

Contrastive Disentanglement for Authorship Attribution

Anonymous ACL submission

Abstract

Authorship Attribution (AA) aims to identify the authorship of texts by analyzing distinctive writing styles. While current AA methods have yielded promising performance, these approaches commonly exhibit suboptimal performance in contexts where the subject matter varies significantly (i.e., topic-shift scenarios). This limitation stems from their inadequacy in differentiating between the topical content and the author’s stylistic signature. Additionally, existing studies predominantly focus on AA at an individual level, thereby neglecting the exploration of regional-level AA, which could reveal common linguistic patterns influenced by cultural and geographical factors. Addressing these gaps, this paper introduces **ContratDistAA**, a novel framework employing contrastive learning coupled with mutual information maximization to segregate content from stylistic features in latent representations for AA tasks. Our comprehensive experimental evaluations reveal that **ContratDistAA** outperforms existing state-of-the-art models in both individual and regional-level AA scenarios. This advancement not only enhances the accuracy of authorship attribution but also expands its applicability to encompass regional linguistic analysis, thus contributing significantly to the broader field of computational linguistics.

1 Introduction

Motivation. Authorship Attribution (AA) is an extensively researched area (Zheng and Jin, 2023). The goal of AA is to identify the author of a piece of text based on distinctive linguistic characteristics inherent in their writing style. Applications of AA span a broad range of domains, including digital forensics (Iqbal et al., 2008) and plagiarism detection (Stamatatos and Koppel, 2011).

Existing methods in AA can be broadly categorized into two groups: traditional stylometric approaches (Seroussi et al., 2011; Bevendorff et al., 2019) and machine learning-based

techniques (Zhang et al., 2018; Saedi and Dras, 2021). Traditional stylometric methods exploit features such as word lengths, sentence lengths, and function words to attribute authorship. Machine learning-based methods, particularly deep learning techniques, were leveraged to capture intricate patterns in writing styles, often surpassing the performance of stylometric methods (Rivera-Soto et al., 2021; Wang et al., 2023).

Despite these advancements, a significant challenge persists in scenarios involving a shift in topics, particularly when the testing phase encompasses topics not present in the training dataset (Sapkota et al., 2014; Hu et al., 2023). This issue primarily arises from the conflation of topic-related content and the author’s unique writing style. Consequently, standard stylistic features employed in AA may inadvertently reflect topical variations rather than the author’s stylistic nuances, leading to inaccuracies in authorship determination based solely on writing style.

Moreover, the majority of existing research in AA predominantly concentrates on the individual author level, thereby overlooking the potential of regional-level AA. Exploring AA at the regional level could reveal distinct linguistic styles shared by authors within the same geographical region, influenced by cultural nuances. For example, in Singapore, the widespread use of English is distinctively marked by local cultural influences and slang, offering a unique dimension essential for effective AA at a regional scale. This warrants further investigation to fully understand and utilize the nuances of regional linguistic variations for the AA task.

Research Objectives. In this paper, we propose **ContratDistAA**, a novel AA approach that leverages contrastive learning and mutual information to disentangle topic and style information in the latent space. This allows us to handle topic shift settings and conduct AA at both individual and regional levels. To facilitate our investigations, we construct

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

084 a new dataset to support the regional-level AA task.
085 We conduct extensive experiments to evaluate **ContrastDistAA**
086 against state-of-the-art baselines on
087 both regional-level and individual-level AA tasks.

088 **Contributions.** Our work makes the following
089 contributions: (i) We introduce a new regional-
090 level AA task and a dataset to support the evalu-
091 ation of AA methods on this new task. (ii) We
092 propose **ContrastDistAA**, which can disentangle
093 content and style information to improve AA per-
094 formance. (iii) We conduct extensive experiments
095 to benchmark **ContrastDistAA** against state-of-the-
096 art AA methods. Our experiment results demon-
097 strate **ContrastDistAA**'s superior performance in
098 both individual-level and regional-level AA tasks.
099 This study not only fills a gap in the AA literature
100 but also sheds light on the intricate interplay be-
101 tween linguistic styles and cultural elements within
102 the realm of AA, offering new perspectives and
103 understanding in the field.

104 2 Related Work

105 2.1 Authorship Attribution

106 AA has been extensively researched, with re-
107 cent surveys providing comprehensive overviews
108 of seminal works and advancements in the field
109 (Zheng and Jin, 2023; Tyo et al., 2022). Re-
110 searchers primarily relied on heuristic and statisti-
111 cal approaches in the nascent stages of AA. These
112 involved the usage of basic stylometric features
113 such as word lengths, sentence lengths, and func-
114 tion words (Neal et al., 2017; Ding et al., 2017).
115 This phase evolved with the training of classical ma-
116 chine learning algorithms as classifiers to link these
117 stylometric features with author identities (Boen-
118 ninghoff et al., 2019b,a; Theóphilo et al., 2019).
119 The emergence of deep learning marked a signif-
120 icant shift in AA, enabling the learning of more
121 complex writing patterns (Shrestha et al., 2017; Hu
122 et al., 2020; Jafariakinabad et al., 2019; Liu et al.,
123 2021). The introduction of pre-trained language
124 models like BERT (Devlin et al., 2019) further rev-
125 olutionized AA, achieving state-of-the-art results
126 through fine-tuning for specific AA tasks (Rivera-
127 Soto et al., 2021; Manolache et al., 2021; Reimers
128 and Gurevych, 2019; El Boukkouri et al., 2020).
129 However, these techniques often performed poorly
130 in topic-shift scenarios, where the topics under eval-
131 uation during the testing phase are not represented
132 in the training data (Altakrori et al., 2021). Our
133 **ContrastDistAA** approach aims to overcome this

challenge by employing contrastive learning and
mutual information to separate content (i.e., topic)
and linguistic style in latent space for AA.

