# StyleLink: User Identity Linkage Across Social Media with Stylometric Representations

**Wenwen Xu[1], Benjamin C. M. Fung[2]**

[1] School of Computer Science, McGill University, Canada
[2] School of Information Studies, McGill University, Canada
wenwen.xu2@mail.mcgill.ca, ben.fung@mcgill.ca

## Abstract

User identity linkage (UIL) is the task of aligning user identities of the same user across different social network platforms. Although existing approaches have explored various aspects such as different user profile attributes and social network structures, the writing styles from user-generated texts, which is commonly known as stylometry, remain relatively underexplored. In this paper, we propose a novel Graph Neural Network (GNN)-based model named StyleLink, which leverages both social network structures and stylometric features derived from user-generated texts to address the UIL problem in an integrated manner. Our model utilizes GNNs to incorporate both stylometric features and the network structure for each social network, effectively embedding the network and enhancing user representation. This is the first work to incorporate stylometric features into GNNs to embed social networks and then conduct UIL between two embedding spaces. Extensive experiments on real-world social network datasets demonstrate the superior performance of StyleLink over existing state-of-the-art methods, achieving higher accuracy in user linkage and improved ranking of identity matches. In addition, we explore the effects of different linguistic characteristics in the identification of user identities and visualizes the effects of applying GNNs for better social network embedding.

## Introduction

With the flourishing Online Social Networks (OSNs), people tend to participate in various social networks to engage in different social activities. According to reports (Pew Research Center 2018, 2021), roughly three-quarters of the public (73%) uses more than one OSN and the median American uses three mainstream social network sites. Each OSN serves different social networking functions in daily life. For instance, users connect with friends on Facebook and Instagram, share updates on X (formerly known as Twitter), and network with colleagues and potential employers on LinkedIn.

As a result, the same individual may have signed up multiple accounts across these diverse platforms, each account reflecting different user attributes, user-generated content (UGC), and behavior patterns, such as follows, likes, etc.

User Identity Linkage (UIL), the process of matching accounts owned by the same individual across multiple social platforms, has both direct and indirect benefits in the for-profit and non-profit areas. From a revenue perspective, UIL helps social media companies create more accurate user profiles, leading to better recommendation systems (Bonhard and Sasse 2006; Mazhari, Fakhrahmad, and Sadeghbeygi 2015; Bok et al. 2016) and data-driven strategies that assist user migration (Kumar, Zafarani, and Liu 2011). Meanwhile, in non-commercial or public areas, UIL assists in detecting and investigating cyber crimes (Zhang et al. 2019; Han et al. 2017), supports user migration between platforms, and strengthens user privacy protection (Fire et al. 2014; Yuan, Chen, and Yu 2010; Shetty et al. 2023; Li, Chen, and Wu 2023) by identifying vulnerabilities linked to account misuse or duplication. Consequently, UIL not only drives profit through better monetization and user engagement but also enhances broader social and ethical initiatives such as public safety and personal data security. With the advancement of network embedding techniques, embedding-based methods have been widely employed to address the user identity linkage (UIL) problem. Existing approaches leverage various dimensional attributes of user identities and can be grouped into three categories: user profile-based, network structure-based, and content-based methods.

User profile-based approaches typically focus on user-provided identifiable information, including username, gender, birthday, email, education, location, etc (Zafarani and Liu 2009; Liu et al. 2013; Ahmad and Ali 2019). While public profile attributes offer valuable insights for identifying users across OSNs, their effectiveness diminishes when applied to large-scale OSNs, where many attributes can be duplicated and easily impersonated.

Network-based approaches aim to link user identities with their network structures, specifically utilizing *topology consistency* (Zhang and Tong 2016). Users who share similar neighborhoods in different networks could be recognized as matched identities. In social networks, social relationships such as follower-followee play a pivotal role in exploring corresponding user identities across different OSNs (Liu et al. 2016; Zhang and Tong 2016; Zhou et al. 2018; Man et al. 2016). However, the assumption of topology consistency is challenged by network heterogeneity. For instance, users may prefer certain platforms, such as favoring Face-

book over graph, leading to active engagement on one network and a subdued presence on another. Additionally, heterogeneity arises from differing semantics of relations, such as those between a career-oriented platform like LinkedIn and a co-authorship network like Google Scholar.

Content-based approaches to user identity linkage have explored various aspects of UGC and behavior patterns. These methods have analyzed tag frequencies (Iofciu et al. 2011), typing patterns (Zafarani and Liu 2013), multi-modal UGC (Chen et al. 2020), and N-gram language modelling (Goga et al. 2013a; Vosoughi, Zhou, and Roy 2015; Zafarani and Liu 2013). However, these approaches still have limitations. By focusing solely on UGC, they overlook the crucial network structure and user connectivity, which are the most typical characteristics of OSNs. They also face challenges with platform-specific content variations and scalability issues with large datasets. Moreover, the exclusive focus on content neglects the fundamental purpose and dynamics of social networking platforms. These limitations highlight the need for a more comprehensive approach that integrates content analysis with network structural information to achieve more robust cross-platform user identity linkage.

In OSNs, user profile attributes and network structure are closely interrelated. For instance, users with similar attributes are more likely to be connected as friends, and groups of users with shared characteristics often form dense communities. Drawing inspiration from the success of applying graph neural networks (GNNs) and attention mechanisms to OSN embeddings (Wang, Ye, and Zhou 2020; Wang et al. 2020), we propose a novel GNN-based model for User Identity Linkage (UIL). This application-driven approach leverages both social network structures and stylometric features derived from UGC to address the aforementioned limitations. To the best of our knowledge, this is the first work to incorporate stylometric features into GNNs to embed social networks and then conduct UIL between two embedding spaces.

- We present a novel methodology that leverages both user-generated content (UGC) and network structure to establish correspondences between user accounts across OSNs. This approach reduces reliance on user-provided identifiable information, which may be inconsistent or deliberately obscured. Instead, we focus on analyzing user activities, including writing styles and social connections, which are harder to impersonate and accumulate over time with consistent social network engagement.

- We introduce StyleLink, an innovative Graph Neural Network (GNN)-based approach to tackle the UIL problem. StyleLink consists of three primary components: a) Stylometric feature extraction, where we identify distinctive linguistic patterns in UGC to capture unique writing styles; b) Network embedding: we employ GNN models to generate user representations that incorporate both stylometric features and network structure; and c) Supervised linkage learning, where we use a Multi-Layer Perceptron (MLP) as the mapping function to learn the embedding transformation between source and target net-

works, thus predicting aligned user identities.

- We validate the effectiveness of StyleLink on real-world social network datasets, specifically X and Foursquare. Experimental results demonstrate that our method significantly outperforms existing baselines in terms of both accuracy and efficiency.

## Related Work

### User Identity Linkage

We categorize the existing approaches on user identity linkage (UIL) into three categories: user profile-based, network-based, and content-based, from the perspectives of features in social networks. Researchers have demonstrated the utility of a wide range of attributes in addressing the UIL problem and leveraging a combination of them often leads to more reliable user identification. Attributes such as username, gender, age, occupation, hometown, location, email address, and profile photo (Motoyama and Varghese 2009; Carmagnola and Cena 2009; Goga et al. 2013b, 2015; Sharma and Dyreson 2018) have all been proven effective in existing approaches for solving the UIL problem. However, the effectiveness of these methods is limited in large-scale social networks due to attribute duplication and user impersonation. Many users also intentionally mask or falsify their personal information, further challenging these approaches.

Network-based approaches aim to link user identities with their network structures by leveraging the *topology consistency*(Zhang and Tong 2016), which means users who share similar neighborhoods in different networks could be recognized as matched ones. In social networks, social relationships, such as follower-followee relations, play a pivotal role in exploring corresponding user identities across different SNs (Liu et al. 2016; Zhou et al. 2016, 2018; Man et al. 2016; Chu et al. 2019). However, they focus solely on the connections between users, ignoring the rich user-specific information contained in profile attributes, behaviors, and content. This can lead to embeddings that fail to capture important aspects of user identity, such as interests, preferences, and linguistic patterns. When users have relatively few connections, i.e. the network connections are sparse, structure-only embeddings can not adequately represent users..

