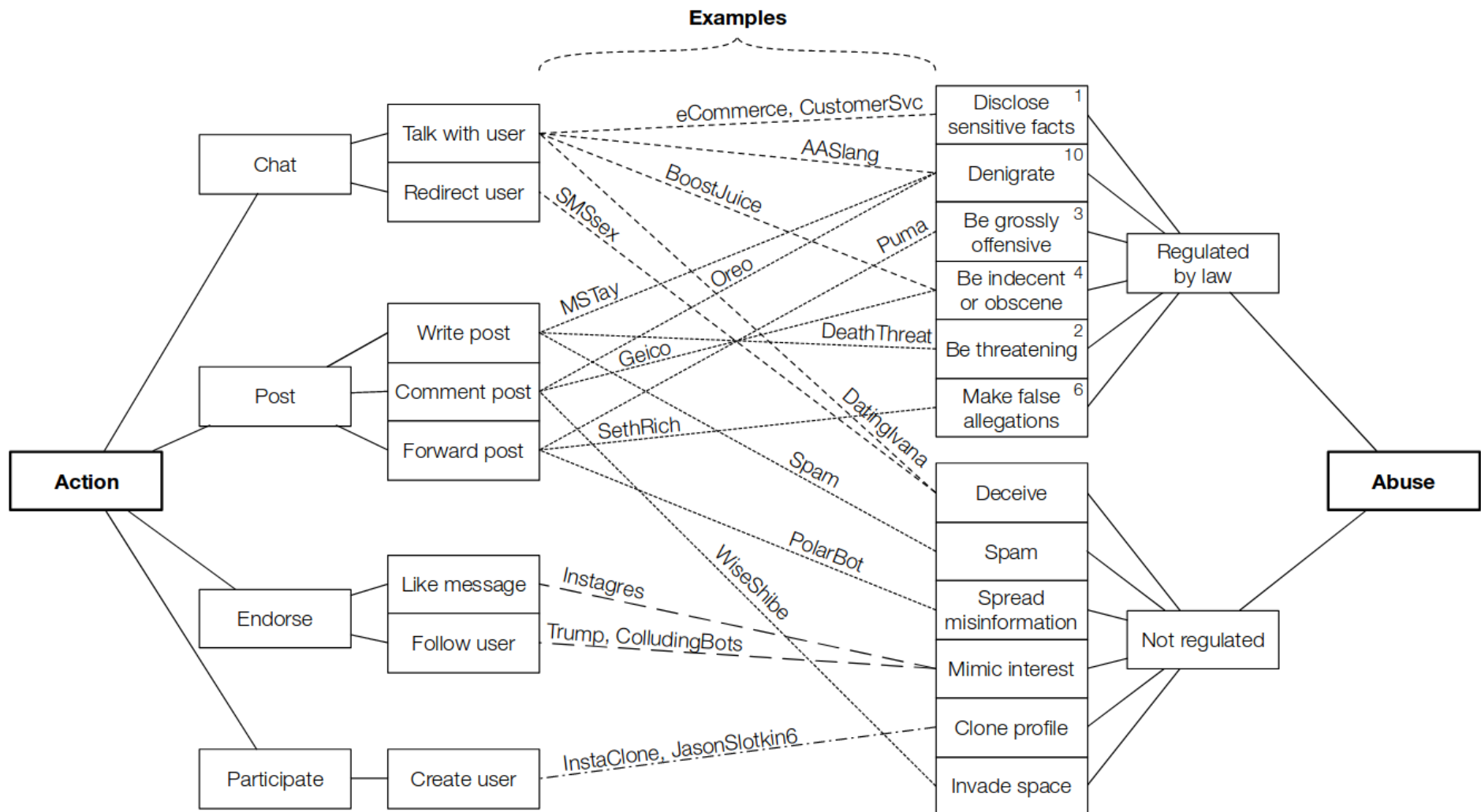




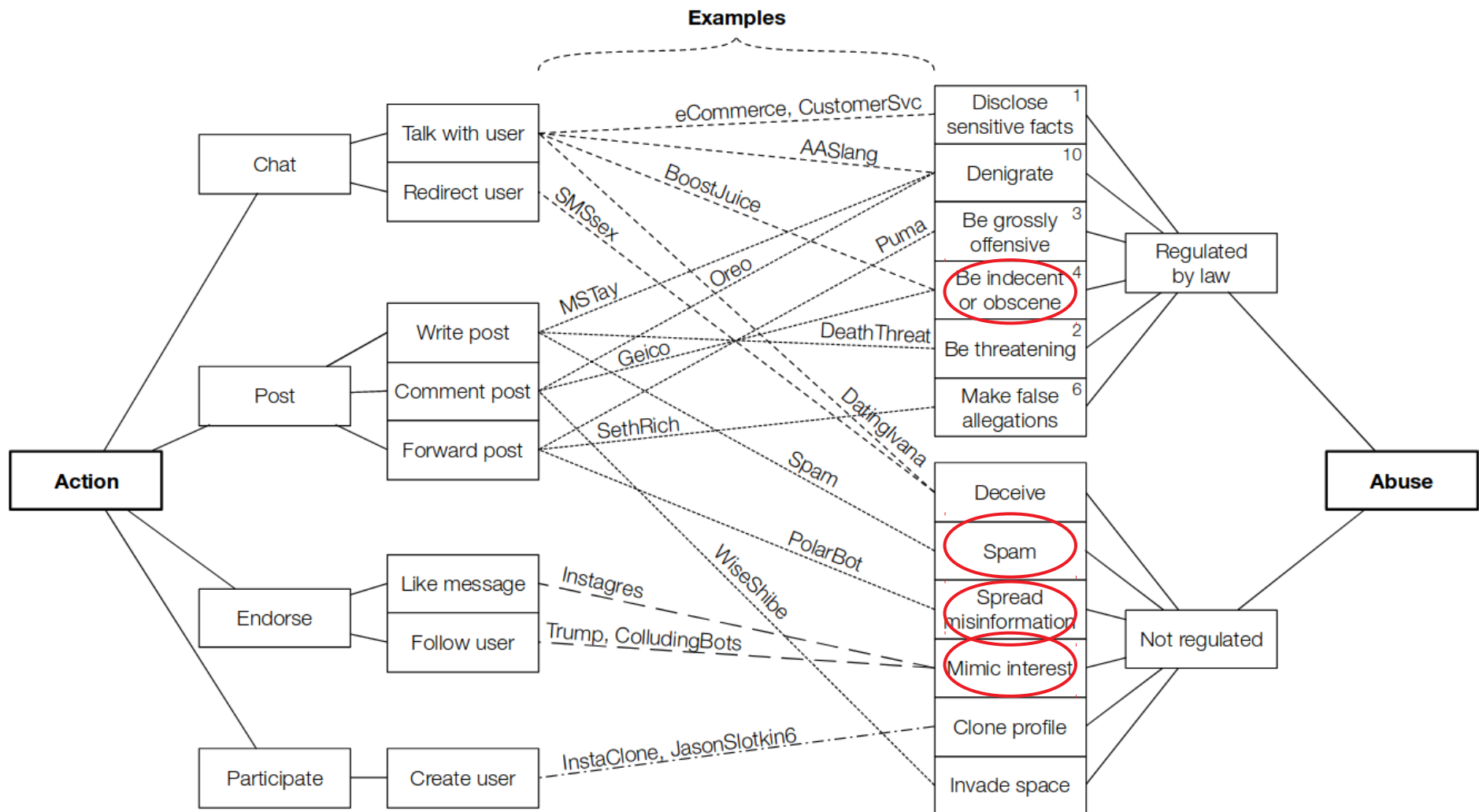
POLITECNICO
MILANO 1863

**Detection and Classification of Harmful Bots in Social
Human-Bot Interaction**

Social bots



Twitter bots

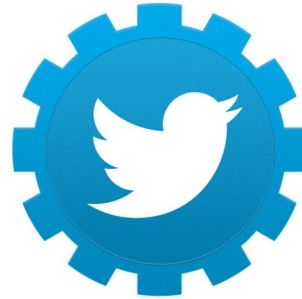


Twitter bots

- **NSFW**
- **News-spreaders**
- **Spambots**
- **Fake-followers**
- ✓ **Genuine**

Tools

→ **Twitter API**



→ **Botometer API**



→ **Hoaxy API**



Visualize the spread of claims and fact checking

Datasets

→ Caverlee-2011

- Bots
- Humans

→ Cresci-2017

- spambots (job offers)
- spambots (mobile app)
- Spambots (Amazon products)
- fake followers
- Humans

→ Varol-2017

- Bots
- Humans

→ BotBlock

- NSFW bots (adult contents)



Data Collection

NSFW

- Get ids from BotBlock list
- Scrape users and tweets information with Twitter API

Spambots

From Cresci dataset:

- traditional spambots 1 (generic Spabots)
- social spambots 2 (mobile app)
- social spambots 3 (Amazon products)

Data Collection

Fake-followers

From Cresci dataset:

→ Fake-followers

followers bought from.

