligned} \quad (15)$$

Subtracting the two terms yields the simplified expression for  $2s_t$ :

$$2s_t = \boldsymbol{\delta}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\delta} + 2\boldsymbol{\delta}^T \boldsymbol{\Sigma}^{-1} \mathbf{z}. \quad (16)$$

Taking the mathematical expectation of the above equation, since  $\mathbf{z}$  is zero-mean noise, the expectation of the linear cross-term is  $\mathbb{E}[\mathbf{z}] = \mathbf{0}$ . Therefore:

$$\mathbb{E}[s_t] = \frac{1}{2} \boldsymbol{\delta}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\delta} = \frac{1}{2} \|\boldsymbol{\delta}\|_{\boldsymbol{\Sigma}^{-1}}^2. \quad (17)$$

Since the covariance matrix  $\boldsymbol{\Sigma}$  is positive definite, its inverse matrix  $\boldsymbol{\Sigma}^{-1}$  is also positive definite. According to the conclusion in Section 1, the systematic offset  $\|\boldsymbol{\delta}\| > 0$ ; thus,  $\mathbb{E}[s_t]$  is strictly positive. ■

This proposition not only proves the validity of the Mahalanobis distance but also reveals the contribution of the semantic collapse described in Section B.2 to detection performance. The expectation of the statistic  $s_t$  is proportional to  $\boldsymbol{\delta}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\delta}$ . The observed shrinkage in the distribution volume of MGT in Figure 5 implies that the eigenvalues of its covariance matrix become smaller, which in turn causes the eigenvalues of its inverse matrix  $\boldsymbol{\Sigma}^{-1}$  to increase. This expansion of the precision matrix essentially acts as a magnifying glass, significantly amplifying the weight of the systematic bias  $\boldsymbol{\delta}$  in the statistic, thereby further enhancing the discriminative signal-to-noise ratio of the model against MGT.

## C Detailed Derivations and Proofs

### C.1 Theoretical Derivation of Supervised Subspace Projection

This section provides mathematical proofs for the recursive optimization problem defined in Section 3.2 of the main text. We derive the closed-form solution for the optimal projection direction, prove the orthogonality of the feature deflation step, and finally elucidate the necessity of introducing multi-dimensional features for distribution estimation.

**Proposition 2.** *Given the centered residual embedding matrix  $\mathbf{E}_{j-1} \in \mathbb{R}^{N \times D}$  and the centered*

*label vector  $\mathbf{y} \in \mathbb{R}^N$ , the optimal solution  $\mathbf{p}_j^*$  to the optimization problem*

$$\mathbf{p}_j^* = \underset{\mathbf{p} \in \mathbb{R}^D}{\operatorname{argmax}} (\operatorname{Cov}(\mathbf{E}_{j-1} \mathbf{p}, \mathbf{y}))^2 \quad (18)$$

$$\text{s.t. } \|\mathbf{p}\|_2 = 1,$$

*is the normalized cross-covariance vector between the current residual features and the labels.*

**Proof:** Since  $\mathbf{E}_{j-1}$  and  $\mathbf{y}$  are centered vectors, their sample covariance is proportional to their inner product. The objective function can be rewritten as:

$$J(\mathbf{p}) = (\mathbf{y}^T \mathbf{E}_{j-1} \mathbf{p})^2 = \mathbf{p}^T (\mathbf{E}_{j-1}^T \mathbf{y} \mathbf{y}^T \mathbf{E}_{j-1}) \mathbf{p}. \quad (19)$$

We introduce the Lagrange multiplier  $\lambda$  to construct the Lagrangian function:

$$\mathcal{L}(\mathbf{p}, \lambda) = \mathbf{p}^T \mathbf{M}_{j-1} \mathbf{p} - \lambda (\mathbf{p}^T \mathbf{p} - 1), \quad (20)$$

where  $\mathbf{M}_{j-1} = (\mathbf{E}_{j-1}^T \mathbf{y})(\mathbf{y}^T \mathbf{E}_{j-1})$  is a rank-1 positive semi-definite matrix. Taking the partial derivative with respect to  $\mathbf{p}$  and setting it to zero yields the eigenvalue equation:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{p}} = 2\mathbf{M}_{j-1} \mathbf{p} - 2\lambda \mathbf{p} = 0 \implies \mathbf{M}_{j-1} \mathbf{p} = \lambda \mathbf{p}. \quad (21)$$

This indicates that  $\mathbf{p}_j^*$  must be the eigenvector corresponding to the maximum eigenvalue of matrix  $\mathbf{M}_{j-1}$ . Let  $\mathbf{u} = \mathbf{E}_{j-1}^T \mathbf{y}$ , then  $\mathbf{M}_{j-1} = \mathbf{u} \mathbf{u}^T$ . Applying matrix  $\mathbf{M}_{j-1}$  to vector  $\mathbf{u}$ :

$$\mathbf{M}_{j-1} \mathbf{u} = (\mathbf{u} \mathbf{u}^T) \mathbf{u} = \mathbf{u} (\mathbf{u}^T \mathbf{u}) = \|\mathbf{u}\|^2 \mathbf{u}. \quad (22)$$

It follows that  $\mathbf{u} = \mathbf{E}_{j-1}^T \mathbf{y}$  is the principal eigenvector corresponding to the maximum eigenvalue  $\lambda_{max} = \|\mathbf{u}\|^2$ . Combining this with the unit norm constraint  $\|\mathbf{p}\|_2 = 1$ , we obtain the closed-form solution:

$$\mathbf{p}_j^* = \frac{\mathbf{E}_{j-1}^T \mathbf{y}}{\|\mathbf{E}_{j-1}^T \mathbf{y}\|_2}. \quad (23)$$

This vector indicates that the optimal projection direction is collinear with the current residual-label cross-covariance direction. ■

**Proposition 3.** *The feature deflation step ensures that the deflated residual matrix  $\mathbf{E}_j$  is orthogonal to the current projection direction  $\mathbf{p}_j^*$ .*

**Proof:** By definition,  $\mathbf{E}_j = \mathbf{E}_{j-1} - s_j(\mathbf{p}_j^*)^T$ , where  $s_j = \mathbf{E}_{j-1}\mathbf{p}_j^*$ . Examining the projection of the residual matrix  $\mathbf{E}_j$  onto the current direction  $\mathbf{p}_j^*$ :

$$\begin{aligned} \mathbf{E}_j\mathbf{p}_j^* &= (\mathbf{E}_{j-1} - s_j(\mathbf{p}_j^*)^T)\mathbf{p}_j^* \\ &= \mathbf{E}_{j-1}\mathbf{p}_j^* - s_j(\mathbf{p}_j^*)^T\mathbf{p}_j^*. \end{aligned} \quad (24)$$

Substituting  $s_j$  and utilizing the unit vector property  $(\mathbf{p}_j^*)^T\mathbf{p}_j^* = 1$ , we obtain:

$$\mathbf{E}_j\mathbf{p}_j^* = s_j - s_j(1) = \mathbf{0}. \quad (25)$$

That is, the column space of the residual matrix  $\mathbf{E}_j$  is orthogonal to the current projection direction  $\mathbf{p}_j^*$ . ■

**Analysis on Multi-dimensional Subspace Construction.** The above propositions establish the mathematical foundation for the iterative solution. It is necessary to clarify why constructing a subspace with  $k > 1$  remains essential in a binary classification task, where  $y$  is a 1D scalar.

According to Proposition 2, the first basis vector  $\mathbf{p}_1^*$  effectively captures all linear mean differences between HWT and MGT. However, the probabilistic model proposed in Section 3.3 relies not only on the mean  $\boldsymbol{\mu}$  but also heavily on the estimation of the covariance matrix  $\boldsymbol{\Sigma}$  to capture precise distributional characteristics. If only  $k = 1$  is selected, the feature space degenerates into a line, and the covariance matrix degenerates into a scalar variance, failing to describe the spatial distribution shape of the data. Through the aforementioned iterative process, subsequent projection vectors  $\mathbf{p}_j^*$  ( $j > 1$ ) continue to extract major structural information of the data in the residual space orthogonal to  $\mathbf{p}_1^*$ . This set of orthogonal bases  $\{\mathbf{p}_1, \dots, \mathbf{p}_k\}$  jointly constitutes a low-dimensional complete feature subspace. This allows us to retain the original discriminative information while accurately estimating the covariance structure of the Gaussian distribution. Mathematically, this aligns with the NIPALS algorithm of partial least squares (PLS) (Wold, 1966), which constructs stable low-dimensional representations through iterative deflation.

## C.2 Theoretical Analysis of Distribution Modeling

### C.2.1 Proof of Asymptotic Normality

This section establishes the distributional properties of the meta-semantic vector  $\mathbf{x}_t$ , defined in Section 3.3 of the main text, as the original embedding dimension  $D \rightarrow \infty$ .

**Definition C.1 (Construction of Meta-Semantic Vectors)** Let  $\mathbf{e}_t, \mathbf{e}_{t+1} \in \mathbb{R}^D$  be the original embedding vectors at time steps  $t$  and  $t + 1$ . Let  $\mathbf{P} \in \mathbb{R}^{D \times k}$  be the column-orthogonal projection matrix obtained in Section 3.2. We define the meta-semantic vector  $\mathbf{x}_t \in \mathbb{R}^{2k}$  as the concatenation of projected features from adjacent time steps:

$$\mathbf{x}_t = \begin{bmatrix} \mathbf{P}^T \mathbf{e}_t \\ \mathbf{P}^T \mathbf{e}_{t+1} \end{bmatrix}. \quad (26)$$

**Assumption 1** (High-dimensional Noise Decomposition). *Following common settings in high-dimensional statistics, we model the original embedding  $\mathbf{e}_t$  as the sum of a deterministic mean signal and random noise:*

$$\mathbf{e}_t = \boldsymbol{\mu}_t + \boldsymbol{\epsilon}_t, \quad (27)$$

where  $\boldsymbol{\mu}_t = \mathbb{E}[\mathbf{e}_t]$  is the intrinsic semantic mean within this context, and  $\boldsymbol{\epsilon}_t$  is a zero-mean random noise vector.

We construct a joint noise vector  $\boldsymbol{\xi} \in \mathbb{R}^{2D}$  by concatenating the components of  $\boldsymbol{\epsilon}_t$  and  $\boldsymbol{\epsilon}_{t+1}$ :

$$\boldsymbol{\xi} = [\epsilon_{t,1}, \dots, \epsilon_{t,D}, \epsilon_{t+1,1}, \dots, \epsilon_{t+1,D}]^T. \quad (28)$$

We assume the components  $\{\xi_j\}_{j=1}^{2D}$  of  $\boldsymbol{\xi}$  satisfy the following conditions:

1. **Zero Mean:**  $\mathbb{E}[\xi_j] = 0$ .
2. **Finite Variance:**  $\text{Var}(\xi_j) = \sigma_j^2 < \infty$ .
3. **Boundedness:** There exists a constant  $M < \infty$  such that  $|\xi_j| \leq M$  holds almost everywhere (this assumption is guaranteed by the normalization mechanism of LayerNorm layers).