137 2.2 Contrastive Learning

138 Contrastive Learning has emerged as a pivotal ap-
139 proach in forming embedding spaces, where it clus-
140 ters similar data points together while distancing
141 dissimilar ones. Its efficacy is particularly evident
142 in computer vision, as seen in the work of Chen
143 et al. (2020) with their data augmentation frame-
144 work, and He et al. (2020) through the Momen-
145 tum Contrast (MoCo) for enhanced representation
146 learning. In Natural Language Processing (NLP),
147 contrastive learning has been instrumental in refin-
148 ing sentence representations, exemplified by the
149 methodologies of Giorgi et al. (2021) and Gao
150 et al. (2022), which utilize contrastive loss for learn-
151 ing textual embeddings. Additionally, significant
152 progress has been made in formulating strategies
153 for generating positive and negative samples, with
154 Robinson et al. (2021) addressing the challenge of
155 hard negatives through user-controlled sampling.

156 2.3 Disentangled Representation Learning

157 Disentangled Representation Learning, a method
158 that isolates distinct attributes of data into separate
159 variables, has significantly influenced various fields.
160 In computer vision, it is exemplified by CycleGAN,
161 which uses latent embeddings for image translation
162 without paired examples (Zhu et al., 2020). In
163 speech processing, this approach involves using
164 mutual information minimization to separate voice
165 style from content (Yuan et al., 2021). In NLP,
166 models like ADNet, which combine motivational
167 and adversarial losses, effectively disentangle style
168 and meaning in text embeddings (Romanov et al.,
169 2019). Notable developments include the multi-
170 decoder model of Fu et al. (2017) for text transfer
171 tasks with limited parallel corpora and Shen et al.
172 (2020)'s use of denoising objectives for sentence
173 reconstruction. Inspired by these advances, our
174 work adopts a similar approach to meticulously
175 disentangle content and style information in textual
176 data for the AA task.

177 3 Methodology

178 This section outlines our proposed model, **ContrastDistAA**,
179 designed to learn a disentangled rep-
180 resentation of writing style for AA. As depicted
181 in Figure 1, **ContrastDistAA** is structured in two

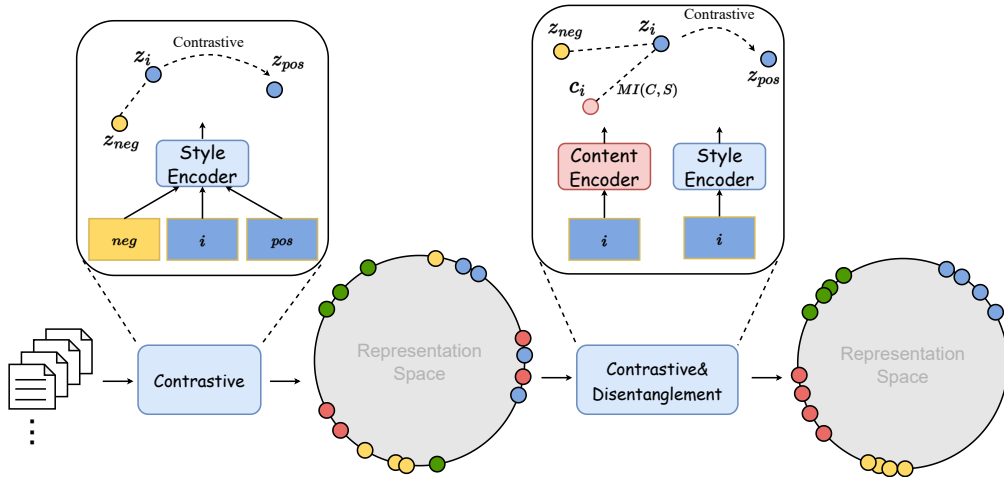


Figure 1: The architecture overview of ContrastDistAA model. The proposed models contains two-stages training process: (i) training using contrastive loss, and (ii) training using both contrastive loss with disentanglement loss.

distinct phases. The initial phase employs supervised contrastive loss to extract key stylistic features from labeled data. However, given the potential for content-related information to be intertwined with style, thus impacting the robustness of AA models, the subsequent phase of ContrastDistAA introduces a mutual information-based approach. This technique aims to separate style and content representations in the latent space, thereby enhancing the effectiveness of contrastive learning by clearly differentiating between style and content-specific attributes, including topical elements.

In subsequent sections, we will first review the contrastive learning component and the associated contrastive losses. This is followed by an introduction to mutual information, which is applied to disentangled representation learning for AA.

3.1 Contrastive Learning

Self-supervised representation learning has seen considerable progress in recent years, largely attributable to the application of contrastive learning (Wu et al., 2018; Hénaff et al., 2020; Oord et al., 2018; Chen et al., 2020). The fundamental mechanism of contrastive learning involves drawing an anchor and a “positive” sample closer in an embedding space, while simultaneously distancing the anchor from multiple “negative” samples, thus yielding meaningful representations. Specifically for AA tasks, we define “positive pair” consists of a text sample authored by the same individual as the anchor within a minibatch. In contrast, “negative pairs” are formed by aligning the anchor with randomly chosen samples from different authors within the same minibatch.

The initial phase of ContrastDistAA involves applying contrastive learning to train a style encoder, which extracts style features from texts authored by individuals or authors from specific regions. We utilize BERT (Devlin et al., 2019), acclaimed for its proficiency in capturing writing styles, as the style encoder. This encoder transforms discrete text into representations within latent space. Following this, supervised contrastive loss is applied to align representations of texts by the same author or from the same region more closely, while simultaneously distinguishing those from different authors or regions. This methodology enhances the style encoder’s ability to discern and learn discriminative style representations.