→ instakipci.com/

→ rantic.com/buy-legit-twitter-followers

Genuine

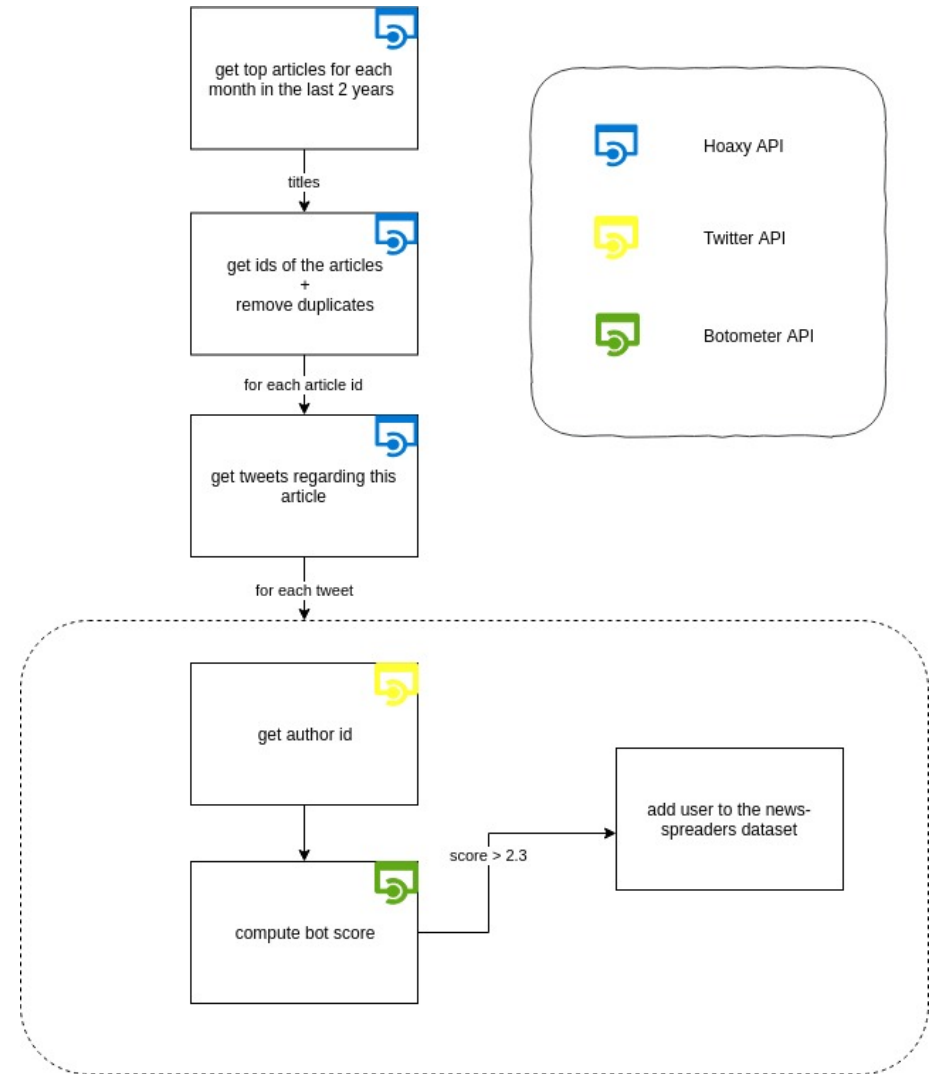
From Cresci dataset:

