# HATE SPEECH DETECTOR
# HATE SPEECH TWEETS CLASSIFIER WITH APPLICATION OF NEURAL NETWORKS

Karol Kowalski (220985), Jakub Kałużny (218159)

## DEFINITION

**Hate speech** denotes verbal or nonverbal attack against a person or a social group. This concept does not concern only the **insults on the background of protected attributes** such as race, religion or national origin. It concerns every **abusive, humiliative or bullying language**.

## DATASET AND CLASSIFIER

English training data were obtained from **Davidson et al.** Github webpage (https://github.com/t-davidson/hate-speech-and-offensive-language) and polish from **Poleval 2019** website (http://2019.poleval.pl/index.php/tasks/task6). English dataset contains 77.43% hate and 22.57% non-hate speech examples. Polish dataset contains respectively 8.92% and 91.08%.

The test data were scrapped from Twitter. Appropriate tweets were searched by hashtag 'Iran' as english data and 'tylkoniemownikomu' as polish. Then it has been performed following analyses. First, which hashtags clearly denote hate-speech classified tweets. Second, what are the percentages of hate- and non-hate-speech tweets denotation for 10 most frequent ambiguous hashtags.
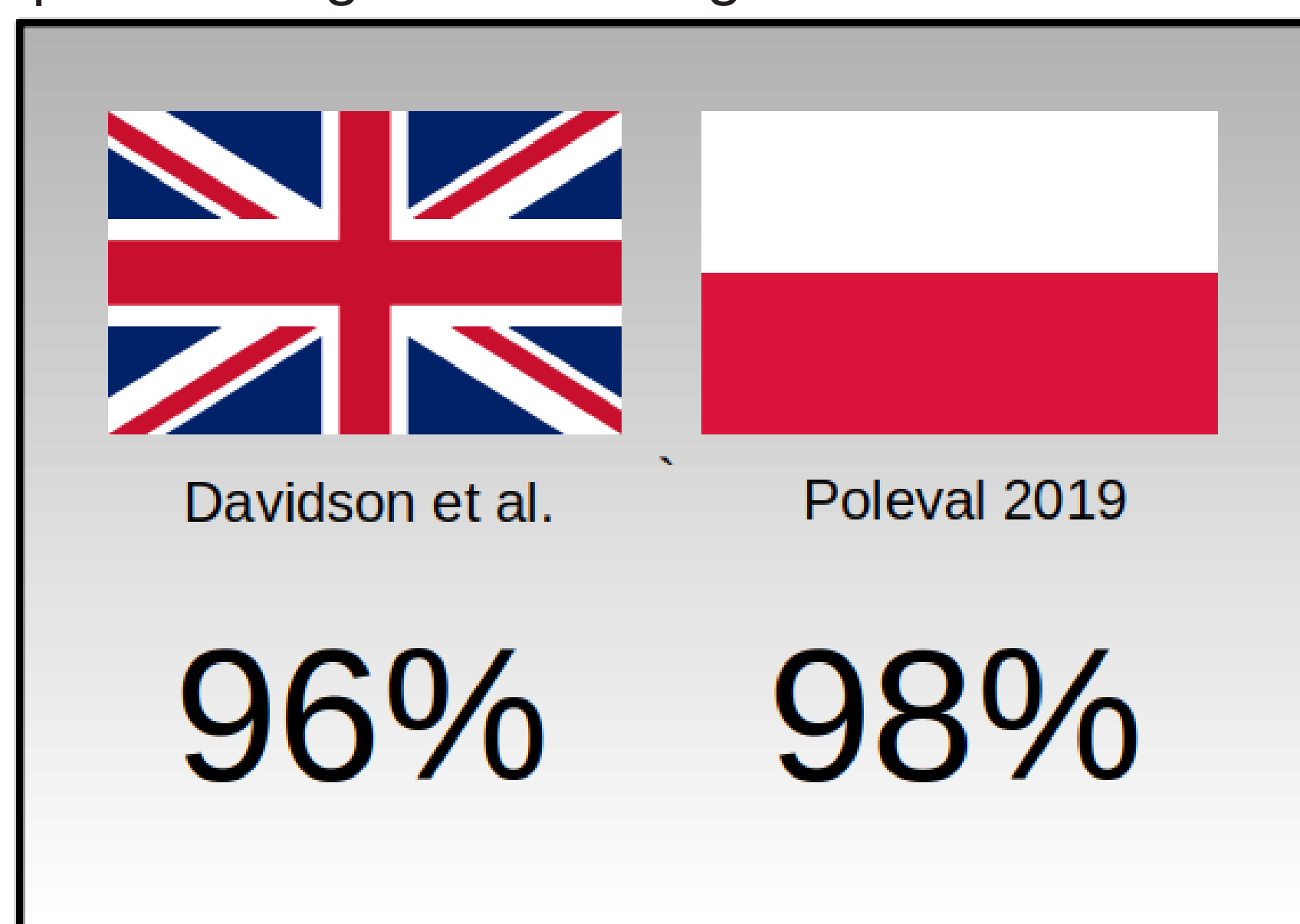
Davidson et al. — 96%    Poleval 2019 — 98%

*Figure 2:* LSTM neural network classifier accuracies for Davidson et al. and Poleval 2019 tweets.

