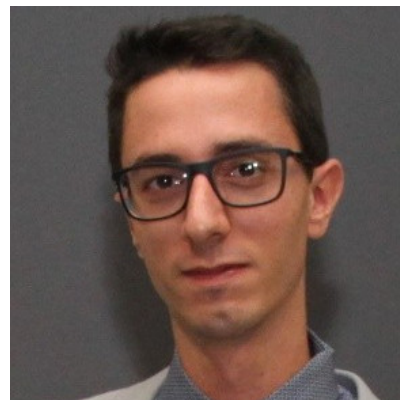


On Twitter Bots Behaving Badly: Empirical Study of Code Patterns on GitHub



POLITECNICO
MILANO 1863

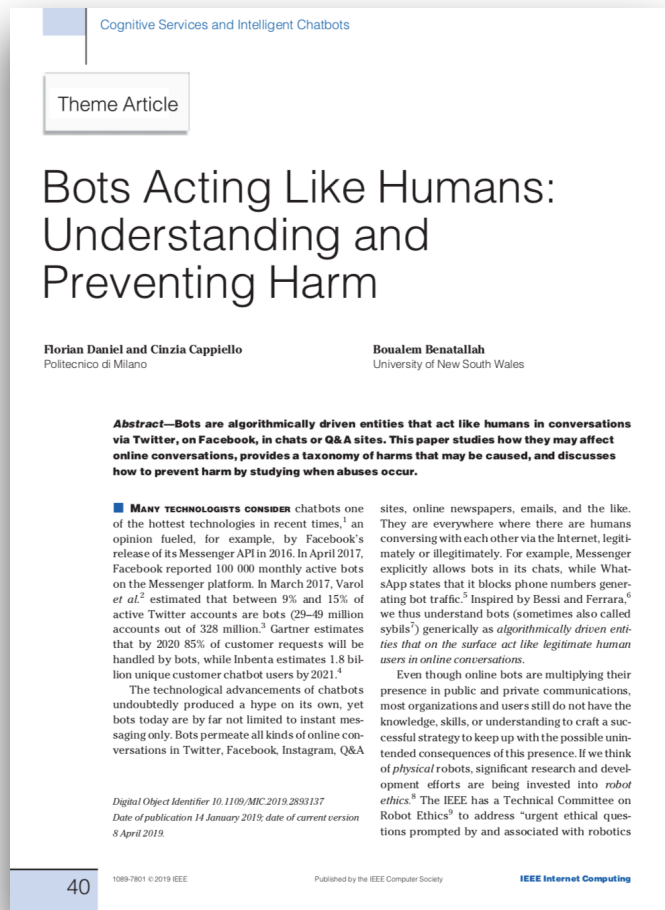
DIPARTIMENTO DI ELETTRONICA
INFORMAZIONE E BIOINGEGNERIA



Andrea
Millimaggi



Florian
Daniel



F. Daniel, C. Cappiello, B. Benatallah.
Bots Acting Like Humans:
Understanding and Preventing
Harm. IEEE Internet Computing
23(2), 2019, Pages 40-49.