→ Genuine

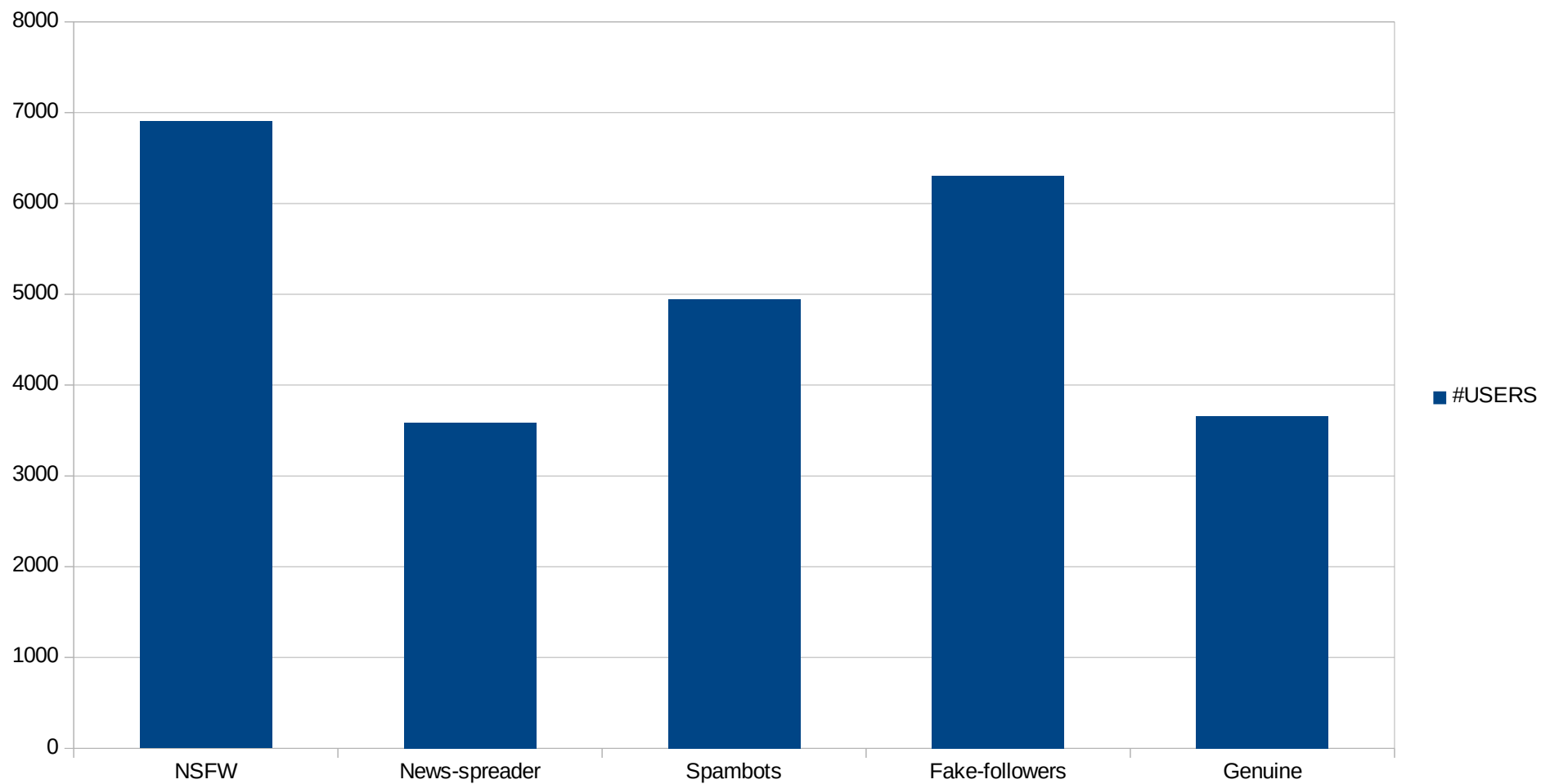
Data Collection

News-spreaders

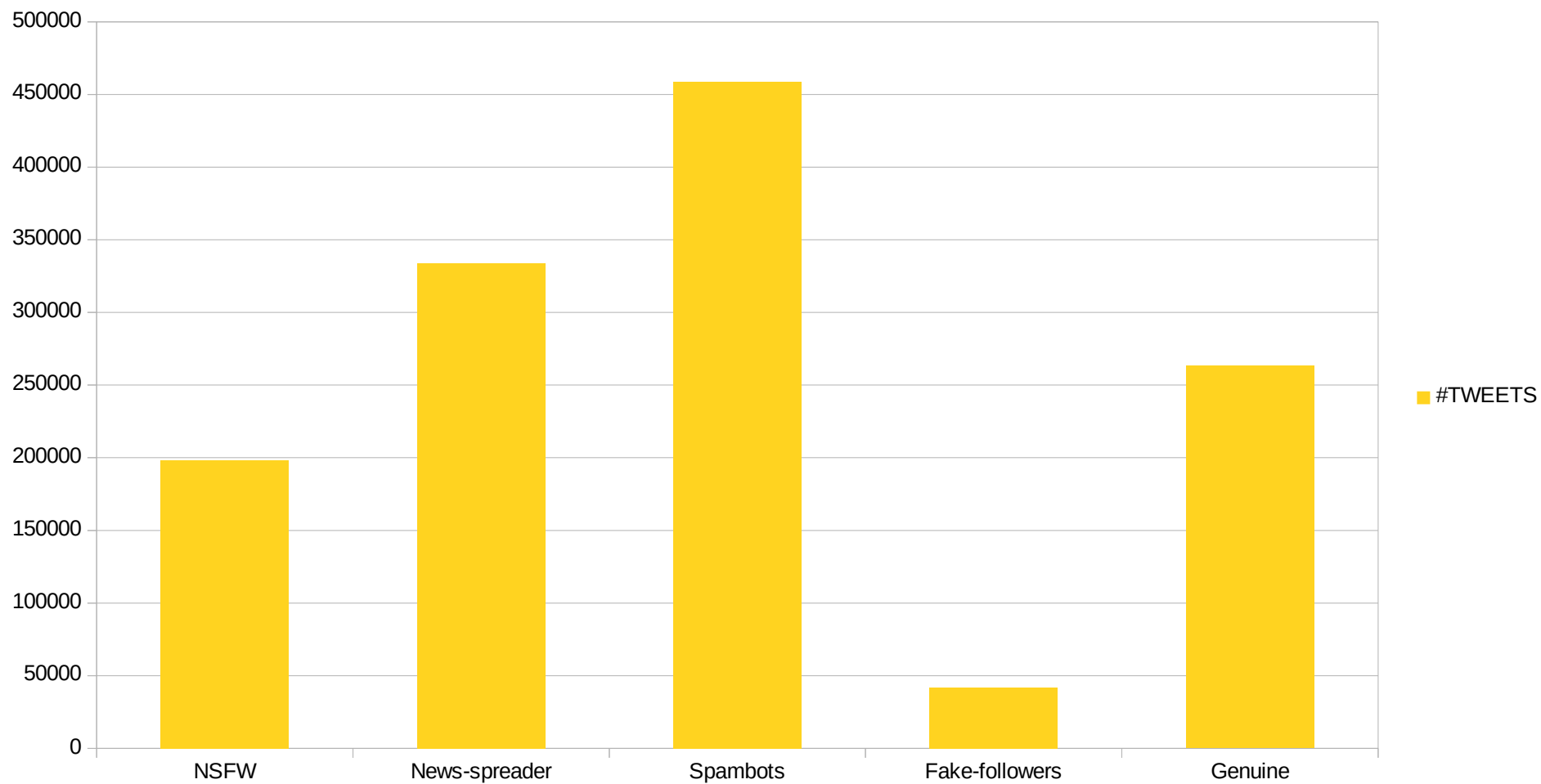
- Find fake-news tweets (Hoaxy)
- Find authors information (Twitter API)
- Check if authors are bots (Botometer)



Final dataset



Final dataset



Original Features

User Features

- Personal data
- Metadata
- Setting preferences

Tweet Features

- Text
- Media attached
- Mentions
- Source

Features Engineering

4 categories of features added

- Descriptive
- Intrinsic
- Extrinsic
- Image

Descriptive Features

- **"Meta features" related to tweets**
- **Synthesis statistics and counters**

Max, min, avg of tweet lengths	percentage of retweets, made by user, over its tweets
Max, min, avg number of retweets	percentage of media content incorporated in tweets
Max, min, avg number of likes	percentage of URL links placed inside tweets
Max, min, avg number of hashtag used	percentage of quotes, made by the user, over its tweets
amount of tweets per day (up to 100)	

Intrinsic Features

- Related to monotony of the user
- euclidean distance from a TF-IDF-encoded centroid

Tweets intradistance	Monotony coefficient about tweet texts
URL intradistance	Monotony coefficient about domain promoted

Extrinsic Features

- **Features related to common behaviours**
- **Creation of 5 different dictionaries**
- **Weighted words in dictionary**
- **Scores computed according to the use of words belonging to dictionaries**

NSFW words score

news spreaders words score

spam bots words score

fake followers words score

genuine words score

Image Features

- Related to profile image and attached media
- Scores computed with a CNN trained to identify NSFW contents

NSFW profile	NSFW evaluation of the profile picture
NSFW average	Average of NSFW scores of the last 10 tweeted media

Features Vector

The final feature vector contains 38 features

- 12 default user features
- 17 descriptive features
- 2 intrinsic features
- 5 extrinsic features
- 2 image features