**Assumption 2** (Feller Condition). *For any unit projection direction, we assume that the projection coefficients do not overly concentrate on any single original dimension. Formally, for linear combination coefficients  $\gamma_j$ , the Feller delocalization condition is satisfied as  $D \rightarrow \infty$ :*

$$\lim_{D \rightarrow \infty} \frac{\max_{1 \leq j \leq 2D} |\gamma_j|}{\sqrt{\sum_{j=1}^{2D} \gamma_j^2 \sigma_j^2}} = 0. \quad (29)$$

**Proposition 4** (Asymptotic Normality of Meta-Semantic Distribution). *Based on Assumptions 1 and 2, as the original dimension  $D \rightarrow \infty$ , the meta-semantic vector  $\mathbf{x}_t$  converges in distribution to a multivariate normal distribution:*

$$\mathbf{x}_t \xrightarrow{d} \mathcal{N}(\boldsymbol{\mu}_x, \boldsymbol{\Sigma}_x), \quad (30)$$

where  $\boldsymbol{\mu}_x = [\boldsymbol{\mu}_t^T \mathbf{P}, \boldsymbol{\mu}_{t+1}^T \mathbf{P}]^T$ .

**Proof:** According to the Cramér-Wold theorem (Cramér and Wold, 1936), a necessary and sufficient condition for the random vector  $\mathbf{x}_t$  to follow a multivariate normal distribution is that for any non-zero constant vector  $\boldsymbol{\lambda} \in \mathbb{R}^{2k}$ , the scalar random variable  $Z_D = \boldsymbol{\lambda}^T \mathbf{x}_t$  follows a univariate normal distribution.

Partition  $\boldsymbol{\lambda}$  into  $\boldsymbol{\lambda} = [\boldsymbol{\alpha}^T, \boldsymbol{\beta}^T]^T$ , where  $\boldsymbol{\alpha}, \boldsymbol{\beta} \in \mathbb{R}^k$ . Substituting the definition of  $\mathbf{x}_t$  into  $Z_D$ :

$$\begin{aligned} Z_D &= \boldsymbol{\alpha}^T \mathbf{P}^T (\boldsymbol{\mu}_t + \boldsymbol{\epsilon}_t) + \boldsymbol{\beta}^T \mathbf{P}^T (\boldsymbol{\mu}_{t+1} + \boldsymbol{\epsilon}_{t+1}) \\ &= \mu_Z + S_D, \end{aligned} \quad (31)$$

where:

$$\begin{aligned} \mu_Z &= \boldsymbol{\alpha}^T \mathbf{P}^T \boldsymbol{\mu}_t + \boldsymbol{\beta}^T \mathbf{P}^T \boldsymbol{\mu}_{t+1} \\ S_D &= \boldsymbol{\alpha}^T \mathbf{P}^T \boldsymbol{\epsilon}_t + \boldsymbol{\beta}^T \mathbf{P}^T \boldsymbol{\epsilon}_{t+1}, \end{aligned}$$

here,  $\mu_Z$  denotes the deterministic mean, and  $S_D$  represents the random fluctuation term. Since  $\mu_Z$  is constant, we only need to prove that the random term  $S_D$  converges in distribution to  $\mathcal{N}(0, \sigma_Z^2)$ .

Define the coefficient vector  $\boldsymbol{\gamma} \in \mathbb{R}^{2D}$  as  $\boldsymbol{\gamma} = [(\mathbf{P}\boldsymbol{\alpha})^T, (\mathbf{P}\boldsymbol{\beta})^T]^T$ . Then  $S_D$  can be rewritten as a weighted sum of joint noise components:

$$S_D = \sum_{j=1}^{2D} \gamma_j \xi_j, \quad (32)$$

Calculate the variance sequence  $B_D^2$  of  $S_D$ :

$$B_D^2 = \text{Var}(S_D) = \sum_{j=1}^{2D} \gamma_j^2 \sigma_j^2, \quad (33)$$

According to the Lindeberg-Feller central limit theorem (Feller, 1991) for sums of independent (or weakly dependent) non-identically distributed random variables, convergence to a normal distribution is guaranteed if the Lindeberg condition holds. We verify that for any  $\tau > 0$ :

$$L = \lim_{D \rightarrow \infty} \frac{1}{B_D^2} \sum_{j=1}^{2D} \mathbb{E} [(\gamma_j \xi_j)^2 \cdot \mathbb{I}(|\gamma_j \xi_j| > \tau B_D)] = 0. \quad (34)$$

From the boundedness in Assumption 1, we know  $|\xi_j| \leq M$ , thus  $|\gamma_j \xi_j| \leq |\gamma_j| M$ . The indicator function  $\mathbb{I}(\cdot)$  is non-zero only if  $|\gamma_j \xi_j| > \tau B_D$ , which implies the necessary condition:

$$|\gamma_j| M > \tau B_D \iff \frac{|\gamma_j|}{B_D} > \frac{\tau}{M}. \quad (35)$$

However, according to Assumption 2:

$$\lim_{D \rightarrow \infty} \frac{\max_j |\gamma_j|}{B_D} = 0. \quad (36)$$

This implies that for sufficiently large  $D$ , given any fixed  $\tau, M > 0$ , there exists no index  $j$  such that  $|\gamma_j| M > \tau B_D$  holds. Therefore, the indicator function  $\mathbb{I}(|\gamma_j \xi_j| > \tau B_D)$  is identically zero. Consequently, the limit  $L = 0$ , and the Lindeberg condition is satisfied.

Since the Lindeberg condition is met, the random variable  $S_D$  converges in distribution to a scaled standard normal distribution:

$$\frac{S_D}{B_D} \xrightarrow{d} \mathcal{N}(0, 1) \implies S_D \xrightarrow{d} \mathcal{N}(0, B_D^2). \quad (37)$$

Substituting back into the expression for  $Z_D$ , we have  $Z_D \xrightarrow{d} \mathcal{N}(\mu_Z, B_D^2)$ . Since  $\boldsymbol{\lambda}$  is arbitrary, by the Cramér-Wold Theorem, the meta-semantic vector  $\mathbf{x}_t$  converges in distribution to a multivariate Gaussian distribution. ■

**Remark:** Although the above theorem is derived based on the asymptotic  $D \rightarrow \infty$ , in our experimental setup, the original semantic space dimension  $D$  (e.g., 1024 for RoBERTa-large) satisfies the requirements for high-dimensional statistical approximation. Furthermore, the projection matrix  $\mathbf{P}$  obtained via supervised subspace learning tends to utilize global semantic information. This results in projection coefficients  $\gamma_j$  exhibiting significant non-sparse characteristics, preventing any single original dimension from dominating the distribution. Therefore, modeling compact semantic features using a multivariate Gaussian distribution in Section 3.3 is not only theoretically grounded but also consistent with the statistical laws of high-dimensional data.

## C.2.2 Argument for Model Distribution Selection based on Maximum Entropy Principle

Although the asymptotic normality of meta-semantic features has been proven above, in practical applications with finite dimensions, observed data may not perfectly follow a Gaussian distribution. This section proves from an information-theoretic perspective that, given observations of only the sample mean and covariance, the Gaussian distribution model is the optimal choice as it introduces the least prior bias.

**Proposition 5.** For a random variable  $\mathbf{x} \in \mathbb{R}^{2k}$ , if we observe only its first moment and second moment from dataset  $\mathcal{D}$ , then among all probability distributions  $p(\mathbf{x})$  satisfying these moment constraints, the multivariate Gaussian distribution has the maximum Differential Entropy.

**Proof:** Our objective is to maximize the differential entropy  $H(p)$ :

$$\begin{aligned} \max_p \quad & - \int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x} \\ \text{s.t.} \quad & \int p(\mathbf{x}) d\mathbf{x} = 1 \\ & \int \mathbf{x} p(\mathbf{x}) d\mathbf{x} = \boldsymbol{\mu} \\ & \int (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T p(\mathbf{x}) d\mathbf{x} = \boldsymbol{\Sigma}. \end{aligned} \quad (38)$$

We construct the Lagrangian function:

$$\begin{aligned} \mathcal{L} = H(p) + \lambda_0 \left( \int p - 1 \right) + \boldsymbol{\lambda}_1^T \left( \int \mathbf{x} p - \boldsymbol{\mu} \right) \\ + \text{Tr} \left( \boldsymbol{\Lambda}_2 \left( \int (\mathbf{x} \mathbf{x}^T) p - (\boldsymbol{\Sigma} + \boldsymbol{\mu} \boldsymbol{\mu}^T) \right) \right). \end{aligned} \quad (39)$$

Taking the derivative with respect to  $p(\mathbf{x})$  using variational methods and setting it to zero, the form of  $p(\mathbf{x})$  must be an exponential family distribution:

$$p^*(\mathbf{x}) = \exp \left( -1 + \lambda_0 + \boldsymbol{\lambda}_1^T \mathbf{x} + \mathbf{x}^T \boldsymbol{\Lambda}_2 \mathbf{x} \right). \quad (40)$$

Substituting the constraints to solve for the Lagrange multipliers, we finally obtain:

$$\begin{aligned} p^*(\mathbf{x}) = \frac{1}{(2\pi)^k |\boldsymbol{\Sigma}|^{1/2}} \\ \times \exp \left( -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right), \end{aligned} \quad (41)$$

which is the multivariate Gaussian distribution. ■

Since we possess a large-scale dataset (e.g., 700,000 meta-structure pairs), by the Law of Large Numbers, we can estimate  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  with high precision. In the absence of prior knowledge regarding higher-order statistics of the data, adopting a Gaussian distribution model is equivalent to making the minimal assumption at the information-theoretic level. Any other distributional assumption would implicitly introduce additional structural information not supported by the data, thereby increasing the risk of model overfitting.