Content-based approaches usually exploit any available multi-modal UGC, including texts (Zheng et al. 2006; Goga et al. 2013a; Srivastava and Roychoudhury 2020), images (Zhang et al. 2019; Huang et al. 2019; Chen et al. 2020) , and activity patterns of users (Vosoughi, Zhou, and Roy 2015; Iofciu et al. 2011). Apart from the aforementioned limitations, without network structure, it is difficult to capture how users are influenced by or influence others in the network, which can be important for understanding user behavior and identity.

### Network Embedding

Network embedding is a crucial technique for learning low-dimensional representations of vertices within networks, in order to capture and preserve the network structure. Most existing network embedding approaches rely on shallow

models, including DeepWalk (Perozzi, Al-Rfou, and Skiena 2014) treats random walks on graphs as sentences and applies the Skip-gram model from natural language processing to learn node representations; Node2vec (Grover and Leskovec 2016), an extension of DeepWalk, introduces a flexible notion of a node's network neighborhood and employs biased random walks to efficiently explore diverse neighborhoods; LINE (Tang et al. 2015) optimizes a meticulously crafted objective function that maintains both local and global network structures. They have been widely used in many UIL approaches when conducting network embedding, for instance, Node2Vec in RLink (Li et al. 2021), random walks in DeepLink, CLF, and CRW(Zhou et al. 2018; Zhang and Philip 2015; Zhan, Zhang, and Yu 2019), etc. However, due to the inherent complexity of network structures, shallow models struggle to effectively capture the complex, highly non-linear structure, leading to suboptimal representations. The second issue is that they cannot generate embeddings for nodes not in the training set and are inherently transductive.

To alleviate these drawbacks, deep representation learning (Wang, Cui, and Zhu 2016) and Graph Neural Networks (GNNs) (Kipf and Welling 2017; Hamilton, Ying, and Leskovec 2018; Xu et al. 2019; Veličković et al. 2018) in recent years leverage multiple layers to capture the intricate, non-linear relationships within networks, leading to more expressive and accurate embeddings. Empowered by the message passing mechanism, GNNs iteratively update node representations by aggregating information from neighboring nodes, enabling them to capture both local and global structural information. The techniques such as Graph Convolutional Networks (GCNs) (Kipf and Welling 2017) and Graph Attention Networks (GATs) (Veličković et al. 2018) have been particularly influential in network embedding and downstream tasks. For instance, DyGNN (Ma et al. 2020), a Dynamic Graph Neural Network model, which can model the dynamic information of a graph as it evolves. AS-GCN (Yu et al. 2021) unifies the neural topic model and GCNs with text-rich network representation. Personality GCN (Wang et al. 2020) can effectively use text information, including TF-IDF and PMI, to detect user personalities. By leveraging this dual-aspect approach, these models capture not only the network's structural properties but also the rich information associated with individual nodes. Drawing inspiration from the success of these techniques, we are innovated to incorporate stylometric features into our GNN-based social network embedding process. This novel integration of linguistic characteristics with network topology allows StyleLink to create more nuanced and informative representations of users within social networks, potentially enhancing the accuracy and robustness of user identity linkage.

## Problem Definition

A social network is a graph $G = \{V, E, X\}$, where $V = \{v_1, v_2, ..., v_N\}$ is a set of nodes representing the users, and $E \in V \times V$ is a set of edges representing the social relationships among users, e.g., follower/followee in Twitter. Each user $v_i$ is associated with a $d$-dimensional stylometric fea-

ture vector $x_i$ (the $i$-th row in $\mathbf{X}$), which is extracted from the text written by the user $v_i$.

Let $G^s = \{V^s, E^s, X^s\}$ and $G^t = \{V^t, E^t, X^t\}$ be the source and target networks, respectively, In these networks, $V^s$ and $V^t$ are the sets of users, $E^s$ and $E^t$ are the sets of edges representing connections between users, and $X^s$ and $X^t$ are the sets of stylometric features. In addition, a set of known anchor nodes $T = \{(v_i^s, u_j^t)|v_i^s \in V^s, u_j^t \in V^t\}$ is provided, where each pair $(v_i^s, u_j^t)$ represents accounts belonging to the same individual between the two networks. In real-life social networks, anchor links naturally exist due to users registering accounts on multiple platforms. Users may explicitly mention or link their other social network accounts in their profiles or posts, providing clear anchor links.

The goal of User Identity Linkage (UIL) is to predict whether a user $v_i^s$ in the source network and a user $v_j^t$ in the target network correspond to the same individual in the real world. Formally, this can be expressed as a function $f(v_i^s, v_j^t|T, G^s, G^t)$, which is defined as:

$$f(v_i^s, v_j^t|T, G^s, G^t) = \begin{cases} 1, & \text{if } v_i^s = v_j^t \\ 0, & \text{otherwise} \end{cases}$$

## Methodology

To solve the problem of User Identity Linkage, we propose a GNN-based model, named StyleLink. As shown in Figure 1, StyleLink consists of three key components: stylometric feature engineering, network embedding via Graph Neural Networks (GNNs), and supervised linkage learning. We will discuss each component in detail.

### (a) Stylometric Feature Extraction

To model users' writing styles, stylometric features like word choice, frequency, punctuation, and sentence length can be easily identified (Sari, Stevenson, and Vlachos 2018) and assembled into sets of representative characteristics. In this paper, we extract and characterize the writing styles of users from the following aspects, following the framework proposed by (Zheng et al. 2006). We evaluated 274 static features including lexical, syntactical, structural features, and idiosyncratic features specially designed for UGC on social networks. The frequency of misspellings can reflect a user's attention to detail, educational background, and language proficiency. The use of abbreviations varies greatly among users, reflecting their communication style, level of formality, and adaptation to platform norms. Users who frequently interact with various topics may exhibit a broader range of abbreviations, reflecting their engagement level on OSNs. Therefore, we choose to incorporate these idiosyncratic features into the stylometric feature set and validate their capability to enhance the accuracy and reliability of UIL. All static stylometric features are listed in Table **??** in the Appendix.

### (b) Graph Neural Networks for Social Network Embedding

In StyleLink, both source and target networks are embedded into low-dimensional spaces, denoted as $Z^s$ and $Z^t$ respectively, and a mapping function $\Phi : Z^s \rightarrow Z^t$, which maps
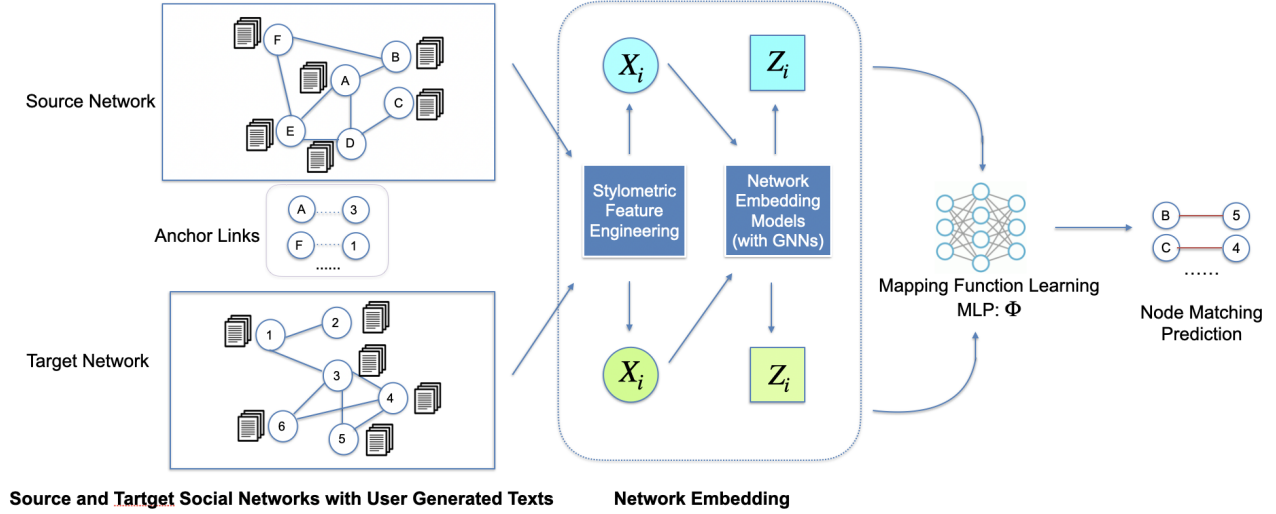
Figure 1: Illustration of the StyleLink model workflow. The process begins with obtaining source and target social network information, including user connections and their publicly posted texts (textual UGC). Next, stylometric features are extracted from the UGC and input into a Graph Neural Network to generate network embeddings that better represent the users. Subsequently, a mapping function is constructed to learn the relationships across the two OSNs. Finally, the user linkage results are produced.
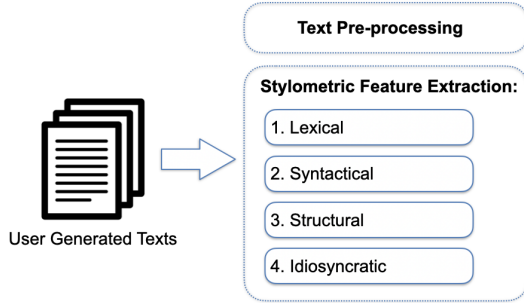


