# AUTOMATED DETECTION OF INAPPROPRIATE COMMENTS IN ONLINE DISCUSSIONS

# moderateIT
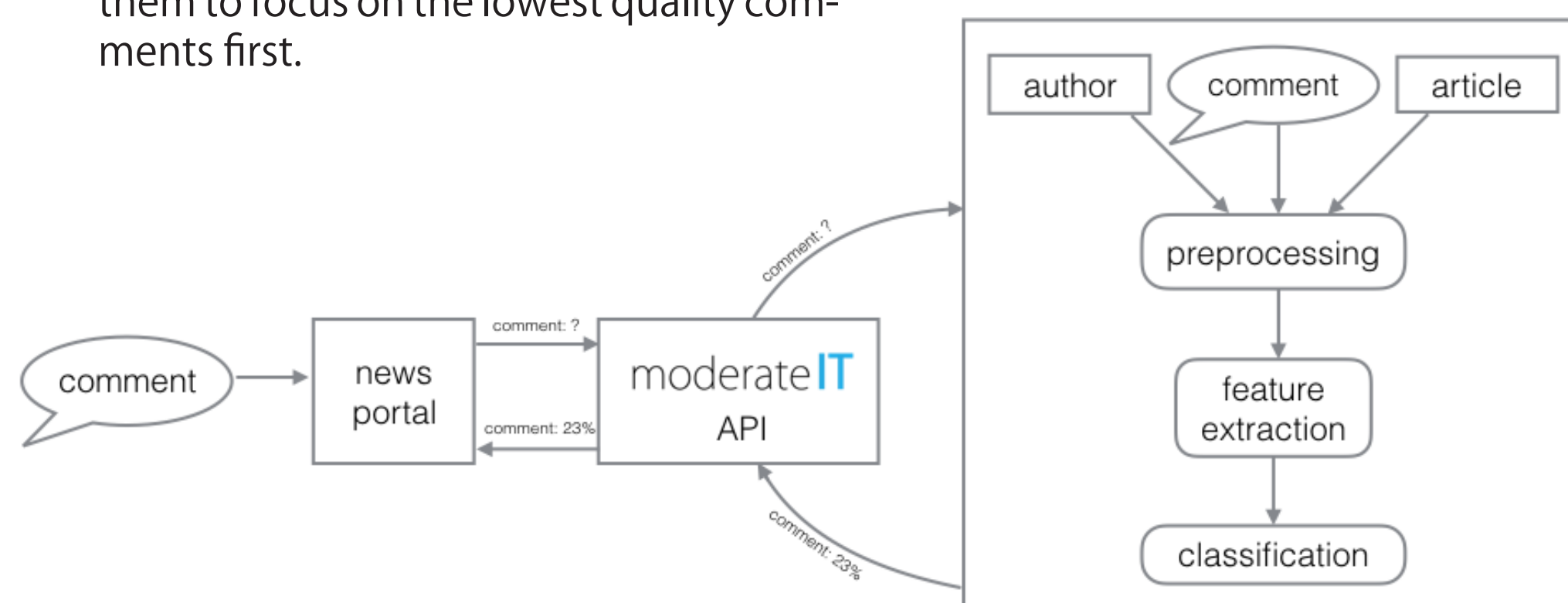## moderate online discussions using IT

Jakub Adam, Monika Filipčiková, Andrej Švec, Filip Vozár
supervisor: Jakub Šimko

## OVERVIEW

Our goal is to help online news portals easily manage discussion sections under their articles, which are very often full of hateful and inappropriate comments. Because of the volume of new comments being posted every day, moderators can't check and determine the quality of every comment.