### C.3 2-Wasserstein Distance for Gaussian Distributions

Let  $P_H^\pi = \mathcal{N}(\boldsymbol{\mu}_H^\pi, \boldsymbol{\Sigma}_H^\pi)$  and  $P_M^\pi = \mathcal{N}(\boldsymbol{\mu}_M^\pi, \boldsymbol{\Sigma}_M^\pi)$  be two Gaussian distributions on  $\mathbb{R}^{2k}$ . The squared

2-Wasserstein distance is defined by the optimal transport problem:

$$W_2^2(P_H^\pi, P_M^\pi) = \inf_{\gamma \in \Pi(P_H^\pi, P_M^\pi)} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \gamma} [\|\mathbf{x} - \mathbf{y}\|_2^2]. \quad (42)$$

By expanding the squared Euclidean norm, we decompose the expectation based on the first and second moments:

$$\begin{aligned} \mathbb{E}[\|\mathbf{x} - \mathbf{y}\|_2^2] = \|\boldsymbol{\mu}_H^\pi - \boldsymbol{\mu}_M^\pi\|_2^2 + \mathbb{E}[\|\bar{\mathbf{x}} - \bar{\mathbf{y}}\|_2^2] \\ + 2(\boldsymbol{\mu}_H^\pi - \boldsymbol{\mu}_M^\pi)^\top \mathbb{E}[\bar{\mathbf{x}} - \bar{\mathbf{y}}], \end{aligned} \quad (43)$$

where  $\bar{\mathbf{x}} = \mathbf{x} - \boldsymbol{\mu}_H^\pi$  and  $\bar{\mathbf{y}} = \mathbf{y} - \boldsymbol{\mu}_M^\pi$  are centered random variables. Since  $\mathbb{E}[\bar{\mathbf{x}}] = \mathbb{E}[\bar{\mathbf{y}}] = \mathbf{0}$ , the cross-term vanishes. The problem reduces to finding the optimal coupling for the centered covariances:

$$W_2^2(P_H^\pi, P_M^\pi) = \|\boldsymbol{\mu}_H^\pi - \boldsymbol{\mu}_M^\pi\|_2^2 + \inf_{\gamma} \mathbb{E}[\|\bar{\mathbf{x}} - \bar{\mathbf{y}}\|_2^2]. \quad (44)$$

Minimizing the covariance term is equivalent to maximizing the correlation trace. According to the properties of optimal transport for Gaussian measures, the optimal value yields the Bures-Wasserstein metric. Combining this with the mean difference, we obtain the final metric  $w_\pi$ :

$$\begin{aligned} w_\pi = \left( \|\boldsymbol{\mu}_H^\pi - \boldsymbol{\mu}_M^\pi\|_2^2 + \right. \\ \left. \text{Tr} \left( \boldsymbol{\Sigma}_H^\pi + \boldsymbol{\Sigma}_M^\pi - 2(\boldsymbol{\Sigma}_H^{\pi \frac{1}{2}} \boldsymbol{\Sigma}_M^\pi \boldsymbol{\Sigma}_H^{\pi \frac{1}{2}})^{\frac{1}{2}} \right) \right)^{1/2}. \end{aligned} \quad (45)$$

### C.4 Proof of Wasserstein Discriminative Bound

This section provides the theoretical proof for the discriminative bound based on the Wasserstein distance discussed in Section 3.3 of the main text. We demonstrate that for any function satisfying the Lipschitz continuity condition, the difference in its expectations over two distributions is strictly bounded by the 2-Wasserstein distance between these distributions.

**Theorem 2** (Wasserstein Discriminative Bound). *Let  $P_H$  and  $P_M$  be two probability measures on  $\mathbb{R}^{2k}$ . For any  $K$ -Lipschitz continuous discriminative function  $f : \mathbb{R}^{2k} \rightarrow \mathbb{R}$ , the difference in expectations between the two distributions satisfies:*

$$|\mathbb{E}_{\mathbf{x} \sim P_H}[f(\mathbf{x})] - \mathbb{E}_{\mathbf{y} \sim P_M}[f(\mathbf{y})]| \leq K \cdot W_2(P_H, P_M). \quad (46)$$

**Proof:** Let  $\Pi(P_H, P_M)$  be the set of all couplings whose marginals are  $P_H$  and  $P_M$ . For any coupling  $\gamma \in \Pi(P_H, P_M)$ , if the random variable pair  $(\mathbf{X}, \mathbf{Y}) \sim \gamma$ , then by the property of marginal distributions,  $\mathbf{X} \sim P_H$  and  $\mathbf{Y} \sim P_M$ .

Consider the absolute difference of the expectations:

$$\begin{aligned} \Delta_f &= |\mathbb{E}_{\mathbf{X} \sim P_H}[f(\mathbf{X})] - \mathbb{E}_{\mathbf{Y} \sim P_M}[f(\mathbf{Y})]| \\ &= \left| \int f(\mathbf{x}) dP_H(\mathbf{x}) - \int f(\mathbf{y}) dP_M(\mathbf{y}) \right| \\ &= \left| \iint (f(\mathbf{x}) - f(\mathbf{y})) d\gamma(\mathbf{x}, \mathbf{y}) \right|. \end{aligned} \quad (47)$$

By the integral triangle inequality, we have:

$$\Delta_f \leq \iint |f(\mathbf{x}) - f(\mathbf{y})| d\gamma(\mathbf{x}, \mathbf{y}). \quad (48)$$

Utilizing the  $K$ -Lipschitz continuity condition of  $f$ :

$$\begin{aligned} \Delta_f &\leq \iint K \|\mathbf{x} - \mathbf{y}\|_2 d\gamma(\mathbf{x}, \mathbf{y}) \\ &= K \cdot \mathbb{E}_{(\mathbf{X}, \mathbf{Y}) \sim \gamma} [\|\mathbf{X} - \mathbf{Y}\|_2]. \end{aligned} \quad (49)$$

According to Lyapunov's inequality, the first moment of a random variable is less than or equal to the square root of its second moment, i.e.,  $\mathbb{E}[Z] \leq (\mathbb{E}[Z^2])^{1/2}$ . Applying this to the Euclidean distance  $\|\mathbf{X} - \mathbf{Y}\|_2$ :

$$\mathbb{E}[\|\mathbf{X} - \mathbf{Y}\|_2] \leq (\mathbb{E}[\|\mathbf{X} - \mathbf{Y}\|_2^2])^{1/2}. \quad (50)$$

Therefore, for any  $\gamma \in \Pi(P_H, P_M)$ , it holds that:

$$\Delta_f \leq K \left( \iint \|\mathbf{x} - \mathbf{y}\|_2^2 d\gamma(\mathbf{x}, \mathbf{y}) \right)^{1/2}. \quad (51)$$

Note that the left side  $\Delta_f$  is a constant dependent only on the marginals  $P_H$  and  $P_M$ , and is independent of the specific coupling  $\gamma$ . Thus, the inequality holds for any  $\gamma$  in  $\Pi(P_H, P_M)$ , and consequently holds for the infimum over  $\gamma$  on the right side:

$$\Delta_f \leq K \cdot \inf_{\gamma \in \Pi} (\mathbb{E}_{\gamma}[\|\mathbf{X} - \mathbf{Y}\|_2^2])^{1/2}. \quad (52)$$

Due to the continuity and monotonicity of the square root function,  $\inf(\mathbb{E}^{1/2}) = (\inf \mathbb{E})^{1/2}$ . Combining this with the definition of the 2-Wasserstein distance  $W_2(P_H, P_M) = (\inf_{\gamma \in \Pi} \mathbb{E}[\|\mathbf{X} - \mathbf{Y}\|_2^2])^{1/2}$ , we finally obtain:

$$|\mathbb{E}_{P_H}[f] - \mathbb{E}_{P_M}[f]| \leq K \cdot W_2(P_H, P_M). \quad (53)$$

This concludes the proof. ■

**Remark:** This theorem provides the theoretical basis for the weighting strategy presented in Section 3.3 of the main text. The term  $W_2(P_H, P_M)$  on the right side measures the geometric separability between human and machine text in the semantic space under a specific syntactic structure  $\pi$ . Meanwhile,  $K$  represents the sensitivity of the discriminator to semantic perturbations (i.e., the smoothness of the classification surface). This bound indicates that the larger the  $W_2$  distance between two distributions, the greater the theoretical expected difference any smooth discriminative function can achieve in distinguishing them.

### C.5 Derivation of the Modified Mahalanobis Distance

This section derives the equivalence relationship between the log-likelihood ratio (LLR) test described in Section 3.4 of the main text and the modified Mahalanobis Distance.

**Proposition 6** (Equivalence of LLR and Mahalanobis Distance). *Let  $P_H(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_H, \boldsymbol{\Sigma}_H)$  and  $P_M(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_M, \boldsymbol{\Sigma}_M)$  be the feature distributions of human and machine text under a specific structure, respectively. Define the modified Mahalanobis distance as  $\mathcal{M}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) + \ln |\boldsymbol{\Sigma}|$ . Then, for a sample  $\mathbf{x}_t$ , its log-likelihood ratio statistic  $s_t$  satisfies:*

$$\begin{aligned} s_t &= \ln \frac{P_M(\mathbf{x}_t)}{P_H(\mathbf{x}_t)} \\ \iff s_t &= \frac{1}{2} \left[ \mathcal{M}(\mathbf{x}_t; \boldsymbol{\mu}_H, \boldsymbol{\Sigma}_H) - \mathcal{M}(\mathbf{x}_t; \boldsymbol{\mu}_M, \boldsymbol{\Sigma}_M) \right]. \end{aligned} \quad (54)$$

**Proof:** For a  $d$ -dimensional multivariate Gaussian distribution  $P(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ , the log-likelihood function is:

$$\begin{aligned} \ln P(\mathbf{x}) &= -\frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln |\boldsymbol{\Sigma}| \\ &\quad - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}). \end{aligned} \quad (55)$$

Expanding the log-likelihood ratio  $s_t$ :

$$\begin{aligned} s_t &= \ln P_M(\mathbf{x}_t) - \ln P_H(\mathbf{x}_t) \\ &= \left[ -\frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln |\boldsymbol{\Sigma}_M| - \frac{1}{2} Q_M(\mathbf{x}_t) \right] \\ &\quad - \left[ -\frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln |\boldsymbol{\Sigma}_H| - \frac{1}{2} Q_H(\mathbf{x}_t) \right], \end{aligned} \quad (56)$$

where  $Q(\mathbf{x}_t) = (\mathbf{x}_t - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_t - \boldsymbol{\mu})$  is the quadratic term.

Canceling the constant term  $-\frac{d}{2} \ln(2\pi)$  and rearranging the terms, we obtain:

$$s_t = \frac{1}{2} (Q_H(\mathbf{x}_t) + \ln |\boldsymbol{\Sigma}_H|) - \frac{1}{2} (Q_M(\mathbf{x}_t) + \ln |\boldsymbol{\Sigma}_M|). \quad (57)$$

According to the definition of the modified Mahalanobis distance  $\mathcal{M}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \triangleq Q(\mathbf{x}) + \ln |\boldsymbol{\Sigma}|$ , substituting this into the equation yields:

$$s_t = \frac{1}{2} \mathcal{M}(\mathbf{x}_t; \boldsymbol{\mu}_H, \boldsymbol{\Sigma}_H) - \frac{1}{2} \mathcal{M}(\mathbf{x}_t; \boldsymbol{\mu}_M, \boldsymbol{\Sigma}_M). \quad (58)$$

This concludes the proof. ■

Theorem 1 requires the discriminative function  $f(\mathbf{x})$  to satisfy the  $K$ -Lipschitz continuity condition. Although the Mahalanobis distance, as a quadratic function, does not possess global Lipschitz properties over the entire  $\mathbb{R}^{2k}$  space, within the framework of our method, the domain of the feature variable  $\mathbf{x}$  is bounded, thereby guaranteeing that it satisfies the Lipschitz condition.