## PHASES OF PROJECT

Phases of Project

5 — Summary — Drawing conlusions from twitter

4 — Scrapping twitter data — Gathering of real twitter data to analyse hate speech data on

3 — Models train — Searching for best deep neural network model by train it on test data.

2 — Searching for training Data — Searching for hate speech examples in polish and english language for twitter social network

1 — Defying hate speech — Find out what hate speech realy mean?

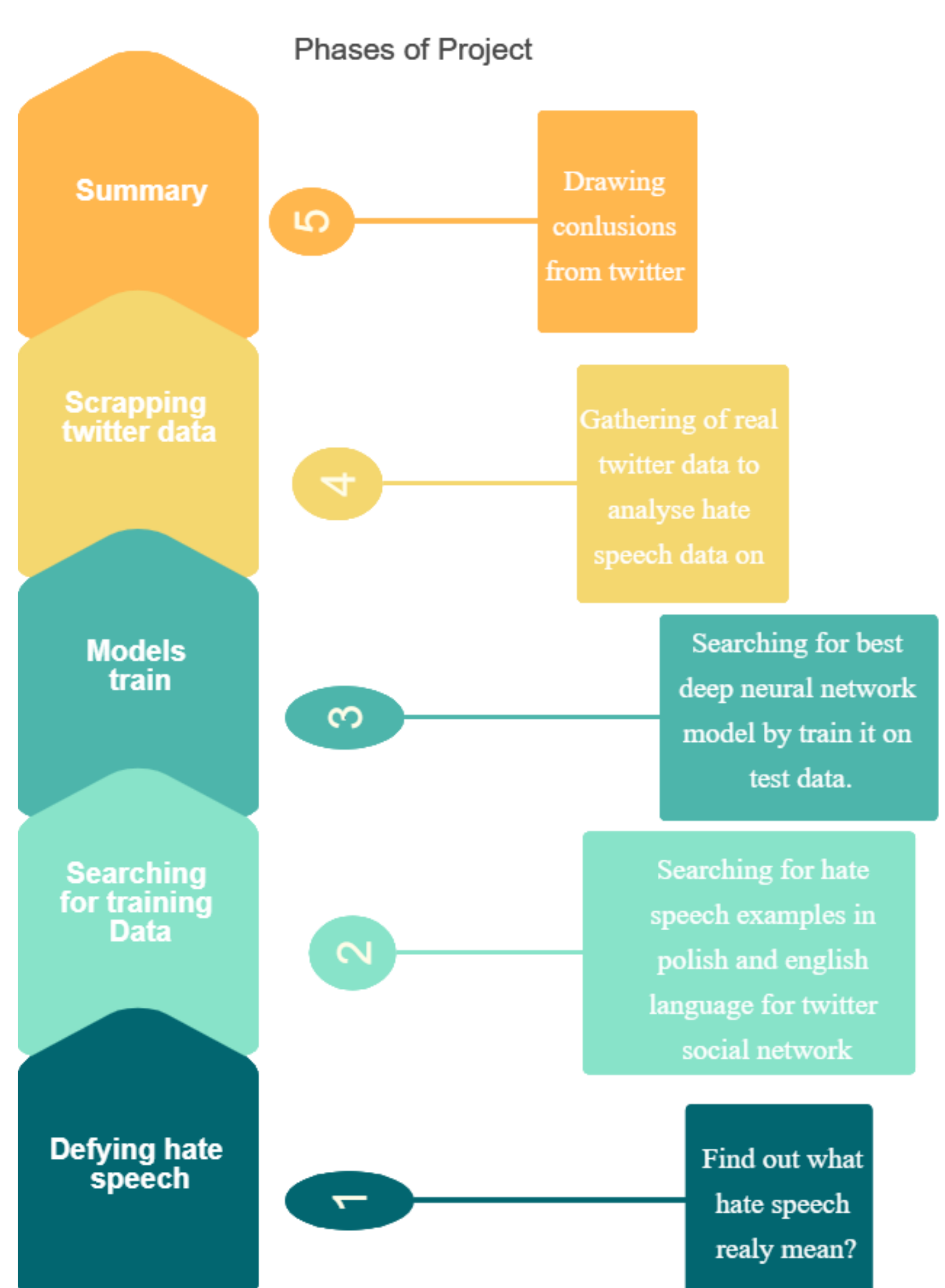*Figure 4:* Main project phases

## HASHTAGS CLOUDS

Hate-speech hashtags ("Iran")