Figure 2: Overview of Stylometric Feature Extraction: User-generated texts undergo text pre-processing followed by extraction of four key stylometric features—Lexical, Syntactical, Structural, and Idiosyncratic—to represent writing styles.

the latent spaces from the source to the target, is learned. Firstly, we apply Graph Convolutional Networks (GCNs) (Kipf and Welling 2017), an effective graph neural network that captures high-order neighborhood information, to embed the source and target networks. For multi-layer GCN, the layers can be mathematically defined as:

$$H^1 = \sigma(D^{-\frac{1}{2}}\widetilde{A}D^{-\frac{1}{2}}XW^0) \qquad (1)$$

$$H^{k+1} = \sigma(D^{-\frac{1}{2}}\widetilde{A}D^{-\frac{1}{2}}H^kW^k) \qquad (2)$$

where $H^k$ is the node embedding matrix at layer $k$, $X$ in Eq.1 is the matrix of stylometric features, and also the initial layer $H^0$, $\sigma(\cdot)$ is a non-linear activation function (e.g.,

ReLU $\sigma(x) = max(0, x)$), $\tilde{A}$ is the graph adjacency matrix with the addition of self-loops, ensuring that each node's own features are included in the aggregation process. $\tilde{D}$ is the graph degree matrix with the addition of self-loops, and $W^l \in R^{d_l \times d_{l+1}}$ is the trainable weight matrix for the $l$-th layer.

The final output of the GCN, $Z = H^k$, represents the network embeddings after $k$ layers, where each row corresponds to a continuous and low-dimensional embedding of a user) in the network. These embeddings incorporate information not only from the node's own features but also from the features of its $k$-hop neighbors, effectively embedding the network. We name the StyleLink variant as StyleLink-GCN if its embeddings are from GCNs.

Incorporating the attention mechanism into the GNN learning process enables a node to gather the most relevant information from its neighbors and update its features using learned attention weights. This allows the model to focus on important nodes or edges in the graph while reducing the impact of noise during message passing in the network (Fountoulakis et al. 2023). Here we applied multi-head mechanism as described in Graph Attention Network (GAT) (Veličković et al. 2018) to be the second variant: StyleLink-GAT. Mathematically, the multi-head attention mechanism can be defined as follows:

Let $h_i$ denote the hidden state of node $i$. The attention coefficient $e_{ij}$ between node $i$ and node $j$ is computed as:

$$e_{ij} = \text{LeakyReLU}(\alpha^T[Wh_i \,\|\, Wh_j]) \qquad (3)$$

where $\alpha$ is the attention vector, $W$ is the weight matrix, and $\|$ denotes concatenation. The normalized attention coefficients $\alpha_{ij}$ are then obtained using the softmax function:

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k \in \mathcal{N}(i)} \exp(e_{ik})} \quad (4)$$

For the multi-head attention GAT layer, the new representation of node $i$ is computed as a weighted sum of its neighbors' representations, taking into account the attention coefficients:

$$\mathbf{h}'_i = \sigma \left( \sum_{j \in \mathcal{N}(i)} \alpha_{ij} \mathbf{W} \mathbf{X}_j \right) \quad (5)$$

where $\sigma$ is a non-linear activation function, and $\mathcal{N}(i)$ denotes the set of neighbors of node $i$.

Thus, the attention mechanism enables the model to focus on important nodes or edges while effectively aggregating the writing styles of neighbors.

### (c) Supervised Linkage Learning

After obtaining the representation $Z^s \in R^{d \times n}, Z^t \in R^{d \times n}$ from GNNs for the source and target network graphs, respectively, the next step is to learn a mapping function $\Phi : Z^s \to Z^t$. This mapping function is supervised using the known anchor links $T$.

In contrast to many existing methods that typically choose cosine similarity for their loss function (Zhou et al. 2018; Zhong et al. 2018; Xie et al. 2018; Chen et al. 2020; Kong, Zhang, and Yu 2013), our proposed model adopts the Triplet Loss for the objective function in supervised linkage learning. The concept of triplet loss was initially developed for facial recognition applications (Schroff, Kalenichenko, and Philbin 2015). It has shown improved ability to distinguish between different items in the embedding space. Unlike cosine similarity, which only looks at pairs of items, triplet loss considers groups of three. This approach pushes the mapping function to position correct matches significantly closer together in the latent space compared to incorrect matches.

We need to minimize the objective function as follows during the mapping function learning:

$$\mathcal{L}_{\text{triplet}} = \arg \min_{W,b} \sum_{(a,p,n) \in \mathcal{T}} \left[ \| \Phi(Z_a^s) - Z_p^t \|_2^2 \right.$$
$$\left. - \| \Phi(Z_a^s) - Z_n^t \|_2^2 + \alpha \right] \quad (6)$$

, where Anchor ($a$) represents a node from the source network; Positive ($p$) represents the corresponding node from the target network, i.e. the same user; and Negative ($n$) is a different node from the target network.

Based on comparisons between linear and non-linear mapping functions in (Man et al. 2016), we also decide to employ Multi-Layer Perceptron (MLP) as our mapping function $\Phi$, which is able to capture the non-linear mapping relationship between the source and target social networks.

The overall complexity of the method can be summarized as:

- GNN embedding: $O(L \cdot (|E_s| + |E_t|) \cdot d)$
- MLP Training: $O(k \cdot d^2 \cdot m \cdot n)$
- Linkage: $O(|V_s| \cdot |V_t| \cdot d)$

while the method is computationally intensive due to multiple steps involving large-scale graph operations, embeddings, and neural network training, the overall complexity is manageable with respect to modern computational resources and can be optimized based on specific use cases and network sizes.

## Experiment

In this section, we evaluate the proposed StyleLink model, including both its GCN and GAT variants, through two experimental tasks. Our investigation aims to address several key aspects of the model's performance and capabilities. Firstly, we explore the effectiveness of StyleLink in predicting user identities across different OSNs. We want to see how our model performs compared to state-of-the-art (SOTA) methods in the field of UIL. In addition, we focus on how these linguistic characteristics contribute to creating more representative and informative network embeddings for OSNs. We examine various aspects of stylometric features to understand their individual and collective impact on both the quality of network embeddings and the overall performance of user identity linkage tasks. All the experiments are carried out on a Windows Server equipped with two Xeon E5-2697 CPUs (36 cores), 384 GB of RAM, and four NVIDIA TITAN XP GPUs.

### Datasets

To verify our approach, we conduct experiments on the real-world partially aligned OSN datasets: X (Twitter) - Foursquare. The statistics information of the datasets is shown as Table **??**. This dataset was originally provided by Zhang et al. 2016, where users of two social networks are partially aligned. The ground truth of 1,609 anchors is publicly provided in their Foursquare profiles. We extended the original datasets with additional UGC scraping until the year of 2023, such that the average number of posts is increased by 30% in Table **??** to avoid stylometric features being too sparse.