The reasoning is as follows based on three conditions. (1) Input boundedness: The meta-semantic vector  $\mathbf{x}_t$  is formed by concatenating adjacent projected features:  $\mathbf{x}_t = [\mathbf{v}_t; \mathbf{v}_{t+1}]$ . The projected feature  $\mathbf{v} = \mathbf{P}^T \mathbf{e}$  is obtained by transforming the original semantic embedding  $\mathbf{e}$  via the probe matrix  $\mathbf{P} = [\mathbf{p}_1, \dots, \mathbf{p}_k]$ . (2) Original space boundedness: The output embeddings  $\mathbf{e}$  of modern LLMs typically undergo layer normalization, ensuring their Euclidean norms are bounded. That is, there exists a constant  $R_e$  such that  $\|\mathbf{e}\|_2 \leq R_e$ . (3) Projection constraint: In the optimization objective of Section 3.2, we explicitly impose a unit norm constraint on the probe vectors:  $\|\mathbf{p}_i\|_2 = 1$ .

By applying the Cauchy-Schwarz inequality (Anderson et al., 1958), the norm of the projected features satisfies:

$$\|\mathbf{v}\|_2 = \|\mathbf{P}^T \mathbf{e}\|_2 \leq \|\mathbf{P}\|_F \|\mathbf{e}\|_2 \leq \sqrt{k} \cdot 1 \cdot R_e. \quad (59)$$

Therefore, the meta-semantic vector lies within a compact subset  $\mathcal{C}$  of the space  $\mathbb{R}^{2k}$ . Consider the gradient of the discriminative function  $f(\mathbf{x})$ :

$$\nabla f(\mathbf{x}) = \boldsymbol{\Sigma}_H^{-1} (\mathbf{x} - \boldsymbol{\mu}_H) - \boldsymbol{\Sigma}_M^{-1} (\mathbf{x} - \boldsymbol{\mu}_M). \quad (60)$$

Since  $\mathbf{x} \in \mathcal{C}$  is bounded and  $\boldsymbol{\Sigma}_H, \boldsymbol{\Sigma}_M$  are positive definite matrices, the gradient norm  $\|\nabla f(\mathbf{x})\|_2$

has an upper bound  $K_{max}$  on the region  $\mathcal{C}$ . According to the Mean Value Theorem, if the gradient of a function is bounded on a convex compact set, then the function is Lipschitz continuous.

## D Benchmark and Baselines

### D.1 Benchmark Construction and Quality Control

To ensure comparability and rigor amidst rapid model iteration, we strictly followed the DetectRL protocol proposed by (Wu et al., 2024). for dataset construction. We also enforced equivalent standards of manual review during the post-processing stage.

**Human-Written Text (HWT) Curation.** We completely replicated the original data configuration, covering four domains prone to misuse: academic writing (ArXiv Abstracts), news reports (XSum), creative writing (Writing Prompts), and social media reviews (Yelp Reviews). To avoid data contamination from LLMs, all selected human samples date from the pre-ChatGPT era. Each domain contains 2,800 filtered high-quality samples, constituting a reliable detection baseline.

**SOTA Model Generation Strategy.** For the construction of MGT, we upgraded the generation sources to current SOTA models: GPT-5.1, Claude-Sonnet-4, and Gemini-3-Flash, Grok-4.1. Abandoning simple unconstrained generation, we strictly adhered to the task-specific conditional generation strategy from the original paper. Carefully designed prompts forced the models to generate content strictly aligned with the length of the corresponding HWT within specific contexts.

**Quality Control & Data Partitioning.** We implemented a multi-stage quality control process to ensure dataset safety and academic rigor. First, we automatically filtered invalid samples with insufficient length during the preprocessing stage. Second, despite the automated generation process, we introduced a manual sampling review mechanism before storage. This aimed to identify and remove samples with potential logical breakdowns, repetition, or offensive content. Regarding data partitioning, we followed the original protocol (Random Seed = 42), splitting the data into a training set (1,800 HWT/MGT pairs) and a test set (1,000 HWT/MGT pairs). Stratified sampling was employed to ensure the proportions of domains and model sources in each subset remained consistent with the original paper.

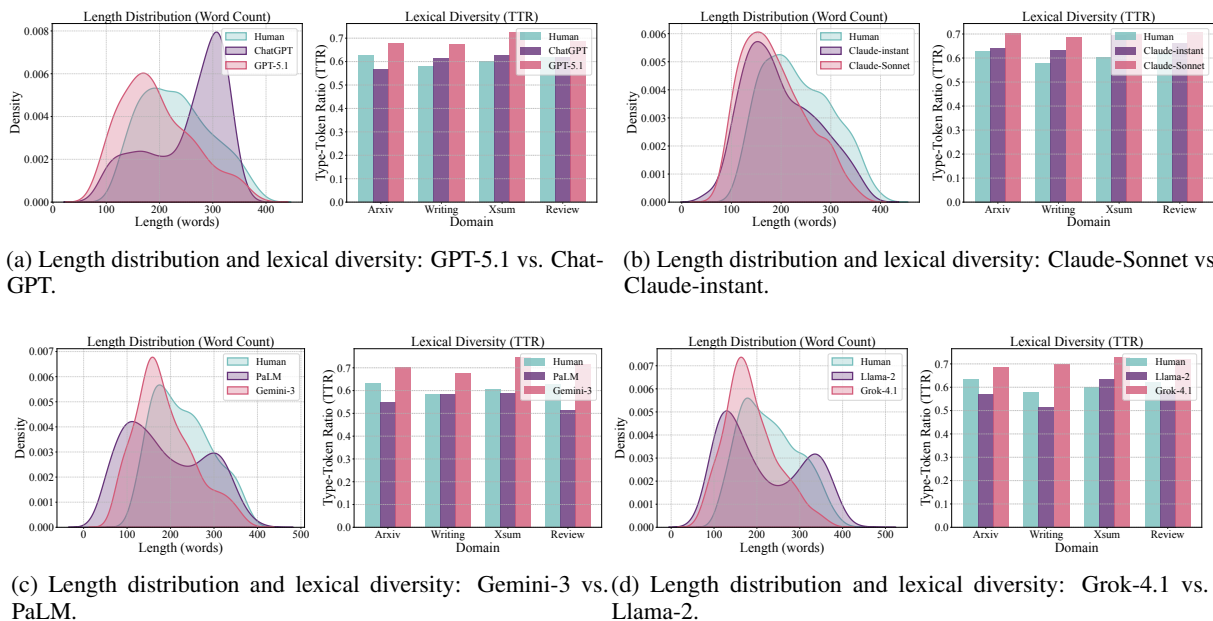


Figure 6: **Overview of length distributions and lexical diversity comparisons.** Each subplot illustrates the comparative analysis of human text, old models, and new models across different datasets, visualizing the kernel density estimation (KDE) for length and type-token ratio (TTR) for diversity.

## D.2 Benchmark Descriptive Statistics

Figure 6 illustrates the comparison of Length distribution and lexical diversity across four domains. It compares our constructed datasets with original model-generated texts and human texts.

First, kernel density estimation (KDE) results for Length Distribution demonstrate that SOTA MGT effectively covers the HWT length range while maintaining a more concentrated distribution. As illustrated, the probability densities for both GPT-5.1 and Grok-4.1 fall primarily between 50 and 400 words, consistent with the HWT span. Notably, in contrast to baseline models which exhibit irregularities such as the abnormal spike at approximately 300 words in ChatGPT or the bimodal distribution in Llama-2, the SOTA models display smoother, unimodal curves. This high degree of distributional alignment effectively diminishes length as a trivial discriminatory feature. It ensures that models accurately simulate human writing patterns across diverse domains, from the conciseness of academic abstracts to the extensiveness of creative writing, thereby guaranteeing the fairness of the detection task.

Second, regarding lexical diversity, type-token ratio (TTR) statistics indicate that current SOTA MGT surpasses the human baseline in vocabulary richness. Observing the TTR bar charts, models such as GPT-5.1 and Claude-Sonnet generally

show higher TTR values than HWT in domains like ArXiv and XSum. Particularly in the Writing Prompts domain, the models do not exhibit the monotony common in traditional generation. Their lexical diversity metrics significantly outperform early models like PaLM or Claude-instant in multiple samples. This data characteristic suggests that advanced models utilize broader vocabulary distributions during generation, rendering simple detection methods based on vocabulary repetition ineffective on this dataset. Furthermore, the slightly higher lexical diversity of generated text does not imply a decline in data quality or semantic divergence. Conversely, this characteristic positively reflects that SOTA models, under generation constraints, have successfully avoided common mode collapse and repetitive degeneration. This indicates that the generated text possesses complex syntactic structures and rich word choices, accurately simulating high-stealth, high-deception MGT in the current environment.

## D.3 Baseline Details

To establish a comprehensive benchmark, we compare ProSSD against a diverse set of state-of-the-art methods, including those published within the last three months, covering both zero-shot detection and training-based approaches. The specific introductions and experimental configurations for each baseline are as follows:

1621	• <b>Log-Likelihood Log-Rank Ratio (LRR)</b> (Su et al., 2023): This method utilizes log-rank information to distinguish between human and machine-generated texts, based on the hypothesis that MGT exhibits specific statistical properties in rank distribution. In our experiments, we employ GPT-Neo-2.7B as the base scoring model.	We directly evaluate the open-source RoBERTa-base-openai-detector checkpoint as a strong baseline.	1670 1671 1672
1622			
1623			
1624			
1625			
1626			
1627			
1628			
1629	• <b>Normalized Perturbed Log-Rank (NPR)</b> (Su et al., 2023): As an enhanced version of LRR, NPR introduces text perturbations to quantify the sensitivity of log-rank scores. Although computationally more expensive, it achieves higher detection precision. We configure GPT-Neo-2.7B as the scoring model and use T5-small to generate perturbed samples, setting 100 perturbations per sample.		1673 1674 1675 1676 1677 1678 1679 1680 1681 1682
1630			
1631			
1632			
1633			
1634			
1635			
1636			
1637			
1638	• <b>DetectGPT</b> (Mitchell et al., 2023): This method operates on the hypothesis that text generated by LLMs tends to reside in the negative curvature regions of the model’s log-probability function. It performs detection by comparing the log-probability differences between the original and perturbed texts. We use GPT-Neo-2.7B as the base model and T5-small as the mask-filling model to generate perturbations, with 100 perturbations per sample.		1683 1684 1685 1686 1687 1688 1689 1690 1691 1692 1693 1694 1695
1639			
1640			
1641			
1642			
1643			
1644			
1645			
1646			
1647			
1648			
1649	• <b>Fast-DetectGPT</b> (Bao et al., 2024): This method utilizes conditional probability curvature and an efficient sampling strategy to replace the perturbation steps in DetectGPT, maintaining high precision while improving inference speed. Following the standard settings of the original paper, we use GPT-Neo-2.7B as the scoring model and GPT-J-6B as the reference model.		1696 1697 1698 1699 1700 1701 1702 1703 1704 1705 1706 1707
1650			
1651			
1652			
1653			
1654			
1655			
1656			
1657			
1658	• <b>Binoculars</b> (Hans et al., 2024): This is a low false-positive rate zero-shot detection method that identifies AI text by contrasting the perplexity of an "Observer" model with the cross-perplexity of a "Performer" model. In our configuration, we select Falcon-7b and Falcon-7b-instruct as the Observer and Performer models, respectively.		1708 1709 1710 1711 1712 1713 1714 1715 1716 1717
1659			
1660			
1661			
1662			
1663			
1664			
1665			
1666	• <b>RoBERTa-base</b> (Solaiman et al., 2019): As a classic supervised baseline, RoBERTa identifies machine text by fine-tuning a binary classifier on specific datasets.		
1667			
1668			
1669			

the OOD setting, the model constructs feature directions based on source domain data.