Harm in human-bot interactions

Psychological harm

Someone's psychological health or well-being
get endangered or injured

Legal harm

Someone becomes subject to law enforcement
or prosecution

Economic harm

Someone incurs a monetary cost or loses time

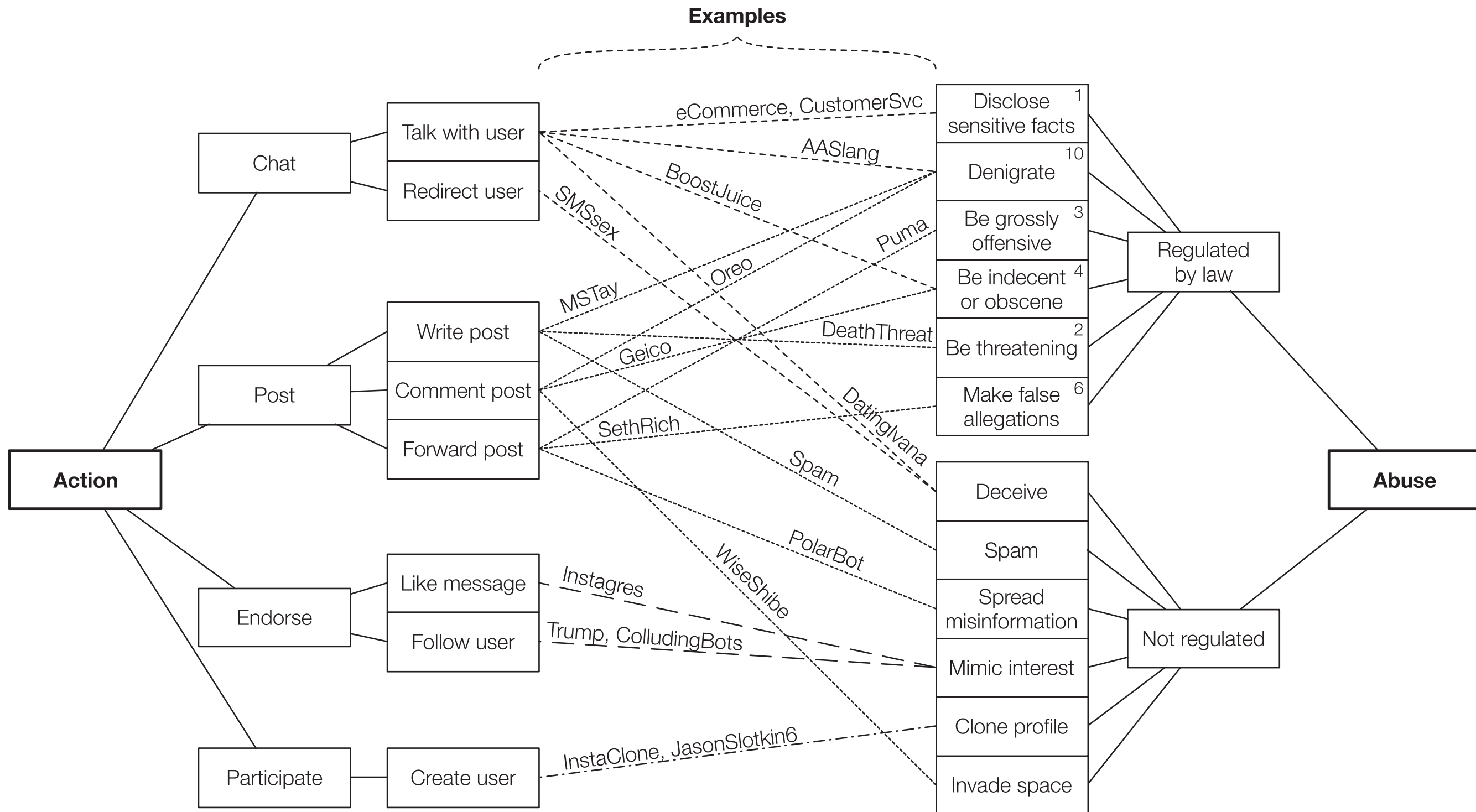
Social harm

Someone's image or standing in a community
gets affected negatively

Democratic harm

Democratic rules and principles are undermined

Harm = consequence of an **action** and an **abuse**



Literature on bots

Bot development (frameworks, APIs, etc.)

Bot detection from externally visible communications

Problem (goal of this paper)

Identify **how** harm is caused by bots

Understand likely underlying **intentions**

>> Abuse-oriented classification of bot **code repositories**
published on GitHub

Before going into the details...

Bots are not negative in general!