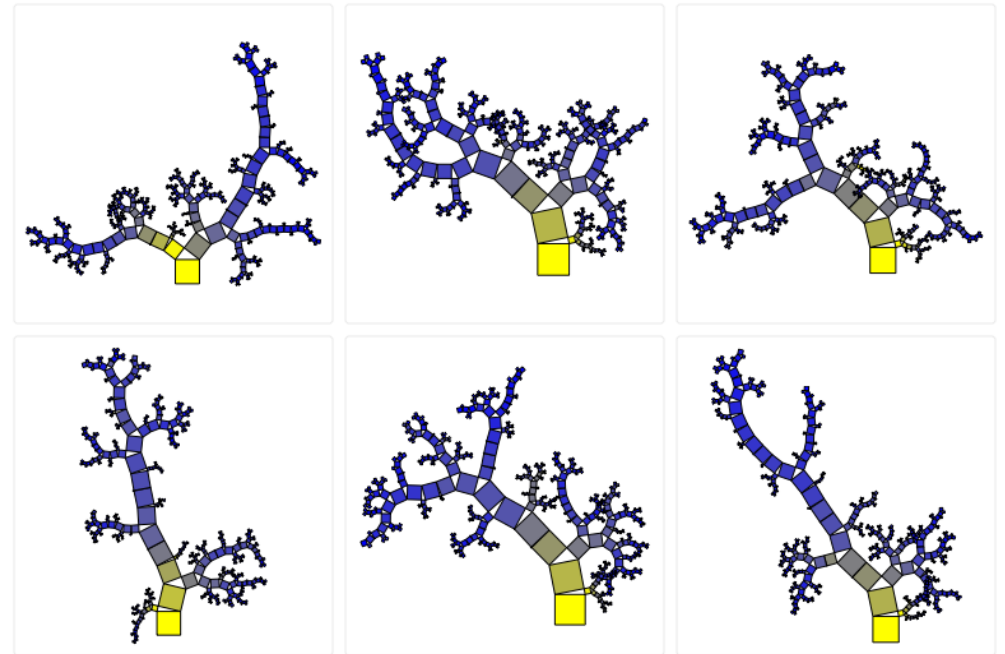
Models

The final solution is composed by a stacking ensemble of three different models

- Random Forest multiclass classifier for users
- Naive Bayes text classifier for tweets
- Random Forest binary classifier (bot or not)

Multiclass classifier

- RandomForestClassifier
 - Considers the whole features vector
 - Parameters tuning performed with Grid Search
 - `max_depth = None`
 - `n_estimators = 250`
 - `criterion = "entropy"`
-
- Cross validation score
 $f1 = 0.946$



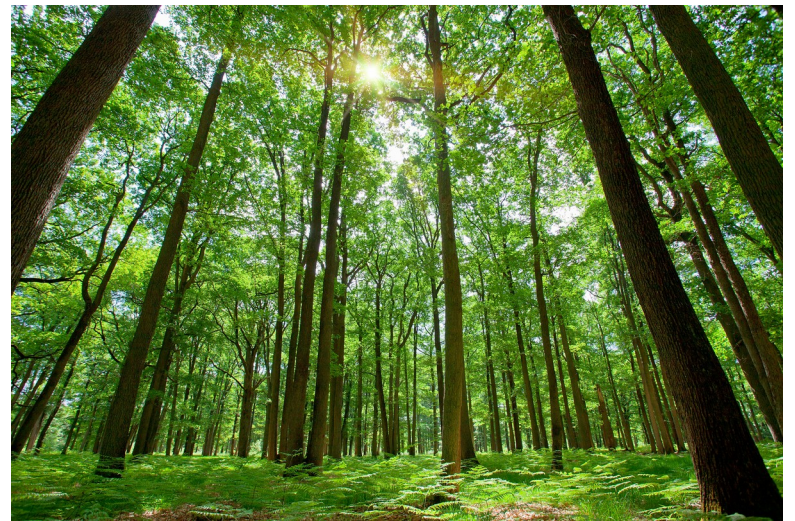
Text classifier

- Every tweet is labeled with the target of its author
- Considers only tweet texts
- Based on unigrams
- Pipeline
 - Stemming
 - TF-IDF encoding
 - MultinomialNB
- Final prediction over user is computed with the average of the predictions of his tweets
- Cross validation score
 - f1 = 0.71