## E More Results

### E.1 Analysis of Supplementary OOD Experiments

**Cross model robustness.** As shown in Table 5, we evaluated the out of distribution performance of the models to examine their generalization capabilities when facing unknown data distributions that differ from the training set. Compared to the in domain results in Table 1 of the main text, although ReprGuard and DetectAnyLLM perform excellently in the source setting, they exhibit significant instability in out of distribution tests. Specifically, when the detector is trained on data generated by Claude instant and directly used to detect text generated by GPT 5.1 (right side of Table 5), the F1 score of ReprGuard drops sharply to 40.23%. In contrast, ProSSD maintains an F1 score of 88.44% and an AUROC of 96.22% under the same setting. Overall, whether trained on Google PaLM or Claude instant, ProSSD maintains an AUROC above 94% on all four target test models (Claude S 4, Gemini 3 F, GPT 5.1, Grok 4.1), demonstrating the best robustness.

**Cross domain adaptability.** Cross domain evaluation examines the ability of the model to handle huge differences in vocabulary and semantic distributions. The task involves training or constructing distributions on dataset A and testing on dataset B. Tables 6 and 7 report the relevant results. As shown in Table 6, although DetectAnyLLM and ReprGuard score slightly higher than ProSSD in individual specific transfer scenarios (such as Writing Prompts  $\rightarrow$  Arxiv and Arxiv  $\rightarrow$  XSum), the gap is minimal. More importantly, judging from the overall scores in Table 6 and Table 7, ProSSD exhibits the lowest performance variance, powerfully proving its strong generalization ability. Combining the above out of distribution analyses, the core advantage of ProSSD lies not in defeating all baselines on every single metric, but in its ability to provide a reliable performance lower bound under different training sources. This characteristic makes it more practically valuable when facing complex and changeable unknown data in the real world.

### E.2 Supplementary Evaluation on Adversarial Attacks

To comprehensively assess the security of the detectors in realistic adversarial environments, we evaluated their performance across four distinct attack scenarios: paraphrase, perturbation, direct prompt, and data mixing. Table 8 presents the comparative results. In general, ProSSD demonstrates the most consistent and superior defense capabilities, achieving the best performance in three out of four scenarios.

As shown in Table 8, while ReprGuard maintains strong competitiveness, particularly achieving the highest AUROC of 99.57% in the perturbation setting, ProSSD exhibits better stability across semantic alterations. Specifically, in the paraphrase scenario, which involves substantial rewriting and semantic restructuring, classic methods like DetectGPT and RoBERTa suffer significant performance degradation, with DetectGPT dropping to an F1 of 22.65%. Although ReprGuard performs well, ProSSD further improves the AUROC to 97.83% and F1 to 93.36%. This suggests that our method, by modeling the joint semantic structural distribution, successfully captures the deep invariant features of machine generation that survive surface level rewriting.

Furthermore, in the data mixing scenario, where machine text is interleaved with human text, ProSSD achieves a dominant AUROC of 99.13%, outperforming the second best method ReprGuard by roughly 1 percentage point and significantly surpassing LRR and NPR. It is worth noting that methods heavily reliant on specific probability curvature or noise sensitivity (such as AdaDetectGPT) exhibit extreme volatility in the perturbation setting, whereas ProSSD maintains a high AUROC of 99.21%. Similar to the findings in the out of distribution analysis, the core advantage of ProSSD lies in its extremely low variance across different attack types. Whether facing simple character level perturbations or complex semantic paraphrasing, ProSSD provides a robust detection boundary, proving its reliability as a secure defense mechanism in dynamic adversarial contexts.

### E.3 Supplementary Explanation of Ablation Studies

To rigorously verify the effectiveness of the proposed framework and the contribution of each core component, we designed three groups of

Table 5: Cross-model generalization results with different training sources. The table evaluates detection performance when the training data is generated by Google-PaLM (Left) and Claude-instant (Right). The detectors are tested on four unseen target LLMs to assess robustness across different generator architectures. Results are reported as AUROC(%) and F1(%). Best and second-best results are highlighted in **bold** and underlined.

Method	Train on Google-PaLM								Train on Claude-instant							
	GPT-5.1		Claude-S-4		Gemini-3-F		Grok-4.1		GPT-5.1		Claude-S-4		Gemini-3-F		Grok-4.1	
	AUR.	F1	AUR.	F1	AUR.	F1	AUR.	F1	AUR.	F1	AUR.	F1	AUR.	F1	AUR.	F1
LRR	68.02	58.08	79.62	71.13	63.30	59.10	62.18	49.50	68.02	58.08	79.62	71.13	63.30	59.10	62.18	49.50
Fast-DetectGPT	67.07	58.59	78.78	72.73	76.13	68.44	52.91	42.29	67.07	58.59	78.78	72.73	76.13	68.44	52.91	42.29
ImBD	<u>90.63</u>	<u>83.26</u>	<u>93.72</u>	84.97	<u>94.90</u>	<u>88.66</u>	<u>85.70</u>	77.83	<u>92.67</u>	<u>84.74</u>	<u>96.62</u>	<u>89.94</u>	<u>95.65</u>	<u>90.29</u>	<u>87.83</u>	<u>81.58</u>
AdaDetectGPT	54.40	55.33	63.59	63.33	61.22	57.07	41.37	66.67	61.47	54.74	73.25	69.41	68.95	59.10	47.08	21.61
DetectAnyLLM	76.96	82.48	81.02	<u>85.14</u>	79.41	85.33	78.34	<u>84.14</u>	83.42	74.66	92.04	84.89	89.86	81.82	80.04	73.50
RepreGuard	69.88	64.32	81.69	75.97	73.04	67.24	69.40	66.63	58.56	40.23	68.66	58.41	62.80	46.15	62.02	39.91
ProSSD	<b>94.87</b>	<b>88.14</b>	<b>97.07</b>	<b>90.31</b>	<b>98.34</b>	<b>93.98</b>	<b>96.21</b>	<b>89.22</b>	<b>96.22</b>	<b>88.44</b>	<b>98.53</b>	<b>93.74</b>	<b>97.89</b>	<b>92.69</b>	<b>97.65</b>	<b>91.60</b>

Table 6: Cross-domain detection performance using different training sources. We report the AUROC (%) and F1 (%) scores when the detector is trained on the Writing-prompt(Left) and Arxiv (Right) domains and evaluated on other unseen text domains. The best and second-best results in each column are highlighted in **bold** and underlined, respectively.

Method	Train on Writing-prompt						Train on Arxiv					
	Arxiv		XSum		Review		Writing		XSum		Review	
	AUR.	F1	AUR.	F1	AUR.	F1	AUR.	F1	AUR.	F1	AUR.	F1
LRR	73.90	67.71	70.14	63.49	84.01	75.11	48.20	14.90	70.14	63.49	84.01	75.11
Fast-DetectGPT	74.81	64.93	57.65	50.42	82.92	75.56	58.88	47.29	57.65	50.42	82.92	75.56
ImBD	97.55	92.42	84.52	77.02	<u>99.36</u>	<u>96.52</u>	90.54	82.18	94.24	87.63	96.98	91.22
AdaDetectGPT	46.40	1.19	38.87	0.00	47.31	66.67	59.73	51.70	62.46	59.31	83.78	75.14
DetectAnyLLM	<b>99.91</b>	<b>98.59</b>	90.62	83.23	97.99	93.72	92.68	86.19	<b>100.00</b>	<b>99.95</b>	97.28	92.09
RepreGuard	<u>99.07</u>	<u>95.53</u>	<u>94.07</u>	<u>86.36</u>	99.20	96.37	<b>99.83</b>	<b>98.75</b>	<u>99.95</u>	<u>99.10</u>	<b>99.87</b>	<b>98.24</b>
ProSSD	97.83	94.03	<b>99.54</b>	<b>97.60</b>	<b>99.92</b>	<b>99.00</b>	<u>97.27</u>	<u>92.34</u>	97.93	93.66	<u>98.85</u>	<u>95.20</u>

comparative experiments as shown in Table 9. First, in the robustness of semantic representations group, we aim to prove that this method does not rely on the semantic embeddings of a specific encoder, but rather is based on the inherent statistical differences between human and machine text. We replaced the default RoBERTa-large base with modern large language model embedding Qwen-3-0.6B-Embedding (Zhang et al., 2025) and randomly initialized embeddings. This setting verifies whether our supervised subspace learning can effectively extract discriminative features in different semantic spaces. Second, the impact of subspace projection group is specifically used to verify the core hypothesis in Section 3.2. We replaced supervised subspace learning with unsupervised PCA, random projection, and a baseline without projection. This comparison aims to prove that extracting label correlated variance is superior to purely maximizing global variance (such as PCA) or preserving random geometric structures. Finally, in the ablation on detection strategies, we explored the

necessity of distribution modeling. We examined the results after removing Wasserstein weights to verify the hypothesis that different syntactic structures contribute unequally to detection; we also evaluated the one class setting that only constructs the human distribution ( $P_H$ ) without utilizing the machine distribution ( $P_M$ ), thereby testing the importance of modeling both HWT and MGT distributions simultaneously.

Table 9 reports the quantitative results on four datasets. To ensure the reliability of statistical results, all experiments were independently repeated 5 times, and we report the mean and standard deviation. Significance tests (t test) confirmed that our method achieved statistically significantly better performance than the comparative ablation experiments on the vast majority of metrics ( $p < 0.05$ ).

**Effectiveness of supervised subspace learning.** The experimental results strongly support our hypothesis in Section 3.2. As shown in the impact of subspace projection part, replacing supervised pro-

Table 7: Cross-domain detection performance using different training sources. We report the AUROC (%) and F1 (%) scores when the detector is trained on the XSum(Left) and Review(Right) domains and evaluated on other unseen text domains. The best and second-best results in each column are highlighted in **bold** and underlined, respectively.

