

Unraveling Interwoven Roles of Large Language Models in Authorship Privacy: Obfuscation, Mimicking, and Verification

Anonymous ACL submission

Abstract

Recent advancements in large language models (LLMs) have been fueled by large-scale training corpora drawn from diverse sources such as websites, news articles, and books. These datasets often contain *explicit* user information, such as person names, addresses, that LLMs may unintentionally reproduce in their generated outputs. Beyond such explicit content, LLMs can also leak identity-revealing cues through *implicit* signals such as distinctive writing styles, raising significant concerns about authorship privacy. There are three major automated tasks in authorship privacy, namely authorship obfuscation (AO), authorship mimicking (AM), and authorship verification (AV). Prior research has studied AO, AM, and AV independently. However, their interplays remain under-explored, which leaves a major research gap, especially in the era of LLMs, where they are profoundly shaping how we curate and share user-generated content, and the distinction between machine-generated and human-authored text is also increasingly blurred. This work then presents the first unified framework for analyzing the dynamic relationships among LLM-enabled AO, AM, and AV in the context of authorship privacy. We quantify how they interact with each other to transform human-authored text, examining effects at a single point in time and iteratively over time. We also examine the role of demographic metadata, such as gender, academic background, in modulating their performances, inter-task dynamics, and privacy risks. All source code will be publicly available.

1 Introduction

Recent advances in LLMs have been extraordinary, driven largely by the massive amounts of training data indiscriminately sourced from diverse online platforms such as websites, news outlets, and books (Brown, 2020; Le Scao et al., 2023; Touvron et al., 2023; Achiam et al., 2023). This training data of-

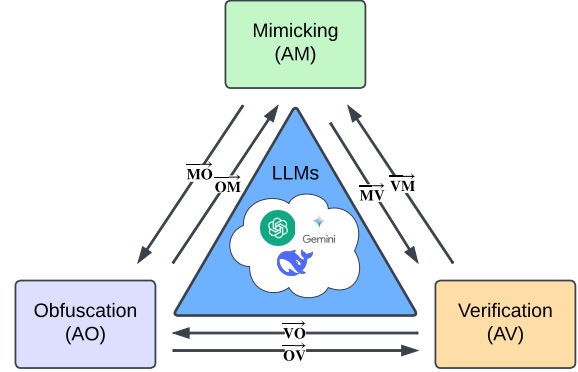


Figure 1: The interactive influence loop between LLMs, obfuscation, mimicking, and verification.

ten includes extensive writing contributions by the same authors, publicly shared across various platforms (Gao et al., 2020; Raffel et al., 2020). These sources frequently contain *explicit* user information, such as names, addresses, and phone numbers, which LLMs can inadvertently expose during their text generation process (Weidinger et al., 2021; Kim et al., 2024). Beyond explicit details, user identification can also be inferred from *implicit* information, such as their distinctive writing styles, that does not immediately give out the authors' identities. Research in human cognitive science and linguistics highlights that individual backgrounds significantly shape writing styles (Zheng et al., 2006; Cheng et al., 2023a; Deshpande et al., 2023; Xing et al., 2024; He et al., 2025), facilitating bidirectional inferences between implicit information (e.g., writing style) and explicit information (e.g., names, ages, or areas of expertise). Recent studies also reveal that text generated by LLMs can also capture human personality traits (Karra et al., 2022; Jiang et al., 2024a,b; Bang et al., 2024; An et al., 2024), and vice versa—i.e., explicit information about specific individuals or groups can be used by LLMs to produce personalized outputs or *mimic individuals' writing styles* (Chen et al., 2024; Salemi et al., 2024).

Although the authorship mimicking (AM) capabilities of LLMs—i.e., their ability to replicate an individual’s writing style, are impressive, this capability could also enable malicious activities, such as impersonating public figures to spread misinformation or commit fraud (Deshpande et al., 2023; Jiang et al., 2024a). For instance, a fraudster could fine-tune an LLM on publicly available texts authored or spoken by a target victim (e.g., social media posts, interviews) and prompt LLMs to generate spam emails or persuasive messages that pretend to be delivered by the victim (Salewski et al., 2023). Contrasting with AM, authorship obfuscation (AO) (Uchendu et al., 2024) aims to conceal an author’s identity by altering stylistic features of text while preserving its original meaning. By masking writing style before public dissemination (e.g., on social media), AO can help protect whistleblowers, such as writers or speakers, from potential anonymity exposure. In addition, authorship verification (AV) is the process of determining the author of a particular piece of writing. AV poses significant privacy risks by enabling the de-anonymization of individuals through their writing style, which can facilitate surveillance, behavioral profiling, and misuse without informed consent.

While AO, AM, and AV have each been studied in isolation, their interactions within a unified framework remain underexplored or limited to only specific pairwise formulations, such as AV and AO in the context of LLM-generated text (Uchendu et al., 2023). In addition, real-world scenarios often involve multiple rounds of text transformation, where content is repeatedly mimicked, obfuscated, and verified—either by different LLMs or within multi-turn dialogue settings where LLMs interact with one another (Duan et al., 2024). To address this gap, our study investigates three key scenarios in which LLMs play triple roles in authorship privacy (Fig. 1), analyzing their individual effects, interdependencies, and collaborative influences. Understanding the interplay among these capabilities is crucial for netizens in today’s LLM era, where users may rely on LLMs to obfuscate their writing style, while others may utilize LLMs to recover or attribute the original authorship. **Our contributions** include: (1) the first unified framework for studying the bidirectional effects among AO, AM, and AV; (2) empirical findings revealing distinct task-specific strengths of various commercial LLMs; (3) detailed analysis showing how demographic and metadata influence these interactions.

Our analysis shows that obfuscation tends to outperform mimicking in interactive settings, effectively disrupting authorial signals. However, mimicking can partially reverse obfuscation over successive cycles, gradually restoring aspects of the original writing style. Furthermore, models with stronger reasoning abilities (e.g., o3-mini, Deepseek) according to the benchmark ¹, excel at verification and concealing authorial traits but are less effective at faithfully replicating an author’s distinctive style.

2 Related Works

Beyond explicit metadata leakage such as names, social security numbers, LLMs’ generations can also reflect *implicit and private authorship signals* such as writing style, tone, or rhetorical structure, many of which are uniquely identifiable to specific individuals (Zheng et al., 2006; Cheng et al., 2023a; Deshpande et al., 2023; Xing et al., 2024; He et al., 2025). Thus, these models may memorize and reproduce identifiable features of authorship through their generated texts, so-called AM, introducing interesting interwoven relationships with LLM-enabled AO and AV.

Authorship Obfuscation (AO) hides the original author’s identity by altering stylistic cues without compromising semantic content. Recent methods include ALISON (Xing et al., 2024), which performs obfuscation by substituting stylistic sequences, and StyleRemix (Fisher et al., 2024), which utilizes AdapterMixup (Nguyen and Le, 2024) to train adapters for various stylistic dimensions and mix them. Different prompting-based approaches using LLMs have also been proposed (Hung et al., 2023; Pape et al., 2024).

Authorship Mimicking (AM) is the reverse of AO, aiming to generate text in the style of a specific author. LLMs excel in this task due to their few-shot and in-context learning capabilities, raising ethical concerns around impersonation, misinformation, and malicious use (Deshpande et al., 2023; Jiang et al., 2024a). Recent work has shown that LLMs can be fine-tuned or prompted to convincingly replicate individual writing styles from publicly available content (Salewski et al., 2023), making these capabilities intersect with privacy risks, such as when the LLMs leak memorized training examples (Carlini et al., 2023; Zhang et al., 2023).

Authorship Verification (AV) seeks to determine

¹<https://www.vals.ai/benchmarks/math500-05-09-2025>

or confirm whether a given text was written by a particular author, based on linguistic cues or stylistic fingerprints (Huang et al., 2025). With the advancement of model size scaling laws, LLMs can now perform AV in few-shot settings (Hung et al., 2023; Huang et al., 2024).