Platform **policies and permissions**



All platforms provide developers with **programmable interfaces**

Typically allow programmatic access to **all functionalities**

Users of the APIs must **authenticate** with the platforms

Almost all platforms impose some kind of **limitation**

Dataset

Focus on **Twitter**

Search by keywords

Collection of code files +
metadata

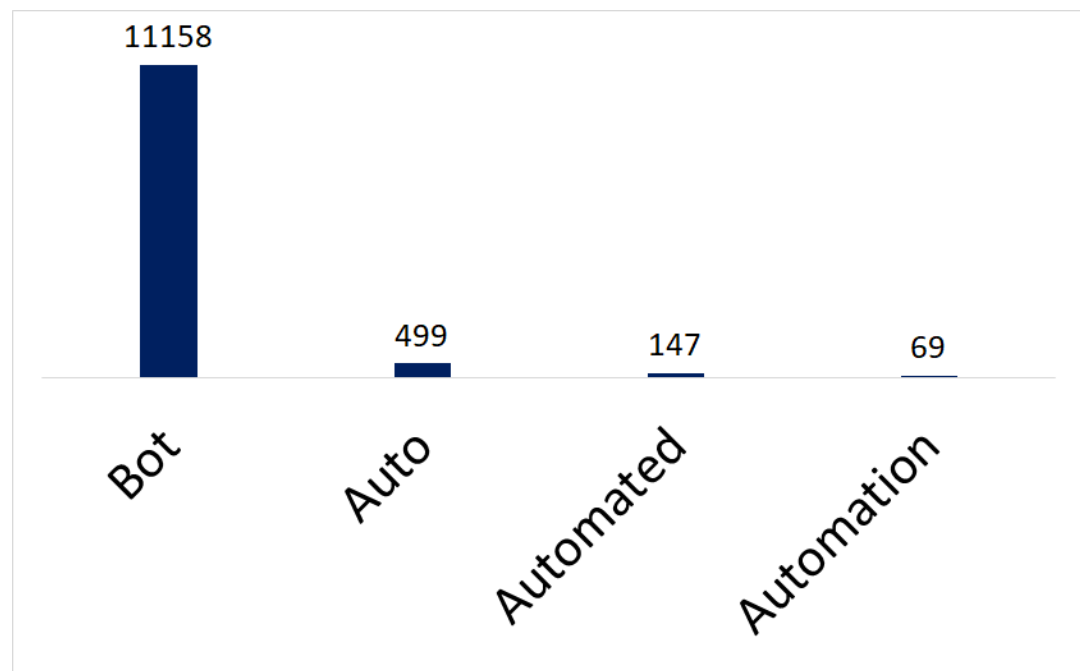


Fig. 1: Distribution of GitHub search results by searched keywords (includes all programming languages).