Method	Train on XSum						Train on Review					
	Arxiv		Writing		Review		Arxiv		Writing		XSum	
	AUR.	F1	AUR.	F1	AUR.	F1	AUR.	F1	AUR.	F1	AUR.	F1
LRR	73.90	67.71	48.20	14.90	84.01	75.11	73.90	67.71	48.20	14.90	70.14	63.49
Fast-DetectGPT	74.81	64.93	58.88	47.29	82.92	75.56	74.81	64.93	58.88	47.29	57.65	50.42
ImBD	<u>98.98</u>	<u>94.98</u>	<b>86.00</b>	<u>77.14</u>	96.66	90.38	97.26	91.26	94.73	87.52	90.90	82.90
AdaDetectGPT	67.23	58.35	55.27	39.20	79.36	72.99	60.77	49.58	52.94	31.95	50.04	12.66
DetectAnyLLM	98.79	<u>95.18</u>	79.39	73.27	<u>97.23</u>	<u>91.35</u>	<b>99.77</b>	<b>98.61</b>	<u>99.74</u>	<u>98.10</u>	<b>99.80</b>	<b>99.50</b>
RepreGuard	97.34	92.58	69.32	66.12	96.27	90.18	<u>99.73</u>	<u>98.29</u>	<b>99.87</b>	<b>98.85</b>	99.54	98.10
ProSSD	<b>99.83</b>	<b>98.55</b>	<b>96.39</b>	<b>91.05</b>	<b>99.25</b>	<b>95.89</b>	98.98	95.51	99.52	97.09	<u>99.72</u>	<u>99.25</u>

Table 8: Adversarial robustness comparison of all detectors. We report detection performance on the original data(Direct Prompt) and against three adversarial attack scenarios(Paraphrase, Perturbation, and Data Mixing). The best results are marked in **bold** and second-best in underlined.

Method	Paraphrase		Perturbation		Direct Prompt		Data Mixing	
	AUR.	F1	AUR.	F1	AUR.	F1	AUR.	F1
LRR	70.39	60.11	97.97	94.66	86.41	78.78	82.78	72.61
NPR	72.96	66.63	98.60	95.69	80.26	72.68	73.42	63.16
RoBERTa	75.73	65.31	66.82	52.95	78.81	64.99	78.78	68.08
DetectGPT	47.58	22.65	88.11	79.40	53.78	41.04	32.63	0.00
Binoculars	81.78	73.11	97.94	94.22	94.93	90.07	93.89	88.43
Fast-DetectGPT	76.33	68.50	93.44	86.72	86.24	79.95	86.52	78.41
ImBD	88.25	79.29	97.48	90.41	93.61	84.48	90.93	80.28
AdaDetectGPT	75.40	64.06	39.20	0.79	84.72	78.44	84.72	76.58
DetectAnyLLM	86.00	76.67	87.38	76.55	86.80	76.40	88.89	76.75
RepreGuard	<u>95.51</u>	<u>89.01</u>	<b>99.57</b>	<b>97.86</b>	<u>97.31</u>	<u>92.52</u>	<u>98.19</u>	<u>94.99</u>
ProSSD	<b>97.83</b>	<b>93.36</b>	<u>99.21</u>	<u>96.35</u>	<b>98.25</b>	<b>93.44</b>	<b>99.13</b>	<b>95.97</b>

jection with PCA ("w/o SSP (PCA)") leads to a significant performance drop, especially on the XSum dataset where F1 drops from 96.96% to 88.02%. This indicates that PCA, which blindly maximizes global variance, is highly likely to retain the semantic noise shared by HWT and MGT, thereby drowning out weak style signals.

**Importance of structured distribution modeling.** The "w/o Wasserstein" ablation experiment exhibits a significant increase in standard deviation (for example  $\pm 1.79$  on Writing) and a decrease in F1 score. This validates Theorem 1, that weighting local decisions according to theoretical discriminability (Wasserstein distance) can build a more robust detector. Furthermore, the "w/o Contrastive" (only constructing human distribution) ablation experiment performed the worst among all valid methods (for example Review F1 dropped to 80.08%). This confirms that MGT are not merely outliers of human text; they possess their own statistical reg-

ularities. Explicitly modeling the likelihood ratio  $P(x|MGT)/P(x|HWT)$  is more effective.

**Robustness across embedding models.** Although the original method uses RoBERTa, the comparative experiment with Qwen-3-0.6B-Embedding also achieved highly competitive results, even slightly outperforming RoBERTa on the Writing dataset (F1 97.27%). This proves that our method has high adaptability for different embedding models. Notably, the "random projection" comparative trial still achieved impressive performance (average F1 > 80%), even surpassing many existing mainstream detection baselines. This indicates that even in the absence of true semantic information, detection with considerable precision can be achieved solely relying on the structured statistical differences captured by our distribution modeling. When high quality semantic embeddings are introduced, performance is further improved by approximately 10% to 15%, illustrating

Table 9: Ablation studies across different domains. We report AUC and F1 scores for ArXiv and Writing, and F1 scores for XSum and Review. Results represent the mean and standard deviation over 5 independent runs. The statistical significance of the performance drop compared to ProSSD is measured by a t-test: \*  $p < 0.05$ , †  $p < 0.01$ , ‡  $p < 0.001$ .

Method	ArXiv		Writing		XSum	Review
	AUC	F1	AUC	F1	F1	F1
<b>ProSSD (RoBERTa + SSP)</b>	<b>99.55±0.26</b>	<b>99.25±0.18</b>	<b>99.34±0.03</b>	96.45±0.14	<b>96.96±0.16</b>	<b>97.23±0.03</b>
<i>Robustness of Semantic Representations</i>						
Qwen-2.5-0.6B-Embed	96.22±0.12‡	91.38±0.34‡	99.16±0.13*	<b>97.27±0.25†</b>	88.93±0.42‡	97.18±0.16
Random Projection	83.50±1.81‡	76.29±2.19‡	85.53±3.37†	77.21±3.58‡	87.36±1.36‡	84.05±2.63‡
<i>Impact of Subspace Projection (SSP)</i>						
w/o SSP (PCA)	99.51±0.22	98.74±0.20†	97.58±0.09‡	92.04±0.19‡	88.02±0.29‡	92.25±0.25‡
w/o SSP (Rand)	92.05±1.52‡	87.97±3.03†	97.84±0.41†	93.23±0.51‡	89.61±0.88‡	93.02±0.54‡
w/o SSP (No-Proj)	95.71±0.83‡	94.60±0.75‡	98.16±0.32†	94.44±0.43‡	90.81±0.28‡	94.17±0.30‡
<i>Ablation on Detection Strategies</i>						
w/o Wasserstein (Uniform)	<b>99.74±0.20</b>	98.07±1.06	96.42±0.88†	90.70±1.79†	<u>91.46±1.56†</u>	92.15±1.06‡
w/o Contrastive (Human-Only)	97.39±0.18‡	93.02±0.30‡	95.73±0.19‡	88.71±0.35‡	86.27±0.57‡	80.08±0.53‡

that our method better combines deep semantic flow with syntactic statistical features.

#### E.4 Visualization Details

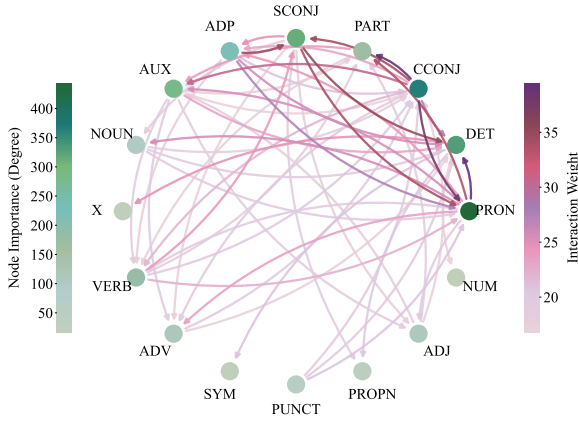
##### Visualization of discriminative meta-structures.

To visualize the relationships between structures, we construct a discriminative meta-structure topology as shown in Figure 7. In this graph, each node represents a part of speech (POS) category acting as a local syntactic anchor, while directed edges explicitly characterize the meta-structure transition patterns  $\pi_t = (pos_t, pos_{t+1})$  between adjacent words. The color intensity and thickness of the edges correspond to the Wasserstein discriminant weights ( $w_\pi$ ) derived in Section 3.3. Darker colors and thicker lines indicate that under this syntactic transition path, the conditional semantic distribution of machine generated text exhibits significant statistical deviation from that of human written text, thus being assigned a higher discriminant weight. Accordingly, the node size reflects its cumulative usage, characterizing the importance of that syntactic structure as a hub for differentiated semantic flow.

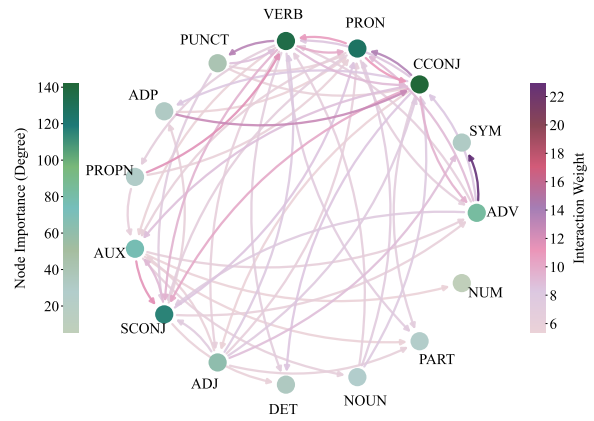
The meta-structure topology across different domains reveals that the structural information generated by machines is not static but demonstrates strong context dependency. As shown in Figure 7, discriminative hotspots present distinctly different distribution results. In formal domains with high syntactic complexity such as news or Wikipedia (Figure 7b), the discriminative network presents a dense interaction structure centered on subordinat-

ing conjunctions (SCONJ) and coordinating conjunctions (CCONJ). This suggests that although large models perform well in local fluency, they still struggle to reproduce the rigorous logic and coherence of human authors in deep syntactic structures when dealing with long complex sentences and logical clauses, leading to significantly elevated Wasserstein distances at these connectives. Conversely, in domains with strong subjectivity or informal contexts like reviews or social media (Figure 7c), the center of discrimination significantly shifts towards pronouns (PRON), adjectives (ADJ), and even interjections (INTJ). In this context, statistical differences mainly stem from microscopic variations in stance expression and emotional coloring. Human written text often exhibits extremely high variance and idiosyncrasy in subjective descriptions and first person narratives, whereas model generated text tends to fall into semantic collapse, manifesting as more mediocre and conservative semantic choices that fail to mimic diverse authentic human emotional states.