To improve the accuracy of our analysis, we pre-process the datasets by removing non-English UGCs first. Then for each valid tweet from X or tip from Foursquare, we remove user mentions (e.g., "@username"), retweets, hashtags (e.g., #topic), and replace URLs uniformly with a specified token. These elements are excluded because they are often generic and repetitively used by many users, making them less useful for represent the unique writing styles of users.

### Evaluation Metrics and Baselines

To quantitatively evaluate the performance of our proposed model, we consider the metrics: *Precision@k (P@k)* and *Mean Average Precision (MAP)*, which are commonly used in previous studies (Zhou et al. 2018; Man et al. 2016). The higher the value of each of the measures, the better the performance of UIL.

*Precision@k (P@k)* is the metric for evaluating the linking accuracy, which is the same as *Recall@k* and $F_1@k$. It

|  | N | V | Avg Degree | Avg Posts | Vocabulary Size |
|---|---|---|---|---|---|
| X | 5,120 | 130,575 | 60.28 | 1,405.5 | 90,661 |
| Foursquare | 5,313 | 54,233 | 26.05 | 270.6 | 480,135 |

Table 1: Summary Statistics of X - Foursquare Dataset, with 1,609 anchors users.

| UIL method | Type | Topology | Attribute |
|---|---|---|---|
| IONE | supervised | Y | N |
| DeepLink | supervised | Y | N |
| PALE | supervised | Y | N |
| MNA | supervised | N | Y |
| RLink | reinforcement | Y | N |

Table 2: Comparison among different baseline UIL methods. Y(es) or N(o) stands for whether their network embedding methods involve topology and attributes.

is defined as:

$$P@k = \sum_{i}^{n} I_i\{success@k\} \tag{7}$$

where $I_i\{success@k\}/n$ measures whether the positive matching identity exists in the *top-k* $(k \leq n)$ list, and $n$ is the number of testing anchor nodes.

*Mean Average Precision (MAP)* is calculated as:

$$MAP@k = \frac{1}{n}(\sum^{n} \frac{1}{ra}) \tag{8}$$

Compared with Precision@k, it is more concerned with the performance of the returned items ranked ahead.

We evaluated the performance of StyleLink by comparing it with the following baselines, and the differences among them are shown in Table **??**. Among these baseline models we choose, network embedding methods are all employed, but in different manners, such that the user latent space is obtained for aligning the user identities.

- **IONE** (Liu et al. 2016): In Input-Output Network Embedding (IONE), a network embedding method is designed to simultaneously learn each user's follower-ship and followee-ship while utilizing the input and output context vectors to maintain the proximity of anchor users.
- **PALE (MLP)** (Man et al. 2016): Predicting Anchor Links via Embedding (PALE) conducts network embedding to capture its major structural regularity. In the matching stage, it learns an MLP mapping function across two low-dimensional latent spaces.
- **DeepLink** (Zhou et al. 2018): DeepLink is a deep reinforcement learning based algorithm that applies unbiased Random Walk to generate embeddings and uses MLP in a dual learning way to map users.
- **MNA** (Kong, Zhang, and Yu 2013): Multi-Network Anchoring (MNA) extracts social features, including spatial, temporal, and text content features (bag-of-words vectors weighted by TF-IDF), and neighborhood-based network features and match user identity pairs.

- **RLink** (Li et al. 2021): RLink applies Node2Vec (Grover and Leskovec 2016) to pre-train the network embedding and concatenates the embeddings of source and target networks to represent network structure information. Specifically, it is the first to consider UIL as a sequence decision problem and proposes a deep reinforcement learning model.
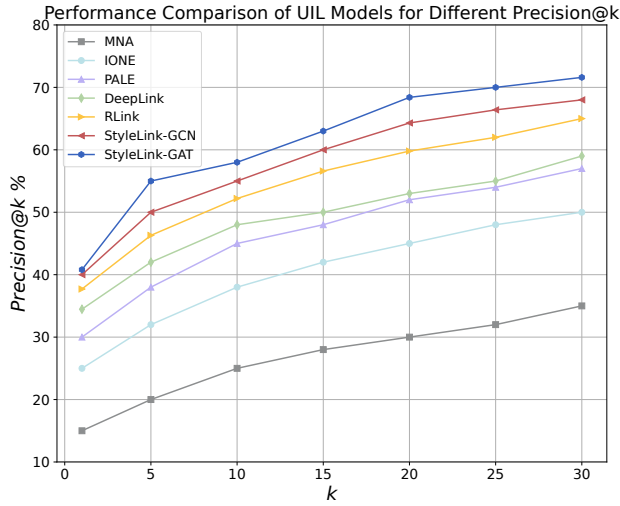
## Experimental Results

First of all, we compare the performances of various approaches by linking precision $P@k$, as presented in Figure 3a. We set the training ratio $\alpha$ to be 0.7 and present the results of different $P@k$. The results in Figure 3a show that both StyleLink-GCN and StyleLink-GAT consistently outperform the other models across all values of $k$, with StyleLink-GAT achieving the highest accuracy. On average, our method of both variants achieves a 9.2% improvement over the baseline model RLink and a 21.2% improvement over DeepLink on the X-Foursquare datasets. We observe that models utilizing deep learning techniques, such as DeepLink, PALE, RLink, and our proposed variants, StyleLink-GCN and StyleLink-GAT, generally achieve higher linking precision compared to models that do not employ neural networks, such as IONE and MNA. Specifically, IONE, DeepLink, PALE, RLink, and the StyleLink variants significantly outperform MNA, which achieves only 36.04% precision at $P@30$, whereas the other models achieve comparable precision at $P@5$. Compared to PALE and DeepLink, both of which use supervised mapping with deep learning methods, StyleLink demonstrates superior performance by integrating writing style features into the network embeddings.
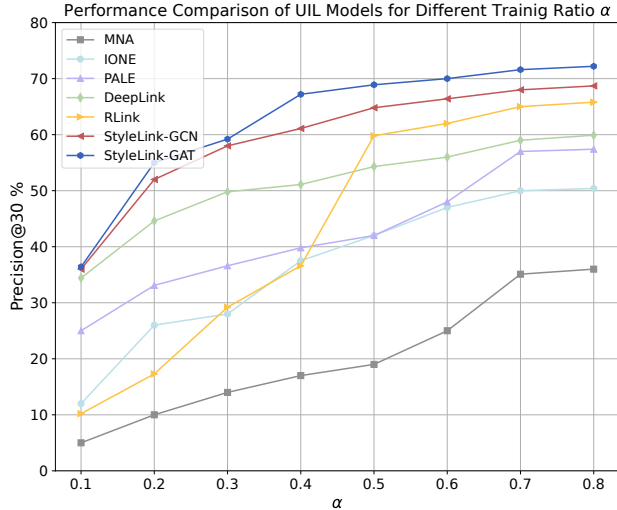
Furthermore, we varied the training ratio settings from 0.1 to 0.8 and evaluated the $P@30$ for each method. The proportion of anchor nodes $T$ used during training significantly impacts the performance of UIL models. While RLink exhibits competitive performance, particularly at higher training ratios, it does not reach the precision levels achieved by the StyleLink models. This suggests that while considering UIL as a sequence decision problem is beneficial, the network embeddings generated via Node2Vec in RLink are not as effective as those produced by our GNN-based embeddings.

To summarize, the above observations demonstrate that our proposed StyleLink models, both StyleLink-GCN and StyleLink-GAT, effectively address the User Identity Linkage (UIL) problem. Compared to other baseline models, StyleLink demonstrates significantly better performance with a lower proportion of training anchor nodes. It can effectively learn meaningful representations and perform well in scenarios where training data is incomplete or imbalanced, which is common for authentic social network datasets. In addition, StyleLink-GAT, which incorporates an

attention mechanism, achieves superior linkage performance over other models.



(a) This figure shows the performance of UIL models for different P@k, on the X-Foursquare dataset.



