

Track Me If You Can: Insights into Profile Interlinking on Social Networks

Sergej Denisov

Bielefeld University of Applied Sciences

Bielefeld, Germany

sergej.denisov@fh-bielefeld.de

Frederik S. Bäumer

Bielefeld University of Applied Sciences

Bielefeld, Germany

frederik.baeumer@fh-bielefeld.de

Michaela Geierhos

Universität der Bundeswehr München

Munich, Germany

michaela.geierhos@unibw.de

Abstract—Social networks shape today’s Web with modern communication capabilities and the ability to share heterogeneous media. The various networks have different focuses and reach different user groups so that users are often registered and active on multiple networks to take advantage of the respective benefits. In the ADRIAN project, we study threats to users on the Web, focusing on the creation of Digital Twins of real users. Here, we investigate the interlinking of user profiles on Social Networks and derive insights that help us model Digital Twins. We discuss the possibilities of using links to find additional profiles and assign them to users. To do this, both the links and the information in the profiles are examined. Only with high-quality data, it is possible to warn users reliably about the dangers of disclosing data.

Index Terms—Social network services, Data privacy

I. INTRODUCTION

People share their opinions and daily experiences on Online Social Networks (OSNs) [1]. However, platforms differ in terms of functionality and user experience. For this reason, users share their posts on the same topic on different OSNs (so-called cross-platform content sharing) [2]. On average, a Web user is expected to have 7.5 social media accounts in 2022 [3]. Moreover, the various OSNs often have different audiences, as well as different rules, therefore the posts are adapted accordingly by the users, especially in language style [2]. Significant differences were identified between the platforms in terms of usage and posting behavior. For example, Twitter is primarily used for information purposes, Twitter and Instagram for social sharing, and Instagram for entertainment purposes [4].

As a result, the available information of individual users also differs from one social network to another. By interlinking users across OSNs, very comprehensive user profiles can be obtained so that, on the one hand, their entire profile and behavior and, on the other hand, their preferences, activities, and friend network can be reconstructed. In the area of cybercrime the digital footprint can be used to track and target such users and create Digital Twins (DTs) [5]. To more effectively combat cyberbullying and identity theft, it is necessary to develop preventive methods that uncover the digital footprints, reveal the links across social media profiles, and thus point out the associated individual exposures.

The research project ADRIAN (*Authority-Dependent Risk Identification and Analysis in online Networks*) focuses on the disclosure of information by individuals in Web 2.0 [5].

This is not a new topic, but has been researched for some time now [6]. With the rise of modern OSNs, the threat to the author through published information became more concrete [7]. One form of this threat is doxing, which is the collection and publication of information about a person. Cyberbullying such as doxing is increasing [8], so automated solutions are being worked on to detect and prevent it [9]. From the user’s point of view, it is often difficult to understand that information distributed across different “places on the Web” can, in combination, pose a threat. We attempt to make this threat visible to users by merging the information and modeling it as DT. However, merging information across OSNs (i.e., profile interlinking) is a data processing challenge and is the subject of current research [10].

In this short paper, we look at the underlying data available for creating link profiles and initializing DTs. For this purpose, we look at sample data from Twitter, YouTube, Facebook, and Instagram and highlight challenges that arise from a data science perspective. In this context, we focus on the possibilities of reliable profile merging and data quality. The structure of this paper is as follows: In Section II, related work is outlined. Building on this, Section III presents our data set. The results and implications for the ADRIAN research project are discussed in Section IV before we conclude our work.

II. RELATED WORK

In the following, we review relevant work on profile interlinking (cf. Section II-A) and DTs (cf. Section II-B).

A. Profile Interlinking

The fundamental problem of profile interlinking is not new. The problem of identity matching was first mentioned by Newcombe in the late 1950s [11], long before the emergence of Web 2.0 and Internet use by the general public. The mathematical foundations for this followed ten years later by Fellegi and Sunter [12]. Since then, research has addressed this topic in the areas of databases, statistics, natural language processing, and data mining [13], among others.