In attempt to solve this situation, we've created moderateIT - a service which can automatically determine comment's quality based on several aspects. Moderators can then sort comments according to our evaluation, allowing them to focus on the lowest quality comments first.
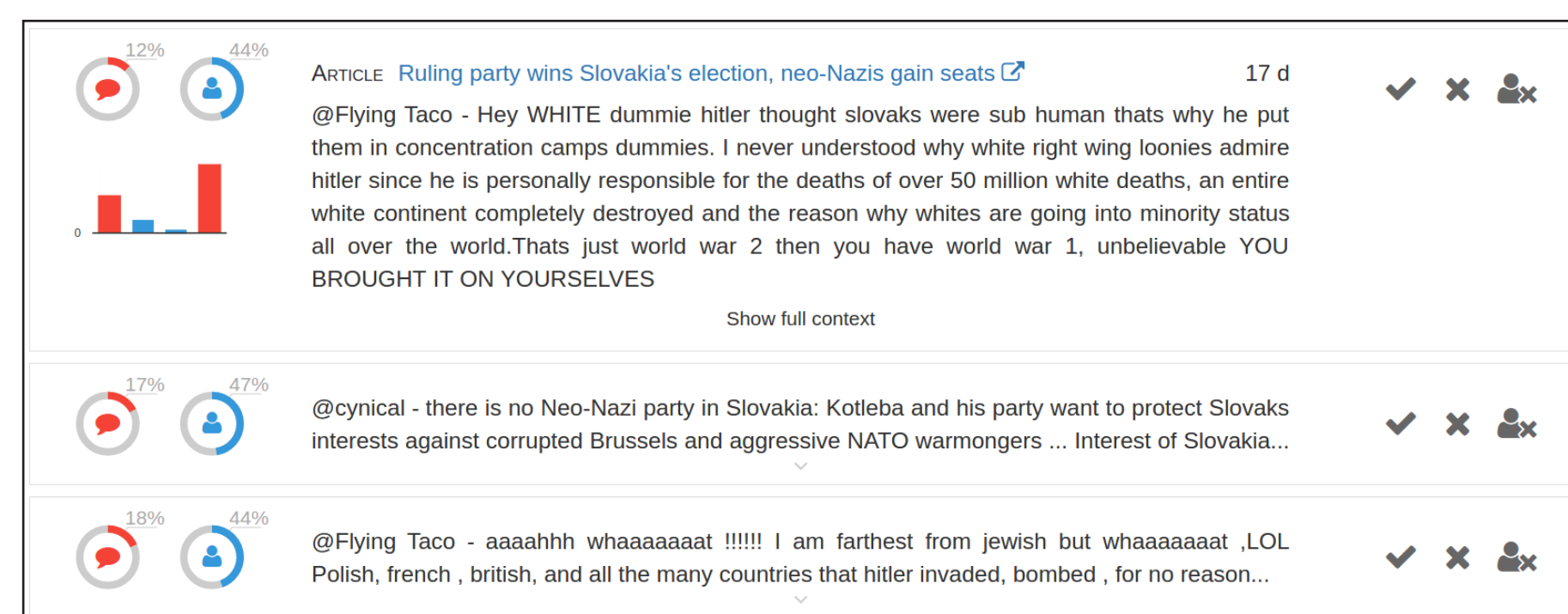
## FRONT-END



Figure 1. Screen with list of comments. They are sorted from worst to best.

We provide a moderation tool for moderators. It is a web front-end built upon our service. This application displays comments on their website along with many features concerning the comment. The moderation task is therefore more efficient. All performed actions are propagated immediately to the news portal. The portal serves also as a presentation tool for Imagine Cup judges or for our future customers.