(b) This figure shows the performance of UIL models for different training ratios $\alpha$, on the X-Foursquare dataset.

Figure 3: Performance comparison between baseline methods and our model as well as its variants on X-Foursquare datasets. Each experiment was repeated 10 times, and the mean evaluation results were recorded.

## Ablation Study

An ablation study was carried out to determine the contribution of different components of stylometric features to the network embedding and UIL performance on OSNs. In Table **??**, stylometric features are divided into 4 categories, from the perspective of different linguistics. Therefore, we conducted experiments on X-Foursquare datasets between different variants of StyleLink, with different category of stylometric features padded zeros respectively. Negative values indicate a decrease in performance when that category of features is padded with zeros.

Overall, each category of stylometric feature types contributes positively to the performance of our model, but to different extents. Lexical features remain the most critical for StyleLink-GCN according to both metrics. For StyleLink-GAT, idiosyncratic features have the largest impact in terms of P@10, while lexical and syntactic features affect MAP more. Compared to StyleLink-GAT, which demonstrates more balanced sensitivity across different feature types, StyleLink-GCN exhibits higher sensitivity to feature ablation, particularly for lexical and syntactic features. Structural features have the least impact on StyleLink-GCN but are more influential for StyleLink-GAT.

## Effectiveness of Social Network Embedding via GCNs

We conduct this experiment to validate the effectiveness of our approach of generating the stylometric features and then applying GCNs to embed the whole network. We present visualizations in Figure 4 to illustrate that applying GNNs in network embedding indeed helps generate more meaningful and distinguishable embeddings.
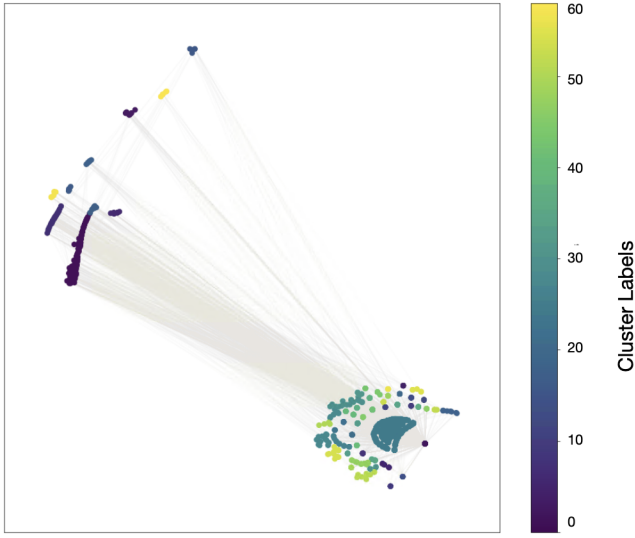
Firstly, the embeddings are reduced to two dimensions using t-SNE (van der Maaten and Hinton 2008) for visualization. Then, we can use some density-based clustering algorithm, for example, we adopted DBSCAN (Ester et al. 1996) for large spatial databases with noise in OSNs, to group similar nodes based on the reduced embeddings. Nodes with a similarity score above a defined threshold are filtered for clarity. Positions for these nodes are determined, and colors are assigned according to their cluster labels. The filtered nodes and their connecting edges are then visualized in a network graph, with a color bar indicating cluster labels, helping to highlight similarities in writing styles.

From the visualization comparison in Figure 4, we observe that representing users solely with stylometric features results in several clustering communities, where users with similar embeddings tend to cluster closely together in the embedding space. However, as seen in Figure 4a, despite the same number of users being represented as in Figure 4b, most nodes overlap significantly, forming extremely dense clusters. This can lead to the incorrect alignment of dissimi-
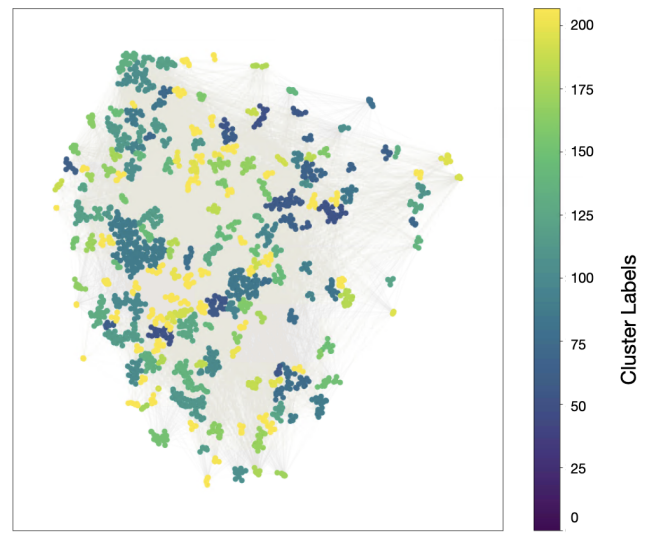
| Features | StyleLink-GCN | | StyleLink-GAT | |
| --- | --- | --- | --- | --- |
| | P@10 | MAP | P@10 | MAP |
| all features | 55.3 | 47.1 | 57.5 | 50.6 |
| (−) Lexical | -1.87 | -2.32 | -1.00 | -0.89 |
| (−) Syntactic | -1.43 | -1.21 | -0.60 | -1.10 |
| (−) Structural | -0.83 | -0.29 | -0.40 | -0.50 |
| (−) Idiosyncratic | -1.02 | -1.29 | -1.40 | -1.00 |

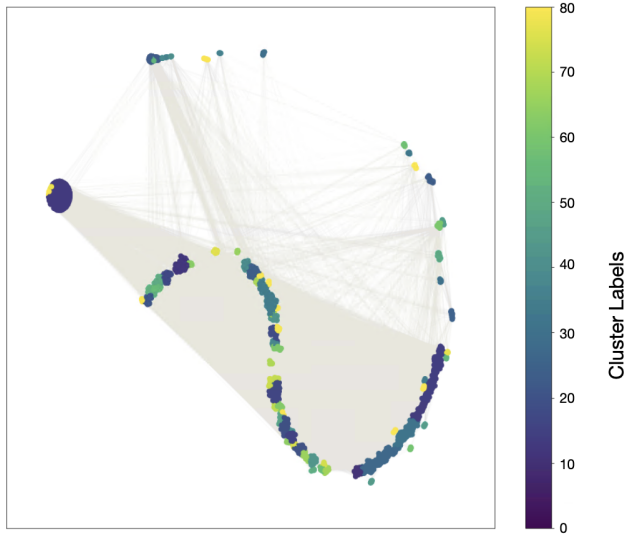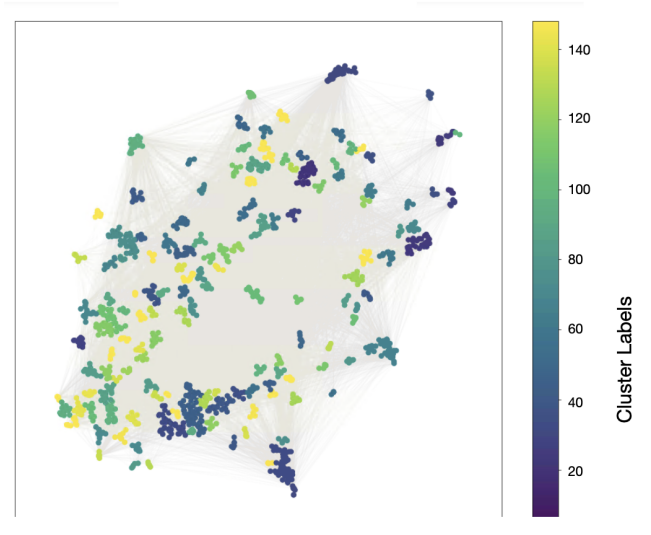Table 3: Stylometric Feature Ablation Results.

(a) This figure visualizes Foursquare embeddings based on stylometric features, before the application of GCNs.

(b) This figure visualizes Foursquare embeddings after applying GCNs.

(c) This figure visualizes X (Twitter) embeddings based on stylometric features, before the application of GCNs.