Interdependency of AO, AM, and AV. Prior research has largely treated AO, AM, and AV in isolation. However, their **pairwise interactions**, especially under the influence of LLMs, remain underexplored and foundational to many practical scenarios. For instance, for AO–AV, users obfuscate their writing style to protect identity, while adversaries re-identify authorship, creating a privacy-versus-attribution dynamic; for AM–AV, attackers mimic a target author’s style to deceive attribution models, challenging the robustness of verification systems; and for AO–AM, one can attempt to reconstruct authorial style from obfuscated text, testing the boundaries of stylistic recovery. Moreover, AO, AM, and AV can also form a closed loop in a **triplet-wise interaction**, reflecting how a text authorship changes under the influence of LLMs overtime. Our work is the first to address all pairwise and triplet-wise interdependencies of LLM-enabled AO, AM, and AV.

3 Research Questions and Formulation

3.1 Research Questions

We propose three **research questions (RQs)** to investigate both isolated and multi-level interdependencies among LLM-enabled authorship privacy tasks AO, AM, and AV, aiming to understand how individual and joint model behaviors influence the privacy and stylometry—i.e., writing styles, in complex authorship pipelines. Practical implications of our RQs are motivated in Appendix. A.1.

RQ1: How effectively can different LLMs perform AO, AM, AV *in isolation*, and which models are best suited for specific goals such as privacy preservation and stylistic imitation?

RQ2: How do LLM-enabled AO, AM, and AV influence one another to transform individuals’ stylometries when used *in conjunction at one point in time*, including their pairwise and triplet interactions?

RQ3: How do LLM-enabled AO, AM, and AV influence one another to transform individuals’ stylometries when used *in conjunction iteratively through time*?

To answer these RQs, we first formally de-

fine the evaluation of AO, AM, and AV of a target LLM $f(\cdot)$. For a given author a , let $\mathcal{D}_a = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ represent a set of a ’s original written documents paired with their corresponding author labels. M_a denotes the metadata associated with author a , such as name, field of study. We define $C_a = \{M_a, \mathcal{D}_a\}$ as the context available to $f(\cdot)$. For example, $f^{AO}(x|C_a)$ denotes the output obfuscation text of LLM $f(\cdot)$ on the input text x given the context C_a . $d(\cdot)$ is a stylometric distance defined on the two sets of input texts.

3.2 Isolation - No Interdependency

We begin by formulating AO, AM, and AV in isolation to evaluate the standalone performance of a specific LLM f . This setting is the most common in prior work, where researchers aim to quantify how well an individual LLM performs on specific authorship privacy tasks (Hung et al., 2023; Huang et al., 2024; Fisher et al., 2024; Pape et al., 2024; Salewski et al., 2023).

AO. To evaluate the effectiveness of AO on an input text x , we compute the distance $d(\cdot)$ between the original authentic texts and the obfuscated one (Eq. 1). The larger the distance, the more divergent the obfuscated text becomes from the original, suggesting more effective obfuscation.

$$AO = d(f^{AO}(x|C_a), \mathcal{D}_a) \quad (1)$$

AM. We evaluate the effectiveness of AM on an input text x by computing the distance between the original texts and the mimicked text (Eq. 2). The smaller stylometric distance $d(\cdot)$, the more similar the mimicked text is to the original, suggesting more effective mimicking.

$$AM = d(f^{AM}(x|C_a), \mathcal{D}_a) \quad (2)$$

AV. We evaluate the effectiveness of AV on an input text x by comparing its binary predictive verification—i.e., whether the text was written by author a or by someone else (Eq. 3). The higher verification accuracy, the more effective $f(\cdot)$ is at correctly identifying the author’s text.

$$AV = \mathbb{I}(f^{AV}(x|C_a) == a) \quad (3)$$

3.3 Pairwise Interdependency

Netizens are increasingly relying on LLMs to refine or disguise their writing through polishing, paraphrasing, or rephrasing, before sharing and publishing their content. These scenarios highlight a growing trend in which multiple LLMs are employed

within a single pipeline: one model generates or modifies text, while another evaluates or attributes authorship. Consequently, the input to these models is not always original author-written text but may already have undergone AI-driven transformation (Uchendu et al., 2023). To better understand these interactions, we conduct *pairwise interdependency* evaluations that measure their bidirectional relationships—i.e., *how one LLM’s capabilities influence the performance of others* (Fig. 1). To reflect the realistic scenario where the users prefer the best models for specific tasks, we designate a “judge” f_{judge} for each task, or the LLM that is selected based on its highest standalone performance in isolation (§ 3.2), for this evaluation.

Influence of Obfuscation. We factorize the influence of AO in the authorship pipeline into (1) how AO influences AM (\overrightarrow{OM}) and (2) how AO influences AV (\overrightarrow{OV}). For \overrightarrow{OM} , we first generate the obfuscated versions of an input text x , denoted x_{obf} , using various LLMs. Each of the obfuscated texts then serve as an input for the mimicking “judge”—a “ground-truth” LLM with the highest AM performance in isolation (§ 3.2), which attempts to reconstruct the original style of input x (Eq. 4). We compare the mimicked outputs to the original, authentic texts. The greater their stylistic divergence is, the more effective the obfuscated input, and hence the more influential the corresponding AO, and vice versa:

$$\overrightarrow{OM} = d(f_{judge}^{AM}(x|x_{obf}), \mathcal{D}_a) \quad (4)$$

For \overrightarrow{OV} , we pass the obfuscated texts x_{obf} to a verification “judge”. We compute verification accuracy on the original input x given the obfuscated texts (Eq. 5). The lower the accuracy, the more effective the obfuscation is; otherwise, it suggests the author’s style remains identifiable. This evaluation provides a practical measure of AO by testing whether others can still attribute the distorted writing to its original author. Such insights are particularly valuable in privacy-sensitive settings—e.g., anonymous investigative journalism or whistleblowing—where safeguarding the author’s identity is paramount:

$$\overrightarrow{OV} = \mathbb{I}(f_{judge}^{AV}(x|x_{obf}) == a) \quad (5)$$

Influence of Mimicking. We factorize the influence of AM in the authorship pipeline into (1) how AM influences AO (\overrightarrow{MO}) and (2) how AM influences AV (\overrightarrow{MV}). For (\overrightarrow{MO}), we first generate mimicking versions of the input text x , denoted

as x_{mimic} , using various LLMs. These mimicked texts then serve as the reference inputs for the obfuscation “judge”. Then, we compare the resulting obfuscated outputs to the original, authentic texts (Eq. 6). Obfuscation style significantly diverging from the originals indicates that the mimicking was effective in replicating the author’s writing style, and vice versa:

$$\overrightarrow{MO} = d(f_{judge}^{AO}(x|x_{mimic}), \mathcal{D}_a) \quad (6)$$

For \overrightarrow{MV} , we feed x_{mimic} into a verification “judge”. We calculate the verification accuracy of the predictive author with x ’s original author a (Eq. 7). A high verification accuracy indicates that the mimicked text effectively replicates the original author’s writing style, whereas a low accuracy suggests poor stylistic imitation:

$$\overrightarrow{MV} = \mathbb{I}(f_{judge}^{AV}(x|x_{mimic}) == a) \quad (7)$$

Influence of Verification. We factorize the influence of AV in the authorship pipeline into (1) how AV influences AO (\overrightarrow{VO}) and (2) how AV influences AM (\overrightarrow{VM}). In other words, AV acts as a filtering process to select only the texts verified as being authored by a as the input *contexts* for AO and AM. Intuitively, AV decides how pure or contaminated C_a is. To do this, we randomly sample n *noisy texts* or documents written by authors different from a , supposedly these are imposter samples. In both settings, we assess AV performance under two conditions: (1) *perfect C_a* : where all input context are genuine samples from the target author, and (2) *noisy \overline{C}_a* : where we introduce imposter samples from other authors that the model nonetheless classifies as the target author. Persistent positive classification of these imposter texts indicates weaker verification robustness. We then compute the distance of mimicking and obfuscation texts on the original input x , with the ground truth samples are all genuine and noisy (Eq. 8, Eq. 9).