Users systematically adapt their profile to the platform-specific standards with regard to language and wording in the profile. In doing so, they distinguish between formal and informal platforms, and even age- and gender-specific differences can be observed [2]. Even posts from the same

person on different OSNs have linguistic variations. This is because users adapt their language style to the platform-specific norms [2]. Not only in terms of language, but also in terms of behavior, differences can be observed among users.

Recently, deep learning have also been applied to enable profile interlinking. Xu et al. [10] proposed an anchor node embedding method based on dual domain adaptation to learn the anchor node representation considering the attributes, topological structure and difference between domains. Users in different OSNs are called anchor nodes, and edges between users are called anchor links. In addition, they developed a node adaptation method based on a domain adaptation by backpropagation to learn the appropriate adaptation function using a backpropagation neural network. Moreover, Wang et al. [14] introduced a system called Fusion Embedding for User Identification (FEUI), in which user-pair graphs were interactively integrated by network structure, node attribute information, and node label. Thereby, the FEUI framework exploited a single-input and dual-output deep neural network to represent complex correlation from different information sources. Furthermore, Guo et al. [15] set up a deep neural tensor network-based model to represent the interactions between entities and extract the relationships between users from a higher dimension.

B. Digital Twin

The term DT is used in several areas of research and in practice. Among others, it appears in mechanical engineering, medicine, and computer science [16]. In artificial intelligence, the term has gained broader usage. In general, “DTs can be defined as (physical and/or virtual) machines or computer-based models that are simulating, emulating, mirroring, or ‘twinning’ the life of a physical entity, which may be an object, a process, a human, or a human-related feature” [16]. Here, we refer to the term as the digital representation of a human being that is created by personal information available on the Web [5]. The DT can never reflect the whole complexity of a real person, but represents characteristics that, alone or in combination with other characteristics, may pose a risk to the real person. Modeling DTs is based on established and freely available standards of the semantic web, such as Schema.org and FOAF (Friend of a Friend). At the same time, the overwhelming number of possible sources of information, the quality of the data, and a multitude of contradictory data make modeling challenging. However, studies [17] show that a large amount of relevant information is knowingly and unknowingly disclosed by users themselves [18].

III. LINK RECORD & PROFILE DATA

For the analysis of user behavior when posting cross-platform links (cf. Section III-A), data from YouTube was obtained as a starting point. In addition, Twitter data was collected based on the links included in YouTube videos to facilitate comparison. This allows us to determine which OSNs are interlinked and which OSNs are suitable entry points for data collection to create DTs. However, the relevance of the

OSNs also comes from the trackable information. Again, what information can be found on multiple OSNs is crucial to ensure data quality through data matching. We therefore compare the different data points provided by the OSNs in Section III-B.

A. Link Record Analysis

The first goal is to analyze the collected data from YouTube and Twitter (cf. Table I).

TABLE I
DESCRIPTIVE STATISTICS: YOUTUBE AND TWITTER DATASET

YouTube	#	Twitter	#
Total Videos	4,605	Total Tweets	345,748
Total Channels	2,841	Total Users	842
Total Links	32,464	Total Links	467,834
Total Videos with Link	4,108	Total Tweets with Link	345,748
Videos/Channel (min)	1	Tweets/User (min)	1
Videos/Channel (mean)	1.62	Tweets/User (mean)	411.12
Videos/Channel (max)	39	Tweets/User (max)	87,343
Links/Video (min)	0	Links/Tweet (min)	1
Links/Video (mean)	6.94	Links/Tweet (mean)	1.35
Links/Video (max)	88	Links/Tweet (max)	8
Links/Channel (min)	1	Links/User (min)	1
Links/Channel (mean)	11.43	Links/User (mean)	607.58
Links/Channel (max)	513	Links/User (max)	174,356

We have collected a total of 4,605 videos belonging to 2,841 channels. On average, a video contains 6.94 links, a tweet has significantly fewer links than that, with 1.35 on average. The Twitter dataset has significantly more links with 607.58 links per user compared to 11.43 links per channel. This is obviously due to the fact that we have more tweets per user than videos per channel. Building on this, we turn to a deeper analysis of the links to better understand the link structure. We focus here on YouTube’s video descriptions, which allow users to insert a large number of links.