3.1.1 Supervised Contrastive Loss for AA

In the ContrastDistAA model, we implement a supervised contrastive loss for AA. Consider a batch consisting of N textual samples from distinct authors. Let $i \in I \equiv \{1, 2, \dots, N\}$ represent an individual sample in the minibatch, and let $A(i) \equiv I \setminus \{i\}$ denote the set of other texts excluding i . The negative samples for anchor i , denoted as $NEG(i) \equiv \{neg \in A(i) : y_{neg} \neq y_i\}$, are those not sharing the same author as i , while $POS(i) \equiv \{pos \in A(i) : y_{pos} = y_i\}$ represents the positive samples, sharing the same author as i . The supervised contrastive loss is particularly effective in scenarios where multiple samples belong to the same class, as it utilizes the available labels (Khosla et al., 2021). The formulation of the supervised contrastive loss for AA tasks is as follows:

$$L^{sup} = \sum_{i \in I} \frac{-1}{|POS(i)|} \sum_{pos \in POS(i)} \log \frac{\exp(z_i \cdot z_{pos} / \tau)}{\sum_{neg \in NEG(i)} \exp(z_i \cdot z_{neg} / \tau)} \quad (1)$$

where $z_i = StyleEncoder(x_i)$, the \cdot symbol denotes the inner product, $\tau \in \mathcal{R}^+$ is a scalar temperature parameter, $POS(i) \equiv \{pos \in A(i) : y_{pos} = y_i\}$ is the set of indices of all positive samples distinct from i , and $|POS(i)|$ is its cardinality.

3.2 Mutual Information for Style-Content Disentanglement

The style encoder, trained using supervised contrastive loss, becomes proficient at extracting representations that encapsulate both style and content attributes. Therefore, to refine the style encoder’s focus on capturing writing style more distinctly, we integrate mutual information with contrastive learning. This synergy aims to separate style and content information within the latent space.

Mutual information, a fundamental concept in information theory, measures the dependence between two random variables. For our model, mutual information between style (z) and content (c) representations is crucial. Its mathematical definition involves the expectation of the logarithm of the ratio of the joint distribution of z and c to their respective marginal distributions, which can be expressed as follows:

$$I(z; c) = \mathbb{E}_{p(z,c)} \left[\log \frac{p(z,c)}{p(z)p(c)} \right] \quad (2)$$

In practice, accurately calculating mutual information is challenging due to the intractability of the integral involved (Chen et al., 2016; Belghazi et al., 2018; Poole et al., 2019). To address this, we employ the Contrastive Log-ratio Upper Bound (CLUB) estimation method (Cheng et al., 2020). This approach is particularly suitable when conditional distributions such as $p(z|c)$ or $p(c|z)$ are not explicitly available. We approximate $p(z|c)$ using a variational distribution $q_\theta(z|c)$, parameterized by θ , leading to the definition of the variational CLUB term (vCLUB) as follows:

In disentangled representation learning, a common objective is to minimize the mutual information between varying types of embeddings, aligning with our training target (Poole et al., 2019). However, determining the exact value of

mutual information presents challenges in practical settings, as the integral in Eq. 2 is often intractable. To overcome this, several mutual information estimation methods have been proposed (Chen et al., 2016; Belghazi et al., 2018; Poole et al., 2019). We employ the estimation method known as the Contrastive Log-ratio Upper Bound (CLUB) (Cheng et al., 2020), which is suitable for the scenario where the conditional distributions $p(z|c)$ or $p(c|z)$ is not provided. A variational distribution $q_\theta(z|c)$ with parameter θ is used to approximate $p(z|c)$. Consequently, a variational CLUB term (vCLUB) is defined as follows:

$$I_{vCLUB}(z; c) := \mathbb{E}_{p(z,c)} [\log q_\theta(z|c)] - \mathbb{E}_{p(z)} \mathbb{E}_{p(c)} [\log q_\theta(z|c)] \quad (3)$$

The unbiased estimator for vCLUB is derived from a set of samples, effectively quantifying the mutual information in a computationally feasible manner. unbiased estimator for vCLUB with sample $\{z_i, c_i\}$ is expressed as follows:

$$\begin{aligned} \hat{I}_{vCLUB} &= \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N [\log q_\theta(z_i|c_i) - \log q_\theta(z_j|c_i)] \\ &= \frac{1}{N} \sum_{i=1}^N [\log q_\theta(z_i|c_i) - \frac{1}{N} \sum_{j=1}^N \log q_\theta(z_j|c_i)]. \end{aligned} \quad (4)$$

In summary, to facilitate style-content disentanglement in ContrastDistAA, we first deploy a content encoder, also a BERT model, to extract content representations, denoted as c . Meanwhile, the pre-trained style encoder from the first stage extracts style representations, denoted as z . Each post i thus has two distinct representations: the content representation c_i and the style representation s_i . Here, we apply the vCLUB estimator to minimize the mutual information between these content and style representations, refining the distinctiveness of each. Concurrently, the supervised contrastive loss continues to enhance the style encoder’s ability to capture writing style nuances. During the evaluation phase, only the style encoder is used to extract style representations from posts authored by individuals or from specific regions. The regional or individual author style representations are then calculated by averaging the post-style representations, facilitating a comprehensive and nuanced assessment of writing styles.

Dataset	#Users	#Train	#Valid	#Test
Regional Tweets	87,836	382,598	42,513	42,513
CCAT50	50	1,766	442	465
Twitter1000	1,000	6,000	2,000	2,000
IMDB62	62	37,200	12,400	12,400

Table 1: Statistics of datasets

4 Experiments

4.1 Experimental Settings

Datasets. To evaluate ContrastDistAA effectively on both individual and regional AA tasks, we utilize four datasets in our experiments. The statistical distributions of the datasets are shown in Table 1.

Regional Tweets: This dataset, aimed at exploring regional writing styles, comprises English tweets from Southeast Asia, collected using the Twitter API from 2021 to 2022. It includes 425,111 tweets from 87,836 users across six regions: Singapore, Kuala Lumpur, Manila, Jakarta, Hanoi, and Bangkok. The selection criteria focused on English tweets with more than three words for better data quality. The dataset is divided into training, validation, and testing sets in an 8:1:1 ratio.

CCAT50: A subset of the Reuters Corpus and a prominent resource in AA research, the CCAT50 dataset (Liu et al., 2012) focuses on the top 50 contributors in the CCAT (corporate/industrial) subtopic. It consists of 5,000 texts (50 per author) divided into distinct training, validation, and testing sets following a 6:2:2 ratio, based on the processed version by (Tyo et al., 2022).