$$\overrightarrow{VO} = d(f_{judge}^{AO}(x|C_a), f_{judge}^{AO}(x|\overline{C}_a)) \quad (8)$$

$$\overrightarrow{VM} = d(f_{judge}^{AM}(x|C_a), f_{judge}^{AM}(x|\overline{C}_a)) \quad (9)$$

3.4 Triplet-wise Interdependency

While previous evaluations identify which models excel at individual tasks and how they are pairwise-interdependent, this section investigates *the authorship pipeline cycle as a whole* (Fig. 1)—i.e., how AO and AM alter verification accuracy and the linguistic distribution of original human texts. By orchestrating multiple LLMs, each deployed

Models	AO		AM		AV
	PPL (\uparrow)	SIM (\downarrow)	PPL (\downarrow)	SIM (\uparrow)	Acc (\uparrow)
<i>4o-mini</i>	0.72	0.12	0.65	0.13	0.45
<i>o3-mini</i>	2.71	0.10	1.57	0.11	0.89
<i>deepseek</i>	1.08	0.11	1.86	0.12	0.74
<i>gemini</i>	<u>0.31</u>	<u>0.12</u>	1.00	<u>0.13</u>	<u>0.39</u>

Table 1: Isolation evaluation on AO, AM, and AV across different models. **Bold** and underline indicate each metric’s best and second-best performance, respectively.

for its strongest capability, whether AO, AM or AV, we evaluate their collective impact on authorship privacy. This integrated perspective mirrors real-world workflows in which texts undergo successive AI-mediated transformations, from iterative edits in anonymous online forums to chained paraphrasing and verify in whistleblowing activities.

4 Experiment Setup

Models. We utilize the well-known commercial LLMs of varying presence of reasoning capability and origins: GPT-4o-mini (Achiam et al., 2023), GPT-o3-mini (Brown, 2020), Gemini-2.0 (Team et al., 2023), and Deepseek-v3 (Liu et al., 2024).

Datasets. We utilize three datasets: *Speech*: US Presidents’ speeches from Fisher et al. (2024), *Quora*: Quora blog posts by diverse users with active online presence that we collect ourself; and *Essay*: writing essays from layperson (Li and Wan, 2025). These corpora vary in text length and author notoriety, descending from *Speech*, *Quora* and *Essay*. They also allow us to evaluate LLMs’ performance on writing by both native and non-native English speakers.

Prompts. Following previous works such as LIP (Huang et al., 2024), we design prompts along four key dimensions: *Context*, *Task*, *Instruction*, and *Output* to characterize open-ended LLMs’ behavior systematically (Cheng et al., 2023b). Specifically, we prompt LLMs to focus on writing style rather than topic or content differences.

Metrics. In our work, authorship privacy depends on identifying linguistic traits that are unique to individuals and can also help differentiate human-authored text from that generated by LLMs. Particularly, we examine how 4 key linguistic features change before and after an authorship task AO, AM and AV is performed. Central to this is word distribution, quantified using TF-IDF similarity (denoted as **SIM**), which is also widely applied in detecting deepfake text by revealing unnatural or overly

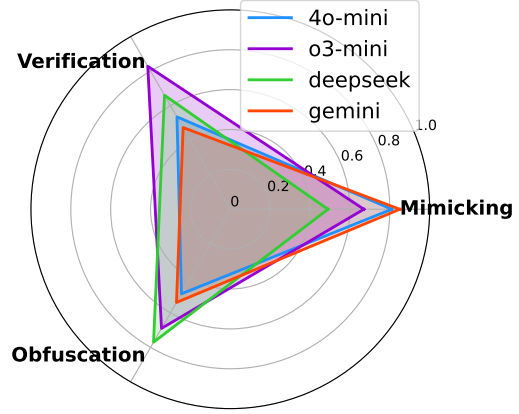


Figure 2: We present an overall pairwise interdependency evaluation of each LLM across the tasks of AO, AM, and AV. For each aspect, the final score is computed as the average across two “judge” evaluations to enable relative comparison.

consistent vocabulary usage (Becker et al., 2023). Additionally, we evaluate language naturalness using perplexity (denoted as **PPL**) and also report the **KL** divergence over the distribution of text PPL scores. This metric is commonly employed to capture the natural writing patterns of individuals and to detect machine-generated text that may appear overly fluent or statistically optimized compared to genuine human writing. In our experimental setup, we conduct evaluations both with metadata $C_a = \{M_a, D_a\}$ and without metadata $C_a = \{D_a\}$. Details of the datasets, prompts, and metrics are provided in the Appendix. A.2, A.3, A.4, respectively.

5 Experiment Results

5.1 Isolation Evaluation (RQ1)

Overall, *o3-mini* performs the best in AO and AV tasks, and *4o-mini* leads in faithful AM (Table 1). Particularly, *o3-mini* achieves the highest perplexity (2.71) and lowest similarity (0.10) in AO, indicating more distinct and less traceable outputs. For AM, *4o-mini* excels with the lowest PPL (0.65) and highest similarity (0.13), reflecting better stylistic imitation of the original texts. For AV, *o3-mini* identifies authorships with the highest accuracy (0.89).

5.2 Pairwise Interdependency (RQ2)

From the isolation evaluation (Sec. 5.1), we select *o3-mini* as both the obfuscation and verification judge, and *4o-mini* as the mimicking judge to assess the interplays among the authorship tasks. Fig. 2 presents a comprehensive comparison of the four models’ influence capabilities across AO, AM,

Models	Speech			Quora			Essay		
	KL	SIM	ACC	KL	SIM	ACC	KL	SIM	ACC
<i>4o-mini</i>	0.14	0.08	0.71	1.21	0.19	0.67	1.99	0.18	0.59
<i>o3-mini</i>	0.91	0.08	0.58	1.96	0.16	0.59	2.15	0.15	0.51
<i>w. meta</i> <i>gemini</i>	0.23	0.09	0.66	1.82	0.18	0.69	1.83	0.15	0.61
<i>deepseek</i>	1.15	0.08	0.53	2.15	0.13	0.57	2.23	0.15	0.51
<i>w.o. meta</i> <i>4o-mini</i>	0.39	0.08	0.63	1.25	0.16	0.62	1.87	0.17	0.62
<i>o3-mini</i>	1.41	0.07	0.59	1.84	0.16	0.63	2.01	0.15	0.53
<i>w.o. meta</i> <i>gemini</i>	0.05	0.08	0.70	1.76	0.17	0.76	1.86	0.15	0.60
<i>deepseek</i>	1.76	0.06	0.52	1.85	0.16	0.62	2.21	0.15	0.53

Table 2: Evaluation on obfuscation. KL (\uparrow), SIM (\downarrow), and Verification Accuracy (ACC) (\downarrow) between the mimicked and original text.

and AV. Overall, *gemini*, *deepseek*, and *o3-mini* are the most influential or effective with mimicking, obfuscation, and verification, respectively. We analyze each authorship task in detail as follows.

Influence of Obfuscation. To quantify AO, we employ a mimicking judge (*4o-mini*) and a verification judge (*o3-mini*). Table 2 reports the KL and SIM between mimicked and original texts and the verification accuracy on original texts when using obfuscated texts as the ground truth. Overall, among all models, *deepseek* consistently demonstrates the strongest obfuscation influence across all datasets, achieving the highest KL and lowest SIM scores. This indicates that its obfuscated outputs deviate the most from the original writing style.

In addition, the results also show that *obfuscation without user metadata generally outperforms the versions that incorporate metadata*. This suggests that metadata may inadvertently constrain the models, making it more difficult to mask the original writing style. In other words, the mimicking judge can utilize the same user metadata to reconstruct the original author’s writing style, making the obfuscation less impactful. Furthermore, the performance gap between the with-metadata and without-metadata settings is most pronounced in the Speech dataset, which features more well-known authors. This gap progressively narrows in the Quora and Essay datasets, reflecting a trend: *it is easier to conceal the identity of less well-known authors, regardless of metadata inclusion*.

