# Understanding the Role of AI And Large Language Models in Catfishing on Dating Sites: A Multi-Dimensional Study

Thanh An Vu
Faculty Mentor: Alexi Brooks

University of Wisconsin-Stout

## INTRODUCTION

The definition of online dating deception, related to the term "catfishing," involves approaching a potential romantic partner online without genuine intentions of creating an in-person relationship.

This study aims to demonstrate algorithms for detecting fake profiles on online dating sites based on a provided dataset. Current online dating sites already have their security layers. However, they only work for low-sophistication attacks, not higher ones. I present two vulnerability concepts that may increase the fallibility of a detection system.

## APPROACH

I explain the security techniques that online dating platforms are applying to protect users. Utilizing the idea of a large language model to detect online catfishing and provide related algorithms to help the investigation.

Many online dating platforms on the market are currently in terms of target audience, appearance design, functions, and popularity in today's world. Followed by Forbes, some online dating services attracted clients such as Tinder, Bumble, Elite Singles, Silver Singles, Plenty of Fish, and Match. These applications let users very their preferences when deciding on an attractive partner whom they want to have a conversation with. The suggestions from online dating sites will recommend a person based on common interests, goals, or personal descriptions that clients put in. In addition, the system will consider a "match" when both users give each other a "like".

In the past, dating websites already have some techniques for detecting fake accounts. For example, Jiayuan (a popular dating platform in China) has a behavioral-based, IP address-based, photograph-based, and text-based detection system to recognize fake accounts. Specifically, scammer accounts will be flagged if they use the same IP address to create multiple accounts, using duplicating images or similar messages to potential partners. However, these techniques only work for low-sophistication attacks when scammers have little understanding of the technology.
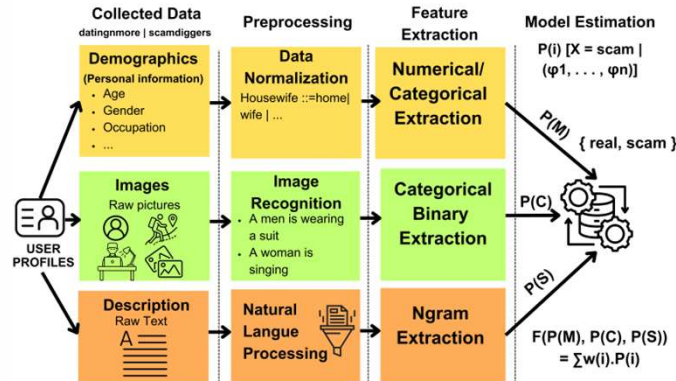

Scam Alert!



**Figure 1:** Machine learning model based on Suarez-Tangil et al (2020).

- For high-sophistication attacks in which most of the profiles are real, the traits are more varied. Those listed above in Figure 1 will work for most available online dating sites, including: demographics, images, and descriptions. Each of the classifiers deals with missing data because those fields are not mandatory for users to complete due to the private preferences policy before calculating the probabilities of $\theta_M$, $\theta_C$, and $\theta_S$. Using only one classifier will lead to less accurate results instead of applying multiple ones.

**Potential Vulnerabilities:**

- **New Users Never Use Dating Services**: There exists the possibility of creating an account on dating sites based on a real profile on other social media accounts that have not used online dating sites before. Hence, increasing the possibility of wrong predictions of a fraudulent account.
- **Replayed Attack:** This is a technique in which hackers can intercept and retransmit data that was recorded in the system by using hacking tools. This can happen when people open online dating platforms while using public free Wi-Fi. Attackers can capture the data transmitted including authentication credentials or encrypted messages and then send it back to trick the server. By taking advantage of the Three-Ways Handshake when a device accesses a website, attackers can take over user account.
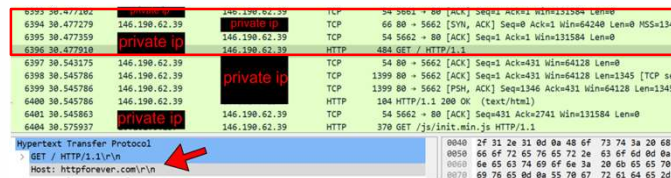


**Figure 2:** Using Wireshark to capture Three Ways Handshake Packets when accessing a website using HTTP protocol: httpforever.com
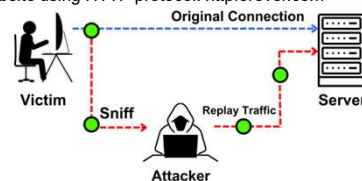


**Figure 3:** Replay attack simulation

- **File Upload Vulnerability:** This vulnerability may affect the system because all users are required to upload at least one profile picture on a dating service. Most web-based applications are written in JavaScript or PHP languages and the format of popular uploaded pictures (JPEG or PNG). Attackers can utilize this to send other web files such as (test.php) or (test.js.png/ test.php.jpeg) to access the local data.
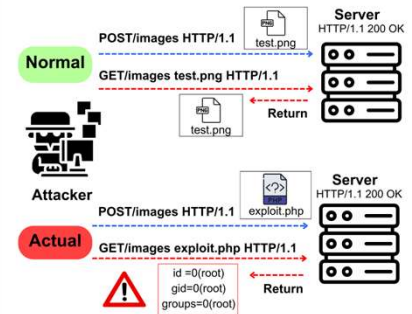


**Figure 4:** File upload vulnerability simulation

## DISCUSSION

- In terms of using AI algorithms to detect fake accounts on dating sites, we also have to mention the importance of the probability of the algorithms for misclassifying known as True vs. False and Positive vs. Negative classifications as fake profiles can be detected as real ones.

- I believe the vulnerabilities I mentioned would affect user experience and reputation of dating service companies. Especially in the case where scammers can steal people's profile information on a specific social media app because we tend to share more information online currently.

- It is ethical to find the vulnerabilities of a system to improve the user's experiences and prevent security threats. Additionally, I did not create an account to avoid the situation that my profile would become a potential partner for real users on dating sites and use it for illegal purposes.

| | Scam | Real |
|---|---|---|
| **Scam** | True Positive (TP): Scam profiles correctly classified as scams. | False Positive (FP): Real profiles misclassified as scams. |
| **Real** | False Negative (FN): Scam profiles misclassified as real. | True Negative (TN): Real profiles correctly classified as real. |

**Figure 5:** True-False Positive and Negative Prediction Model

## REFERENCES

Suarez-Tangil, G., Edwards, M., Peersman, C., Stringhini, G., Rashid, A., & Whitty, M. (2020). Automatically Dismantling Online Dating Fraud. *IEEE Transactions on Information Forensics and Security, 15*, 1128–1137. https://doi.org/10.1109/TIFS.2019.2930479

Huang, J., Stringhini, G., & Yong, P. (2015). Quit playing games with my heart: Understanding online dating scams. *Detection of Intrusions and Malware, and Vulnerability Assessment*, 216–236. https://doi.org/10.1007/978-3-319-20550-2_12