Preliminary analysis
of actions implemented
in repositories

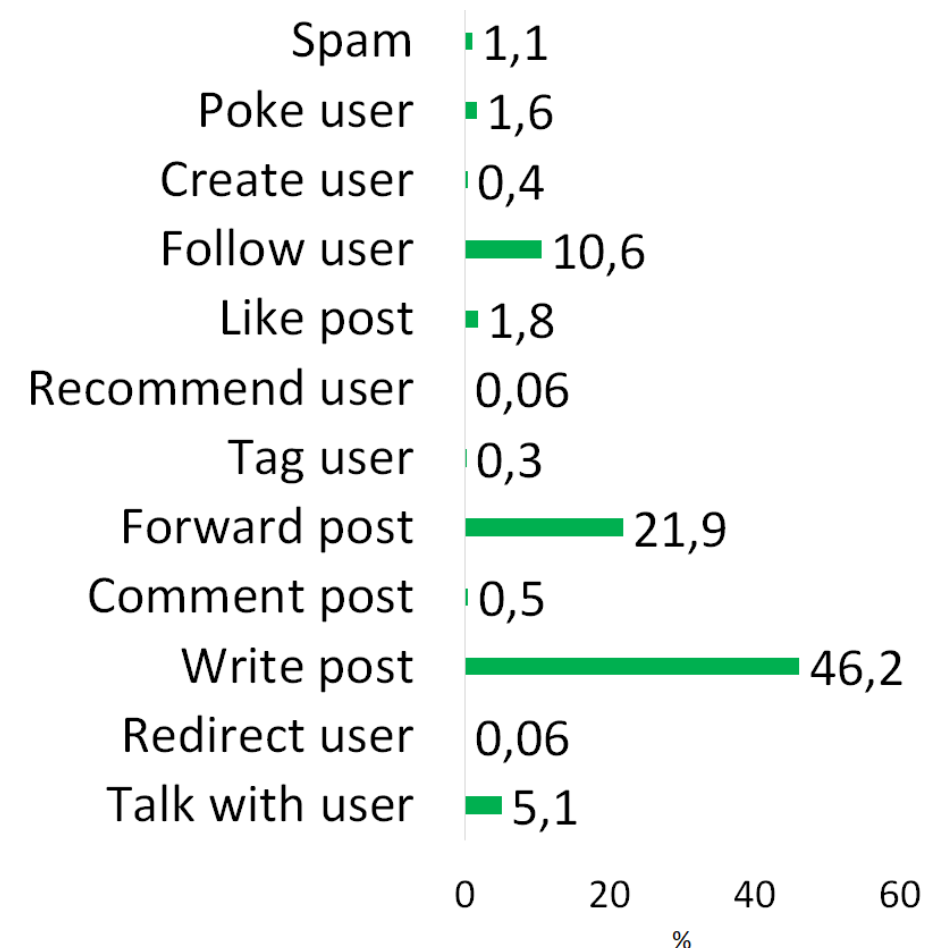


Fig. 2: Labels of repositories.

Starting point Final dataset

Selection criteria:

Programming language = Python

Exclusion of repositories that are out of scope

5 best repositories for each of the most used actions

5 random repositories from the rest

10 best repositories we could not classify

10 random repositories from the rest

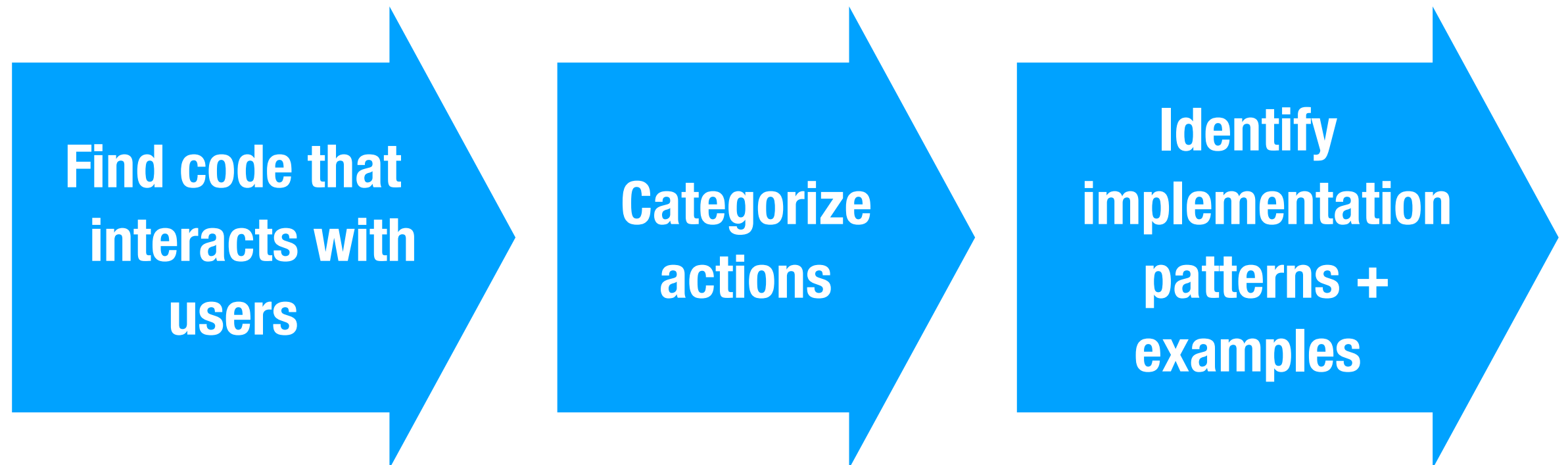
= 60 GitHub Twitter bot repositories programmed in Python

average number of files per repository: 3

average number of lines of code: 192

average size of repository: 21.39 KBytes

Methodology



Systematic, **manual** code review!