Influence of Mimicking. To quantify AM, we evaluate the mimicked texts using two distinct judges: an obfuscation judge (*o3-mini*) and a verification judge (*o3-mini*). Table 3 reveals several consistent trends across datasets. Gemini achieves the strongest overall performance in text obfuscation and verification, followed by *4o-mini*, with Gem-

Models	Speech			Quora			Essay		
	KL	SIM	ACC	KL	SIM	ACC	KL	SIM	ACC
<i>4o-mini</i>	3.25	0.05	0.73	2.51	0.17	0.78	2.32	0.20	0.68
<i>o3-mini</i>	2.95	0.06	0.70	2.30	0.19	0.73	2.14	0.19	0.65
<i>w. meta</i> <i>gemini</i>	3.29	0.05	0.87	3.20	0.15	0.89	2.98	0.18	0.71
<i>deepseek</i>	2.95	0.07	0.65	2.18	0.18	0.82	1.97	0.21	0.67
<i>w.o. meta</i> <i>4o-mini</i>	3.32	0.06	0.70	2.13	0.16	0.79	2.16	0.20	0.63
<i>o3-mini</i>	3.26	0.06	0.62	2.24	0.19	0.64	1.98	0.22	0.60
<i>w.o. meta</i> <i>gemini</i>	3.28	0.05	0.82	2.48	0.15	0.87	2.79	0.19	0.69
<i>deepseek</i>	2.58	0.07	0.59	2.37	0.17	0.81	2.03	0.22	0.62

Table 3: Evaluation on mimicking. KL (\uparrow), SIM (\downarrow), and Verification Accuracy (ACC) (\uparrow) between the obfuscation and original text.

ini leading in most KL (\uparrow), SIM (\downarrow), and ACC (\uparrow) metrics. *Contrast with previous AO evaluation, incorporating user metadata to AM significantly enhances verification quality specially on Speech data*. Notably, the performance gap between settings with and without metadata narrows from well-known to lesser-known authors, suggesting that metadata plays a more critical role in capturing and disguising distinctive writing styles. Specifically, in the Speech dataset, the gap in KL divergence and SIM metrics between the metadata and without-metadata settings is substantially larger for AO than for AM. This implies that metadata is more influential in AO or that AO is generally more effective than AM. One possible explanation is that the input text contains many identifiable linguistic patterns, making it easier to alter (for obfuscation) than to replicate (for mimicking).

Influence of Verification. We construct noisy samples \bar{C}_a by doing AV across the 4 models, which then serve as inputs for obfuscation and mimicking judge. Overall, *o3-mini* achieves the highest precision and recall, with *deepseek* showing strong recall, while *4o-mini* and *gemini* perform less effectively in AV. We refer to Appendix. A.5 for detailed setup and results.

Table 4 reports how AV influences AO and AM when feeding AV with perfect (C_a) and noisy samples (\bar{C}_a). Overall, models with higher precision, indicating fewer *false positives* in \bar{C}_a (Eq. 8, Eq. 9) and reduced noise in the few-shot ground truth, exhibiting smaller divergence between obfuscation texts generated with perfect and imperfect samples. This suggests that cleaner sample ground truth examples make the obfuscation texts more indistinguishable. Moreover, removing metadata during obfuscation amplifies the divergence between obfuscated texts, potentially because the obfuscation

Models	\overrightarrow{VO}				\overrightarrow{VM}			
	Speech		Quora		Speech		Quora	
	KL	SIM	KL	SIM	KL	SIM	KL	SIM
<i>w. meta</i>								
4o-mini	1.47	0.24	1.89	0.19	0.21	0.33	0.39	0.26
o3-mini	1.08	0.27	1.57	0.24	0.19	0.34	0.30	0.28
gemini	1.65	0.22	1.80	0.18	0.22	0.30	0.40	0.25
deepseek	1.21	0.24	1.74	0.21	0.20	0.33	0.35	0.26
<i>w.o. meta</i>								
4o-mini	1.72	0.22	1.91	0.17	0.34	0.29	0.41	0.25
o3-mini	1.24	0.24	1.60	0.23	0.24	0.31	0.36	0.27
gemini	1.71	0.18	1.83	0.17	0.33	0.28	0.43	0.25
deepseek	1.45	0.21	1.72	0.20	0.29	0.31	0.38	0.26

Table 4: Evaluation on verification. KL (\downarrow) and SIM (\uparrow) measure similarity between two obfuscated texts. Full results are shown in Table A5.

judge can utilize the metadata to force the obfuscated texts to be similar. Lastly, across datasets, the gap in KL and SIM becomes narrower as the author becomes less well-known, reflecting the *diminishing influence of author-specific features in obfuscation*.

In terms of \overrightarrow{VM} , overall, mimicked texts derived from ground-truth examples of LLMs with higher precision exhibit lower divergence, reflected by smaller KL and higher SIM, because higher precision reduces false positives and thus introduces less noise during the mimicking process. Additionally, *AV’s access to metadata consistently improves the AM judge’s ability to perform accurate text mimicking compared to settings without metadata, although this benefit diminishes as the authors become less well-known*. The reason might be LLMs’ familiarity with famous people, and hence able to effectively utilize metadata.

5.3 Triplet-wise Interdependency (RQ2, RQ3)

This section analyzes five *iterative cycles* of AO, AM, and AV to evaluate how LLMs progressively shape stylometric patterns over time. Without loss of generality, we begin with mimicking followed by obfuscation, as their outputs are iteratively used as inputs for the subsequent task throughout the evaluation process. An interesting observation is the emergence of zig-zag patterns in all plots in Fig. 3, suggesting an ongoing “tug-of-war” between mimicking and obfuscation. Obfuscation appears to be more dominant, though the nature of this interplay varies depending on (1) the dataset and (2) the presence or absence of metadata.

Authorship Verification. Overall, mimicking demonstrates the ability to recover the original text to some extent (first plot in Fig. 3). However, its

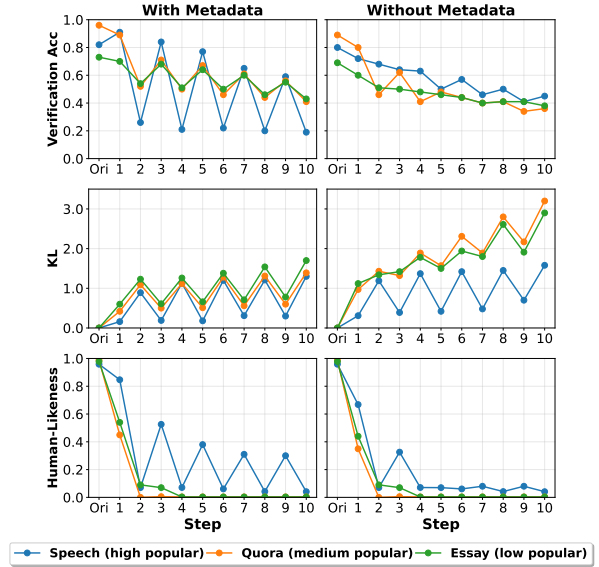


Figure 3: Verification accuracy (\uparrow), KL (\downarrow), and Human-likeness scores of mimicked and obfuscated texts compared to original texts across datasets, both with and without metadata. The x-axis represents the step order, ranging from 1 to 10 for 5 iterations *alternating between* $AM \rightarrow AO \rightarrow AM \rightarrow \dots \rightarrow AO$. AV is used as an intermediate step after AO and does not generate any texts, so we hide it for clarity. We refer to Table A6 for the detailed results.

effectiveness diminishes over successive iterations, due to the cumulative noise introduced by repeated obfuscation steps, which makes it increasingly difficult for the mimicker to reconstruct the original content. This degradation is particularly evident in the Quora and Essay datasets, where mimicking accuracy drops sharply after the first iteration. In terms of obfuscation, we observe a substantial reduction in verification accuracy for the Speech dataset compared to Quora and Essay. This suggests that *obfuscation is more effective when author identity is strongly encoded in the text*, as is the case for public figures whose speech styles are easily recognizable. Notably, *removing metadata from AO/AM consistently decreases verification accuracy across all datasets and iterations*, further demonstrating the value of auxiliary information in authorship verification.