Twitter1000: Derived from a larger Twitter dataset used in AA research (Shrestha et al., 2017; Schwartz et al., 2013), Twitter1000 includes tweets from the top 1,000 authors by volume, with 100 tweets randomly selected from each. The dataset is organized into training, validation, and testing subsets, also following a 6:2:2 ratio.

IMDB62: Recognized for long-text AA studies (Seroussi et al., 2014), the IMDB62 dataset includes contributions from 62 authors, each providing 1,000 texts. Similar to the others, this dataset is partitioned into training, validation, and testing sets in a 6:2:2 ratio.

Evaluation Metrics. Following existing AA studies, we adopt Macro-F1 and Micro-F1 as the evaluation metrics in our experiments.

4.2 Baselines

We benchmark our model against commonly used and state-of-the-art AA models. These baselines are trained or fine-tuned to perform both the regional-level and individual-level AA tasks.

LR-Stylo: This logistic regression model, leveraging stylometric features as inputs, is grounded in prior research (Sari, 2018; Aborisade and Anwar, 2018). Based on (Fabien et al., 2020), it uses ten different stylometric features like text length and word count for classification.

LR-TF-IDF: Employing Term Frequency - Inverse Document Frequency (TF-IDF) at the word level, this logistic regression classifier follows the approach of (Fabien et al., 2020). Pre-processing includes stemming and stop-word removal before constructing the TF-IDF features.

LR-Char: This model uses character N-gram-based features, shown to be effective in AA (Bischoff et al., 2020; Shrestha et al., 2017; Altakrori et al., 2021). Following (Tyo et al., 2022), the logistic regression classifier is trained with a mix of character N-gram, part-of-speech N-gram, and summary statistics.

LSTM: An LSTM model, inspired by recent studies (Oliva et al., 2022), incorporates a dense layer followed by a max pooling layer. It focuses on the hidden states of the LSTM for AA tasks.

BertAA: Utilizing a pre-trained BERT language model, BertAA (Fabien et al., 2020) is fine-tuned specifically for AA, integrating a dense layer and softmax activation function for AA classification.

DistilBert: Known for its efficiency as a compact language model, DistilBERT (Sanh et al., 2019) is fine-tuned for AA tasks.

Roberta: Employing the Roberta model (Liu et al., 2019), we follow the original hyperparameters and fine-tune it on AA datasets over a specific number of epochs.

4.3 Implementation.

Our experiments were carried out on a system operating on Ubuntu 20.04.3 LTS, equipped with robust hardware specifications including 24 CPU cores, 128 GB of RAM, and a base clock speed of 2.9 GHz. To facilitate efficient training of the pre-trained models, Nvidia GTX 3090 graphics cards were utilized. BERT, with its pre-trained weights, served dual roles as both the style and content encoders in our experiment, which was divided into two distinct stages.

Method	Regional Tweet		CCAT50		Twitter1000		IMDB62	
	Macro F1	Micro F1	Macro F1	Micro F1	Macro F1	Micro F1	Macro F1	Micro F1
LR-Stylo	0.176	0.251	0.013	0.037	0.019	0.035	0.013	0.037
LR-TF-IDF	0.402	0.446	0.554	0.554	0.566	0.566	0.554	0.554
LR-Char	0.252	0.308	0.180	0.209	0.077	0.128	0.503	0.503
LSTM	0.186	0.290	0.244	0.274	0.124	0.126	0.307	0.326
BertAA	0.433	0.472	0.518	0.512	0.226	0.249	0.627	0.654
DistilBERT	0.407	0.449	0.453	0.447	0.213	0.242	0.402	0.441
Roberta	0.476	0.522	0.466	0.497	0.622	0.626	0.735	0.749
ContrastDistAA	0.510	0.550	0.578	0.584	0.960	0.961	0.813	0.816
ContrastDistAA (w/o dist)	0.505	0.508	0.552	0.566	0.960	0.916	0.803	0.816

Table 2: Macro and Micro F1 scores for baselines and ContrastDistAA on four benchmark datasets.