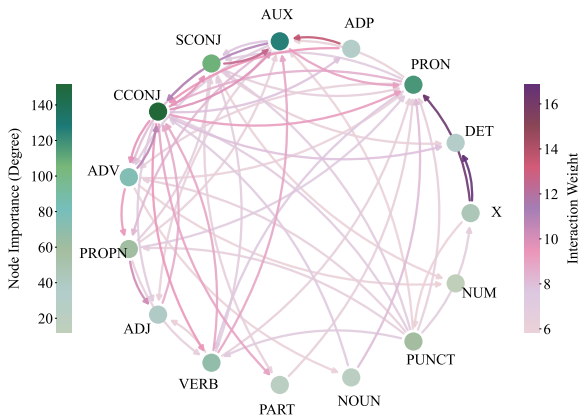
**Distribution of detection scores.** To intuitively evaluate the ability of the detector to distinguish between human and machine text, we visualized the score distribution of all samples in the test set (Figure 8). The green and red curves represent the probability density distributions of scores for human written text and large language model text, respectively. We calculated the overlap rate, defined as the intersection area of the two probability distribution functions, to quantify the degree of confusion. A lower overlap rate implies a greater



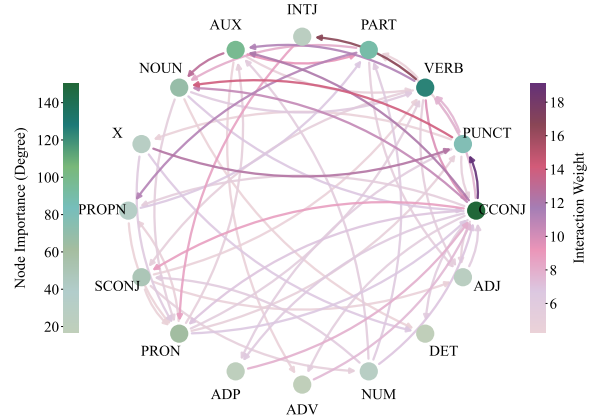
(a) Discriminative meta-structures in academic texts.



(b) Discriminative meta-structures in news summaries.



(c) Discriminative meta-structures in subjective reviews.



(d) Discriminative meta-structures in creative writing.

Figure 7: **Visualization of Discriminative meta-structure Topology across different domains.** Each node represents a POS tag, and edges denote structural transitions. The edge thickness and color intensity correspond to the Wasserstein discriminative weight ( $w_\pi$ ), highlighting where MGT deviates most from HWT.

1963 distance between the two classes of text in the fea-  
 1964 ture space and a clearer decision boundary.

1965 Comparing the baseline RepreGuard with our  
 1966 ProSSD reveals the significant advantage of  
 1967 ProSSD in distinctiveness. First, ProSSD demon-  
 1968 strates extremely low distribution overlap. While  
 1969 RepreGuard shows an overlap rate exceeding 14%  
 1970 on multiple models (such as 14.4% on GPT 5.1),  
 1971 posing a higher risk of misjudgment, ProSSD com-  
 1972 presses the overlap rate to an extremely low level  
 1973 (such as only 1.0% on Gemini 3 Flash), effectively  
 1974 partitioning the two types of text into distinct dis-  
 1975 tribution regions. Second, ProSSD establishes a  
 1976 wider safety margin in score values. Observing the  
 1977 horizontal axis, the scores of RepreGuard are con-  
 1978 centrated in a narrow interval from 0 to 3, whereas  
 1979 the score span of ProSSD significantly expands to  
 1980 between -80 and +20. This broad numerical differ-  
 1981 ence builds a sufficient safety buffer between hu-  
 1982 man and machine text, proving that this method can

1983 extract features with greater discriminative power  
 1984 and achieve more robust sample separation in the  
 1985 numerical space.

1986 **Word level discriminative analysis.** We con-  
 1987 ducted a sample analysis to visualize the discrimi-  
 1988 native contribution of each token via color coding,  
 1989 as shown in Figures 9. As indicated by the topol-  
 1990 ogy in Figure 7, while attending to text content,  
 1991 our method keenly captures interactions between  
 1992 part of speech structures. The nodes for pronouns  
 1993 (PRON), coordinating conjunctions (CCONJ), ad-  
 1994 positions (ADP), and determiners (DET), along  
 1995 with their connecting edges, exhibit the darkest col-  
 1996 ors, indicating that they are key to distinguishing  
 1997 between human and AI text.

1998 As shown in Figure 9b, this text is classified  
 1999 as LLM generated. For instance, when handling  
 2000 clause introductions in long complex sentences  
 2001 (such as "whose") and parallel structures (such as  
 2002 "and"), the semantic transitions surrounding these

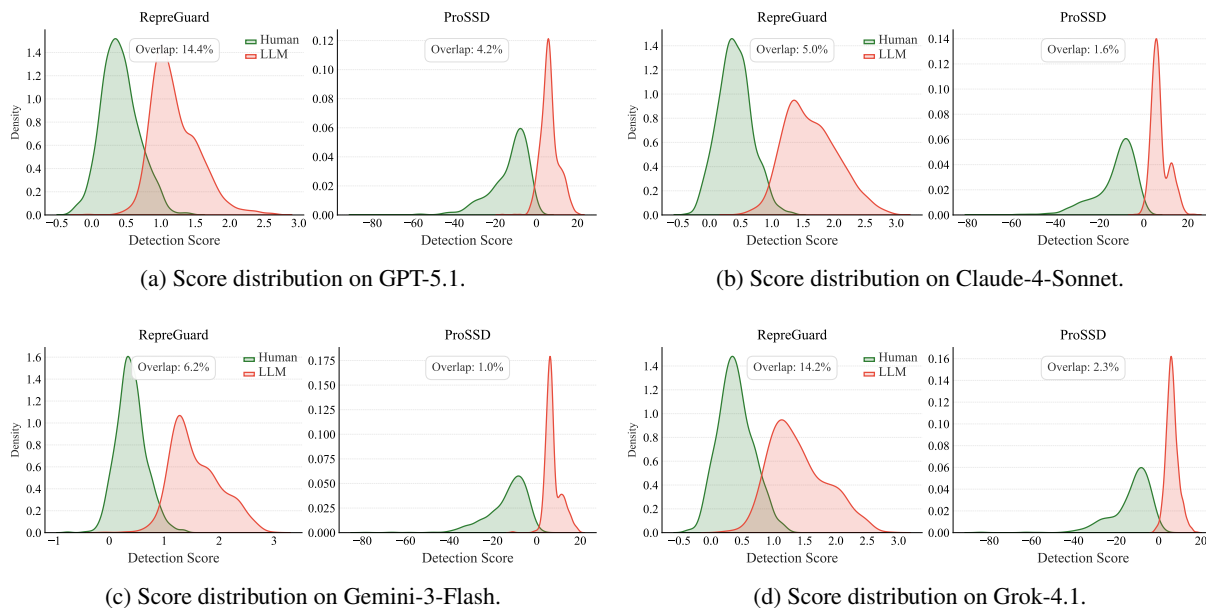


Figure 8: **Visualization of detection score distributions across four LLMs.** The green and red curves represent the probability density functions of human-written and machine-generated texts, respectively. We compare our method, **ProSSD**, with the state-of-the-art baseline, **RepreGuard**.

tokens are extremely smooth and lack variation. Their distribution biases towards the center of our constructed AI Gaussian distribution, exhibiting extremely low semantic variance.

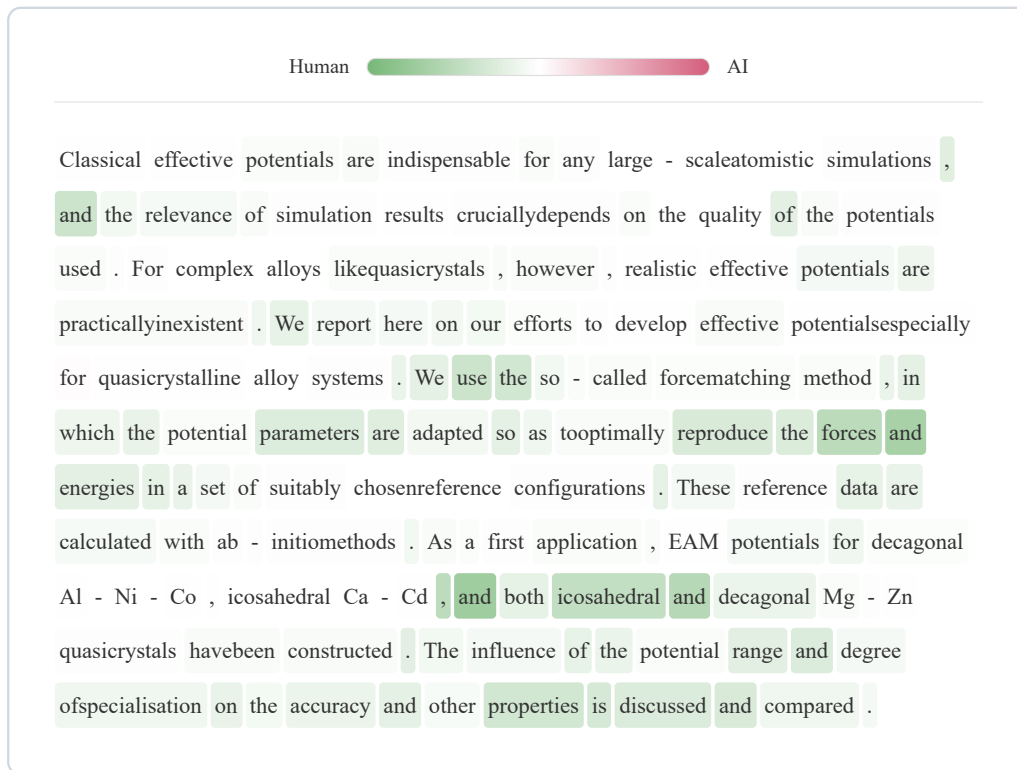
As shown in Figure 9a, this text is classified as human written. Consider the placement of "and", particularly the tight combination of CCONJ and DET in "and both". When describing complex chemical structures like icosahedral and decagonal quasicrystals, the human author employs the nested parallel structure "and both... and...". This usage exhibits strong specificity in the semantic vector space. Such representation lies closer to the human distribution  $\mu_H$  and further from the AI distribution. Consequently, our method classifies it as human text.

## E.5 Extended Evaluation on MIRAGE Benchmark

To further verify the generalization ability of the proposed method in complex real world scenarios, we extended the evaluation scope to the latest MIRAGE benchmark (Fu et al., 2025) proposed within three months. MIRAGE is a comprehensive evaluation framework covering 10 different corpora across 5 text domains, utilizing 17 advanced LLMs to construct diverse adversarial samples. The benchmark includes two settings: disjoint input generation (DIG) and shared input generation (SIG), aiming to test the performance of detectors

under different input output correspondences. To optimize the spatial layout of Table 10, we categorized the evaluation tasks into "Gen." (generation) and "Rev." (revision). Specifically, the "Rev." category is the weighted average result of the "polish" and "rewrite" tasks in the original benchmark, as both tasks involve modifications to human written text rather than generating content from scratch. The results of all baseline methods, including the previous state of the art methods DetectAnyLLM and RepreGuard, are partly cited directly from the original MIRAGE paper to ensure fairness and consistency of comparison.

Table 10 presents the quantitative comparison results between our method and existing baseline methods. In the "Gen." task, ProSSD demonstrates significant performance advantages, outperforming the strongest baseline DetectAnyLLM in all metrics under both DIG and SIG settings. Notably, in the MIRAGE DIG generation subset, ProSSD achieved a TPR@5% of 88.62%, realizing a significant improvement of approximately 10.9 percentage points compared to DetectAnyLLM (77.70%). This indicates that our method constructs a more robust discriminative boundary for completely machine generated content. In the "Rev." task, ProSSD remains highly competitive, achieving an AUROC score slightly higher than current state of the art methods, for example 92.70% versus 92.64% under the DIG setting. Although



(a) Human-written text sample.