Results **Actions** for participation

Action	Description
<i>Search</i>	Search users or tweets using names, keywords, hashtags, ids or similar or by navigating social network relationships (e.g., friends of friends, followers of friends, friends of followers, followers of followers)
<i>Follow</i>	Follow users to establish social relationships
<i>Like</i>	Like tweets by other users to endorse them
<i>Tweet</i>	Post a new tweet to communicate content
<i>Mention</i>	Mention other users in tweets using @ to attract attention
<i>Retweet</i>	Re-post tweets by other users to endorse them
<i>Talk to</i>	Send direct messages to users to converse with them
<i>Pause</i>	Pause the conversation flow of the bot
<i>Store</i>	Store content retrieved from the social network for later use

Synthesis of online communication actions implemented by Twitter bots

Table 2: Taxonomy of code patterns used for the implementation of actions.

Action	Pattern	Description
Search	<i>User search</i>	Search user account by name, keyword, id or similar
	<i>Tweet search</i>	Search tweets by keyword or hashtag
	<i>Trend search</i>	Search trending topics or hashtags by location
Follow	<i>Indiscriminate follow</i>	Follow users without checking suitability of users, user-names or content shared
	<i>Whitelist-based follow</i>	Follow only users whose attributes or tweets match some element of a given whitelist
	<i>Blacklist-based follow</i>	Don't follow users whose attributes or tweets satisfy one or more criteria specified in a blacklist
	<i>Phantom follow</i>	Follow users and unfollow them as soon as a given condition is satisfied, e.g, a limit of friends reached or being followed back
Like	<i>Indiscriminate like</i>	Like tweets without checking suitability of content, user or username
	<i>Whitelist-based like</i>	Like only tweets by users whose attributes or content match some element of a whitelist
	<i>Blacklist-based like</i>	Don't like tweets whose attributes or users match an element of a blacklist
	<i>Mass like</i>	Aggressively like tweets of given users
Tweet	<i>Fixed-content tweet</i>	The content of the tweet is taken from a fixed, static collection of predefined messages
	<i>AI-generated tweet</i>	The text of the tweet is automatically generated using AI/NLP tools
	<i>Trusted source tweet</i>	The content of the tweet is taken from a source that can be considered trusted
	<i>Tweet with opt-in</i>	Tweets are sent only to people who ask to interact with the bot, sending it a message or mentioning it in a tweet
Mention	<i>Indiscriminate mention</i>	Mention other users without checking suitability of username or content shared
	<i>Targeted mention</i>	Classify users on the basis of their tweets and mention them in targeted messages
	<i>Whitelist-based mention</i>	Mention only users whose attributes match some element of a whitelist
	<i>Blacklist-based mention</i>	Don't mention users whose attributes match elements of a blacklist
Retweet	<i>Indiscriminate retweet</i>	Retweet tweets without checking content or username for suitability
	<i>Whitelist-based retweet</i>	Retweet content only from users whose attributes match some element of a whitelist
	<i>Blacklist-based retweet</i>	Don't retweet tweets whose attributes or users satisfy some condition expressed in a blacklist
	<i>Mass retweet</i>	Aggressively retweet multiple tweets by selected users
Talk to	<i>Indiscriminate talk</i>	Send direct, instant messages to users without checking their suitability
	<i>Talk with opt-in</i>	Reply only to messages sent to the bot (passive behavior)
	<i>AI-generated talk</i>	Generate messages using AI/NLP tools
	<i>Fixed-content talk</i>	Take message from a fixed list of predefined phrases
	<i>Targeted talk</i>	Classify users based on their tweets or attributes and target message accordingly
Pause	<i>Mimic human</i>	Use pauses in instant messages to deliver human-like conversation experience to other humans
	<i>Satisfy API constraints</i>	Use as short as possible pauses just to avoid being blocked by API usage limitations
Store	<i>Store persistently</i>	Store retrieved content or user information for later use