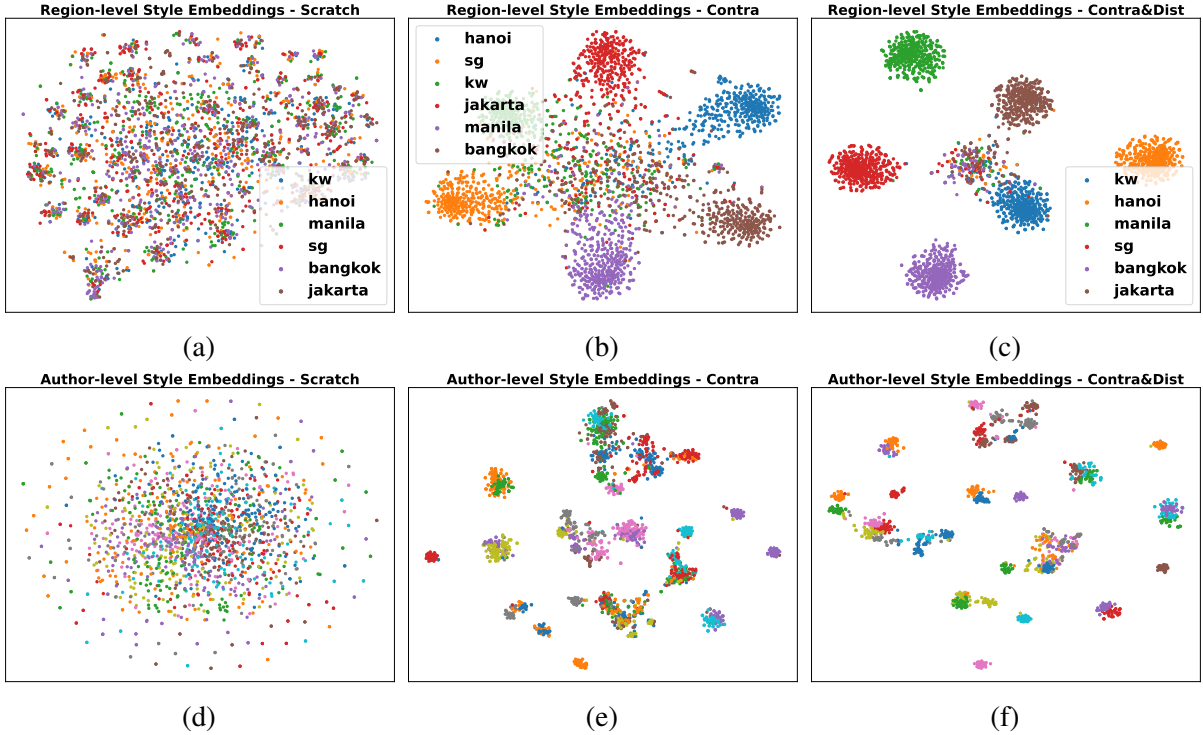


Figure 2: t-SNE visualization of posts from Regional Tweets and CCAT50. Specifically, we select 100 posts from each region in the Regional Tweets dataset and 50 posts from each author in CCAT50. The top three visualizations display the posts from Regional Tweets, while the bottom three pertain to CCAT50.

We train ContrastDistAA in two stages. In the first stage, the style encoder was the sole focus, trained using supervised contrastive loss over 30 epochs. The subsequent stage marked the joint training of both the content and style encoders. This phase, extending for an additional 20 epochs, employed supervised contrastive loss alongside a mutual information estimator. The implementation of the mutual information estimator was based on the source code¹ provided by (Cheng et al., 2020). Consistency in training parameters was maintained throughout, with a learning rate set at 1e-3 and a batch size of 32 for both stages. This setup ensured a balanced and rigorous training process for the

¹<https://github.com/Linear95/CLUB>

ContrastDistAA model.

4.4 Experiment Results

In our study, the efficacy of the ContrastDistAA model was thoroughly assessed on both regional and individual-level datasets, with its performance benchmarked against a range of established baseline models. The comparative results, evaluated using F1 scores, are detailed in Table 2.

The ContrastDistAA model consistently exhibited superior performance across these datasets. For instance, within the Regional Tweets dataset, it attained a Micro F1 score of 0.55, representing a notable 7% improvement compared to the next closest model, BertAA. In the context of the CCAT50 dataset, ContrastDistAA surpassed all

452 baselines in every evaluated metric, achieving a
453 significant 16% improvement in Micro F1 scores.
454 The model also demonstrated exceptional perfor-
455 mance on the Twitter1000 dataset, registering a
456 substantial 29% increase in F1 scores. Further-
457 more, on the IMDB62 dataset, ContrastDistAA
458 achieved a 6.7% improvement in performance, in-
459 dicative of its robustness even in the presence of
460 textual complexity. These results collectively af-
461 firm the ContrastDistAA model’s capability in ef-
462 fectively discerning writing styles at both regional
463 and individual levels, thereby establishing it as a
464 state-of-the-art benchmark in the AA tasks.

465 Interestingly, we also noted that the models’ F1
466 scores are generally lower for the Regional Tweets
467 dataset, suggesting the difficulty of the region-level
468 AA task. The individual authors typically have
469 more distinct and consistent writing styles com-
470 pared to a group of authors from a region. This
471 uniqueness in individual writing styles makes it
472 easier for models to attribute authorship accurately,
473 leading to higher F1 scores. In contrast, regional-
474 level AA deals with broader, less distinct writing
475 styles shared by a group, which can be more chal-
476 lenging to differentiate.

477 4.5 Ablation Study

478 We also conduct an ablation study, which aimed to
479 assess the impact of the dual-stage training process
480 on ContrastDistAA. This study involved compar-
481 ing the model’s performance after the initial train-
482 ing phase, which utilized solely contrastive loss,
483 against its performance following the second train-
484 ing stage that integrated both contrastive and disen-
485 tangle-ment losses. The results, detailed in last two
486 rows of Table 2, emphasize the significant contri-
487 bution of representation disentanglement learning
488 to the model’s efficacy.

489 Crucially, the findings reveal that ContrastDis-
490 tAA demonstrates an improvement in F1 scores
491 when the disentanglement loss is incorporated in
492 the second training stage, compared to the model
493 trained only with contrastive loss. This improve-
494 ment underscores the value of the second training
495 stage in enhancing the model’s capability. By ef-
496 fectively separating content-related elements from
497 style-related information in the training process,
498 the model becomes more adept at isolating and re-
499 cognizing distinctive stylistic features inherent to
500 different regional writings. This separation is key
501 to the improved performance, illustrating the effec-
502 tiveness of the comprehensive two-stage training

approach in ContrastDistAA. 503

504 4.6 Qualitative Analysis

505 To demonstrate the efficacy of ContrastDistAA,
506 we employed the t-SNE algorithm (Van der Maaten
507 and Hinton, 2008) to visually represent post style
508 embeddings in two-dimensional space. This vi-
509 sualization aimed to show how different training
510 methodologies influence the distribution of post
511 representations. We selected 100 posts from each
512 region in the Regional Tweets dataset and 50 posts
513 per author from the CCAT50 dataset, extracting
514 their latent representations using three approaches:
515 (i) BERT in its basic form, (ii) a style encoder
516 trained with contrastive loss, and (iii) a style en-
517 coder trained using both contrastive loss and mu-
518 tual information.

519 Figure 2 presents these representations. The first
520 three visualizations pertain to posts from the Re-
521 gional Tweets dataset, while the latter three focus
522 on the CCAT50 dataset. Notably, with the applica-
523 tion of contrastive loss, distinct clusters emerge, in-
524 dicating the style encoder’s ability to capture style
525 information effectively. However, challenges are
526 evident, such as the central clustering in Figure 2
527 (b), reflecting the limitations of contrastive learning
528 with complex samples. The incorporation of mu-
529 tual information for disentangling content and style
530 in latent space results in more distinct clustering
531 patterns, as seen in Figure 2 (c). This suggests that
532 the integration of both contrastive and disentanglement
533 learning notably enhances the style encoder’s
534 capability to discern style information, thereby im-
535 proving its application in AA tasks.

