

# StyloThai: A Scalable Framework for Stylometric Authorship Identification of Thai Documents

RAHEEM SARWAR and THANASARN PORTHAVEEPONG, VISTEC  
 ATTAPOL RUTHERFORD, Chulalongkorn University  
 THANAWIN RAKTHANMANON, Kasetsart University and VISTEC  
 SARANA NUTANONG, VISTEC

Authorship identification helps to identify the true author of a given anonymous document from a set of candidate authors. The applications of this task can be found in several domains, such as law enforcement agencies and information retrieval. These application domains are not limited to a specific language, community, or ethnicity. However, most of the existing solutions are designed for English, and a little attention has been paid to Thai. These existing solutions are not directly applicable to Thai due to the linguistic differences between these two languages. Moreover, the existing solution designed for Thai is unable to (i) handle outliers in the dataset, (ii) scale when the size of the candidate authors set increases, and (iii) perform well when the number of writing samples for each candidate author is low. We identify a stylometric feature space for the Thai authorship identification task. Based on our feature space, we present an authorship identification solution that uses the probabilistic  $k$  nearest neighbors classifier by transforming each document into a collection of point sets. Specifically, this document transformation allows us to (i) use set distance measures associated with an outlier handling mechanism, (ii) capture stylistic variations within a document, and (iii) produce multiple predictions for a query document. We create a new Thai authorship identification corpus containing 547 documents from 200 authors, which is significantly larger than the corpus used by the existing study (an increase of 32 folds in terms of the number of candidate authors). The experimental results show that our solution can overcome the limitations of the existing solution and outperforms all competitors with an accuracy level of 91.02%. Moreover, we investigate the effectiveness of each stylometric features category with the help of an ablation study. We found that combining all categories of the stylometric features outperforms the other combinations. Finally, we cross compare the feature spaces and classification methods of all solutions. We found that (i) our solution can scale as the number of candidate authors increases, (ii) our method outperforms all the competitors, and (iii) our feature space provides better performance than the feature space used by the existing study.

CCS Concepts: • **Computing methodologies** → **Language resources**; **Supervised learning by classification**; *Classification and regression trees*; • **Information systems** → *Content analysis and feature selection*; *Information extraction*; • **Applied computing** → *Investigation techniques*; *Evidence collection, storage and analysis*;

The research was partially supported by the Digital Economy Promotion Agency (project# MP-62-0003), and the Thailand Research Fund and Office of the Higher Education Commission (MRG6180266).

Authors' addresses: R. Sarwar, T. Porthaveepong, and S. Nutanong (corresponding author), School of Information Science and Technology, VISTEC, Wangchan Valley 555 Moo 1 Payupnai, Wangchan, Rayong, Thailand, 21210; emails: {raheem.s, thanasarn.p, snutanon}@vistec.ac.th; A. Rutherford, Department of Linguistics at Faculty of Arts Chulalongkorn University, Bangkok, Thailand; email: attapolrutherford@gmail.com; T. Rakthanmanon, Department of Computer Engineering, Kasetsart University, Thailand, and School of Information Science and Technology, VISTEC, Wangchan Valley 555 Moo 1 Payupnai, Wangchan, Rayong, Thailand, 21210; email: thanawin.r@ku.ac.th.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](https://permissions.acm.org).

© 2020 Association for Computing Machinery.

2375-4699/2020/01-ART36 \$15.00

<https://doi.org/10.1145/3365832>

Additional Key Words and Phrases: Authorship analysis, stylometry, similarity search, Thai authorship identification

#### ACM Reference format:

Raheem Sarwar, Thanasarn Porthaveepong, Attapol Rutherford, Thanawin Rakthanmanon, and Sarana Nutanong. 2020. *StyloThai: A Scalable Framework for Stylometric Authorship Identification of Thai Documents*. *ACM Trans. Asian Low-Resour. Lang. Inf. Process.* 19, 3, Article 36 (January 2020), 15 pages. <https://doi.org/10.1145/3365832>

## 1 INTRODUCTION

Stylometry is the science of measuring the writing style of an author [1, 7, 14]. Stylometry relies on the observation that each author has a unique writing style that can help differentiate among the documents written by different authors [29]. This concept is known as the *authorial fingerprint* [23, 29, 32]. Coulthard [4] describes that every author has his or her own form of the language, which is known as *idiolect*. An authors' idiolect manifests itself distinctive and cumulatively unique rule-governed choices for the written communication. Specifically, every author has stored a large set of vocabulary built up over many years. The vocabulary set of an author may differ considerably or slightly from the vocabulary set of all other authors. The difference among the vocabulary sets of authors occurs in terms of the stored vocabulary items, passive vocabulary items, and, most importantly, their preferences for selecting and combining these items for written communication. These differences can help identify authors. One prominent task performed by using stylometry is authorship identification, which identifies the original author of a given anonymous document, and is formally defined as follows.

*Definition 1.1 (Authorship Identification).* “Given an anonymous/disputed text  $x$ , a set of candidate authors  $Y$ , and their writing samples  $X$ , identify the most likely author of  $x$  in  $Y$  by analyzing the writing samples in  $X$  and comparing them with  $x$ ” [23].

Applications of the authorship identification task span across several areas, such as *intelligence agencies work*, where authorship identification can help in linking the intercepted messages to known enemies [1, 21, 29]; *criminal law*, where authorship identification can help in identifying the true author of harassing letters and ransom notes [29]; and *plagiarism detection*, where an authorship identification solution can help in identifying the true authors of student submissions [23]. Moreover, these days, managing large text repositories has become a major challenge and has received significant attention by researchers from several areas, such as web information management [25], natural language processing [34], and information retrieval [38].

The application domains of authorship identification are not limited to a specific language, community, or ethnicity. Thai is a member of the Kra-Dai languages family. The Kra-Dai languages<sup>1</sup> include Thai, Lao, and the tonal languages spoken in Southeast Asia, Northeast India, and Southern China. More than 90 million people speak Kra-Dai languages, and Thai is the most widely spoken Kra-Dai language. However, most of the existing authorship identification solutions are designed for English [15, 16, 18]. These solutions are not directly applicable to Thai due to linguistic differences between English and Thai. For example, unlike English [15, 16, 18], (i) Thai has 44 consonants and 4 tone marks; (ii) Thai has 18 vowel symbols that create many compound vowels, and few special symbols; (iii) Thai text samples do not have word/sentence boundaries; and (iv) the first person pronunciation in Thai can be gender specific and gender-neutral (i.e., men tend to use the former and women tend to use the latter) [18]. These characteristics of Thai make

<sup>1</sup> Available from [https://en.wikipedia.org/wiki/Kra-Dai\\_languages](https://en.wikipedia.org/wiki/Kra-Dai_languages).

the stylometric features extraction process noisier in comparison to English and require a solution associated with outlier handling techniques (see Section 4.1 for more details). To the best of our knowledge, only one solution [18] has been proposed to perform authorship identification for Thai, and we call it *SVM-CWPEST* for short (see Section 2.1.3 for more details about this existing solution). However, there are several limitations to this solution, which are described next.

