# Identity and User Behavior in Online Communities

**Cheng Guo**
Clemson University
Clemson, SC 29630, USA
chengg@clemson.edu

## Abstract

In online communities, people share and discuss information at all levels of topic sensitivity. Identity policies within these communities range from real names to anonymity. The amount of user engagement, the quality of the information, disinformation behavior (e.g., trolling) may differ under different types of identity, which is currently unclear. Most of these online communities have a mechanism of content moderation. The relationship between identity and moderation is also unclear. Finally, yet little is known about how and why people make decisions of self-disclosure in online communities. My dissertation research aims to deepen our understanding of identity and user behavior in online communities. My research will benefit privacy researchers, online social network designers, policymakers, and researchers in the field of Human-Computer Interaction who study online identity and social media.

## Author Keywords

anonymity, online identity, online communities, privacy, trolling, moderation, information quality

## CCS Concepts

•**Security and privacy** → **Social aspects of security and privacy;** •**Human-centered computing** → **Empirical studies in HCI;**
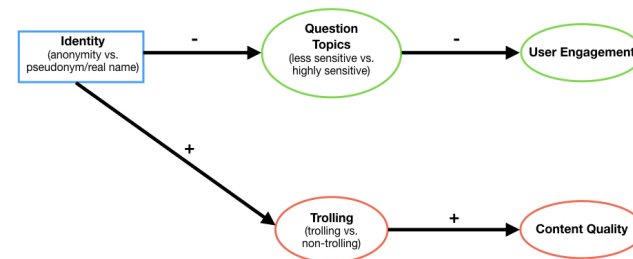
## Introduction and Research Questions

People use different types of identity in online communities ranging from real names, to pseudonyms, and even anonymity. The quality of the information in online communities varies from excellent to poor [1]. Some users may be more engaged than other users to participate in online information sharing and discussion [8]. Deviant behavior such as trolling is also very common in online communities [3]. One factor that may associate with user engagement, information quality, deviant behavior (e.g., trolling) and currently unclear is the online identity which may vary from site to site. Moreover, most of the online communities use content moderation mechanisms to keep their communities healthy using Human moderators and computer bots. However, the moderated data in online communities is often missing in research studies [4] due to inaccessibility and thus, the relationship between identity and content moderation is unclear. Finally, little is known about how and why people make identity self-disclosure under different situations in online communities. In my dissertation, I aim to deepen our understanding of the relationship between identity and user behavior in online communities by a set of empirical studies.

## Research in Progress

To answer my research questions and address my research goals, my dissertation will use both qualitative and quantitative research methods in a series of studies.

*Identity, User Engagement, Information Quality, and Trolling*
Focusing on the relationship between different types of identity and a group of user behaviors, my first part of the dissertation conducts an empirical study on three Q&A



**Figure 1:** A concept model with identity, question topics, user engagement, trolling and content quality.

sites[1][2][3]. These three Q&A sites have different identity policies. Quora has a restricted real name policy that enforces users to use their real names. Yahoo makes users use any username their desire which I considered as pseudonyms. Zhihu recommends users to use real names but it's not mandatory. As a result, users use both real names and pseudonyms on Zhihu. Moreover, these three sites all provide users a privacy feature that allows users to be anatomized whenever they'd like to. Thus, there are three types of identity on these three sites: real name, pseudonym, and anonymity. I first collected ten highly sensitive topics and ten less sensitive topics from the literature [2, 6, 7]. Then, I collected 50 questions and all the answers to these questions for each topic on Quora, Yahoo, and Zhihu for further analysis. I found that anonymity is indeed associated with different user behaviors such as asking more sensitive questions, lowering user engagement and leading to more deviant behavior, but not lowering the content quality (See Figure 1). This study is completed and I am revising this manuscript.

---

[1]https://quora.com
[2]https://answers.yahoo.com/
[3]https://zhihu.com

*Identity and Content Moderation*

One question still left out in the first study and also a limitation for almost all the studies using online forum data is caused by content moderation [4] take by online communities. Online forum data are usually heavily moderated and the raw data is usually inaccessible for researchers. Thus, the first study was conducted on the data that has already been moderated. This may lead to inaccurate research results. For example, the trolling posts may be already moderated at the time I collected the data. However, this is a limitation that very difficult or even impossible for researchers to solve because online forums usually don't provide any APIs. Even if some online forums do provide APIs, the data extracted via APIs are always moderated. In this study, I hope to explore the relationship between identity and content moderation. I settled down using Reddit[4] as my research platform. I collected raw (unmoderated) data via publicmodlogs[5], an account that makes moderation logs on Reddit publicly. Certain subreddits can voluntarily make their moderation logs public by inviting publicamodlogs to be one of their moderators. Besides collecting moderated data, I also collected unmoderated data using Using the Reddit's APIs. In a pilot study, I collected moderated and unmoderated data from three highly sensitive subreddits (operationalized as not suitable for work (NSFW) subreddits) and three less sensitive subreddits (operationalized as suitable for work (SFW) subreddits) from publicmodlogs. I found that most of the moderation was still done by Human moderators on Reddit. I also found that content posted anonymously (using throwaway accounts [5]) did not lead to more content deletion by moderation. This pilot study is completed and currently under revision. In the next step, I plan to expand the number of subreddits and collect more data. Then, I plan to do both qualitative and quantitative