536 4.7 Case Studies for Regional AA

537 To highlight the unique writing styles prevalent in
538 different regions, we conducted a linguistic analy-
539 sis of posts from these areas. This involved select-
540 ing three posts from each region and calculating the
541 cosine similarity between their representations and
542 the corresponding regional style representations,
543 providing insights into how closely these posts
544 align with predominant regional writing styles.

545 Our analysis revealed distinct linguistic features
546 characteristic of each region, often embodied in
547 specific words or expressions that encapsulate re-
548 gional nuances and evoke emotional responses. For
549 instance, authors from Bangkok frequently con-
550 clude sentences with unique words such as “kub”,
551 “naka”, “krub”, or “na” adding an expressive and
552 emotive quality to their writing. In Jakarta, au-

Regions	Sentences	Similarity
Bangkok	1. @USER thank u naa	0.825
	2. @USER You're very welcome I feel honored and very happy . ka pleading_face two_hearts	0.977
	3. @USER You make all of us lazy people feel ashamed on a Sunday morning na krub .	0.990
Hanoi	1. isit Indonesian #Booth in Ly Thao To Park , DATE	0.995
	2. Those light is fierce ! #welldone @USER Trang Tien Plaza HTTPURL	0.996
	3. try some coconut coffee hot_beverage USER Cong Caphe HTTPURL	0.996
Jakarta	1. @USER Serem amat :loudly_crying_face:	0.991
	2. @USER batman who laughs lumayan lah atleast	0.951
	3. Mantul the babbies nyusul the daddies	0.992
Manila	1. Salamat sa live selling at unboxing ! Lol char . Love you bestie ! Congratulations ! HTTPURL	0.996
	2. Wow , salamat po sa Dios To God be the Glory sparkles Are Your Prayers Heard #PureDoc-trinesOfChrist HTTPURL	0.996
	3. DATE nabudol ako sa film life . Excited for youuuuuu . @USER Stay Broke , Shoot Film . HTTPUR . HTTPURL	0.959
Singapore	1. STOp . the tarot card readings gotta STOOOOOOOop pls lah	0.857
	2. So much things on my mind rn ! Inshallah all goes well	0.651
	3. @USER i no have scandal leh u my one and only	0.984
Kuala Lumpur	1. Pusing lah kot mana pun , no one else is calling it democratic . Except PN of course	0.908
	2. Say goodbye to grainy spycam footage . Tak main lah video quality Nokia	0.501
	3. adut saya order 138 utk pastikan bontot staff saya 8p m 5pm tak ke Pavilion	0.964

Table 3: Examples showcasing the unique writing expressions (highlighted in yellow) from each region. The similarity score is the cosine similarity between the post representation and the region style embedding.

553 thors use expressions like "lumayan" to indicate a
554 moderate experience, "seem amat" for excitement,
555 and "mantul" to denote something extraordinary,
556 showcasing the rich and diverse writing style of
557 this region. Hanoi's writing style, influenced by
558 the modern Latin script and its use of diacritical
559 marks, often features Vietnamese words without
560 these marks. This use reflects a blend of traditional
561 and contemporary linguistic practices, allowing for
562 effective communication while honoring the lin-
563 guistic heritage and subtleties of the region. These
564 findings underscore the distinct linguistic identities
565 of each region, as mirrored in their writing styles.

566 5 Conclusion

567 In this study, we introduced ContrastDistAA, a
568 model designed to effectively separate content and
569 style information, thereby enhancing AA perfor-
570 mance. A significant contribution of our research
571 is the introduction of the regional-level AA task,
572 along with a dedicated dataset to evaluate AA meth-
573 ods in this new context. Through comprehensive
574 experiments, ContrastDistAA was benchmarked
575 against state-of-the-art AA techniques, demonstrat-

ing its superior performance in both individual-
576 level and regional-level AA tasks. 577

The results from our case studies indicate that
578 ContrastDistAA is adept at identifying unique lin-
579 guistic features indicative of regional writing styles.
580 Specifically, the contrastive learning and represen-
581 tation disentanglement approach have helped to
582 effectively segregate content from stylistic features
583 for AA tasks. This capability is crucial for un-
584 derstanding how linguistic styles and cultural in-
585 fluences interplay in AA. Our research addresses
586 a previously unexplored aspect of AA and offers
587 fresh perspectives on the relationship between lin-
588 guistic styles and cultural elements. 589

For future work, we will focus on further explor-
590 ing regional and cultural writing styles. We aim
591 to include a broader range of cultural characteris-
592 tics and regional diversity, thereby enhancing the
593 understanding of AA in diverse linguistic and cul-
594 tural contexts. This ongoing research will continue
595 to expand the horizons of AA, contributing to a
596 deeper understanding of the intricate relationship
597 between authorship, language, and culture. 598

6 Limitations

This study makes noteworthy contributions to the field of Authorship Attribution (AA), but it also acknowledges two key limitations. The first limitation pertains to the methodology of obtaining style representations for regions and authors, which is based on averaging post representations. This approach, while practical, is susceptible to the cluster center shift problem, especially when outliers are included in the calculations. Outliers can significantly skew the average, leading to potential misrepresentations of the typical writing style of a region or an author.

The second limitation is the geographical scope of the dataset used. The Regional Tweets dataset is confined to six regions within Southeast Asian countries, which, while providing valuable regional insights, limits the broader applicability and generalizability of the study’s findings. To enhance the scope and robustness of future research in AA, it would be beneficial to include more diverse regions from various countries and cultural areas. This expansion would offer a more comprehensive understanding of the diverse linguistic and stylistic nuances that characterize writing styles globally, and contribute to the development of AA methods that are universally relevant and sensitive to regional and cultural variations.

References

Opeyemi Aborisade and Mohd Anwar. 2018. Classification for authorship of tweets by comparing logistic regression and naive bayes classifiers. In *2018 IEEE International Conference on Information Reuse and Integration (IRI)*, pages 269–276. IEEE.