$$p(C_k|\mathbf{x}) = \frac{p(\mathbf{x}|C_k)p(C_k)}{p(\mathbf{x})}$$

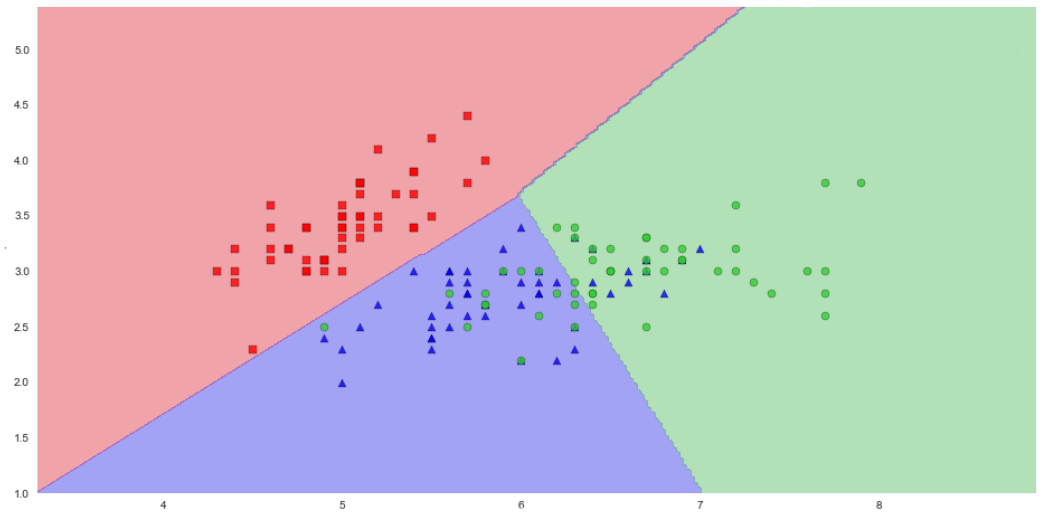
Binary classifier

- RandomForestClassifier
- Considers the whole features vector except image and extrinsic features
- Parameters tuning performed with Grid Search
- $\text{max_depth} = 26$
- $\text{n_estimators} = 175$
- $\text{criterion} = \text{"entropy"}$
- Cross validation score
 $\text{AUC} = 0.936$

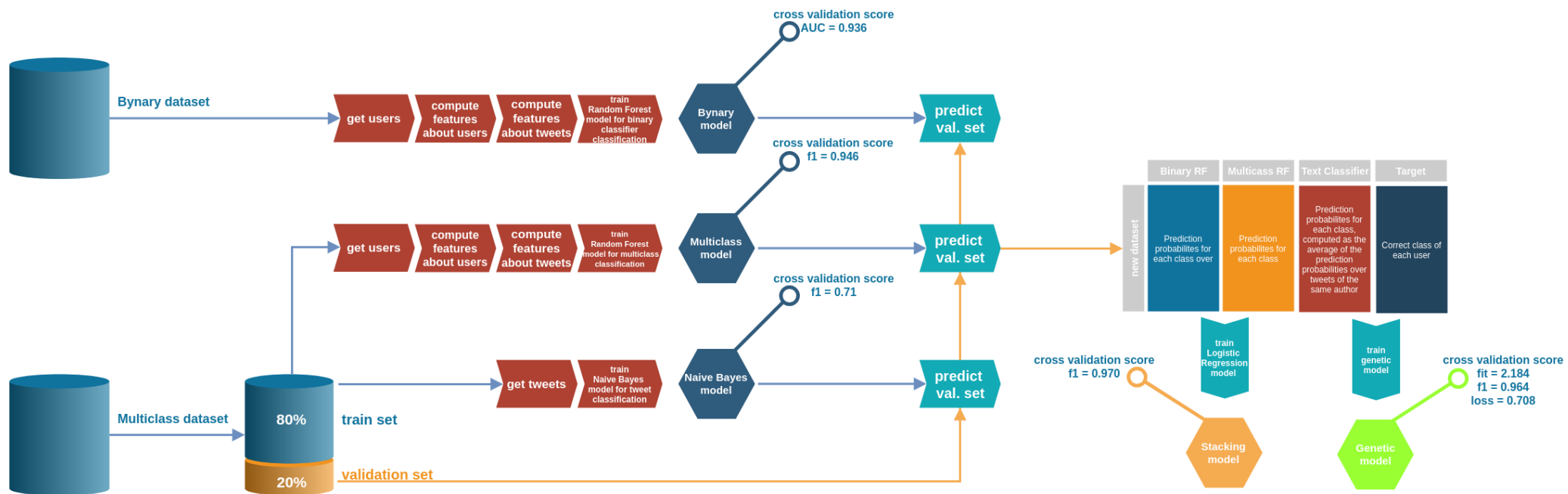


Ensemble

- Stacking with LogisticRegression
 - Input data consist in the output probabilities of the other models
 - `max_iter = 100`
 - `solver = "saga"`
 - `class_weight = "balanced"`
 - `multi_class = "multinomial"`
-
- cross validation score
 $f1 = 0.9734$