TABLE II
NUMBER OF LINKS REFERRING TO A SPECIFIC DOMAIN

Domain in YouTube Videos	# of Links	Domain in Tweets	# of Links
YouTube	4,647	Twitter	277,058
Bit	4,285	Screammov	87,054
Instagram	4,258	Trib	20,772
Twitter	2,334	Independent	11,061
Amazon	1,441	Bit	7,634
Facebook	1,415	TheGuardian	5,904
TikTok	903	LiverpoolEcho	4,881
Twitich	826	WioNews	4,567
Discord	540	FoxNews	2,983
Lnk	447	YouTube	2,925

Table II shows the domain distribution for links from YouTube videos and tweets. Most of the links in the YouTube videos lead to Instagram, Twitter, and Facebook. This correlates with the social media platforms relevant to our use case. For this reason, we analyze these three OSNs in more detail. In this context, it is important to determine how many links to the various platforms are contained in all videos or per channel.

TABLE III
ANALYSIS OF THE LINKS INCLUDED IN YOUTUBE VIDEOS

# of Links per Video/Channel	on Twitter	on Facebook	on Instagram
= 1	1,044	785	1,611
> 1	261	76	436
> 2	37	13	237
> 3	28	6	150
> 4	24	4	106
> 5	17	2	71
> 6	14	2	51
> 7	6	1	38
> 8	5	1	28
> 9	3	1	22
> 10	2	1	19

As shown in Table III, Instagram has the most links within embedded videos from YouTube in the descriptions, followed by Twitter and Facebook. Also, for any number between 2 and 10 links, Instagram achieves the highest number. It is also important to determine how many links overlap. A YouTube channel that contains links to all three OSNs (cf. Figure 1) is particularly relevant for creating a cross-platform profile.

B. Profile Data Overview

We divide the data into the following categories: Channel/User Identity, Channel/User Information, Content Information, Links to Images or Videos, External Links, Location Information, and Channel/User Metrics. The categories represent different aspects of user profiles, which generally require separate methods for correlation. The first category is used to establish the identity of a person. For this purpose, the name and username provided by the OSNs are used. The channel title or username appears in all OSNs, while the real name only appears on Instagram and Twitter. The challenge here is to determine, first, whether the user has provided the real name, and second, whether the username includes, for example, the user’s first or last name. The next category is mainly about textual information for the channel/user. From the description of a profile in connection with the identity, it can be derived whether it is a person or an organization. In the case of persons, for example, the profession, interests, or further external links are revealed. The next category represents the content published by the channel or user. Tweets and Facebook posts are textual content enriched with images or videos, while YouTube and Instagram are images or videos with a description. The tags can provide initial information about the content and then form the basis for analyses over a period of time that provide insights into the behavior of the channel/user. The next two categories can provide very strong evidence for correlation. Images are suitable for a direct comparison of user profiles in different OSNs. The same applies to external links, e.g., when the profiles link to each other. Location data is available in all OSNs. Usually, except for Twitter, only the country or city is given as text. In the case of Twitter, many other types of representation are available, such as latitude and longitude coordinates, a bounding box, or

even automatic extraction of locations from published content. Finally, the reach of a channel/user is also important. The number of followers or following can be used to determine user activity. In the case of Twitter, followers can be retrieved directly and direct connections between users can be analyzed.

TABLE IV
PROFILE DATA OVERVIEW FOR DIFFERENT OSNs

YouTube	Instagram	Facebook	Twitter
user_id	user_id	user_id	user_id
title	username	username	username
X	fullName	X	name
description	biography	about	description
publishedAt	X	X	created_at
privacyStatus	private	X	X
X	verified	X	verified
topicCategories	X	category	X
X	X	type	X
post_id	post_id	post_id	post_id
publishTime	timestamp	timestamp	created_at
title	caption	X	X
description	X	text	text
defaultLanguage	X	X	lang
madeForKids	X	X	possibly_sensitive
tags	hashtags	X	hashtags
viewCount	X	X	retweet_count
commentCount	commentsCount	X	reply_count
likeCount	likesCount	likes	like_count
kind	type	X	X
url	images	images	X
externalUrl	displayUrl	X	profile_image_url
X	profilePicUrl	X	media
X	externalUrl	link	url
X	X	X	entities.urls
X	facebookPage	X	X
X	X	X	location
X	X	X	annotations
X	X	X	coordinates
country	X	X	country_code
X	X	X	place_type
X	locationName	places_lived	full_name
subscriberCount	followersCount	followers	followers_count
X	followsCount	following	following_count
videoCount	postCount	X	tweet_count
viewCount	X	X	X