*Limitations of the existing study.*

- (1) *Low accuracy of language processing tools.* The aforementioned unique characteristics of Thai, such as how Thai text does not contain any word/sentence boundary, makes it harder for Thai language processing tools such as TLTK<sup>2</sup> to yield a high accuracy [5, 6, 16, 18]. Consequently, the stylometric features extraction process for Thai is noisier in comparison to English. However, the existing solution (*SVM-CWPEST*) [18] is not associated with any noise handling mechanism. We aim at designing an authorship identification solution for Thai that can mitigate the effect of outliers in the dataset to improve the authorship identification accuracy.
- (2) *Small number of writing samples per author.* The existing solution (*SVM-CWPEST*) [18] is unable to perform well when the average number of writing samples for a candidate author is between 2 and 3 (see Section 5.2 for more details). Using a large number of different documents (i.e., writing samples) enhances the ability of capturing the stylistic variation information of an author. For example, the existing study [18] uses 25 writing samples for each candidate author. However, such a large number of writing samples may not be available for each candidate author in real-world scenarios. Thus, we aim at designing an authorship identification solution for Thai that can perform well (i.e., achieve more than 90% accuracy) when the average number of writing samples for each candidate author is low (i.e., less than 3).
- (3) *Large size of candidate author set.* The existing solution (*SVM-CWPEST*) [18] drastically drops the accuracy when the size of candidate authors set increases (see Section 5.2 for more details). Moreover, the existing study [18] is limited to six candidate authors only. However, in a real-world scenario, such as plagiarism detection in student submissions, there can be hundreds of candidate authors. We aim at designing an authorship identification solution for Thai that can handle a large number of candidate authors.

In this investigation we identify a stylometric feature space (LSS) to perform authorship identification on Thai. Specifically, our feature space (LSS) consists of 46 stylometric features that can be organized into three main categories, including 27 lexical features (L), 17 syntactic features (S), and 2 structural features (S). These features are explained in Section 4.1. Our feature space (LSS) is better than the feature space used in existing work (CWPEST) [18]. This is because, unlike CWPEST, the LSS feature space contains syntactic features (i.e., part-of-speech (POS)-based features), which can play an important role in distinguishing between documents written by different authors [14, 23, 32] (see Section 5.2 for experimental results).

*Research questions.* In addition to addressing the aforementioned limitations of the existing study [18], we answer the following research questions in this article:

- *Research question 1.* Recall that, unlike the feature space used by existing work (CWPEST) [18] that does not contain syntactic features, our stylometric feature space (LSS) contains syntactic features in addition to lexical and structural features. Thus, we investigate

<sup>2</sup><https://pypi.org/project/tltk/>.

this question: How important is it to use syntactic features for the Thai authorship identification process? In addition to this, we also investigate the importance of each category of the stylometric features to perform Thai authorship identification with the help of an ablation study.

- *Research question 2.* How important is it to use all three categories of the stylometric features in the authorship identification process?
- *Research question 3.* How important is it to use set similarity measures associated with outlier handling mechanisms in comparison to the standard set similarity measure (i.e., without an outlier handling mechanism) in the Thai authorship identification process? The set similarity measures are discussed in Section 4.2.

As for the classification method, we adopt the *probabilistic k nearest neighbors* (PkNN) classifier to perform scalable authorship identification with a limited number of writing samples per candidate author [11]. However, the PkNN is sensitive to noise in the dataset [11]. To address this issue, we use a document transformation model that relies on set similarity search [33] such that the stylistic variations between the text samples can be computed as a set distance [12]. By using a corpus of 547 Thai documents from 200 authors, which is significantly larger than the existing study (an increase of 32 folds in terms of the number of candidate authors), we perform experimental studies to show that our solution can (i) mitigate the effect of outliers in the dataset; (ii) handle a large number of candidate authors; (iii) perform well when the number of writing samples per candidate author is low; and (iv) achieve an accuracy level of 91.02%, which is higher than that of all competitors.

*Summary of our contributions.* The contributions of this work includes the following:

- (1) We formulate an effective stylometric features space (LSS) for the Thai authorship identification task. Based on LSS, we present an authorship identification solution for Thai that can overcome the limitations of existing study and achieve an accuracy level of 91.02%, which is higher than that of all competitors.
- (2) We create a new significantly larger Thai authorship identification corpus than the existing study (i.e., an increase of 32 folds in terms of the number of candidate authors).
- (3) We summarize the findings of our studies here to compare the performance of our solution against (i) SVM-CWPEST, the only existing authorship identification solution for Thai, and (ii) four extensively used classifiers in authorship identification studies in different settings.

The rest of the article is organized as follows. Section 2 reviews existing studies on authorship identification. Section 3 illustrates our corpus. Section 4 describes our solution. Section 5 presents the experimental results. Section 6 contains the concluding remarks.

## 2 LITERATURE REVIEW

Authorship identification is generally performed in two steps. The first step is related to the stylometric features extraction from the true documents of the candidate authors. The stylometric features are the writing style markers that can help distinguish among the documents from different authors. Stylometric features can be organized into three main categories: lexical features, syntactic features, and structural features [20, 23, 24].

- (1) *Lexical features* are the statistical measures of character-based and word-based lexical variations in a document, such as vocabulary richness [29] and word length distributions [23].
- (2) POS tags and function words are examples of *syntactic features* [20].

Table 1. Implementations of the Classification Algorithms and Their Parameters

Method	Implementation	Parameters Changed from the Default
Support vector machines (SVM)	*.functions.LibSVM	—
Naïve Bayes (NB)	*.bayes.NaïveBayes	kernel: Radial Basis
Decision trees (DT)	*.trees.J48	—
Random forests (RF)	*.trees.RandomForests	—

\* Available under WEKA.Classifiers.