Non-hate-speech hashtags ("Iran")

Hate-speech hashtags ("tylkoniemownikomu")

Non-hate-speech hashtags ("tylkoniemownikomu")

*Figure 1:* Examples of hashtags which appeared in hate- and non-hate-speech for english and polish tweets. Thanks to these clouds it could be found out which hashtags are strongly corelated with hate or non-hate speech. Judging by above analyses the **polish users of twitter are less prone to use hate speech than english users**. Whereas considering the hashtags contained in non-hate-speech classified polish tweets the conclusion is that the **model might have been tested on irrepresentative data** because some of the hashtags rather denote hate-speech tweets.

## HASHTAGS ANALYSIS

English hashtags analysis (topic: "Iran")
Precentages of: hate, non-hate-speech and ambiguous tweet hashtags distibutions

Clearcut hashtags

#irannewspaper #peterabaum #fromourcontributorsandmembers #shameonyoujavadzarifliarfakeminister #mini — 4.55

#fighterpilot #combat #shortreads #militaryaction #sirya — 95.45

hate-speech hashtags
non-hate-speech hashtags

Ambiguous hashtags
#iran 12.36 | #iranprotests 10.66 | #iraq 12.71 | #soleimani 8.23 | #wwiii 10.41 | #trump 9.38 | #worldwar3 8.65 | #internetiran 10.14 | #us 10.81 | #usa 8.07

Polish hashtags analysis (topic: "tylkoniemownikomu")
Precentages of: hate, non-hate-speech and ambiguous tweet hashtags distibutions