Malik Altakrori, Jackie Chi Kit Cheung, and Benjamin C. M. Fung. 2021. [The topic confusion task: A novel evaluation scenario for authorship attribution](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 4242–4256, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeshwar, Sherjil Ozair, Yoshua Bengio, Aaron Courville, and Devon Hjelm. 2018. Mutual information neural estimation. In *International conference on machine learning*, pages 531–540. PMLR.

Janek Bevendorff, Benno Stein, Matthias Hagen, and Martin Potthast. 2019. Generalizing unmasking for short texts. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language*

Technologies, Volume 1 (Long and Short Papers), pages 654–659. 650
651

Sebastian Bischoff, Niklas Deckers, Marcel Schliebs, Ben Thies, Matthias Hagen, Efstathios Stamatatos, Benno Stein, and Martin Potthast. 2020. The importance of suppressing domain style in authorship analysis. *arXiv preprint arXiv:2005.14714*. 652
653
654
655
656

Benedikt Boenninghoff, Steffen Hessler, Dorothea Kolossa, and Robert M Nickel. 2019a. Explainable authorship verification in social media via attention-based similarity learning. In *2019 IEEE International Conference on Big Data (Big Data)*, pages 36–45. IEEE. 657
658
659
660
661
662

Benedikt Boenninghoff, Robert M Nickel, Steffen Zeiler, and Dorothea Kolossa. 2019b. Similarity learning for authorship verification in social media. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2457–2461. IEEE. 663
664
665
666
667
668

Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR. 669
670
671
672
673

Xi Chen, Yan Duan, Rein Houthoofd, John Schulman, Ilya Sutskever, and Pieter Abbeel. 2016. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. *Advances in neural information processing systems*, 29. 674
675
676
677
678

Pengyu Cheng, Weituo Hao, Shuyang Dai, Jiachang Liu, Zhe Gan, and Lawrence Carin. 2020. Club: A contrastive log-ratio upper bound of mutual information. In *International conference on machine learning*, pages 1779–1788. PMLR. 679
680
681
682
683

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186. 684
685
686
687
688
689
690
691

Steven HH Ding, Benjamin CM Fung, Farkhund Iqbal, and William K Cheung. 2017. Learning stylometric representations for authorship analysis. *IEEE transactions on cybernetics*, 49(1):107–121. 692
693
694
695

Hicham El Boukkouri, Olivier Ferret, Thomas Lavergne, Hiroshi Noji, Pierre Zweigenbaum, and Jun’ichi Tsujii. 2020. [CharacterBERT: Reconciling ELMo and BERT for word-level open-vocabulary representations from characters](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6903–6915, Barcelona, Spain (Online). International Committee on Computational Linguistics. 696
697
698
699
700
701
702
703
704

705	Maël Fabien, Esau Villatoro-Tello, Petr Motliceck, and Shantipriya Parida. 2020. BertAA : BERT fine-tuning for authorship attribution . In <i>Proceedings of the 17th International Conference on Natural Language Processing (ICON)</i> , pages 127–137, Indian Institute of Technology Patna, Patna, India. NLP Association of India (NLP AI).	26–29, 2021, <i>Proceedings, Part I 22</i> , pages 403–411. Springer.	760 761
712	Zhenxin Fu, Xiaoye Tan, Nanyun Peng, Dongyan Zhao, and Rui Yan. 2017. Style transfer in text: Exploration and evaluation .	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. <i>arXiv preprint arXiv:1907.11692</i> .	762 763 764 765 766
715	Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2022. Simcse: Simple contrastive learning of sentence embeddings .	Zhi Liu, Zongkai Yang, Sanya Liu, and Wenting Meng. 2012. A novel random subspace method for online writeprint identification. <i>J. Comput.</i> , 7(12):2997–3004.	767 768 769 770
718	John Giorgi, Osvald Nitski, Bo Wang, and Gary Bader. 2021. Declutr: Deep contrastive learning for unsupervised textual representations .	Andrei Manolache, Florin Brad, Elena Burceanu, Antonio Barbalau, Radu Ionescu, and Marius Popescu. 2021. Transferring bert-like transformers’ knowledge for authorship verification. <i>arXiv preprint arXiv:2112.05125</i> .	771 772 773 774 775
721	Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. 2020. Momentum contrast for unsupervised visual representation learning . In <i>2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 9726–9735.	Tempestt Neal, Kalaivani Sundararajan, Aneez Fatima, Yiming Yan, Yingfei Xiang, and Damon Woodard. 2017. Surveying stylometry techniques and applications. <i>ACM Computing Surveys (CSuR)</i> , 50(6):1–36.	776 777 778 779
726	Xinyu Hu, Weihang Ou, Sudipta Acharya, Steven HH Ding, Ryan D’Gama, and Hanbo Yu. 2023. Tdrlm: Stylometric learning for authorship verification by topic-debiasing. <i>Expert Systems with Applications</i> , 233:120745.	Christian Oliva, Santiago Palmero Muñoz, Luis F Lago-Fernández, and David Arroyo. 2022. Improving lstms’ under-performance in authorship attribution for short texts. In <i>Proceedings of the 2022 European Interdisciplinary Cybersecurity Conference</i> , pages 99–101.	780 781 782 783 784 785
731	Zhiqiang Hu, Roy Ka-Wei Lee, Lei Wang, Ee-peng Lim, and Bo Dai. 2020. Deepstyle: User style embedding for authorship attribution of short texts. In <i>Web and Big Data: 4th International Joint Conference, APWeb-WAIM 2020, Tianjin, China, September 18-20, 2020, Proceedings, Part II 4</i> , pages 221–229. Springer.	Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation learning with contrastive predictive coding. <i>arXiv preprint arXiv:1807.03748</i> .	786 787 788
738	Olivier J. Hénaff, Aravind Srinivas, Jeffrey De Fauw, Ali Razavi, Carl Doersch, S. M. Ali Eslami, and Aaron van den Oord. 2020. Data-efficient image recognition with contrastive predictive coding .	Ben Poole, Sherjil Ozair, Aaron Van Den Oord, Alex Alemi, and George Tucker. 2019. On variational bounds of mutual information. In <i>International Conference on Machine Learning</i> , pages 5171–5180. PMLR.	789 790 791 792 793
742	Farkhund Iqbal, Rachid Hadjidj, Benjamin CM Fung, and Mourad Debbabi. 2008. A novel approach of mining write-prints for authorship attribution in e-mail forensics. <i>digital investigation</i> , 5:S42–S51.	Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. <i>arXiv preprint arXiv:1908.10084</i> .	794 795 796
746	Fereshteh Jafariakinabad, Sansiri Tarnpradab, and Kien A Hua. 2019. Syntactic recurrent neural network for authorship attribution. <i>arXiv preprint arXiv:1902.09723</i> .	Rafael A. Rivera-Soto, Olivia Elizabeth Miano, Juanita Ordonez, Barry Y. Chen, Aleem Khan, Marcus Bishop, and Nicholas Andrews. 2021. Learning universal authorship representations . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 913–919, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	797 798 799 800 801 802 803 804
750	Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. 2021. Supervised contrastive learning .	Joshua David Robinson, Ching-Yao Chuang, Suvrit Sra, and Stefanie Jegelka. 2021. Contrastive learning with hard negative samples . In <i>International Conference on Learning Representations</i> .	805 806 807 808
754	Jianbo Liu, Zhiqiang Hu, Jiasheng Zhang, Roy Ka-Wei Lee, and Jie Shao. 2021. A syntax-aware encoder for authorship attribution. In <i>Web Information Systems Engineering–WISE 2021: 22nd International Conference on Web Information Systems Engineering, WISE 2021, Melbourne, VIC, Australia, October</i>	Alexey Romanov, Anna Rumshisky, Anna Rogers, and David Donahue. 2019. Adversarial decomposition of text representation .	809 810 811
759		Chakaveh Saedi and Mark Dras. 2021. Siamese networks for large-scale author identification. <i>Computer Speech & Language</i> , 70:101241.	812 813 814