(d) This figure visualized X (Twitter) embeddings after applying GCNs.

Figure 4: We present embedding visualizations for the X - Foursquare datasets, comparing the representations before and after applying GCNs. Although StyleLink-GAT could achieve better performance, we choose to visualize with Style-GCN for the scalability and simplicity. High dimensional embeddings of $V$, $X$ and $Z$, are projected to 2D dimensional space and the light grey lines represent the edges $E$ from the network graphs. To enhance clarity and improve visualization quality, we filtered out nodes with similarity scores below a certain threshold. These filtered nodes, colored in dark purple, contribute to visual clutter if not removed. After applying this filtering process, 2,004 Twitter users and 2,369 Foursquare users remain, which were used to generate the visualizations shown above.

lar users due to similar embeddings.

In contrast, after applying GCNs to generate network embeddings, we notice that users with similar embeddings (i.e., users with similar colors) still form clusters, but these clusters are more dispersed and better by clearer boundaries. This observation indicates that GCNs can create more distinct clusters, which helps differentiate between various users.

These observations are also evident in Figures 4c and 4d, where similar patterns in the X dataset can be observed.

## Discussion

This study aims to develop and evaluate the StyleLink model, integrating both GCN and GAT variants, to improve user identity linkage (UIL) across various online social networks (OSNs). Our main objective was to evaluate whether incorporating stylometric features, linguistics characteristics inherent in users' writing styles, into Graph Neural Networks could enhance the linking precision and quality of network embeddings for UIL tasks. We compare StyleLink's performance against leading UIL methods, examining how these stylometric features contribute to the model's ability to generate superior social network embeddings and accurately link user identities across different OSNs. Our results align with previous studies highlighting the importance of linguistic features in user identification (Zheng et al. 2006; Goga et al. 2013a; Srivastava and Roychoudhury 2020). However, our work extends these findings by demonstrating the effectiveness of GNNs in leveraging these features for cross-platform identity linkage. While our study demonstrates the effectiveness of StyleLink, our experiments were conducted on one benchmark dataset, however with rich information on textual UGCs. The model's performance is expected to be assessed on additional OSNs with varying platform functionalities and user behaviors.

## Conclusions and Future Work

In this study, we introduced and evaluated the StyleLink model, incorporating GCN and GAT variants, for user identity linkage (UIL) across different online social networks (OSNs). While our evaluation was conducted on a single, well-structured X-Foursquare dataset with both textual and network features, the results remain robust, demonstrating the effectiveness of combining stylometric features with graph neural networks for user identity linkage. The results demonstrate that StyleLink, particularly the GAT variant, significantly outperforms state-of-the-art UIL methods in precision, especially at a lower training ratio. This performance improvement highlights the effectiveness of integrating stylometric features into GNN models, providing a more effective and representative embedding of user identities across OSNs. The proposed StyleLink method not only introduces a novel approach to UIL, but also provides practical insights for practical applications. There are several directions worth investigating in the future. As UGC is also associated with timestamps, we aim to explore whether temporal writing style evolution plays a significant role in user identity linkage on social networks. We also intend to ex-

plore the application of more advanced Graph Neural Network architectures, such as Graph Transformers (Yun et al. 2019; Hu et al. 2020) or GraphSAGE (Hamilton, Ying, and Leskovec 2018), to potentially enhance the capability in capturing complex and large-scale network structures. Moving forward, we will also explore how this approach can be further tailored for heterogeneous platforms and extensive datasets, which is the persistent challenge in UIL (Shu et al. 2017; Senette, Siino, and Tesconi 2024), aiming to expand its practical applicability in real-world scenarios.

## Acknowledgements

## References

Ahmad, W.; and Ali, R. 2019. Social account matching in online social media using cross-linked posts. *Procedia Computer Science*, 152: 222–229.

Bok, K.; Lim, J.; Yang, H.; and Yoo, J. 2016. Social group recommendation based on dynamic profiles and collaborative filtering. *Neurocomputing*, 209: 3–13.

Bonhard, P.; and Sasse, M. A. 2006. 'Knowing me, knowing you'—Using profiles and social networking to improve recommender systems. *BT Technology Journal*, 24(3): 84–98.

Carmagnola, F.; and Cena, F. 2009. User identification for cross-system personalisation. *Information Sciences*, 179(1): 16–32.

Chen, X.; Song, X.; Cui, S.; Gan, T.; Cheng, Z.; and Nie, L. 2020. User identity linkage across social media via attentive time-aware user modeling. *IEEE Transactions on Multimedia*, 23: 3957–3967.

Chu, X.; Fan, X.; Yao, D.; Zhu, Z.; Huang, J.; and Bi, J. 2019. Cross-Network Embedding for Multi-Network Alignment. In *The World Wide Web Conference*, WWW '19, 273–284. New York, NY, USA: Association for Computing Machinery. ISBN 9781450366748.

Ester, M.; Kriegel, H.-P.; Sander, J.; and Xu, X. 1996. A density-based algorithm for discovering clusters in large spatial databases with noise. In *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining*, KDD'96, 226–231. AAAI Press.

Fire, M.; Kagan, D.; Elyashar, A.; and Elovici, Y. 2014. Friend or foe? Fake profile identification in online social networks. *Social Network Analysis and Mining*, 4: 1–23.

FORCE11. 2020. The FAIR Data principles. https://force11.org/info/the-fair-data-principles/.

Fountoulakis, K.; Levi, A.; Yang, S.; Baranwal, A.; and Jagannath, A. 2023. Graph attention retrospective. *Journal of Machine Learning Research*, 24(246): 1–52.

Gebru, T.; Morgenstern, J.; Vecchione, B.; Vaughan, J. W.; Wallach, H.; Iii, H. D.; and Crawford, K. 2021. Datasheets for datasets. *Communications of the ACM*, 64(12): 86–92.

Goga, O.; Lei, H.; Parthasarathi, S. H. K.; Friedland, G.; Sommer, R.; and Teixeira, R. 2013a. Exploiting innocuous activity for correlating users across sites. In *Proceedings of the 22nd international conference on World Wide Web*, 447–458.

Goga, O.; Loiseau, P.; Sommer, R.; Teixeira, R.; and Gummadi, K. P. 2015. On the reliability of profile matching across large online social networks. In *Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*, 1799–1808.

Goga, O.; Perito, D.; Lei, H.; Teixeira, R.; and Sommer, R. 2013b. Large-scale correlation of accounts across social networks. *University of California at Berkeley, Berkeley, California, Tech. Rep. TR-13-002*.

Grover, A.; and Leskovec, J. 2016. node2vec: Scalable feature learning for networks. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, 855–864.

Hamilton, W. L.; Ying, R.; and Leskovec, J. 2018. Inductive Representation Learning on Large Graphs. arXiv:1706.02216.

Han, X.; Wang, L.; Cui, C.; Ma, J.; and Zhang, S. 2017. Linking Multiple Online Identities in Criminal Investigations: A Spectral Co-Clustering Framework. *IEEE Transactions on Information Forensics and Security*, 12(9): 2242–2255.

Hu, Z.; Dong, Y.; Wang, K.; and Sun, Y. 2020. Heterogeneous graph transformer. In *Proceedings of the web conference 2020*, 2704–2710.

Huang, F.; Zhang, X.; Xu, J.; Li, C.; and Li, Z. 2019. Network embedding by fusing multimodal contents and links. *Knowledge-Based Systems*, 171: 44–55.

Iofciu, T.; Fankhauser, P.; Abel, F.; and Bischoff, K. 2011. Identifying users across social tagging systems. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 5, 522–525.

Khairutdinov, R. R.; Mukhametzyanova, F. G.; and Gaysina, A. R. 2017. Socio-psychological characteristics of the subject use of slang and abbreviations in English-speaking social networks. *Turk. Online J. Des. Art Commun*.

Kipf, T. N.; and Welling, M. 2017. Semi-Supervised Classification with Graph Convolutional Networks. In *International Conference on Learning Representations*.

Kong, X.; Zhang, J.; and Yu, P. S. 2013. Inferring anchor links across multiple heterogeneous social networks. In *Proceedings of the 22nd ACM international conference on Information & Knowledge Management*, 179–188.