- (3) *Structural features* are associated with the organization of the document, such as the average number of words in a sentence or a paragraph [23].

The second step is related to learning a classification model to predict the true author of the anonymous document.

## 2.1 Authorship Identification Methods

**2.1.1 Deep Learning–Based Methods to Authorship Identification.** Recently, deep learning methods have received significant attention by researchers. Specifically, deep learning methods do not require manual features engineering, which makes them more effective over traditional techniques. This is because a right set of features is required to achieve state-of-the-art accuracy [9, 13, 22, 27, 35, 37]. Nevertheless, a tremendous amount of data is required to train deep learning models. In other words, in a convolutional neural network (CNN), the implicit data representation is learned in hidden layers, and based on this learned data representation, the classification is performed at the output layer. Given the huge amount of training data, the learned data representation is better in comparison to handcrafted features and provides better accuracy.

Several existing studies focused on English used deep learning to perform the authorship identification task. For example, Solorio et al. [36] performed authorship identification using a three-layer CNN model based on character bi-grams. They reported an accuracy level of 76.1% using a corpus written by 50 authors where each author has 1,000 samples. Moreover, Ge et al. [8] performed authorship identification using feed-forward neural networks. Ge et al. [8] reported 95% accuracy and noted that this task is too easy to perform due to the availability of huge training data.

*Comparison with our work.* Deep learning methods may achieve high accuracy for the authorship identification task only when a large amount of training data is available. However, in this investigation, we aim at designing an authorship identification solution for Thai that can perform well in data-poor conditions where the average number of writing samples for each candidate author is between 2 and 3.

**2.1.2 Machine Learning Methods to Authorship Identification.** The well-known machine learning methods for authorship identification include random forests (RF), decision trees (DT), naïve bayes (NB), and support vector machines (SVM) [1, 18, 29, 33]. In this work, we compare the accuracy of our method against these well-known, extensively used methods by varying the size of the candidate author set. In addition to directly comparing our solution against these competitors, we cross compare the feature spaces and methods by formulating different solutions (see Section 5.2 for more details). The implementation details of these methods are given in Table 1. Specifically, we use WEKA’s implementation as given in Table 1. Among these methods, LibSVM is not available directly in WEKA, and we included it manually.

**2.1.3 Thai Authorship Identification (SVM-CWPEST).** We note that most existing solutions are designed for English. To the best of our knowledge, there is only one study that is focused on



authorship identification of Thai text. This study is performed on a corpus from six authors where each author has 25 text samples [18]. This study extracts 53 features called *CWPEST* from each writing sample, applies the SVM and DT (Weka's J48) classifiers to predict the true author of the anonymous text, and shows that SVM yields better accuracy. We call this Thai authorship identification solution *SVM-CWPEST* for short.

*Comparison with our method.* Note that unlike our method that represents each document as a collection of point sets, the SVM-CWPEST method represents each document as one single data point in a multidimensional space (see Section 4 for more details). As a result, the SVM-CWPEST method is unable to (i) capture the writing style variations within the same document, (ii) produce multiple predictions for the same document, and (iii) apply set distance measures to handle outliers in the dataset. Moreover, unlike the feature space used by existing method (CWPEST), our feature space contains POS-based features in addition to lexical and structural features. Moreover SVM-CWPEST was applied on short text samples. The applications of short-text authorship identification can be found in the social media domain, such as author identification of controversial posts by virtual identities on social media, authorship analysis on Facebook posts, Twitter status, chat conversations, and short message service (SMS) messages [1]. The applications of long-text authorship identification can be found in several areas associated with managing large text repositories and plagiarism detection in student theses. Specifically, retrieving and categorizing documents with respect to their authors and plagiarism detection have been receiving a significant attention by researchers in several areas such as natural language processing [2, 3, 34], web information management [25], and information retrieval [10, 26, 28, 30, 31, 38].

### 3 DATA COLLECTION

There are two main issues associated with authorship identification corpora: (i) the number of publicly available corpora is limited, and (ii) the size of the publicly available corpora is small in terms of the number of candidate authors. To the best of our knowledge, there is no benchmark corpora available for the Thai authorship identification task. To perform experiments, we created a new Thai authorship identification corpus extracted from an online *Dek-D*<sup>3</sup> repository. Our scraper is written in Python and extracts the data in two steps: (i) retrieve all URLs of each author, and (ii) based on the retrieved URLs, extract the documents of each author from the website. Our corpus contains 547 Thai documents from 200 authors where the average length of documents is 25,334 tokens. Moreover, our corpus is significantly larger than the existing study [18] (i.e., an increase of 32 folds in terms of the number of candidate authors). Furthermore, on average, there are 2.73 samples per class (author), which is a more realistic scenario where a large number of writing samples per author may not be available.

### 4 METHODOLOGY

We explain our solution with the help of Figure 1. Our solution consists of four main parts: (i) preprocessing, (ii) set similarity search, (iii) PkNN classification, and (iv) prediction aggregation.

#### 4.1 Preprocessing

The preprocessing part of our solution transforms each document into a collection of fragments (i.e., collection of point sets) using a three-step process [33]: (i) partition each document into fixed-size fragments, (ii) partition each fragment obtained from the first step into fixed-size chunks,<sup>4</sup> and

<sup>3</sup><https://www.dek-d.com>.

<sup>4</sup>A chunk is a collection of tokens.