Language Naturalness. Overall, KL divergence increases over iterations, mirroring verification trends and signaling growing linguistic drift from the original text (second plot in Fig. 3). Mimicking degrades over time, especially without metadata, while obfuscation consistently drives text away from its original form. Mimicking works best on shorter, structured texts like Speech, whereas obfuscation excels on longer, more variable texts like

Quora and Essay due to richer linguistic features for distortion.

Anthropomorphism Analysis. We investigate whether generated text becomes more human-written or machine-generated through successive iterations of AM→AO (Cheng et al., 2025). To quantify this, we employ GPTZero², one of the most popular commercial deepfake text detectors, to assess the degree to which a given text *resembles human writing*. Fig. 3 reports the human-likeness score- GPTZero’s estimated probability that a given text is written by a human. The first mimicked texts often appear most human-like, especially on the Speech dataset, while obfuscated texts consistently score low. Mimicking after obfuscation can partially restore human-like style, but this effect *fades over time as the text becomes increasingly machine-generated*. For Quora and Essay, texts generated after the second iteration are generally classified as machine-generated. This may be attributed to the lower popularity and variability in writing styles within these datasets, making it harder for mimicking models to recover stylistic patterns. Without metadata, this effect intensifies across all datasets, texts quickly adopt machine-like traits after two iterations, with minimal recovery by AM even in the Speech dataset.

Topic Distribution. We analyze how mimicking and obfuscation alter topic distributions using LDA (Blei et al., 2003) and find that iterative authorship tasks gradually shift texts away from their original themes. For instance, in the Speech dataset, the initial texts cover topics such as *politics, elections, health/life, war/terror*, and *economy/jobs* are replaced by more generic, repetitive content over time. This degradation may result from the *compounding effects of generation*, as LLMs tend to produce less specific and more repetitive content (Holtzman et al., 2020). Detailed topic trends are in Appendix A.8.

6 Discussions

Relationship between Authors’ Popularity and Metadata’s Effectiveness. Including metadata significantly boosts AV effectiveness, especially for well-known individuals, heightening privacy risks through easier re-identification or impersonation. Otherwise, lesser-known authors are less affected, indicating that popularity increases identifiability. While obfuscation helps, it does not reliably ensure

anonymity. These results carry important implications for LLM providers like OpenAI, Google: (1) LLMs may unintentionally erode user privacy by leveraging publicly available or leaked metadata; second, (2) incorporating privacy-preserving mechanisms into authoring and editing tools; (3) providing transparency and safeguards around how metadata is used or inferred in LLM-driven authorship tasks.

The Double-edged Sword of LLMs: Empowering Privacy or Enabling Threats? LLMs are double-edged tools. On one hand, users can utilize LLMs for privacy-preserving purposes. For instance, whistleblowers or vulnerable individuals may rely on LLM-powered obfuscation tools to share sensitive content anonymously. On the other hand, the same technology can be misused for impersonation or misinformation. Our results show that LLMs can convincingly mimic writing styles, especially when metadata such as demographics is available, opening the door for social engineering attacks or deepfake text generation. Therefore, individuals must be aware that their public user-generated content, even absent explicit identifiers, can leave behind implicit rich digital traces. This raises an urgent need for tools that proactively evaluate and adjust online writings to minimize their digital traces.

Impersonation and Misuse at Scale. The interplay between AO and AM reveals that obfuscated text can still be reverse-engineered by powerful LLMs, especially with demographic cues. This poses real risks: malicious actors could impersonate public figures or institutions at scale to spread misinformation. As a result, stronger authorship detection tools are essential to identify AI-generated impersonations and trace their origins.

7 Conclusion

In this work, we present a unified framework to evaluate how LLMs interact across authorship obfuscation, mimicking, and verification, highlighting task-specific strengths and the role of demographic metadata. Our analysis quantifies interdependencies among these tasks and shows that obfuscation generally dominates mimicking in disrupting authorial signals, though mimicking can partially recover stylistic traits over time. Notably, models with stronger reasoning excel at verification and style concealment but struggle to faithfully replicate an author’s unique voice.

²<https://gptzero.me/>

Acknowledgements

We thank ChatGPT and Grammarly for their assistance in typo correction.

Limitation

Despite presenting a comprehensive evaluation framework for the three core authorship privacy tasks—authorship verification, obfuscation, and mimicking—using diverse linguistic metrics across a range of real-world datasets, our study is limited by the absence of human-centered evaluation. While automated metrics offer scalability and consistency, incorporating human judgment would provide valuable insights into the perceived naturalness, fluency, and effectiveness of obfuscated or mimicked text. This is especially important in assessing whether generated text truly conceals authorship or convincingly imitates another writing style from a human perspective. Future work could benefit from human-in-the-loop studies to better align evaluation with real-world perceptions and practical usability.

Broader Impacts and Ethics Statement

This work raises important ethical considerations in authorship privacy. While our framework helps evaluate and improve privacy-preserving techniques, it also reveals how LLMs can deanonymize writers or impersonate them, posing risks to vulnerable individuals and enabling potential misuse, such as spreading misinformation. We urge the development of safeguards, such as tools that warn users of identifiability risks and stronger detection systems for AI-generated content. All data used are publicly available and handled following ethical research standards.

References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv*.

Jiafu An, Difang Huang, Chen Lin, and Mingzhu Tai. 2024. Measuring gender and racial biases in large language models. *arXiv*.

Yejin Bang, Delong Chen, Nayeon Lee, and Pascale Fung. 2024. Measuring political bias in large language models: What is said and how it is said. *arXiv*.

Jonas Becker, Jan Philip Wahle, Terry Ruas, and Bela Gipp. 2023. Paraphrase detection: Human vs. machine content. *arXiv*.

David M Blei, Andrew Y Ng, and Michael I Jordan. 2003. Latent dirichlet allocation. *JMLR*.

Tom B Brown. 2020. Language models are few-shot learners. *NeurIPS*.

Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and Chiyuan Zhang. 2023. Quantifying memorization across neural language models. *ICLR*.

Jin Chen, Zheng Liu, Xu Huang, Chenwang Wu, Qi Liu, Gangwei Jiang, Yuanhao Pu, Yuxuan Lei, Xiaolong Chen, Xingmei Wang, et al. 2024. When large language models meet personalization: Perspectives of challenges and opportunities. *WWW*.

Myra Cheng, Su Lin Blodgett, Alicia DeVrio, Lisa Egede, and Alexandra Olteanu. 2025. Dehumanizing machines: Mitigating anthropomorphic behaviors in text generation systems. *arXiv*.

Myra Cheng, Esin Durmus, and Dan Jurafsky. 2023a. Marked personas: Using natural language prompts to measure stereotypes in language models. *ACL*.

Myra Cheng, Tiziano Piccardi, and Diyi Yang. 2023b. Compost: Characterizing and evaluating caricature in llm simulations. *EMNLP*.

Ameet Deshpande, Vishvak Murahari, Tanmay Rajpurohit, Ashwin Kalyan, and Karthik Narasimhan. 2023. Toxicity in chatgpt: Analyzing persona-assigned language models. *EMNLP*.

Haodong Duan, Jueqi Wei, Chonghua Wang, Hongwei Liu, Yixiao Fang, Songyang Zhang, Dahua Lin, and Kai Chen. 2024. Botchat: Evaluating llms’ capabilities of having multi-turn dialogues. *NAACL*.

Jillian Fisher, Skyler Hallinan, Ximing Lu, Mitchell Gordon, Zaid Harchaoui, and Yejin Choi. 2024. Styleremix: Interpretable authorship obfuscation via distillation and perturbation of style elements. *EMNLP*.

Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. 2020. The pile: An 800gb dataset of diverse text for language modeling. *arXiv*.

Jerry Zhi-Yang He, Sashrika Pandey, Mariah L Schrum, and Anca Dragan. 2025. Cos: Enhancing personalization and mitigating bias with context steering. *ICLR*.

Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The curious case of neural text degeneration. *ICLR*.

Baixiang Huang, Canyu Chen, and Kai Shu. 2024. Can large language models identify authorship? *EMNLP*.

776	Baixiang Huang, Canyu Chen, and Kai Shu. 2025. Authorship attribution in the era of llms: Problems, methodologies, and challenges. <i>SIGKDD</i> .	830
777		831
778		832
779	Chia-Yu Hung, Zhiqiang Hu, Yujia Hu, and Roy Ka-Wei Lee. 2023. Who wrote it and why? prompting large-language models for authorship verification. <i>EMNLP</i> .	833
780		834
781		835
782		836
783	Guangyuan Jiang, Manjie Xu, Song-Chun Zhu, Wenjuan Han, Chi Zhang, and Yixin Zhu. 2024a. Evaluating and inducing personality in pre-trained language models. <i>NeurIPS</i> .	837
784		838
785		839
786		840
787	Hang Jiang, Xiajie Zhang, Xubo Cao, Cynthia Breazeal, Deb Roy, and Jad Kabbara. 2024b. Personallm: Investigating the ability of large language models to express personality traits. <i>NAACL</i> .	841
788		842
789		843
790		844
791	Saketh Reddy Karra, Son The Nguyen, and Theja Tula-bandhula. 2022. Estimating the personality of white-box language models. <i>arXiv</i> .	845
792		846
793		847
794	Siwon Kim, Sangdoo Yun, Hwaran Lee, Martin Gubri, Sungroh Yoon, and Seong Joon Oh. 2024. Propile: Probing privacy leakage in large language models. <i>NeurIPS</i> .	848
795		849
796		850
797		851
798	Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. 2023. Bloom: A 176b-parameter open-access multilingual language model. <i>arXiv</i> .	852
799		853
800		854
801		855
802		856
803		857
804	Jiatao Li and Xiaojun Wan. 2025. Who writes what: Unveiling the impact of author roles on ai-generated text detection. <i>arXiv</i> .	858
805		859
806		860
807	Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. 2024. Deepseek-v3 technical report. <i>arXiv</i> .	861
808		862
809		863
810		864
811	Tuc Nguyen and Thai Le. 2024. Adapters mixup: Mixing parameter-efficient adapters to enhance the adversarial robustness of fine-tuned pre-trained text classifiers. <i>EMNLP</i> .	865
812		866
813		867
814		
815	David Pape, Sina Mavali, Thorsten Eisenhofer, and Lea Schönherr. 2024. Prompt obfuscation for large language models. <i>arXiv</i> .	
816		
817		
818	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. <i>OpenAI blog</i> .	
819		
820		
821		
822	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>JMLR</i> .	
823		
824		
825		
826		
827	Alireza Salemi, Sheshera Mysore, Michael Bendersky, and Hamed Zamani. 2024. Lamp: When large language models meet personalization. <i>ACL</i> .	
828		
829		
	Leonard Salewski, Stephan Alaniz, Isabel Rio-Torto, Eric Schulz, and Zeynep Akata. 2023. In-context impersonation reveals large language models’ strengths and biases. <i>NeurIPS</i> .	
	Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. 2023. Gemini: a family of highly capable multimodal models. <i>arXiv</i> .	
	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. <i>arXiv</i> .	
	Adaku Uchendu, Thai Le, and Dongwon Lee. 2023. Attribution and obfuscation of neural text authorship: A data mining perspective. <i>SIGKDD</i> .	
	Adaku Uchendu, Saranya Venkatraman, Thai Le, and Dongwon Lee. 2024. Catch me if you gpt: Tutorial on deepfake texts. In <i>NAACL: Human Language Technologies (Volume 5: Tutorial Abstracts)</i> .	
	Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, et al. 2021. Ethical and social risks of harm from language models. <i>arXiv</i> .	
	Eric Xing, Saranya Venkatraman, Thai Le, and Dongwon Lee. 2024. Alison: Fast and effective stylometric authorship obfuscation. In <i>AAAI</i> .	
	Chiyuan Zhang, Daphne Ippolito, Katherine Lee, Matthew Jagielski, Florian Tramèr, and Nicholas Carlini. 2023. Counterfactual memorization in neural language models. <i>NeurIPS</i> .	
	Rong Zheng, Jiexun Li, Hsinchun Chen, and Zan Huang. 2006. A framework for authorship identification of online messages: Writing-style features and classification techniques. <i>Journal of the American society for information science and technology</i> .	

A Appendix

A.1 Practical applications on authorship privacy

In this section, we present some applications of our research questions related to real-world authorship privacy.

RQ1: A practical application of this research question in authorship privacy is enabling users, such as whistleblowers, activists, or social media participants, to select the most suitable LLM for their goals. For instance, if a user seeks to mask their identity when writing sensitive content, the analysis can guide them toward models with strong authorship obfuscation (AO) performance. Conversely, a journalist or researcher aiming to emulate a public figure’s writing style might benefit from models that excel in authorship mimicking (AM). Similarly, platforms concerned with detecting AI-generated or impersonated text can rely on models with high authorship verification (AV) accuracy. Thus, understanding isolated LLM performance informs the deployment of tailored models in real-world authorship privacy scenarios.

RQ2: A practical application of this research question in authorship privacy lies in improving the design and security of multi-step text processing pipelines used in sensitive communications. Specifically, in scenarios like anonymous online forums, whistleblower disclosures, or secure messaging, texts often undergo multiple transformations—generation, obfuscation, and verification—each performed by different LLMs. Understanding how these models influence one another and the interdependencies that arise helps identify potential privacy risks, such as:

1. whether obfuscation techniques are truly effective in concealing an author’s style. For instance, whistleblowers and journalists who rely on textual obfuscation to anonymize their writing may still be at risk if LLMs can reverse-engineer their original style, allowing adversaries to trace the obfuscated text back to them.
2. anonymizing sensitive documents, e.g. legal testimonies or medical records, where ensuring that downstream mimicking models cannot recover the original author’s style is critical for privacy protection.
3. evaluating the potential misuse of LLMs in impersonation attacks, such as forging stylistically similar content for deception or misinformation.

4. forensic investigations, where reliable verification must distinguish genuine statements from adversarially altered or mimicked texts. Additionally, content moderation systems can leverage these insights to detect and flag deceptive or impersonated content, enhancing online platform safety and trust.

RQ3: A practical application of this question is in developing robust authorship privacy tools that account for real-world scenarios where text undergoes multiple rounds of transformation. For instance, in environments like anonymous publishing platforms or secure communication channels, text might be repeatedly mimicked, obfuscated, and verified using different LLMs. Understanding how these iterative cycles influence each other helps identify how privacy can degrade or be preserved over successive edits. This knowledge allows designers to build more effective multi-stage pipelines that maintain author anonymity, prevent unintended leakage of writing style, and improve the reliability of verification methods, ultimately enhancing the security and trustworthiness of authorship privacy systems.

A.2 Additional statistics on evaluation dataset

We present the statistics on the evaluation dataset in Table A1 and A2.

Dataset	# Exam	Avg length	Avg doc. length	Avg #sen. per doc.	# Authors
Speech	5,172	58.20	17.44	3.34	3
Quora	9,899	294.62	18.83	15.64	5
Essays	154	225.87	9.43	7.24	3

Table A1: Statistics of the evaluation datasets.

A.3 Prompt Construction

Author identification can be generated based on the attributes of each learner, including sex, academic background, level of English proficiency, and country of origin, to build a more targeted background persona. For example: *The author is female. Her academic background is in the Humanities. Her English proficiency level is CEFR B1 (lower). She is from Singapore, an ESL environment (English as a Second Language).* The prompt construction for mimicking, attribution, and obfuscation are written in Table A4.

Attribute	Value	Count
CEFR	B1_1	914
	B1_2	881
	A2_0	470
	B2_0	231
	XX_0	73
Acad. Genre	Sciences & Tech.	1,034
	Social Sciences	762
	Humanities	674
	Life Sciences	99
Lang. Env.	EFL	1,886
	ESL	610
	NS	73
Sex	F	1,430
	M	1,139

Table A2: Distribution of author attributes across 2,569 learners.