Search

Follow

Like

Tweet

Mention

Retweet

Talk to

Pause

Store

Action	Pattern	Description
Search	<i>User search</i>	Search user account by name, keyword, id or similar
	<i>Tweet search</i>	Search tweets by keyword or hashtag
	<i>Trend search</i>	Search trending topics or hashtags by location
Follow	<i>Indiscriminate follow</i>	Follow users without checking suitability of users, usernames or content shared
	<i>Whitelist-based follow</i>	Follow only users whose attributes or tweets match some element of a given whitelist
	<i>Blacklist-based follow</i>	Don't follow users whose attributes or tweets satisfy one or more criteria specified in a blacklist
	<i>Phantom follow</i>	Follow users and unfollow them as soon as a given condition is satisfied, e.g, a limit of friends reached or being followed back
Like	<i>Indiscriminate like</i>	Like tweets without checking suitability of content, user or username
	<i>Whitelist-based like</i>	Like only tweets by users whose attributes or content match some element of a whitelist
	<i>Blacklist-based like</i>	Don't like tweets whose attributes or users match an element of a blacklist
	<i>Mass like</i>	Aggressively like tweets of given users
Tweet	<i>Fixed-content tweet</i>	The content of the tweet is taken from a fixed, static collection of predefined messages
	<i>AI-generated tweet</i>	The text of the tweet is automatically generated using AI/NLP tools
	<i>Trusted source tweet</i>	The content of the tweet is taken from a source that can be considered trusted
	<i>Tweet with opt-in</i>	Tweets are sent only to people who ask to interact with the bot, sending it a message or mentioning it in a tweet

Example

Library

Action 1: search

```
for tweet in tweepy.Cursor(api.search, q=QUERY).items():  
    tweet.user.follow()
```

Action 2: follow

= search users pattern + indiscriminate follow pattern

Function definition

```
def mentions(count, max_seconds_ago, id_blacklist) :  
    return [mention for mention in api.mentions_timeline(count=count)  
            if not mention.id in id_blacklist]
```