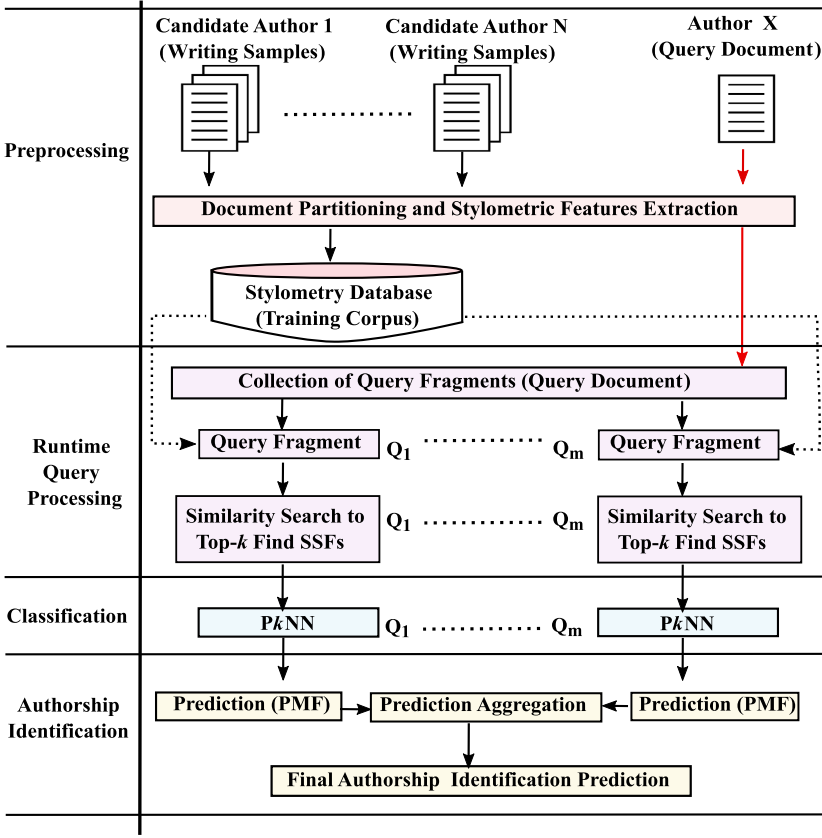


Fig. 1. Overview of the StyloThai framework [33].

(iii) extract the 46 stylometric features (i.e., writing style markers) from each chunk (see Table 2). To obtain reliable stylometric statistics from each fragment and the chunk, we fix their sizes to 7,000 and 700 tokens,<sup>5</sup> respectively. As a result, each chunk is transformed into a point, each fragment is transformed into a point set, and each document is transformed into a collection of point sets in 46-dimensional space. There are three advantages of transforming each document into a collection of point sets. First, we can compute the stylistic variations between text samples as a set distance. Specifically, we can use those set distance measures that are capable of mitigating the effect of outliers in the data such as partial Hausdorff distance (PHD) [12]. Second, we can capture the stylistic variation of an author within a document. This is because each authorship identification prediction is produced by multiple points rather than one single point. Third, we can produce multiple predictions for a query document, which allows us to use only the most certain  $\gamma\%$  predictions of a query document for the prediction aggregation process, as explained later in this section. Once we complete the feature extraction process, we store the feature values into the stylometry database. Our 46-dimensional stylometric feature space can be organized into three categories as shown in Table 3: (i) lexical features (from #1 to #27), (ii) syntactic features (from #28 to #44), and (iii) structural features (from #45 to #46). Note that in our stylometric feature space given in

<sup>5</sup>A token is the content of text separated by a white space character.

Table 2. Example of Stylometric Features Extraction from a Chunk

Thai
(๐๐ พ.ศ. ๒๕๖๒) 07.00 น. @ศูนย์วิทยุ 191 ได้รับแจ้งว่าพบวัตถุต้องสงสัยเป็นระเบิดพร้อมปืนM60จำนวนมากฝังอยู่ใต้ต้นข่อย หลังจากเขาได้ยินเสียงดังตูม บริเวณบ้านไม่มีเลขที่ ในหมู่บ้านเฉลิมพระเกียรติ จึงรีบรถไปตรวจสอบบริเวณดังกล่าว พร้อมกับเจ้าหน้าที่ชุดเก็บกู้
English (Word Based Translation)
(30 MAY 2019) 7 o'clock @ Radio Center 191 was informed that the founded suspected object was a bomb and many M60 guns buried under Tooth brush tree. After he heard loud noises around the unnumbered house in the Chaloem Phra Kiat village, thus hurried to check the area along either the EOD staff.

Table 3, 32 out of 46 features (i.e., except the character-based features from #14 to #27) require the word tokenization and sentence identification, which is a challenging task to perform considering that there are no word/sentence boundaries in Thai. These characteristics of Thai make it harder for Thai language processing tools such as TLTK<sup>6</sup> to yield a high accuracy [5, 6, 16, 18]. Consequently, the stylometric features extraction process for Thai is noisier in comparison to English. Hence, our Thai authorship identification solution is associated with outlier handling mechanisms, as discussed later in this section.

In Table 3,  $N$  represents the count of words,  $V$  represents the count of distinct words,  $V_i$  represents the count of words that occur  $i$  times, and  $C$  represents the total number of characters. As for the lexical features, we identified 13 word-based (from #1 to #13) and 15 character-based (from #14 to #27) stylometric features and computed them from each chunk. These 13 word-based lexical variation features can be considered as language independent [29]. For the rest of the 15 character-based lexical features, 5 features are specific to Thai (i.e., features #16, #17, #21, #23, and #24). As for the syntactic features, we identified 17 features (from #28 to #44) based on the relative frequency of POS categories and computed them from each chunk. Finally, we identified 2 structural features (from #45 to #46) based on the text organization, such as the average number of words per sentence and total number of sentences in a chunk. The word tokenization is performed using DeepCut.<sup>7</sup> The rest of the features are calculated using the Thai Language Toolkit (TLTK).<sup>8</sup> We provide an example of a Thai text sample in Table 2 and show the computed stylometric feature values in Table 3. We measure the effectiveness of each features category for the authorship identification process (see Section 5.2 for experimental results).

## 4.2 Set Similarity Search

While processing a given query document ( $Q$ ), we first apply the preprocessing step of our solution on  $Q$  that transforms it into a collection of point sets. We then execute an independent set similarity query for each query fragment ( $Q$ ) in  $Q$  to retrieve top- $k$  SSFs from the corpus. Note that we execute an individual set similarity query for each  $Q$  in  $Q$ —that is, if a query document results in  $m$  query fragments (point sets), we execute  $m$  independent set similarity queries (see Figure 1). While retrieving the top- $k$  SSFs, we tried three set similarity measures, including (i) standard Hausdorff Distance (SHD), (ii) PHD [12], and modified Hausdorff Distance (MHD) [17] as a proximity measure between two point sets. The SHD between two points sets  $Q$  and  $F$  can be calculated as

<sup>6</sup><https://pypi.org/project/tltk/>.

<sup>7</sup><https://github.com/rkcosmos/deepcut>.

<sup>8</sup><https://pypi.org/project/tltk/>.