Stacking



Web App

Web service that allows users to classify twitter accounts

Frontend development

→ AngularJS



Backend development

→ Flask



Flask

BotBuster

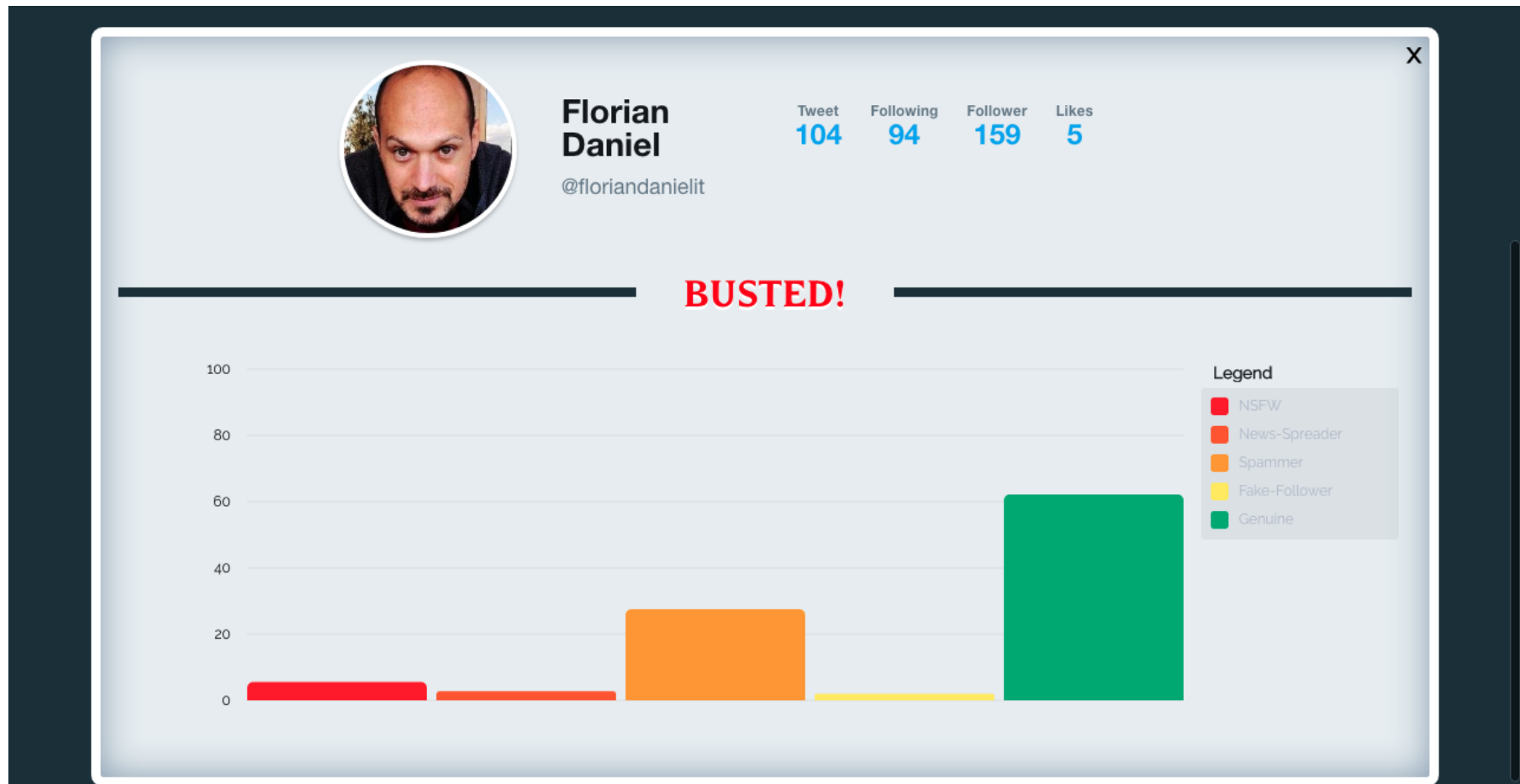


BOTBUSTER

Enter Twitter username

BUST IT

BotBuster



Thanks for your attention