---

[4]https://www.reddit.com
[5]https://www.reddit.com/user/publicmodlogs

analysis of what content was moderated on Reddit. I also want to explore the difference in content (moderated vs. unmoderated) posted by ordinary accounts vs. anonymous (throwaway) accounts. This follow-up study is scheduled to be completed for Fall 2019.

*Identity and Self-disclosure*

My final part of the dissertation is to link the user's online identity and its associated behaviors together through an empirical study. I plan to recruit active online community users from different sites (such as the sites I mentioned earlier, Yahoo, Quora, Zhihu, and Reddit) and have them voluntarily provide me access to their accounts or provide me copies of the content they posted online. With this information, I will be able to collect the content they posted online using ordinary accounts (could be real-name accounts or pseudonym accounts) vs. anonymous accounts. By analyzing these posts quantitatively, I plan to gain an understanding of what extent people post content using different levels of identity. By analyzing these posts qualitatively, I hope to gain an understanding of how people make the decision of self-disclosure and identity choices in online communities under different situations/topics. I also plan to follow up on these participants with interviews to find out why they make such decisions and what factors affect their decisions. This final part of the dissertation is scheduled to be completed in Spring/Summer 2020.

## Expected Contributions

My dissertation contributes to the field of social computing and privacy by extending our knowledge of data sensitivity, user identity, self-disclosure, and user behavior in online communities. My work will benefit both online community practitioners and technology developers with practical insights and guidelines to design high-quality communities with respect to user privacy. My research may also be use-

ful for online community users to make better privacy decisions.

## Expected Benefits

My dissertation proposal defense is scheduled for early 2020. I anticipate by that time I will have already completed the follow-up study of identity and content moderation, and propose the final study of identity and self-disclosure. I hope to get thorough feedback and guidelines from the Group DC committee members since the topics of my work are their expertise. Eventually, the results of my work will be applied to online social networks/online communities. Thus, the suggestions and comments coming from senior GROUP researchers will be very valuable to me and my dissertation.

## Acknowledgments

I would like to thank my awesome advisor Dr. Kelly Caine for all of her consistent help and support on my dissertation and Ph.D. study. I also want to thank the members of the Human and Technology Lab (www.hatlab.org) at Clemson University for their comments and encouragement.

## REFERENCES

[1] Eugene Agichtein, Carlos Castillo, Debora Donato, Aristides Gionis, and Gilad Mishne. 2008. Finding high-quality content in social media. In *Proceedings of the 2008 international conference on web search and data mining*. ACM, 183–194.

[2] Jeremy Birnholtz, Nicholas Aaron Ross Merola, and Arindam Paul. 2015. Is it weird to still be a virgin: Anonymous, locally targeted questions on facebook confession boards. In *Proceedings of the 33rd annual ACM conference on human factors in computing systems*. ACM, 2613–2622.

[3] Justin Cheng, Michael Bernstein, Cristian Danescu-Niculescu-Mizil, and Jure Leskovec. 2017. Anyone can become a troll: Causes of trolling behavior in online discussions. In *Proceedings of the 2017 ACM conference on computer supported cooperative work and social computing*. ACM, 1217–1230.

[4] Devin Gaffney and J Nathan Matias. 2018. Caveat emptor, computational social science: Large-scale missing data in a widely-published Reddit corpus. *PloS one* 13, 7 (2018), e0200162.

[5] Alex Leavitt. 2015. This is a throwaway account: Temporary technical identities and perceptions of anonymity in a massive online community. In *Proceedings of the 18th ACM Conference on Computer Supported Cooperative Work & Social Computing*. ACM, 317–327.

[6] Sai Teja Peddinti, Aleksandra Korolova, Elie Bursztein, and Geetanjali Sampemane. 2014. Cloak and swagger: Understanding data sensitivity through the lens of user anonymity. In *2014 IEEE Symposium on Security and Privacy*. IEEE, 493–508.

[7] Dan Pelleg, Elad Yom-Tov, and Yoelle Maarek. 2012. Can you believe an anonymous contributor? On truthfulness in Yahoo! Answers. In *2012 International Conference on Privacy, Security, Risk and Trust and 2012 International Confernece on Social Computing*. IEEE, 411–420.

[8] Chirag Shah, Jung Sun Oh, and Sanghee Oh. 2008. Exploring characteristics and effects of user participation in online social Q&A sites. *First Monday* 13, 9 (2008).