Blacklist inclusion check

= blacklist-based mention pattern

Results Effects

Patterns may

Enable an abuse

↔ Logic that by design performs an abuse

Prevent an abuse

↔ Logic that prevents the bot from performing an abuse

Be vulnerable to **content abuse**

↔ Interactions with users and/or content that may be inappropriate

Be vulnerable to **trust abuse**

↔ Forward, store or analyze content retrieved from users

















Action	Pattern	Abuse	Disclose sensitive facts	Denigrate	Be grossly offensive	Be indecent or obscene	Be threatening	Make false allegations	Deceive	Spam	Spread misinformation	Mimic interest	Clone profile	Invalidate space
Follow	Indiscriminate follow	—	⚠	⚠	⚠	—	—	—	—	—	GO	—	—	—
	Whitelist-based follow	—	🚫	🚫	🚫	—	—	—	—	—	GO	—	—	—
	Blacklist-based follow	—	🚫	🚫	🚫	—	—	—	—	—	GO	—	—	—
	Phantom follow	—	⚠	⚠	⚠	—	—	—	—	—	GO	—	—	—
Like	Indiscriminate like	—	⚠	⚠	⚠	—	—	—	—	—	GO	—	—	—
	Whitelist-based like	—	🚫	🚫	🚫	—	—	—	—	—	GO	—	—	—
	Blacklist-based like	—	🚫	🚫	🚫	—	—	—	—	—	GO	—	—	—
	Mass like	—	⚠	⚠	⚠	—	—	GO	—	—	GO	—	—	—
Tweet	Fixed-content tweet	🚫	🚫	🚫	🚫	🚫	🚫	—	GO	—	—	—	—	—
	AI-generated tweet	⚠	⚠	⚠	⚠	⚠	⚠	—	GO	⚠	—	GO	—	—
	Trusted source tweet	🚫	🚫	🚫	🚫	🚫	🚫	—	GO	—	—	—	—	—
Mention	Indiscriminate mention	—	GO	GO	GO	—	—	—	GO	—	GO	—	—	—
	Opt-in mention	—	GO	GO	GO	—	—	—	🚫	—	—	—	—	—
	Targeted mention	—	⚠	—	—	—	⚠	—	—	—	—	—	—	—
	Whitelist-based mention	—	🚫	🚫	🚫	—	—	—	—	—	GO	—	—	—
	Blacklist-based mention	—	🚫	🚫	🚫	—	—	—	—	—	GO	—	—	—
Retweet	Indiscriminate retweet	—	⚠	⚠	⚠	⚠	⚠	⚠	—	⚠	GO	—	—	—
	Whitelist-based retweet	—	🚫	🚫	🚫	🚫	🚫	🚫	—	🚫	GO	—	—	—
	Blacklist-based retweet	—	🚫	🚫	🚫	🚫	🚫	🚫	—	🚫	GO	—	—	—
	Mass retweet	—	⚠	⚠	⚠	⚠	⚠	⚠	GO	⚠	GO	—	—	—
Talk to	Indiscriminate talk	—	—	—	—	—	—	—	GO	—	GO	—	—	—
	Fixed-content talk	🚫	🚫	🚫	🚫	🚫	🚫	—	GO	—	GO	—	—	—
	AI-generated talk	⚠	⚠	⚠	⚠	⚠	⚠	—	GO	⚠	GO	GO	—	—
	Talk with opt-in	—	—	—	—	—	—	—	🚫	—	🚫	—	—	—
	Targeted talk	—	—	—	—	—	—	—	GO	—	GO	—	—	—
Pause	Mimic human	—	—	—	—	—	—	GO	—	—	—	—	—	—
	Satisfy API constraints	—	—	—	—	—	—	—	GO	—	—	—	—	—
Store	Store persistently	⚠	—	—	—	—	—	—	—	—	—	⚠	—	—

GO Enables
 🚫 Prevents
 ⚠ Vulnerable to content abuse
 ⚠ Vulnerable to trust abuse

Results

Pattern-effect matrix =
 potential effects
 of patterns

Zoom into **Follow** patterns

Action	Pattern	Abuse												
		Disclose sensitive facts	Denigrate	Be grossly offensive	Be indecent or obscene	Be threatening	Make false allegations	Deceive	Spam	Spread misinformation	Mimic interest	Clone profile	Invade space	
Follow	Indiscriminate follow	—				—	—	—	—	—		—	—	
	Whitelist-based follow	—				—	—	—	—	—		—	—	
	Blacklist-based follow	—				—	—	—	—	—		—	—	
	Phantom follow	—				—	—	—	—	—		—	—	

