# **Police Bot**: Enhancing Social Media Governance with Policing Bots

**Milestone 6 Presentation**

# Group Members:

## Students:

- Gabriel Silva
- Cody Manning
- Liam Dumbell
- Nickolas Falco

## Faculty Advisor / Project Client:

- Khaled Slhoub

## Computer Science Project Instructor:

- Philip Chan

# Overview:

- Discussion of Task Completion:
  - Code Improvements
  - Testing Metrics
  - Decide Module
  - Demo/Commercial
  - Update on Poster and Ebook Page
- Milestone Completion Task Matrix
- Advisor Feedback
- Lessons Learned

# Testing Method

Using List of Known Humans:

- Compiled using the moderator lists of the top subreddits (where bot moderators are declared as bots, eg. AutoModerator)
- Used to test the percentage of false positives (Humans identified as Bots) List of Known Bots

Using List of Known Bots:

- Compiled using various public list of Reddit Bot accounts
- Used to test the percentage of false negatives (Bots identified as Humans)

# Testing Metrics

Positive Assessments (Known Bot List)

- Objective: 70%-80% correct bot detection
- Results: 320/424 (bots) - - - - - 76% Bots identified as Bots

Negative Assessments (Known Human List)

- Objective: <20% false positives
- Results: 102/651 (humans) - - 15% of Humans identified as Bots

# Code Improvements

- Enhanced organization

- Enhanced maintainability

- Enhanced reusability

- Enhanced Scalability

# Decide Module

Decide to report or not a scanned bot:

- Scamming Bots or Harassing Bots (Bad)
  - Are provided with a link to the report page and a tutorial.
- Auto Declared Bots or Known Bots (Good)
  - Ask the user whether to offer the option to report good bots.

# Decide Module

○

```python
1 from colors import cyan, reset
2
3 link = "https://www.reddit.com/user/"
4 def decide(account, ignore_list):
5     if 'exiting' in ignore_list:
6         print(f"url: {cyan}{link}{account.name}{reset}")
7         return True
8     if 'all' in ignore_list or account.good_bot:
9         return False
10    if account.good_bot == '0' and not 'good' in ignore_list:
11        print(f"Autodeclared bot usually not harmful.")
12        x = input("Want to report? (y?) (0 to ignore all)")
13        if x.rstrip().lower() == 'y':
14            print(f"url: {cyan}{link}{account.name}{reset}")
15            return True
16    else:
17        if len(account.reasons) > 0:
18            print(f"url: {cyan}{link}{account.name}{reset}")
19            return True
20        elif not 'inconclusive' in ignore_list:
21            print(f"No decisive reason to report this bot.")
22            z = input("Want to report? (y?)")
23            if z.rstrip().lower() == 'y':
24                print(f"url: {cyan}{link}{account.name}{reset}")
25    return False
```

# Demo/Commercial



https://youtu.be/LVJUa5dlSaI

# EBook Page

| Project Name | Framework to Analyze Behavior of Social Media Bots |
|---|---|
| Team Lead: | Cody Manning |
| Team Member(s): | Gabriel Silva, Liam Dumbell, Cody Manning, Nickolas Falco |
| Faculty Advisor(s): | Dr. Khaled A. Slhoub, Department of Electrical Engineering and Computer Science, Florida Institute of Technology |

**\*\*Do not change font size or text color above this message/delete this before completion. The category will be put in by Staff after submission \*\***

## Project Description:

Social media has become a driving force in many people's lives. Some people have created bots that serve malicious purposes. These bots act like real people, and may be used to steal information or annoy users who unknowingly interact with them. Our framework is created for the purpose of being able to detect these bots, and possibly differentiate them from the bots that are created for beneficial purposes. The framework was created to work on the 'Reddit' social media platform, but the backbone and ideas of the project could be extended to other social media platforms, with some tweaking depending on the features of the social media it is being adapted to.

## Features:

The user will be able to deploy the framework on Reddit using a specific subreddit (which is a collection of topics created by users). The user will also be able to select a specific user (by typing in the suspected users Reddit username). If the user wants to search a specific subreddit, the framework will scan through top posts, newest posts, or the posts that are most popular in a short timeframe. It then asks the user how many posts they would like to search for, and how deep in the posts (how many users) it would like to evaluate. When this is done, it will print out all of the users in the posts it grabbed and give a score based on the likelihood of them being a real human being, or a bot. The framework will also give the user insights on whether the given bot is a 'good' bot (one made to help) or a 'bad' bot (one made to harm).

## Evaluation:

When designing this framework, accuracy was the key for our measurement of success. It doesn't matter much if the results come quickly if they are wrong. To achieve this, we measure against a master list of known bots (which were scraped from several sources, GitHub and Reddit itself in particular). We were shooting for about an 80% accuracy method in detecting whether a user was a bot. When using our known bot list, and a list of known real human accounts, the accuracy rating was well within our desired output. Our timing desire was no more than about 10 seconds per account lookup, and this was unfortunately not really feasible within the context of how the program functioned. It really came down to speed or accuracy, and the team decided that accuracy was what we wanted to focus on.

## Major Challenges:

There were a lot of challenges we encountered and overcame during the process of this project. Particularly detecting bots and distinguishing the good from the bad bots. This is still a widely researched topic in Computer Science, so we were working blind for a lot of this. One notable observation was the inclusion of links. We couldn't find a good reason for bots to direct you outside of the Reddit platform, so it was immediately flagged as suspicious if they wanted you to leave the site and go somewhere else (especially if the outside link was obscured with a link shortener). There is no perfect science for this project, so it is something that needs to be built on more in the future.

# Poster

# Milestone Completion Task Matrix

| Task | Cody | Gabriel | Liam | Falco | To Do |
|---|---|---|---|---|---|
| Finalize the detection algorithms | 15% | 10% | 50% | 25% | N/A |
| Create the decide module | 10% | 40% | 30% | 20% | N/A |
| Finalize the maliciousness algorithms | 5% | 40% | 20% | 5% | N/A |
| Test the framework as a whole | 40% | 40% | 10% | 10% | N/A |
| Create developer / User Manual | 70% | 10% | 20% | 0% | N/A |
| Final Demo | 10% | 50% | 15% | 25% | N/A |

# Advisor Feedback

- Satisfied with our current progress.
  - As aforementioned in early presentations the project's goal was to build from the ground up a framework that polices a social media looking for bots
- Invited us to continue working on this project even after we are done with the class.
- Mentioned writing and publishing a paper with our findings,
- Or making the project available to students in the future.

# Lessons Learned

This project has been quite an ordeal, as aforementioned in early presentations the project's goal was to build from the base up a framework that can be improved in the future.

- Have a clear plan, and stick to it

- Understand API and Platform Limitations

- Start early

- Test early

This concludes our presentation, Thank You