Table 3. List of Stylometric Features

Lexical Features			
Stylometric Features	Values	Stylometric Features	Values
1. <b>N</b> : Total number of words	72	2. <b>V</b> : Total number of distinct words	59
3. Average word length	3.53	4. S.D. of word lengths	22.18
5. $\frac{V}{N}$	0.82	6. $VR(K) = \frac{10^4(\sum i^2 V_i - N)}{N^2}$	223.7654321
7. $VR(R) = \frac{V}{\sqrt{N}}$	6.95	8. $VR(C) = \frac{\log V}{\log N}$	0.95
9. $VR(H) = \frac{(100 \log N)}{(1 - V_1)/V}$	-467.26	10. $VR(S) = \frac{V_2}{V}$	0.050847458
11. $VR(k) = \frac{\log V}{\log(\log N)}$	2.81	12. $VR(LN) = \frac{(1 - V^2)}{V^2(\log N)}$	-0.23
13. Entropy of word freq. ditri.	872.03	14. <b>C</b> : Total # chars	254
15. Freq. of alpha chars	1	16. Freq. of Thai chars	219
17. Freq. of Thai numeric chars	6	18. Freq. of Arabic numeric chars	9
19. Freq. of special chars	7	20. Freq. of white spaces	12
21. Freq. of vowel and tone marks	86	22. $\frac{\text{Freq. of alpha char}}{C}$	0.0039
23. $\frac{\text{Freq. of Thai char}}{C}$	0.85	24. $\frac{\text{Freq. of Thai numeric char}}{C}$	0.024
25. $\frac{\text{Freq. of Arabic numeric char}}{C}$	0.035	26. $\frac{\text{Freq. of special char}}{C}$	0.028
27. $\frac{\text{Freq. of white spaces}}{C}$	0.047		
Syntactic Features			
Stylometric Features	Values	Stylometric Features	Values
28. $\frac{\text{Freq. of adjectives}}{N}$	0.014	29. $\frac{\text{Freq. of adpositions}}{N}$	0.042
30. $\frac{\text{Freq. of adverbs}}{N}$	0.028	31. $\frac{\text{Freq. of auxiliaries}}{N}$	0.014
32. $\frac{\text{Freq. of coordinating conjunctions}}{N}$	0.014	33. $\frac{\text{Freq. of determiners}}{N}$	0.014
34. $\frac{\text{Freq. of interjections}}{N}$	0.014	35. $\frac{\text{Freq. of nouns}}{N}$	0.264
36. $\frac{\text{Freq. of numerals}}{N}$	0.042	37. $\frac{\text{Freq. of particles}}{N}$	0.014
38. $\frac{\text{Freq. of pronouns}}{N}$	0.013	39. $\frac{\text{Freq. of proper nouns}}{N}$	0.014
40. $\frac{\text{Freq. of punctuation}}{N}$	0.19	41. $\frac{\text{Freq. of subconjunction}}{N}$	0.042
42. $\frac{\text{Freq. of symbols}}{N}$	0.013	43. $\frac{\text{Freq. of verbs}}{N}$	0.181
44. $\frac{\text{Freq. of other POS}}{N}$	0.014		
Structural Features			
Stylometric Features	Values	Stylometric Features	Values
45. Total number of sentences	11	46. Avg. number of words per sentence	0.153

Note: N represents the count of words, V represents the count of distinct words,  $V_i$  represents the count of words that occur  $i$  times, and C represents the total number of characters.

$$h(Q, F) = \max_{q_i \in Q} \min_{f_j \in F} d(q_i - f_j).$$

In other words, SHD can be calculated by (i) ranking all data points in a query fragment  $Q$  in accordance with the minimum distance to the fragment  $F$  and (ii) selecting the maximum of the minimum distances. Researchers have argued that SHD is sensitive to the noise in the data [12, 17]. To mitigate the noise (outlier) sensitivity issue associated with SHD, researchers formulated two variants of SHD: MHD [17] and PHD [12]. Specifically, the MHD and PHD measures average out the effect of the outlier over the minimum distances falling into a specified range—for

Table 4. Example of the Prediction Aggregation Process

Query Fragment ( $Q$ )	Query Fragment Prediction (PMF)	Entropy
$Q_1$	[ <i>Author A</i> : 0.33, <i>Author B</i> : 0.34, <i>Author C</i> : 0.33]	1.5848
$Q_2^*$	[ <i>Author A</i> : 0.36, <i>Author B</i> : 0.32, <i>Author C</i> : 0.32]	1.5827
$Q_3^*$	[ <i>Author A</i> : 0.32, <i>Author B</i> : 0.35, <i>Author C</i> : 0.33]	1.5840
$Q_4$	[ <i>Author A</i> : 0.33, <i>Author B</i> : 0.34, <i>Author C</i> : 0.33]	1.5848
Final Prediction	[ <i>Author A</i> : 0.34, <i>Author B</i> : 0.335, <i>Author C</i> : 0.325]	—

\*Top most certain  $\gamma$ 50% predictions.

instance, [50%, 100%] (for MHD, the second parameter value is always 100%). The experimental results regarding set distance measures are reported in Section 5.2.

### 4.3 PkNN Classification

We apply PkNN [11] to the retrieved top- $k$  SSFs to make a probabilistic prediction for each query fragment in a query document. Unlike the simple  $k$ NN classifier where the output is one single class (author), the PkNN classifier produces a probability mass function (PMF) over all classes (candidate authors) associated to the retrieved SSFs. We apply the PkNN [11] that utilizes the distance values of the  $k$ NNs (SSFs in this case) to weight the distribution of the probability. An exponential function is used to smooth the distance-probability mapping [29]. The advantages of using PkNN [11] over other classifiers can be summarized as follows. Little or no training is required for classification [19]. Consequently, there is no information loss associated with generalization [11, 29]. This classifier is capable of performing classification with a limited set of samples [29]. Moreover, it allowed us to apply set distance measures capable of mitigating the effect of outliers in the dataset [12].

### 4.4 Prediction Aggregation

The final step of our solution is to merge all of the fragment probabilistic predictions such that one single authorship identification prediction can be produced for the entire  $Q$  [33]. To do so, one can simply compute the average of all fragment probabilistic predictions. However, all of the fragment probabilistic predictions (one for each  $Q$ ) of a query document  $Q$  are not equally useful—for instance, there can be highly uncertain predictions, and including them into the prediction aggregation process may damage the overall accuracy [33]. At this stage, we apply entropy as an uncertainty measure to find the uncertain fragment predictions and eliminate them from the prediction aggregation process [33]. The final probabilistic prediction of the entire  $Q$  is computed as the average PMF of most certain  $\gamma\%$  prediction. An example of this process is given in Table 4. Assume that the value of  $\gamma$  is 50. The top 50% most certain predictions belong to  $Q_2$  and  $Q_3$  as indicated with an asterisk (\*) the low entropy values. The final prediction of the entire query document  $Q$  is calculated as the average PMF of  $Q_2$  and  $Q_3$ .