Models	Speech		Quora		Essay	
	Precision	Recall	Precision	Recall	Precision	Recall
<i>4o-mini</i>	0.36	0.50	0.36	0.50	0.33	0.50
<i>o3-mini</i>	0.67	0.80	0.54	0.70	0.50	0.70
<i>gemini</i>	0.36	0.50	0.33	0.50	0.27	0.40
<i>deepseek</i>	<u>0.62</u>	0.80	0.54	0.70	<u>0.43</u>	<u>0.60</u>

Table A3: Authorship verification precision and recall of the four LLMs on the Speech, Quora, and Essay datasets.

A.4 Evaluation Metrics

Word Distribution. We employ TF-IDF to quantify each word’s significance within a document relative to the entire corpus. TF counts word occurrences, while IDF down-weights common terms. We extract TF-IDF vectors from our text sources and compute cosine similarity to assess stylistic and thematic alignment.

Language Naturality. Perplexity (PPL) evaluates how well a language model predicts a given text, with lower PPL reflecting greater confidence and closer adherence to learned linguistic patterns. Since LMs capture typical language structures from large corpora, PPL is a proxy for naturalness. Here, we fine-tune GPT-2 (Radford et al., 2019) on the original corpus and compute text-level PPL for both human-written and generated texts.

A.5 Detailed results on Precision and Recall

We construct the imperfect ground truth examples \bar{x}_p by sampling 20 examples from the original texts, including 10 from the author and 10 from others. The target model will be used to verify authorship. All the examples classified as correct verification

will be used as the ground truth for the obfuscation and mimicking processes. Table A3 shows detailed results on Precision and Recall.

A.6 Additional results for VO and VM

We present detailed evaluation results of VO and VM in Table A5.

A.7 Additional evaluation results on triplet-wise interdependency

We present detailed evaluation results on triplet-wise interdependency in Table A6.

A.8 Detailed results on topic distribution analysis

From Table A7 to A17, we show detailed results on topic distribution analysis on the mimicking and obfuscation process.

Task	Prompt
Verification	<ul style="list-style-type: none"> • System Prompt: You are a judge designed to verify the attribution of a human-author written text. • Instruction: You are given sample texts including 5 writings from the author and 5 writings from others. Analyze the writing styles of the input text, disregarding the differences in topic and content. Reasoning based on linguistic features such as phrasal verbs, modal verbs, punctuation, rare words, affixes, quantities, humor, sarcasm, typographical errors, and misspellings. Your task is to verify if the input text was written by <i>{author name}</i>. As output, exclusively return yes or no without any accompanying explanations or comments. • Context: Here is some information about the author: <i>{author identification}</i>. The 10 sample writings: <i>{sample text}</i>. • Task: The input text is: <i>{input text}</i>.
Mimicking	<ul style="list-style-type: none"> • System Prompt: You are an emulator designed to replicate the writing style of a human author. • Instruction: You are given 5 sample writings from the author. The goal of this task is to mimic the author’s writing style while paying meticulous attention to lexical richness and diversity, sentence structure, punctuation style, special character style, expressions and idioms, overall tone, emotion, and mood, or any other relevant aspect of writing style established by the author. Your task is to generate a <i>{avg}</i>-word continuation that seamlessly blends with the provided input text. Ensure that the continuation is indistinguishable from both the input text and the 5 sample writings by the author. As output, exclusively return the text completion without any accompanying explanations or comments. • Context: Here is some information about the author: <i>{author identification}</i>. The 5 sample writings from an author: <i>{sample text}</i>. • Task: The input text is: <i>{input text}</i>.
Obfuscation	<ul style="list-style-type: none"> • System Prompt: You are an emulator designed to hide the writing style of a human author. • Instruction: You are given 5 sample writings from an author. The goal of this task is to conceal the author’s writing style by carefully modifying lexical richness and diversity, sentence structure, punctuation patterns, special character usage, expressions and idioms, overall tone, emotion, mood, and any other distinguishing stylistic elements. Your task is to generate <i>{avg}</i>-word continuation that has writing style significantly different from the provided input text. Strive to make the rewritten text distinguishable from both the input text and the 5 sample writings by the author. As output, exclusively return the text completion without any accompanying explanations or comments. • Context: Here is some information about the author: <i>{author identification}</i>. The 5 sample writings from an author: <i>{sample text}</i>. • Task: The input text is: <i>{input text}</i>.

Table A4: Prompt construction for the 3 tasks to evaluate LLMs ability.

Models		VO						VM					
		Speech		Quora		Essay		Speech		Quora		Essay	
		KL	SIM	KL	SIM	KL	SIM	KL	SIM	KL	SIM	KL	SIM
w. meta	4o-mini	1.47	0.24	1.89	0.19	1.34	0.28	0.21	0.33	0.39	0.26	0.69	0.18
	o3-mini	1.08	0.27	1.57	0.24	1.26	0.31	0.19	0.34	0.30	0.28	0.52	0.19
	gemini	1.65	0.22	1.80	0.18	1.51	0.28	0.22	0.30	0.40	0.25	0.63	0.17
	deepseek	1.21	0.24	1.74	0.21	1.32	0.29	0.20	0.33	0.35	0.26	0.64	0.17
wo. meta	4o-mini	1.72	0.22	1.91	0.17	1.35	0.27	0.34	0.29	0.41	0.25	0.66	0.17
	o3-mini	1.24	0.24	1.60	0.23	1.32	0.29	0.24	0.31	0.36	0.27	0.54	0.18
	gemini	1.71	0.18	1.83	0.17	1.49	0.28	0.33	0.28	0.43	0.25	0.63	0.18
	deepseek	1.45	0.21	1.72	0.20	1.34	0.28	0.29	0.31	0.38	0.26	0.65	0.17

Table A5: Merged results from both evaluations: **Verification Obfuscation** and **Verification Mimicking**. KL (\downarrow) and SIM (\uparrow) measure similarity between two obfuscated texts. **Bold** and underline indicate best and second-best performance per category.

		Verification										KL										
		Original	1		2		3		4		5		1		2		3		4		5	
			AM	AO	AM	AO	AM	AO	AM	AO	AM	AO	AM	AO	AM	AO	AM	AO	AM	AO	AM	AO
w meta	Speech	0.82	0.91	0.26	0.84	0.21	0.77	0.22	0.65	0.20	0.59	0.19	0.16	0.89	0.19	1.13	0.18	1.20	0.31	1.21	0.30	1.30
	Quora	0.96	0.89	0.52	0.71	0.50	0.67	0.46	0.61	0.44	0.56	0.41	0.42	1.09	0.50	1.12	0.51	1.28	0.56	1.31	0.60	1.39
	Essay	0.73	0.60	0.54	0.58	0.51	0.50	0.49	0.46	0.45	0.46	0.43	0.60	1.23	0.61	1.26	0.66	1.38	0.71	1.54	0.78	1.70
wo meta	Speech	0.80	0.72	0.68	0.64	0.63	0.50	0.57	0.46	0.50	0.41	0.45	0.31	1.19	0.39	1.37	0.42	1.42	0.48	1.45	0.70	1.58
	Quora	0.89	0.80	0.46	0.62	0.41	0.48	0.44	0.40	0.41	0.34	0.36	0.97	1.43	1.32	1.89	1.57	2.31	1.89	2.80	2.17	3.20
	Essay	0.69	0.60	0.51	0.50	0.48	0.46	0.44	0.40	0.41	0.41	0.38	1.12	1.34	1.42	1.78	1.50	1.98	1.80	2.61	1.91	2.90

Table A6: Performance analysis across 5 iterations (AM: mimicking, AO: obfuscation) for Verification and KL Divergence metrics.