IV. DISCUSSION

With limited resources and time, it is not possible for attackers to monitor OSNs live or to compare all existing profiles and keep the findings up to date. For this reason, in the ADRIAN project, we rely on OSN users to leave digital traces that point us to additional profiles and information that can be used to create DTs. However, this is not a disadvantage, but exactly the goal of the project. The point is not to create DTs on a daily basis, but to show that traces on the Web make this possible to some extent and enable threats like doxing. In this short paper, we were interested in finding out what traces can be found based on links and whether there are actually enough links to jump from one user profile to another to gather information and the intersections are significant.

Out of 2,841 profiles we identified using YouTube videos, we were able to infer three other OSNs in 507 cases using the links in the video descriptions (cf. Figure 1).

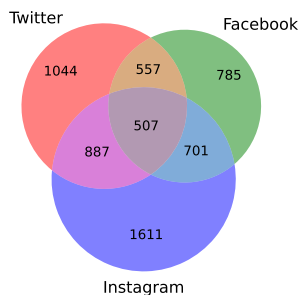


Fig. 1. Profile distribution across different OSNs

In 2,145 cases, we were able to find links to at least one other profile. The link behaviour of users is already a first data point used for identification. Beyond that, however, the large amount of different information about the respective OSNs is also worth mentioning (cf. Table IV). Not all information is available on all OSNs and in the same format, but it is often possible to merge them. In particular, user names, locations, geospatial information, and names or parts of names are often very good clues that help to interlink profiles.

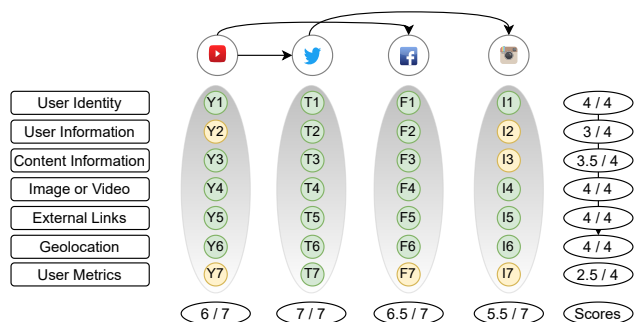


Fig. 2. Different profiles of a user in OSNs

We would like to discuss this with a real-world example – a random YouTube user who maintains multiple links in his video descriptions (cf. Figure 2): The YouTube channel contains a link to Facebook and Twitter profiles. For example, on Facebook, we can determine whether the profile belongs to an individual or an organization. This is important because the creation of the DT initially focuses on individuals. Moreover, the user is active on Twitter and posts, for example, sports activities that lead directly to Strava, an OSN for tracking physical activities that also includes social networking features. Since the source of his tweets is Instagram, the user also indicates another social media account that he uses. On Instagram, he posts photos that contain a lot of private information. As this example shows, the consolidation of the different profiles in OSNs is of considerable importance for the creation of a DT. On the one hand, it enables the validation of information and, on the other hand, information gaps can be closed. As the number of profiles in different OSNs increases, it can be assumed that a higher overall quality of the DT can be achieved.

V. CONCLUSION

We highlighted here that it is possible to infer from user profiles to other profiles. We showed that YouTube is particularly well suited as an entry portal for data acquisition, since it is possible to include many links to other OSNs in descriptions and this is also done regularly by users. In future work, we will match information from different profiles to initialize DTs.

ACKNOWLEDGMENT

This research is funded by dtec.bw – Digitalization and Technology Research Center of the Bundeswehr.