## 5 PERFORMANCE EVALUATION

### 5.1 Experimental Setup

*Evaluation measures.* Recall that we represent each document as a collection of fragments. Hence, we compute the authorship identification accuracy at two levels—fragment level and document level—as follows:

- *Fragment accuracy:* A fragment authorship prediction is considered correct if the true author of the query document is identified as the most likely author.

Table 5. Default Parameter Values of Our Method

$k$	MHD	PHD	$L$	$l$	$\gamma$
5	(50%, 100%]	(50%, 75%]	7,000 tokens	700 tokens	90%

Table 6. StyloThai Document Accuracy:  
Effect of Feature Types

Lexical	Syntactic	Structural	Accuracy
✓		✓	61.21%
	✓	✓	70.23%
✓	✓		79.83%
✓	✓	✓	<b>91.02%</b>

- *Document accuracy*: An aggregated final authorship prediction of the query document is considered correct if the true author of the query document is identified as the most likely author.

*Parameter setting.* Although we have not shown here, we tried different values for each parameter, and the parameter values given in Table 5 resulted in the best accuracy. The  $k$  value denotes the number of closest stylistically similar fragments identified as a result of the set similarity search query to use for PkNN. The values—(50%,100%] and (50%,75%]—denote the MHD and PHD ranges, respectively. The  $L$  and  $l$  denote the sizes of each fragment and each chunk, respectively. The  $\gamma$  denotes the percentage of predictions that we consider for the prediction aggregation process illustrated in Section 4.4.

*Evaluation strategy.* To evaluate the accuracy of all methods in this investigation, we use five-fold cross validation. Recall that, as for our method, each document is represented as a collection of fragments. To avoid test-train set contamination in the evaluation process of our method, we ensure that when a document is used for testing, it is purely used for testing.

## 5.2 Experimental Results

In this section, we report results from our experimental studies. Note that all experiments are performed using corpora containing a limited number of writing samples per candidate author (i.e., between two and three).

*An ablation study of different features (effect of feature types).* This study provides the answers to the first two questions mentioned in Section 1. As can be seen from Table 6, (i) including syntactic features into lexical + structural increases the authorship identification accuracy from 61.21% to 91.02%; and (ii) combining all categories of stylometric features (i.e., lexical + syntactic + structural) outperforms the other combinations (i.e., (a) lexical + structural, (b) syntactic + structural, and (c) lexical + syntactic). These results indicate that the stylometric information captured by different feature categories is complementary and orthogonal. Consequently, combining all feature categories improves the performance of the authorship identification process. Hence, we confine the rest of the experimental studies to combined categories of the stylometric features only (i.e., lexical + syntactic + structural).

*Effect of set distance measures and the prediction aggregation process.* In this study, by using a corpus containing 547 Thai documents from 200 authors, which is significantly larger than the corpus used by the existing study (an increase of 32 folds in terms of the number of candidate authors), we show that our method can (i) mitigate the effect of outliers in the dataset, (ii) handle

Table 7. Proposed Method Only (StyloThai): The Effect of Set Distance Measures and  $\gamma$  Value

Distance Measure	Accuracy				
	Fragment	Document ( $\gamma = 50\%$ )	Document ( $\gamma = 70\%$ )	Document ( $\gamma = 90\%$ )	Document ( $\gamma = 100\%$ )
SHD	73.22%	78.41%	79.09%	80.12%	78.34%
MHD	84.55%	88.66%	90.21%	90.43%	89.03%
<b>PHD</b>	<b>85.89%</b>	89.15%	90.47%	<b>91.02%</b>	89.14%

Table 8. Document Accuracy: Effect of Candidate Author Set Size

Method	Effect of Number of Candidate Authors			
	50	100	150	200
<b>StyloThai-LSS (Our Method)</b>	<b>92.05%</b>	<b>92.13%</b>	<b>91.49%</b>	<b>91.02%</b>
StyloThai-CWPEST	77.67%	72.36%	69.04%	62.47%
SVM-LSS	44.63%	34.47%	25.42%	18.58%
SVM-CWPEST (Competitive Method)	33.84%	21.61%	12.91%	08.97%
RF-LSS	39.27%	25.39%	19.93%	17.34%
RF-CWPEST	34.91%	20.73%	11.82%	07.84%
NB-LSS	36.43%	24.29%	17.24%	15.97%
NB-CWPEST	29.78%	22.43%	13.79%	09.69%
DT-LSS	35.09%	21.96%	17.56%	16.88%
DT-CWPEST	27.05%	16.44%	14.75%	10.35%

a large number of candidate authors, and (iii) perform well in extreme data-poor conditions. The experimental results obtained using our method are shown in Table 7. As mentioned in Section 4.2, unlike MHD and PHD distance measures, SHD is not associated with an outlier handling mechanism. The fact that SHD is significantly outperformed by MHD and PHD shows that our dataset in fact has noise (outliers) to be handled. Moreover, the results show that PHD has a better outlier handling mechanism than MHD. Due to the obvious accuracy gaps, we only adopt the PHD measure in the rest of the experimental studies. Moreover, the experimental results show that instead of using 100% fragment predictions in the prediction aggregation process, using 90% most certain fragment predictions of each query document provides better accuracy. Due to the obvious accuracy gap between the fragment and document accuracies, we only report document accuracy ( $\gamma = 90\%$ ) in rest of the experimental studies.

*Effect of the candidate author set size.* In this study, we provide the performance comparison between our solution and the competitive solutions by varying the size of the candidate author set from 50 to 200. In addition to directly comparing our solution against the competitors, we cross compare the feature spaces and methods by formulating the following solutions:

- (1) *StyloThai-LSS*: Our feature space (LSS) applied to our method (StyloThai) (proposed solution)
- (2) *StyloThai-CWPEST*: Feature space used by the existing study (CWPEST) [18] applied to StyloThai
- (3) *SVM-LSS*: LSS applied to the SVM method
- (4) *SVM-CWPEST*: CWPEST applied to SVM (existing solution for Thai [18])
- (5) *RF-LSS*: LSS applied to the RF method
- (6) *RF-CWPEST*: CWPEST applied to RF

- (7) *NB-LSS*: LSS applied to the NB method
- (8) *NB-CWPEST*: CWPEST applied to the NB method
- (9) *DT-LSS*: LSS applied to the DT method
- (10) *DT-CWPEST*: CWPEST applied to the DT method.