815	Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. <i>arXiv preprint arXiv:1910.01108</i> .	869
816		870
817		871
818		872
819	Upendra Sapkota, Thamar Solorio, Manuel Montes, Steven Bethard, and Paolo Rosso. 2014. Cross-topic authorship attribution: Will out-of-topic data help? In <i>Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers</i> , pages 1228–1237.	873
820		874
821		875
822		876
823		877
824		878
825	Yunita Sari. 2018. <i>Neural and Non-neural Approaches to Authorship Attribution</i> . Ph.D. thesis, University of Sheffield.	879
826		
827		
828	Roy Schwartz, Oren Tsur, Ari Rappoport, and Moshe Koppel. 2013. Authorship attribution of micro-messages. In <i>Proceedings of the 2013 Conference on empirical methods in natural language processing</i> , pages 1880–1891.	880
829		881
830		882
831		883
832		
833	Yanir Seroussi, Ingrid Zukerman, and Fabian Bohnert. 2011. Authorship attribution with latent dirichlet allocation. In <i>Proceedings of the fifteenth conference on computational natural language learning</i> , pages 181–189.	884
834		885
835		886
836		887
837		888
838	Yanir Seroussi, Ingrid Zukerman, and Fabian Bohnert. 2014. Authorship attribution with topic models. <i>Computational Linguistics</i> , 40(2):269–310.	889
839		890
840		891
841	Tianxiao Shen, Jonas Mueller, Regina Barzilay, and Tommi Jaakkola. 2020. Educating text autoencoders: Latent representation guidance via denoising.	892
842		893
843		
844	Prasha Shrestha, Sebastian Sierra, Fabio A González, Manuel Montes, Paolo Rosso, and Thamar Solorio. 2017. Convolutional neural networks for authorship attribution of short texts. In <i>Proceedings of the 15th conference of the European chapter of the association for computational linguistics: Volume 2, short papers</i> , pages 669–674.	894
845		895
846		896
847		
848		
849		
850		
851	Efstathios Stamatatos and Moshe Koppel. 2011. Plagiarism and authorship analysis: introduction to the special issue. <i>Language Resources and Evaluation</i> , 45:1–4.	
852		
853		
854		
855	Antônio Theóphilo, Luís AM Pereira, and Anderson Rocha. 2019. A needle in a haystack? harnessing onomatopoeia and user-specific stylometrics for authorship attribution of micro-messages. In <i>ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 2692–2696. IEEE.	
856		
857		
858		
859		
860		
861		
862	Jacob Tyo, Bhuwan Dhingra, and Zachary C Lipton. 2022. On the state of the art in authorship attribution and authorship verification. <i>arXiv preprint arXiv:2209.06869</i> .	
863		
864		
865		
866	Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-sne. <i>Journal of machine learning research</i> , 9(11).	
867		
868		
	Andrew Wang, Cristina Aggazzotti, Rebecca Kotula, Rafael Rivera Soto, Marcus Bishop, and Nicholas Andrews. 2023. Can authorship representation learning capture stylistic features? <i>Transactions of the Association for Computational Linguistics</i> , 11:1416–1431.	
	Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. 2018. Unsupervised feature learning via non-parametric instance discrimination. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 3733–3742.	
	Siyang Yuan, Pengyu Cheng, Ruiyi Zhang, Weituo Hao, Zhe Gan, and Lawrence Carin. 2021. Improving zero-shot voice style transfer via disentangled representation learning.	
	Richong Zhang, Zhiyuan Hu, Hongyu Guo, and Yongyi Mao. 2018. Syntax encoding with application in authorship attribution. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 2742–2753, Brussels, Belgium. Association for Computational Linguistics.	
	Wanwan Zheng and Mingzhe Jin. 2023. A review on authorship attribution in text mining. <i>Wiley Interdisciplinary Reviews: Computational Statistics</i> , 15(2):e1584.	
	Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A. Efros. 2020. Unpaired image-to-image translation using cycle-consistent adversarial networks.	