Clearcut hashtags

#smutnaprawda #marcinowerecenzje #rachon #francesco #strajknauczcycieli — 0.33

#wieszwięcej #ipptv #zbrodniekościoła #polexit #lgbt — 99.67

hate-speech hashtags
non-hate-speech hashtags

Ambiguous hashtags
#tylkoniemównikomu 0.82 | #tylkoniemownikomu 0.62 | #pedofiliawkościele 0.59 | #pedofilia 0.83 | #kościół 0.63 | #stopkonkordat 0.33 | #pis 1.01 | #sekielski 0.70 | #kler 0.39 | #pedofiliawkościele 0.90
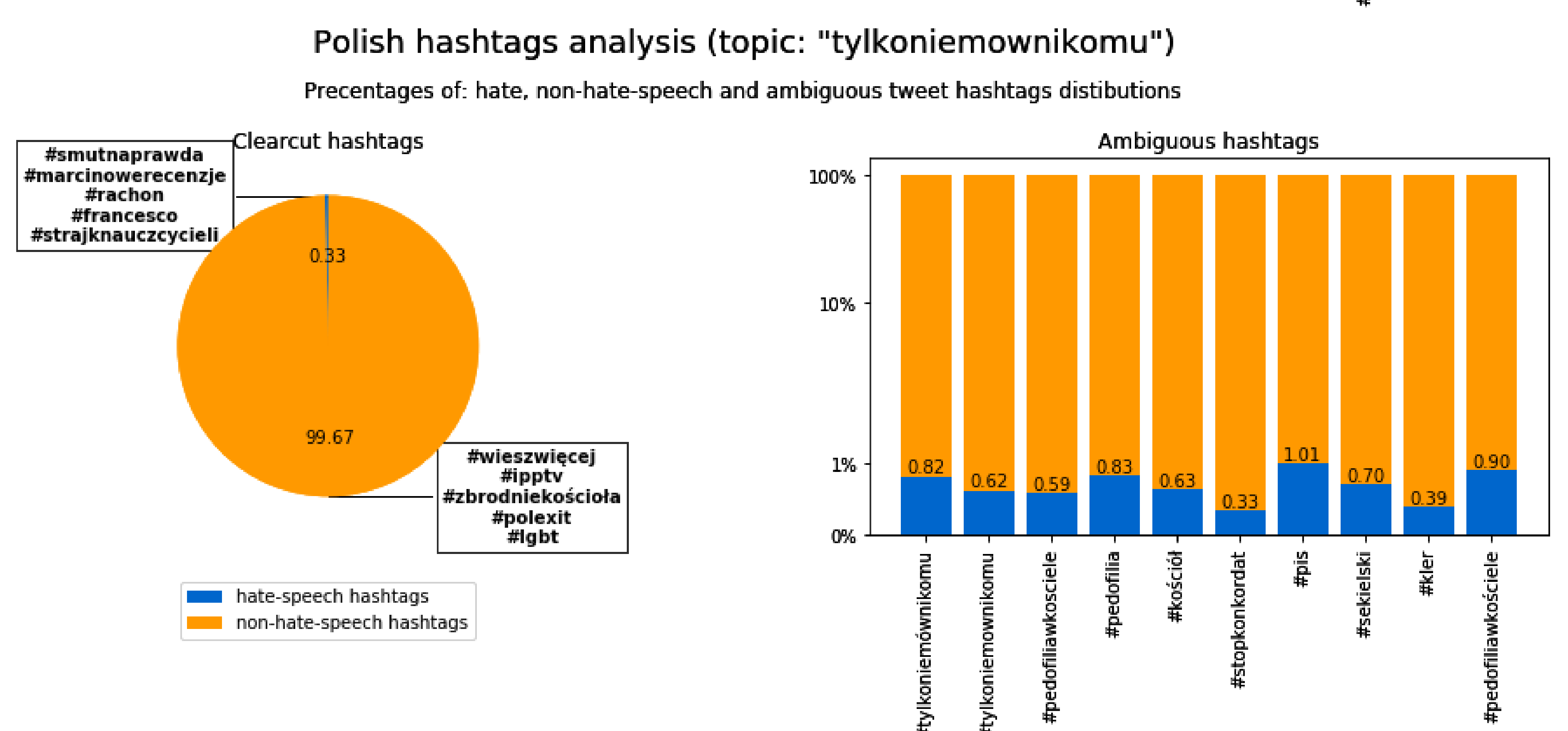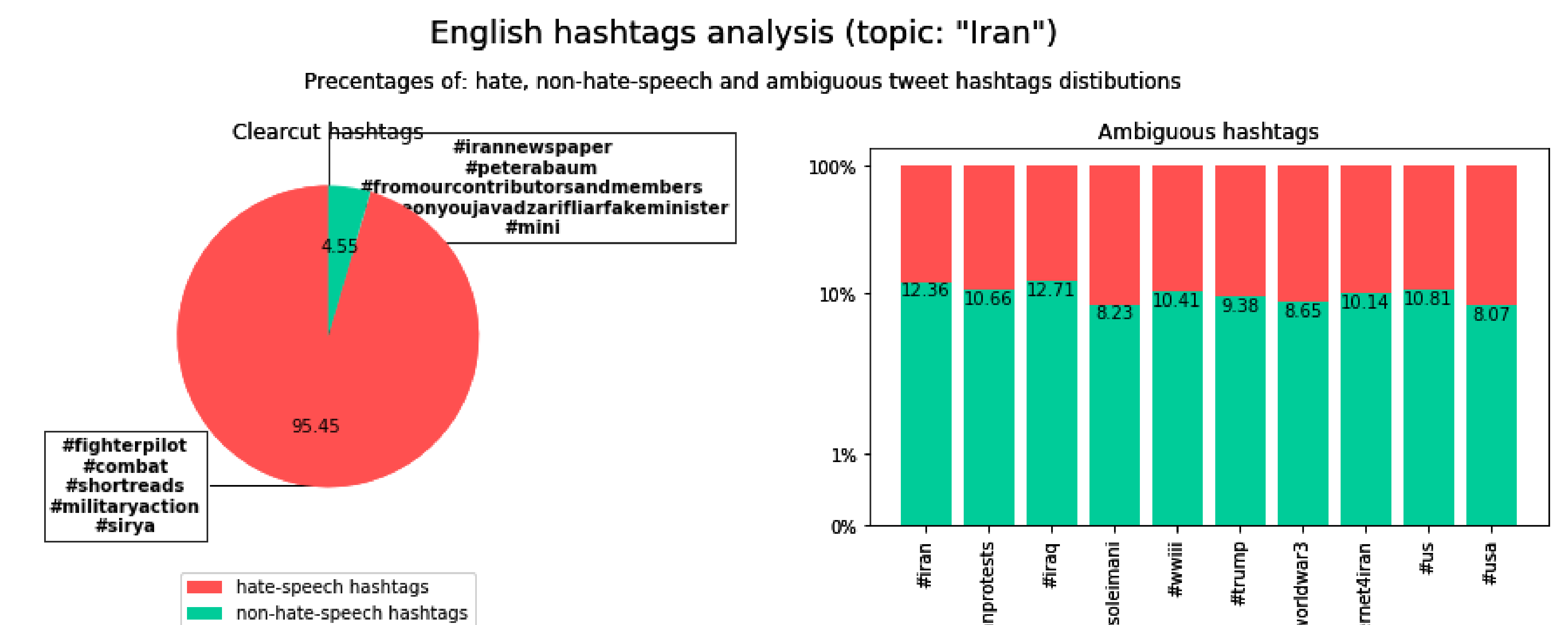
*Figure 3:* Ambiguous and clearcut (hate- or non-hate-speech) hashtags analysis

The above hashtag amount analysis shows that (just like by wordclouds) **english users are more prone to use hate speech in tweet than polish users**. These results however might also be **biased with imbalanced or irrepresentative tweet data**.