The experimental results given in Table 8 show that (i) our solution (StyloThai-LSS) can scale as the number of candidate authors increases, (ii) our method (StyloThai) outperforms all of the competitors, and (iii) our feature space (LSS) provides better performance than the competitive feature space (CWPEST). Moreover, regardless of the feature space, there is a significant accuracy gap between our method (StyloThai) and other methods that are not associated with outlier handling mechanisms (i.e., SVM, RF, NB, and DT).

## 6 CONCLUSION

This article presents a scalable solution for authorship identification of Thai documents. The existing solutions designed for English are not directly applicable to Thai due to the linguistic differences between them. Moreover, the existing solution designed for Thai is (i) not associated with any outlier handling mechanism, (ii) unable to scale when the size of the candidate authors set increases, and (iii) cannot perform well in data-poor conditions. By using a corpus of 547 documents written in Thai from 200 authors, which is significantly larger than the corpus used by the existing study, we perform extensive experimental studies to show that our solution can (i) mitigate the effect of outliers in the dataset; (ii) handle a large number of candidate authors in extreme data-poor conditions; and (iii) achieve an accuracy level of 91.02%, which is significantly higher than all competitors. In addition to addressing the aforementioned limitations of the existing study, we answer the following three research questions in this article: (i) How important is it to use syntactic stylometric features in the Thai authorship identification process? (ii) How important is it to use all three categories of stylometric features in the authorship identification process? (iii) How important is it to use set similarity measures associated with outlier handling mechanisms in comparison to the standard set similarity measure (i.e., without an outlier handling mechanism) in the Thai authorship identification process? We found that (i) including the syntactic features into lexical + structural features increases the authorship identification accuracy from 61.21% to 91.02%, and (ii) combining all categories of stylometric features (i.e., lexical + syntactic + structural) outperforms the other combinations (i.e., (a) lexical + structural, (b) syntactic + structural, and (c) lexical + syntactic). These results indicate that the stylometric information captured by different features categories is complementary and orthogonal. Consequently, combining all feature categories improves the performance of the authorship identification process, and (iii) using PHD, which is associated with an outlier handling mechanism, outperforms the SHD by 10.9 percentage points. This article has laid the foundation for future work in the Thai authorship identification task. We hope that this investigation has opened the door for future work on Thai to keep up with the work in other languages.

## REFERENCES

- [1] Malik H. Altakrori, Farkhund Iqbal, Benjamin C. M. Fung, Steven H. H. Ding, and Abdallah Tubaishat. 2019. Arabic authorship attribution: An extensive study on Twitter posts. *ACM Transactions on Asian and Low-Resource Language Information Processing* 18, 1 (2019), Article 5, 51 pages.
- [2] Sophia Ananiadou, Paul Thompson, and Raheel Nawaz. 2013. Enhancing search: Events and their discourse context. In *Proceedings of the International Conference on Intelligent Text Processing and Computational Linguistics*. 318–334.
- [3] Riza Theresa Batista-Navarro, Georgios Kontonatsios, Claudiu Mihăilă, Paul Thompson, Rafal Rak, Raheel Nawaz, Ioannis Korkontzelos, and Sophia Ananiadou. 2013. Facilitating the analysis of discourse phenomena in an interoperable NLP platform. In *Proceedings of the International Conference on Intelligent Text Processing and Computational Linguistics*. 559–571.