Topic	Top Words
0	day, election, going, people, help, votes, working, could, got, better
1	weapons, tax, best, let, people, made, could, plan, give, think
2	people, country, time, right, look, together, one, border, want, believe
3	iraq, health, costs, people, team, looking, war, year, care, working
4	people, jobs, american, time, america, states, think, right, work, put
5	new, nation, america, american, years, right, peace, workers, great, drug
6	want, people, terrorists, important, college, enforcement, asking, terror
7	one, security, people, country, war, life, let, never, america, american
8	going, government, economy, world, america, afghanistan, iraq, getting, history, go
9	want, going, people, americans, think, true, test, save, health, support

Table A7: Top 10 words for each LDA topic on the original Speech dataset

Topic	Top Words
0	america, nation, people, great, believe, world, together, continue, better, means
1	people, states, world, new, energy, afghan, united, best, take, working
2	weapons, people, country, know, america, act, got, work, tough
3	1st, country, good, could, china, american, always, quality, going, people
4	nation, iraq, america, united, security, safe, people, states, choose, issue
5	people, going, new, americans, great, way, thank, american, want
6	people, going, know, right, american, one, policy, get, great
7	people, economy, going, world, american, country, families, great, challenges, nation
8	want, good, working, people, continue, america, get, let, need
9	america, going, know, day, country, people, give, future, nation

Table A8: (Round1 Step1: Mimicking) Top 10 words for each LDA topic

Topic	Top Words
0	energy, progress, remains, across, people, iraq, together, let, ensure
1	built, probability, plane, boeing, airbus, children, life
2	ensuring, remain, nation, costs, moving, people, yet, accountability, financial
3	one, ensuring, future, fostering, ensure, secure, american, communities,
4	ensuring, essential, fostering, sustainable, future, progress, growth, efforts, economic
5	one, world, unity, life, fostering, resilience, children, wage
6	innovation, progress, challenges, yet, ensuring, fostering, remains, future, essential
7	commitment, last, crucial, energy, year, world, legal, principles, stability
8	future, progress, let, together, challenges, forward, shared, innovation, resilience
9	people, time, seemed, yet, dreams, distant, whispers, future, hope

Table A9: (Round1 Step2: Obfuscation) Top 10 words for each LDA topic

Topic	Top Words
0	people, america, time, one, great, care, future, american, could, talking
1	world, going, future, let, great, build, america, job, look
2	challenges, america, world, new, innovation, good, always, moment, americans, embracing
3	people, america, let, country, opportunity, world, time, essential, american
4	support, people, power, future, let, progress, collective, respect, together
5	commitment, people, principles, trade, essential, future, nation, ensuring, dedication
6	work, future, people, challenges, let, requires, together, vote, commitment
7	america, together, need, nation, challenges, states, total, let, open
8	people, going, country, great, america, let, bad, win, things
9	nation, right, progress, values, back, believe, going, time, security

Table A10: (Round2 Step1: Mimicking) Top 10 words for each LDA topic

Topic	Top Words
0	let, time, together, future, essential, ensure, yet, progress, resilience
1	progress, forward, let, together, path, future, ensuring, yet, test
2	progress, shared, challenges, future, together, unity, resilience, yet, innovation
3	innovation, future, could, change, time, progress, challenges, resilience, ensuring
4	future, let, progress, even, fostering, together, challenges, hope, path
5	day, time, yet, relentless, lies, collective, life, fostering, decisions
6	, offer, energy, greater, yet, dialogue, fostering, often, today, security
7	people, fostering, progress, solutions, ensuring, future, innovation, challenges, efforts
8	progress, future, collective, resilience, let, together, fostering, solutions, efforts
9	sustainable, fostering, growth, future, economic, innovation, ensuring, essential, together

Table A11: (Round2 Step2: Obfuscation) Top 10 words for each LDA topic

Topic	Top Words
0	jobs, future, american, people, look, time, let, lot, need
1	future, nation, together, values, people, going, build, vital, need, american
2	commitment, jobs, right, americans, unwavering, challenges, economy, fill, good, better
3	world, economy, good, fight, prevail, got, forces, iraq
4	together, security, great, strategy, even, america, need, economic, made
5	people, future, work, america, nation, together, values, commitment, american
6	get, progress, ahead, future, nation, true, americans, shared, tomorrow, day
7	going, people, really, see, great, know, coming, thing, election
8	challenges, together, face, world, nation, resolve, forward, future, let
9	always, people, believe, right, nation, freedom, country, bad, working

Table A12: (Round3 Step1: Mimicking) Top 10 words for each LDA topic

Topic	Top Words
0	future, energy, collaboration, fostering, growth, demands, resilience, sustainable, economic
1	ensuring, forward, innovation, together, progress, demands, stability, clear, economy
2	together, progress, future, let, shared, collaboration, innovation, challenges, collective
3	let, progress, yet, together, hope, future, resilience, world, time
4	progress, path, let, future, yet, forward, together, demands, ahead
5	future, together, hope, ensuring, efforts, values, unity, resilience, let
6	quiet, world, shared, becomes, yet, month, america, key, let
7	progress, future, challenges, forward, ensuring, innovation, resilience, let, time
8	time, global, need, one, north, right, let, people
9	future, fairness, ensuring, ensure, economic, together, everyone, fostering, collaboration

Table A13: (Round3 Step2: Obfuscation) Top 10 words for each LDA topic

Topic	Top Words
0	get, want, future, time, let, understand, ensuring, clear, situation
1	world, people, let, american, america, challenges, values, stand, together
2	people, right, help, innovation, prosperity, free, thing, economic, families
3	going, future, people, need, change, country, long, better, day
4	people, see, things, going, something, action, stand, disaster, let
5	people, bad, great, ', nation, america, time, see, going, country
6	progress, journey, shaped, nation, spirit, something, always, human, going
7	america, together, people, future, let, world, requires, true, get
8	people, support, economy, know, work, great, world, time, done
9	world, america, ahead, future, yet, let, resolve, hope, freedom

Table A14: (Round4 Step1: Mimicking) Top 10 words for each LDA topic

Topic	Top Words
0	guide, future, forward, commitment, essential, wage, next, protect
1	path, time, something, forward, energy, address, economic, achieving
2	let, solutions, forward, ?, progress, ahead, time, innovation, change
3	yet, forward, time, small, step, care, health, forged, energy
4	together, future, world, commitment, yet, challenges, let, across
5	shared, peace, resilience, path, collaboration, forward, progress, fostering
6	future, let, forward, progress, path, time, challenges, innovation, keep
7	progress, together, yet, one, let, change, time, vision, path
8	progress, together, resilience, challenges, essential, forward, future, collaboration, ensuring
9	future, together, let, progress, everyone, fostering, ensuring, innovation, build

Table A15: (Round4 Step2: Obfuscation) Top 10 words for each LDA topic

Topic	Top Words
0	together, better, great, world, path, win, easy, opportunity, need
1	let, america, future, country, time, get, believe
2	get, jobs, let, american, america, need, means, investing
3	american, world, let, unwavering, ahead, work, opportunity, people, everyone
4	let, done, people, job, forward, time, get, citizens, keep
5	people, america, let, get, know, country, great, right, american
6	got, ta, room, forward, doubt, ahead, fight, open, freedom, stay
7	people, american, work, always, america, ', time, act, future, nation
8	nation, let, best, got, stay, get, folks, sure
9	going, people, know, work, really, together, believe, want, world

Table A16: (Round5 Step1: Mimicking) Top 10 words for each LDA topic

Topic	Top Words
0	progress, together, shared, resilience, future, challenges, let, fostering, ensuring
1	lost, time, one, momentum, path, progress, america, became
2	future, fostering, innovation, progress, together, embracing, let, resilience, ensuring
3	future, progress, ensuring, let, together, challenges, innovation, collaboration, vision
4	future, progress, together, yet, forward, resilience, remains, shared, collective
5	let, progress, together, forward, step, commitment, action, ensure, future
6	together, future, progress, let, ensuring, path, unity, shared, commitment
7	ensuring, across, future, challenges, resilience, essential, without, forward, together
8	fostering, let, remains, future, progress, approach, financial, essential, together
9	progress, future, shared, forward, together, collaboration, challenges, yet, innovation

Table A17: (Round5 Step2: Obfuscation) Top 10 words for each LDA topic