(b) Machine-generated text sample.

Figure 9: **Visualization of token-level discriminative contribution.** Red shading indicates tokens that the model identifies as more characteristic of LLM generation, while green shading indicates tokens more typical of human writing. The top pane displays a human-written sample, and the bottom pane displays an AI-generated sample.

Table 10: Performance comparison (%) on MIRAGE-DIG and MIRAGE-SIG. Tasks are categorized into **Gen.** (Generation) and **Rev.** (Revision, weighted average of Polish & Rewrite). **ProSSD** achieves superior performance in generation tasks and competitive results in revision tasks.

Method	MIRAGE-DIG								MIRAGE-SIG							
	Gen.				Rev.				Gen.				Rev.			
	AUC	Acc	MCC	TPR	AUC	Acc	MCC	TPR	AUC	Acc	MCC	TPR	AUC	Acc	MCC	TPR
Likelihood	49.36	50.91	1.83	1.47	44.90	50.00	0.00	1.80	49.68	52.07	1.96	1.45	44.55	50.01	0.15	1.70
LogRank	49.92	51.28	2.60	2.20	43.64	50.00	0.00	1.62	50.08	51.83	1.82	1.86	43.41	50.00	0.00	1.63
Entropy	65.22	61.50	25.43	10.99	56.78	54.94	14.55	10.75	64.42	61.23	15.92	10.74	57.52	55.45	7.23	10.76
RoBERTa-Base	55.23	53.97	14.34	12.50	49.42	50.30	1.94	5.16	53.68	53.92	5.29	11.01	49.26	50.68	1.37	5.36
RoBERTa-Large	47.16	52.17	8.42	8.71	53.77	52.72	6.10	7.68	47.03	52.36	4.17	9.10	53.70	52.96	3.46	7.33
LRR	52.15	53.41	7.77	7.01	40.03	50.00	0.00	1.94	52.14	53.11	3.14	6.57	40.25	50.00	0.00	2.05
NPR	61.20	61.40	26.04	1.91	48.85	52.83	8.60	2.71	60.88	61.70	15.71	1.85	49.01	52.39	4.72	2.33
DetectGPT	64.02	62.58	27.58	2.75	52.58	53.94	10.69	3.18	63.53	62.41	17.19	1.93	52.51	53.83	5.46	2.73
Fast-DetectGPT	77.68	72.34	46.28	43.10	55.83	54.99	11.50	11.04	77.06	71.93	20.78	42.00	56.00	55.55	5.65	11.65
ImBD	85.97	77.38	54.97	40.65	78.55	71.07	42.17	28.35	86.12	77.91	55.99	41.83	78.18	70.55	41.85	29.49
DetectAnyLLM	<u>95.25</u>	<u>89.88</u>	<u>79.75</u>	<u>77.70</u>	<u>92.64</u>	<b>87.18</b>	<b>74.66</b>	<b>77.67</b>	<u>95.26</u>	<u>90.59</u>	<u>81.19</u>	<u>77.22</u>	<u>92.34</u>	<b>86.90</b>	<b>73.99</b>	<b>76.73</b>
<b>ProSSD (Ours)</b>	<b>97.23</b>	<b>91.79</b>	<b>83.99</b>	<b>88.62</b>	<b>92.70</b>	84.84	69.89	71.16	<b>97.37</b>	<b>91.88</b>	<b>83.81</b>	<b>88.03</b>	<b>92.52</b>	84.99	70.06	69.94

DetectAnyLLM maintains a slight advantage in accuracy and MCC metrics for revision tasks, the excellent AUROC performance of ProSSD indicates that our method can effectively rank machine revised text, even if specific decision thresholds require further calibration. Overall, these experimental results confirm that our method can effectively capture critical semantic structural distribution differences, a capability that is particularly prominent in generation type tasks.

## F Method Details and Code Statement

### F.1 Algorithm and Experimental Settings

We formally describe our proposed ProSSD discrimination framework in Algorithm 1. The method first constructs a supervised subspace projection matrix. It then performs semantic structural distribution modeling. After obtaining the distribution sets for HWT and MGT, it scores and detects new input text via a detection function based on modified Mahalanobis distance.

The hyperparameter settings of our method include the following aspects. Regarding part of speech tagging, we employ the latest `en_core_web_sm_3.8.0` model released by the `spacy` library. This model covers over 100 part of speech tags. Regarding semantic embeddings, all experiments excluding ablation studies utilize the RoBERTa-large model. We use the 1024 dimensional vector from the final layer output as the word level embedding. All experiments and calculations were completed on a CPU and a single NVIDIA A800 SXM4 80G GPU. For the hyper-

parameters of our method, we set the projection dimension  $k = 4$  in both comparative and ablation experiments. The window size for constructing meta semantics and meta-structures is set to 2, with a stride of 1. Regarding training data volume, our comparisons utilize 1400 HWT and MGT samples extracted from the training set. For general experiments, the random seed is set to 42. In ablation studies, we conduct 5 independent runs with random seeds set to 42, 43, 44, 45, and 46.

In the sensitivity analysis, we primarily conducted two experiments. The first experiment discusses the optimal value of the projection dimension  $k$ . We tested  $k$  in the range of  $\{1, 2, 3, 4, 5, 6, 8, 10, 12, 16, 32\}$ . In this case, the semantic embedding model is controlled as RoBERTa-large, with 1400 HWT and MGT training samples. The second experiment discusses the required training data size. The data volume  $N$  ranges from  $\{50, 100, 200, 400, 600, 800, 1000, 1200, 1400, 1600, 1800\}$  HWT and MGT samples. In this experiment, the projection dimension is  $k = 4$ , and the semantic embedding model is RoBERTa-large.

### F.2 Code Description in Submission Materials

We have submitted the code for the modeling and detection steps of our method. We also included the training and test data for our newly constructed datasets under domain and model testing configurations.

---

**Algorithm 1** Semantic-Structural Distribution Modeling via Subspace Projection
 

---

**Require:** Human corpus  $\mathcal{D}_H$ , AI corpus  $\mathcal{D}_{AI}$ .

**Require:** Embedding Model  $\mathcal{F}$ , Projection dim  $k$ , Window size  $w$ .

**Ensure:** Projection Matrix  $\mathbf{P}$ , Distribution Library  $\mathbb{L}$ .

```

1: procedure LEARNSUBSPACE( $\mathcal{D}_H, \mathcal{D}_{AI}, k$ )
2:   Extract embeddings  $\mathbf{E}_H, \mathbf{E}_{AI}$  using  $\mathcal{F}$ .
3:    $\mathbf{E} \leftarrow [\mathbf{E}_H; \mathbf{E}_{AI}]$ ,  $\mathbf{y} \leftarrow [0 \dots 0; 1 \dots 1]$ 
4:   Initialize  $\mathbf{P} \leftarrow \emptyset$ ,  $\mathbf{E}_{res} \leftarrow \text{Center}(\mathbf{E})$ 
5:   for  $j = 1 \rightarrow k$  do ▷ Sec 3.2: Supervised Projection
6:      $\mathbf{p}_j^* \leftarrow \operatorname{argmax}_{\mathbf{p}: \|\mathbf{p}\|=1} (\operatorname{Cov}(\mathbf{E}_{res}\mathbf{p}, \mathbf{y}))^2$ 
7:      $\mathbf{E}_{res} \leftarrow \mathbf{E}_{res} - (\mathbf{E}_{res}\mathbf{p}_j^*)(\mathbf{p}_j^*)^T$  ▷ Sec 3.2: Feature Deflation
8:     Append  $\mathbf{p}_j^*$  to  $\mathbf{P}$ 
9:   end for return  $\mathbf{P}$ 
10: end procedure

11: procedure BUILDDISTRIBUTIONLIBRARY( $\mathcal{D}_H, \mathcal{D}_{AI}, \mathbf{P}$ )
12:    $\mathbb{L} \leftarrow \emptyset$ 
13:   for  $S \in \{\mathcal{D}_H, \mathcal{D}_{AI}\}$  do
14:      $\mathbf{V} \leftarrow \mathcal{F}(S)\mathbf{P}$  ▷ Sec 3.2: Project to  $k$ -dim semantic space
15:     Construct Meta-Semantics  $\mathcal{X} = \{(\mathbf{x}_t, \pi_t)\}_{t=1}^T$ 
16:     where  $\mathbf{x}_t = [\mathbf{v}_t; \mathbf{v}_{t+1}]$  and  $\pi_t = (\text{post}_t, \text{post}_{t+1})$ 
17:   end for
18:   for each unique structure  $\pi \in \Pi$  do
19:     Estimate  $\mathcal{N}_H^\pi(\boldsymbol{\mu}_H, \boldsymbol{\Sigma}_H)$  and  $\mathcal{N}_M^\pi(\boldsymbol{\mu}_M, \boldsymbol{\Sigma}_M)$  via MLE
20:     Calculate weight  $w_\pi \leftarrow W_2(\mathcal{N}_H^\pi, \mathcal{N}_M^\pi)$  ▷ 3.3: Wasserstein Dist
21:      $\mathbb{L}[\pi] \leftarrow \{\mathcal{N}_H^\pi, \mathcal{N}_M^\pi, w_\pi\}$ 
22:   end for return  $\mathbb{L}$ 
23: end procedure

24: function DETECT(Text  $T, \mathbf{P}, \mathbb{L}$ )
25:    $\mathbf{V} \leftarrow \mathcal{F}(T)\mathbf{P}$ 
26:   Parse structure sequence  $\pi_{1\dots M}$  and features  $\mathbf{x}_{1\dots M}$ 
27:    $\text{ScoreSum} \leftarrow 0$ ,  $\text{WeightSum} \leftarrow 0$ 
28:   for  $t = 1 \rightarrow M$  do
29:     if  $\pi_t \in \mathbb{L}$  then
30:       Retrieve parameters  $(\boldsymbol{\mu}_H, \boldsymbol{\Sigma}_H), (\boldsymbol{\mu}_M, \boldsymbol{\Sigma}_M), w_\pi$  from  $\mathbb{L}$ 
31:       Calculate modified Mahalanobis distances  $\mathcal{M}_H, \mathcal{M}_M$ 
32:        $s_t \leftarrow \frac{1}{2}(\mathcal{M}_H(\mathbf{x}_t) - \mathcal{M}_M(\mathbf{x}_t))$  ▷ Sec 3.4: Ditection
33:        $\text{ScoreSum} \leftarrow \text{ScoreSum} + w_\pi \cdot s_t$ 
34:        $\text{WeightSum} \leftarrow \text{WeightSum} + w_\pi$ 
35:     end if
36:   end for return  $\mathcal{S}(T) \leftarrow \text{ScoreSum}/\text{WeightSum}$ 
37: end function

```

---