- [4] Malcolm Coulthard. 2012. On admissible linguistic evidence. *Journal of Law and Policy* 21 (2012), 441.
- [5] Boonyarit Deewattananon and Usa Sammapun. 2017. Analyzing user reviews in Thai language toward aspects in mobile applications. In *Proceedings of the 14th International Joint Conference on Computer Science and Software Engineering (JCSSE'17)*. 1–6.
- [6] Chenchen Ding, Masao Utiyama, and Eiichiro Sumita. 2019. NOVA: A feasible and flexible annotation system for joint tokenization and part-of-speech tagging. *ACM Transactions on Asian and Low-Resource Language Information Processing* 18, 2 (2019), Article 17, 18 pages.
- [7] Heba El-Fiqi, Eleni Petraki, and Hussein A. Abbass. 2016. Pairwise comparative classification for translator stylometric analysis. *ACM Transactions on Asian and Low-Resource Language Information Processing* 16, 1 (2016), Article 2, 26 pages.
- [8] Zhenhao Ge, Yufang Sun, and Mark J. T. Smith. 2016. Authorship attribution using a neural network language model. In *Proceedings of the 30th AAAI Conference on Artificial Intelligence* 4212–4213.
- [9] Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. 2016. *Deep Learning*. Vol. 1. MIT Press, Cambridge, MA.
- [10] Saeed-Ul Hassan, Raheem Sarwar, and Amina Muazzam. 2016. Tapping into intra-and international collaborations of the organization of Islamic cooperation states across science and technology disciplines. *Science and Public Policy* 43, 5 (2016), 690–701.
- [11] C. C. Holmes and N. M. Adams. 2002. A probabilistic nearest neighbour method for statistical pattern recognition. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 64, 2 (2002), 295–306.
- [12] Daniel P. Huttenlocher, Gregory A. Klanderman, and William Rucklidge. 1993. Comparing images using the Hausdorff distance. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 15, 9 (1993), 850–863.
- [13] Maham Jahangir, Hammad Afzal, Mehreen Ahmed, Khawar Khurshid, and Raheel Nawaz. 2017. An expert system for diabetes prediction using auto tuned multi-layer perceptron. In *Proceedings of the 2017 Intelligent Systems Conference (IntelliSys'17)*. IEEE, Los Alamitos, CA, 722–728.
- [14] Patrick Juola, George K. Mikros, and Sean Vinsick. 2019. A comparative assessment of the difficulty of authorship attribution in Greek and in English. *Journal of the Association for Information Science and Technology* 70, 1 (2019), 61–70.
- [15] Natthawut Kertkeidkachorn, Proadpran Punyabukkana, and Atiwong Suchato. 2016. Acoustic features for hidden conditional random fields-based Thai tone classification. *ACM Transactions on Asian and Low-Resource Language Information Processing* 15, 2 (2016), Article 9, 26 pages.
- [16] Nongnuch Ketui, Thanaruk Theeramunkong, and Chutamanee Onsuwan. 2015. An EDU-based approach for Thai multi-document summarization and its application. *ACM Transactions on Asian and Low-Resource Language Information Processing* 14, 1 (2015), Article 4, 26 pages.
- [17] Rajalida Lipikorn, Akinobu Shimizu, and Hidefumi Kobatake. 1994. A modified Hausdorff distance for object matching. In *Proceedings of the Conference on Pattern Recognition*, Vol. 1. 566–568.
- [18] Rangsan Marukatat, Robroo Somkiadcharoen, Ratthanant Nalintasnai, and Tappasarn Aramboonpong. 2014. Authorship attribution analysis of thai online messages. In *Proceedings of the IEEE International Conference on Information Science and Applications (ICISA'14)*. 1–4.
- [19] Thomas M. Mitchell. 1997. *Machine Learning*. McGraw-Hill, New York, NY.
- [20] Frederick Mosteller and David Wallace. 1964. *Inference and Disputed Authorship: The Federalist*. Addison-Wesley, Reading, MA.
- [21] Arvind Narayanan, Hristo Paskov, Neil Zhenqiang Gong, John Bethencourt, Emil Stefanov, Eui Chul Richard Shin, and Dawn Song. 2012. On the feasibility of Internet-scale author identification. In *Proceedings of the 2012 IEEE Symposium on Security and Privacy (SP'12)*. IEEE, Los Alamitos, CA, 300–314.
- [22] Raheel Nawaz, Paul Thompson, and Sophia Ananiadou. 2012. Identification of manner in bio-events. In *Proceedings of the 8th International Conference on Language Resources and Evaluation (LREC'12)*. 3505–3510.
- [23] Sarana Nutanong, Chenyun Yu, Raheem Sarwar, Peter Xu, and Dickson Chow. 2016. A scalable framework for stylometric analysis query processing. In *Proceedings of the IEEE 16th International Conference on Data Mining (ICDM'16)*. 1125–1130.
- [24] Mathias Payer, Ling Huang, Neil Zhenqiang Gong, Kevin Borgolte, and Mario Frank. 2015. What you submit is who you are: A multimodal approach for deanonymizing scientific publications. *IEEE Transactions on Information Forensics and Security* 10, 1 (2015), 200–212.
- [25] Fuchun Peng, Dale Schuurmans, Shaojun Wang, and Vlado Keselj. 2003. Language independent authorship attribution using character level language models. In *Proceedings of the 10th Conference of the European Chapter of the Association for Computational Linguistics—Volume 1*. 267–274.
- [26] Fahad Sabah, Saeed-Ul Hassan, Amina Muazzam, Sehrish Iqbal, Saira Hanif Soroya, and Raheem Sarwar. 2019. Scientific collaboration networks in Pakistan and their impact on institutional research performance. *Library Hi Tech* 37, 1 (2019), 19–29.



- [27] Ahmad Al Sallab, Ramy Baly, Hazem M. Hajj, Khaled Bashir Shaban, Wassim El-Hajj, and Gilbert Badaro. 2017. AROMA: A recursive deep learning model for opinion mining in Arabic as a low resource language. *ACM Transactions on Asian and Low-Resource Language Information Processing* 16, 4 (2017), Article 25, 20 pages.
- [28] Raheem Sarwar and Saeed-Ul Hassan. 2015. A bibliometric assessment of scientific productivity and international collaboration of the Islamic world in science and technology (S&T) areas. *Scientometrics* 105, 2 (2015), 1059–1077.
- [29] Raheem Sarwar, Qing Li, Thanawin Rakthanmanon, and Sarana Nutanong. 2018. A scalable framework for cross-lingual authorship identification. *Information Sciences* 465 (2018), 323–339.
- [30] Raheem Sarwar and Sarana Nutanong. 2016. The key factors and their influence in authorship attribution. *Research in Computing Science* 110 (2016), 139–150.
- [31] Raheem Sarwar, Saira Hanif Soroya, Amina Muazzam, Fahad Sabah, Sehrish Iqbal, and Saeed-Ul Hassan. 2019. A bibliometric perspective on technology-driven innovation in the Gulf Cooperation Council (GCC) countries in relation to its transformative impact on international business. In *Technology-Driven Innovation in Gulf Cooperation Council (GCC) Countries: Emerging Research and Opportunities*. IGI Global, 49–66.
- [32] Raheem Sarwar, Chenyun Yu, Sarana Nutanong, Norawit Uraileertprasert, Nattapol Vannaboot, and Thanawin Rakthanmanon. 2018. A scalable framework for stylometric analysis of multi-author documents. In *Proceedings of the 23rd International Conference on Database Systems for Advanced Applications (DASFAA'18), Part I*. 813–829.
- [33] Raheem Sarwar, Chenyun Yu, Ninad Tungare, Kanatip Chitavisutthivong, Sukrit Sriratanawilai, Yaohai Xu, Dickson Chow, Thanawin Rakthanmanon, and Sarana Nutanong. 2018. An effective and scalable framework for authorship attribution query processing. *IEEE Access* 6 (2018), 50030–50048.
- [34] Fabrizio Sebastiani. 2006. Classification of text, automatic. *Encyclopedia of Language and Linguistics* 14 (2006), 457–462.
- [35] Matthew Shardlow, Riza Batista-Navarro, Paul Thompson, Raheel Nawaz, John McNaught, and Sophia Ananiadou. 2018. Identification of research hypotheses and new knowledge from scientific literature. *BMC Medical Informatics and Decision Making* 18, 1 (2018), 46.
- [36] Thamar Solorio, Paolo Rosso, Manuel Montes-y-Gómez, Prasha Shrestha, Sebastián Sierra, and Fabio A. González. 2017. Convolutional neural networks for authorship attribution of short texts. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics (EACL'17)—Volume 2: Short Papers*. 669–674.
- [37] Paul Thompson, Raheel Nawaz, John McNaught, and Sophia Ananiadou. 2017. Enriching news events with meta-knowledge information. *Language Resources and Evaluation* 51, 2 (2017), 409–438.
- [38] Ying Zhao and Justin Zobel. 2007. Searching with style: Authorship attribution in classic literature. In *Proceedings of the 30th Australasian Computer Science Conference (ACSC'07)*. 59–68.

Received August 2019; revised October 2019; accepted October 2019