Enables Prevents Vulnerable to content abuse Vulnerable to trust abuse

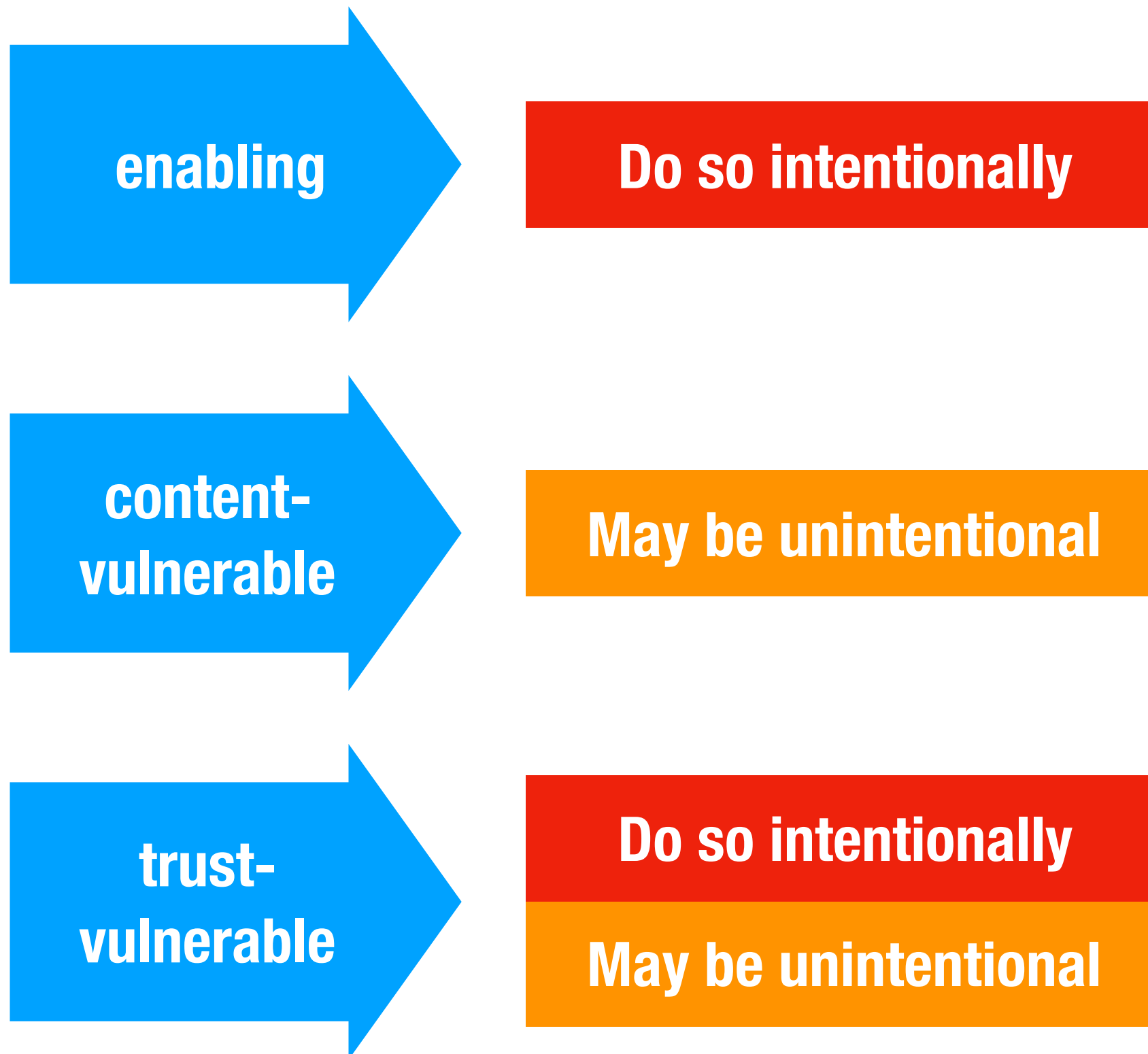
Zoom into **Tweet patterns**

Action	Pattern	Abuse												
		Disclose sensitive facts	Denigrate	Be grossly offensive	Be indecent or obscene	Be threatening	Make false allegations	Deceive	Spam	Spread misinformation	Mimic interest	Clone profile	Invade space	
Tweet	Fixed-content tweet							—		—	—	—	—	
	AI-generated tweet							—			—		—	
	Trusted source tweet							—		—	—	—	—	

Enables Prevents Vulnerable to content abuse Vulnerable to trust abuse

Results

Coming back to the “**why**” question... and using some technical considerations on the nature of patterns...



Summing up

Original perspective on bots for online communication: code

Contributions to state of the art:

1. Identified **31 patterns** and 9 actions from 60 repositories (~ 80 hours of manual code review + x of discussion)
2. Discussed **effects** of patterns and mapped patterns to potential abuses
3. Technical interpretation of **intentionality** underlying bot implementations

Next: formal language for action patterns + patterns search engine for automated pattern retrieval from all collected repositories