REFERENCES

- [1] P. Zhang, H. Zhu, T. Lu, H. Gu, W. Huang, and N. Gu, "Understanding relationship overlapping on social network sites: A case study of weibo and douban," *Proc. ACM Hum. Comput. Interact.*, vol. 1, no. CSCW, pp. 120:1–120:18, 2017.
- [2] C. Zhong, H. Chang, D. Karamshuk, D. Lee, and N. Sastry, "Wearing many (social) hats: How different are your different social network personae?" in *Proc. of the 11th Intl. Conf. on Web and Social Media, ICWSM 2017*. AAAI Press, 2017, pp. 397–406.
- [3] Data Portal, January 2022. [Online]. Available: <https://datareportal.com/reports/digital-2022-global-overview-report> (Accessed 2022-04-01).
- [4] M. Pelletier, A. Krallman, F. Adams, and T. Hancock, "One size doesn't fit all: a uses and gratifications analysis of social media platforms," *JRIM*, vol. 14, no. 2, pp. 269–284, 2020.
- [5] F. S. Bäumer, S. Denisov, Y. Su Lee, and M. Geierhos, "Towards Authority-Dependent Risk Identification and Analysis in Online Networks," in *Proc. of the IST-190 Research Symposium (RSY) on AI, ML and BD for Hybrid Military Operations (AI4HMO)*, A. Halimi and E. Ayday, Eds., October 2021.
- [6] A. Petit, S. Ben Mokhtar, L. Brunie, and H. Kosch, "Towards Efficient and Accurate Privacy Preserving Web Search," in *Proc. of the 9th Workshop on Middleware for Next Generation Internet Computing*, 2014, pp. 1–6.
- [7] M. Fire, R. Goldschmidt, and Y. Elovici, "Online Social Networks: Threats and Solutions," *IEEE Communications Surveys & Tutorials*, vol. 16, no. 4, pp. 2019–2036, 2014.
- [8] M. Chen, A. S. Y. Cheung, and K. L. Chan, "Doxing: What adolescents look for and their intentions," *IJERPH*, vol. 16, no. 2, p. 218, 2019.
- [9] Y. Karimi, A. Squicciarini, and S. Wilson, "Automated detection of doxing on twitter," *arXiv preprint arXiv:2202.00879*, 2022.
- [10] B. Xu, Y. Kou, G. Wang, D. Shen, and T. Nie, "Dualink: Dual domain adaptation for user identity linkage across social networks," in *Web Information Systems and Applications*, C. Xing, X. Fu, Y. Zhang, G. Zhang, and C. Borjigin, Eds. Cham: Springer, 2021, pp. 16–27.
- [11] H. B. Newcombe, J. M. Kennedy, S. J. Axford, and A. P. James, "Automatic linkage of vital records," *Science*, vol. 130, no. 3381, pp. 954–959, 1959.
- [12] I. P. Fellegi and A. B. Sunter, "A theory for record linkage," *J. Am. Stat. Assoc.*, vol. 64, no. 328, p. 1183, 1969.
- [13] A. K. Elmagarmid, P. G. Ipeirotis, and V. S. Verykios, "Duplicate record detection: A survey," *IEEE Trans. Knowl. Data Eng.*, vol. 19, no. 1, pp. 1–16, 2007.
- [14] L. Wang, Y. Zhang, and K. Hu, "FEUI: Fusion Embedding for User Identification across social networks," *APIN*, pp. 1–17, 2021.
- [15] X. Guo, Y. Liu, X. Meng, and L. Liu, "User Identity Linkage Across Social Networks Based on Neural Tensor Network," in *Lecture Notes of the Institute for Computer Sciences, Social-Informatics and Telecommunications Engineering, LNICST*, 01 2021, pp. 162–171.
- [16] B. R. Barricelli, E. Casiraghi, and D. Fogli, "A survey on digital twin: Definitions, characteristics, applications, and design implications," *IEEE Access*, vol. 7, pp. 167 653–167 671, 2019.
- [17] F. S. Bäumer, J. Kersting, M. Orlikowski, and M. Geierhos, "Towards a multi-stage approach to detect privacy breaches in physician reviews." in *SEMANTICS Posters&Demos*, 2018.
- [18] F. S. Bäumer, N. Grote, J. Kersting, and M. Geierhos, "Privacy matters: detecting nocuous patient data exposure in online physician reviews," in *ICIST*. Springer, 2017, pp. 77–89.