Kumar, S.; Zafarani, R.; and Liu, H. 2011. Understanding User Migration Patterns in Social Media. *Proceedings of the AAAI Conference on Artificial Intelligence*, 25(1): 1204–1209.

Li, X.; Cao, Y.; Li, Q.; Shang, Y.; Li, Y.; Liu, Y.; and Xu, G. 2021. RLINK: Deep reinforcement learning for user identity linkage. *World Wide Web*, 24: 85–103.

Li, X.; Chen, L.; and Wu, D. 2023. Adversary for Social Good: Leveraging Adversarial Attacks to Protect Personal Attribute Privacy. *ACM Trans. Knowl. Discov. Data*, 18(2).

Liu, J.; Zhang, F.; Song, X.; Song, Y.-I.; Lin, C.-Y.; and Hon, H.-W. 2013. What's in a name? An unsupervised approach to link users across communities. In *Proceedings of the sixth ACM international conference on Web search and data mining*, 495–504.

Liu, L.; Cheung, W. K.; Li, X.; and Liao, L. 2016. Aligning Users across Social Networks Using Network Embedding. In *IJCAI*, volume 16, 1774–1780.

Ma, Y.; Guo, Z.; Ren, Z.; Tang, J.; and Yin, D. 2020. Streaming Graph Neural Networks. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '20, 719–728. New York, NY, USA: Association for Computing Machinery. ISBN 9781450380164.

Man, T.; Shen, H.; Liu, S.; Jin, X.; and Cheng, X. 2016. Predict anchor links across social networks via an embedding approach. In *Ijcai*, volume 16, 1823–1829.

Mazhari, S.; Fakhrahmad, S. M.; and Sadeghbeygi, H. 2015. A user-profile-based friendship recommendation solution in social networks. *J. Inf. Sci.*, 41(3): 284–295.

Motoyama, M.; and Varghese, G. 2009. I seek you: searching and matching individuals in social networks. In *Proceedings of the eleventh international workshop on Web information and data management*, 67–75.

Perozzi, B.; Al-Rfou, R.; and Skiena, S. 2014. Deepwalk: Online learning of social representations. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, 701–710.

Pew Research Center. 2018. Social Media Use in 2018. https://www.pewresearch.org/internet/2018/03/01/social-media-use-in-2018/. Accessed: 2018.

Pew Research Center. 2021. Social Media Use in 2021. https://www.pewresearch.org/internet/2021/04/07/social-media-use-in-2021/. Accessed: 2021.

Pratiwi, I. D.; and Marlina, L. 2020. An Analysis of Abbreviation in Twitter Status of Hollywood Pop Singers. *English Language and Literature*, 9(1): 127–133.

Sari, Y.; Stevenson, M.; and Vlachos, A. 2018. Topic or Style? Exploring the Most Useful Features for Authorship Attribution. In Bender, E. M.; Derczynski, L.; and Isabelle, P., eds., *Proceedings of the 27th International Conference on Computational Linguistics*, 343–353. Santa Fe, New Mexico, USA: Association for Computational Linguistics.

Schroff, F.; Kalenichenko, D.; and Philbin, J. 2015. FaceNet: A unified embedding for face recognition and clustering. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE.

Senette, C.; Siino, M.; and Tesconi, M. 2024. User Identity Linkage on Social Networks: A Review of Modern Techniques and Applications. *IEEE Access*, 12: 171241–171268.

Sharma, V.; and Dyreson, C. 2018. LINKSOCIAL: Linking user profiles across multiple social media platforms. In *2018 IEEE International Conference on Big Knowledge (ICBK)*, 260–267. IEEE.

Shetty, N. P.; Muniyal, B.; Dokania, A.; Datta, S.; Gandluri, M. S.; Maben, L. M.; Priyanshu, A.; and Rezai, A. 2023. Guarding Your Social Circle: Strategies to Protect Key Connections and Edge Importance. 2023.

Shu, K.; Wang, S.; Tang, J.; Zafarani, R.; and Liu, H. 2017. User identity linkage across online social networks: A review. *Acm Sigkdd Explorations Newsletter*, 18(2): 5–17.

Srivastava, D. K.; and Roychoudhury, B. 2020. Words are important: A textual content based identity resolution scheme across multiple online social networks. *Knowledge-Based Systems*, 195: 105624.

Tang, J.; Qu, M.; Wang, M.; Zhang, M.; Yan, J.; and Mei, Q. 2015. Line: Large-scale information network embedding. In *Proceedings of the 24th international conference on world wide web*, 1067–1077.

van der Maaten, L.; and Hinton, G. 2008. Visualizing Data using t-SNE. *Journal of Machine Learning Research*, 9(86): 2579–2605.

Veličković, P.; Cucurull, G.; Casanova, A.; Romero, A.; Liò, P.; and Bengio, Y. 2018. Graph Attention Networks. arXiv:1710.10903.

Vosoughi, S.; Zhou, H.; and Roy, D. 2015. Digital stylometry: Linking profiles across social networks. In *Social Informatics: 7th International Conference, SocInfo 2015, Beijing, China, December 9-12, 2015, Proceedings 7*, 164–177. Springer.

Wang, D.; Cui, P.; and Zhu, W. 2016. Structural Deep Network Embedding. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '16, 1225–1234. New York, NY, USA: Association for Computing Machinery. ISBN 9781450342322.

Wang, Z.; Wu, C.-H.; Li, Q.-B.; Yan, B.; and Zheng, K.-F. 2020. Encoding Text Information with Graph Convolutional Networks for Personality Recognition. *Applied Sciences*, 10(12).

Wang, Z.; Ye, C.; and Zhou, H. 2020. Geolocation using GAT with Multiview Learning. In *2020 IEEE International Conference on Smart Data Services (SMDS)*, 81–88.

Xie, W.; Mu, X.; Lee, R. K.-W.; Zhu, F.; and Lim, E.-P. 2018. Unsupervised user identity linkage via factoid embedding. In *2018 IEEE International Conference on Data Mining (ICDM)*, 1338–1343. IEEE.

Xu, K.; Hu, W.; Leskovec, J.; and Jegelka, S. 2019. How Powerful are Graph Neural Networks? arXiv:1810.00826.

Yu, Z.; Jin, D.; Liu, Z.; He, D.; Wang, X.; Tong, H.; and Han, J. 2021. AS-GCN: Adaptive Semantic Architecture of Graph Convolutional Networks for Text-Rich Networks. In *2021 IEEE International Conference on Data Mining (ICDM)*, 837–846.

Yuan, M.; Chen, L.; and Yu, P. S. 2010. Personalized privacy protection in social networks. 4(2): 141–150.

Yun, S.; Jeong, M.; Kim, R.; Kang, J.; and Kim, H. J. 2019. Graph transformer networks. *Advances in neural information processing systems*, 32.

Zafarani, R.; and Liu, H. 2009. Connecting corresponding identities across communities. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 3, 354–357.

Zafarani, R.; and Liu, H. 2013. Connecting Users across Social Media Sites: A Behavioral-Modeling Approach. In *Proceedings of the 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '13, 41–49. New York, NY, USA: Association for Computing Machinery. ISBN 9781450321747.

Zhan, Q.; Zhang, J.; and Yu, P. S. 2019. Integrated anchor and social link predictions across multiple social networks. *Knowledge and Information Systems*, 60: 303–326.

Zhang, J.; and Philip, S. Y. 2015. Integrated anchor and social link predictions across social networks. In *Twenty-fourth international joint conference on artificial intelligence*.

Zhang, J.; and Yu, P. S. 2016. PCT: Partial Co-Alignment of Social Networks. In *Proceedings of the 25th International Conference on World Wide Web*, WWW '16, 749–759. Republic and Canton of Geneva, CHE: International World Wide Web Conferences Steering Committee. ISBN 9781450341431.

Zhang, S.; and Tong, H. 2016. FINAL: Fast Attributed Network Alignment. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '16, 1345–1354. New York, NY, USA: Association for Computing Machinery. ISBN 9781450342322.

Zhang, Y.; Fan, Y.; Song, W.; Hou, S.; Ye, Y.; Li, X.; Zhao, L.; Shi, C.; Wang, J.; and Xiong, Q. 2019. Your Style Your Identity: Leveraging Writing and Photography Styles for Drug Trafficker Identification in Darknet Markets over Attributed Heterogeneous Information Network. In *The World Wide Web Conference*, WWW '19, 3448–3454. New York, NY, USA: Association for Computing Machinery. ISBN 9781450366748.

Zheng, R.; Li, J.; Chen, H.; and Huang, Z. 2006. A framework for authorship identification of online messages: Writing-style features and classification techniques. *Journal of the American society for information science and technology*, 57(3): 378–393.

Zhong, Z.; Cao, Y.; Cao, Y.; Guo, M.; Guo, M.; Nie, Z.; and Nie, Z. 2018. CoLink: An Unsupervised Framework for User Identity Linkage. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1).

Zhou, F.; Liu, L.; Zhang, K.; Trajcevski, G.; Wu, J.; and Zhong, T. 2018. Deeplink: A deep learning approach for user identity linkage. In *IEEE INFOCOM 2018-IEEE conference on computer communications*, 1313–1321. IEEE.

Zhou, X.; Liang, X.; Zhang, H.; and Ma, Y. 2016. Cross-Platform Identification of Anonymous Identical Users in Multiple Social Media Networks. *IEEE Transactions on Knowledge and Data Engineering*, 28(2): 411–424.

## Paper Checklist

1. For most authors...

   (a) Would answering this research question advance science without violating social contracts, such as violating privacy norms, perpetuating unfair profiling, exacerbating the socio-economic divide, or implying disrespect to societies or cultures? Yes, our research aligns with the core principles of ethical research and social responsibility, ensuring that the methodologies and datasets used are applied with consideration of fairness, privacy, and the potential societal impact. The datasets are obtained from the public domain.

   (b) Do your main claims in the abstract and introduction accurately reflect the paper's contributions and scope? Yes

   (c) Do you clarify how the proposed methodological approach is appropriate for the claims made? Yes.

   (d) Do you clarify what are possible artifacts in the data used, given population-specific distributions? The dataset used contains only English texts and does not account for population-specific distributions or demographic variation.

   (e) Did you describe the limitations of your work? Yes. Limitations are described in the Discussion section.

   (f) Did you discuss any potential negative societal impacts of your work? No, this research does not foresee any significant negative societal impacts, as it proposes a neutral methodology. Solving UIL problem helps with various downstream tasks, which might bring positive or negative social impacts depending on the intention of users.

   (g) Did you discuss any potential misuse of your work? No, potential misuse of the neutral methodology was not specifically discussed. However, we acknowledge that the application of any UIL models in downstream tasks could, in extreme cases, be misused to suppress freedom of speech or infringe on privacy.

   (h) Did you describe steps taken to prevent or mitigate potential negative outcomes of the research, such as data and model documentation, data anonymization, responsible release, access control, and the reproducibility of findings? Yes. We adhere to standard anonymization protocols for the social media data.

   (i) Have you read the ethics review guidelines and ensured that your paper conforms to them? Yes

2. Additionally, if your study involves hypotheses testing...

   (a) Did you clearly state the assumptions underlying all theoretical results? NA

   (b) Have you provided justifications for all theoretical results? NA

   (c) Did you discuss competing hypotheses or theories that might challenge or complement your theoretical results? NA

   (d) Have you considered alternative mechanisms or explanations that might account for the same outcomes observed in your study? NA

   (e) Did you address potential biases or limitations in your theoretical framework? NA

   (f) Have you related your theoretical results to the existing literature in social science? NA

   (g) Did you discuss the implications of your theoretical results for policy, practice, or further research in the social science domain? NA

3. Additionally, if you are including theoretical proofs...

   (a) Did you state the full set of assumptions of all theoretical results? NA

   (b) Did you include complete proofs of all theoretical results? NA

4. Additionally, if you ran machine learning experiments...

   (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? The feature engineering and algorithm implementation are straightforward and described in detail.

   (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? Yes, in the Experiment Section.

   (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? No, because the standard variances are minimal, and we have omitted them in line with previous studies.

   (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? Yes. We stated in the Experiment Section.

   (e) Do you justify how the proposed evaluation is sufficient and appropriate to the claims made? Yes

   (f) Do you discuss what is "the cost" of misclassification and fault (in)tolerance? No

5. Additionally, if you are using existing assets (e.g., code, data, models) or curating/releasing new assets, **without compromising anonymity**...

   (a) If your work uses existing assets, did you cite the creators? Yes

   (b) Did you mention the license of the assets? Yes. Licenses, where applicable, are mentioned in the cited sources.

   (c) Did you include any new assets in the supplemental material or as a URL? NA

   (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? NA

   (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? Yes. The original raw data includes personally identifiable information, which has been anonymized.

   (f) If you are curating or releasing new datasets, did you discuss how you intend to make your datasets FAIR (see FORCE11 (2020))? NA

(g) If you are curating or releasing new datasets, did you create a Datasheet for the Dataset (see Gebru et al. (2021))? NA

6. Additionally, if you used crowdsourcing or conducted research with human subjects, **without compromising anonymity**...

(a) Did you include the full text of instructions given to participants and screenshots? NA

(b) Did you describe any potential participant risks, with mentions of Institutional Review Board (IRB) approvals? NA

(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? NA

(d) Did you discuss how data is stored, shared, and deidentified? NA

# Appendix

| Abbr. | Meaning | Abbr. | Meaning |
|---|---|---|---|
| AFAIK | As far as I know | AFK | Away from keyboard |
| ASAP | As soon as possible | BC, B/C | Because |
| BFF | Best Friend Forever | BRB | Be right back |
| BTW | By the way | DM | Direct message |
| FYI | For your information | IDK | I don't know |
| IMO | In my opinion | RN | Right now |
| JK | Just kidding | LMK | Let me know |
| LMAO | Laughing my ass off | LOL | Laugh out loud |
| NB | Not bad | NP | No problem |
| NVM | Never mind | OFC | Of course |
| OMG | Oh my God | OMW | On my way |
| PM | Private Message | TBH | To be honest |
| TMI | Too much information | HBD | Happy Birthday |
| TY | Thank You | WTF | What the f*** |
| YW | You're welcome | XOXO | A term to convey affection |

Table 4: Common Social Media Abbreviations and Their Meanings. We chose these 30 words as they are widely used and representative across various OSNs (Khairutdinov, Mukhametzyanova, and Gaysina 2017; Pratiwi and Marlina 2020).

| Categories | Examples |
|---|---|
| Lexical Features F1 | *Character-based features:* <br> 1. Total number of characters(C) <br> 2. Ratio of alphabetic characters/C <br> 3. Ratio of upper-case characters/C <br> 4. Ratio of digits/C <br> 5. Ratio of tabs/C <br> 6–31. Frequency of letters, ignoring case (26 features: A to Z) <br> 32–53. Frequency of special characters (22 features: ()<>%—{} []/#˜+-*=$^& ) <br> *Word-based features:* <br> 54. Total number of words (M) <br> 55. Ratio of short words (less than four characters)/M <br> 56. Total number of characters in words/C <br> 57. Average word length (in characters) <br> 58. Average sentence length (in characters) <br> 59. Average sentence length (in words) <br> 60. Total different words/M <br> 61. Yule's K measure* (A vocabulary richness measure defined by Yule) <br> 62–81. Word length frequency distribution / M (20 features) Frequency of words in different lengths |
| Syntactic Features F2 | 82–89. Frequency of punctuations (8 features, including " , . ? ! : ; ' <br> 90–239. Frequency of function words (150 features) (Zheng et al. 2006) |
| Structural Features F3 | 240. Total number of sentences <br> 241. Average sentences per post <br> 242. Average URL per post |
| Idiosyncratic Features F4 | 243. Average Misspelled words per post <br> 244-273. Abbreviation Frequency <br> 274. Average Abbreviation Diversity |

Table 5: List of Stylometric Features