## COMMENT QUALITY

The analysis process is split into multiple steps. The program runs multiple detectors which identify behavior of each and every comment. Their output then serves as input for classifier and based on these information our program evaluates whether the comment is appropriate or not.
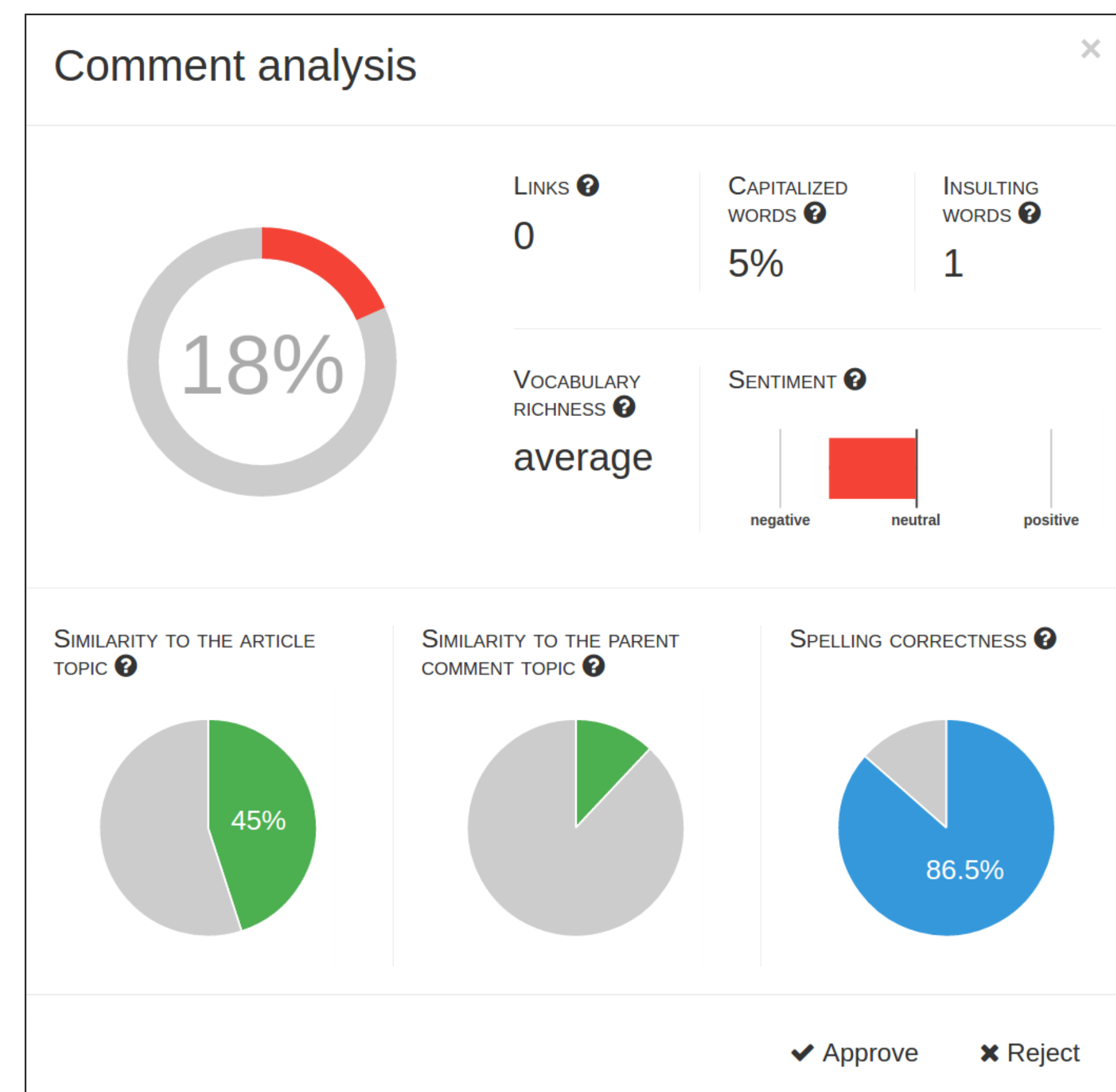


Figure 2. Screen with detailed comment analysis results.

### PREPROCESSING

- language detection
- diacritic for slovak text
- lemmatization
- keywords extraction (RAKE - Rapid Automated Keywords Extraction)

### DETECTORS

- insulting words detector
- sentiment detector
- hyperlinks detector
- off-topic detectors
- spellcheck correctness
- vocabulary richness
- emoticons or other non-word characters
- case sensitivity detector
- detector of sequences of the same characters
- formating detector

## AUTHOR'S REPUTATION

User modeling is used for personalisation purpose to improve accuracy of services and products. Most of the discussants have more than one comment in discussions, it is therefore possible to evaluate their history. From long-term point of view trolls write worse comments than usual participants do. We can define values from previous contributions useful for classifier.
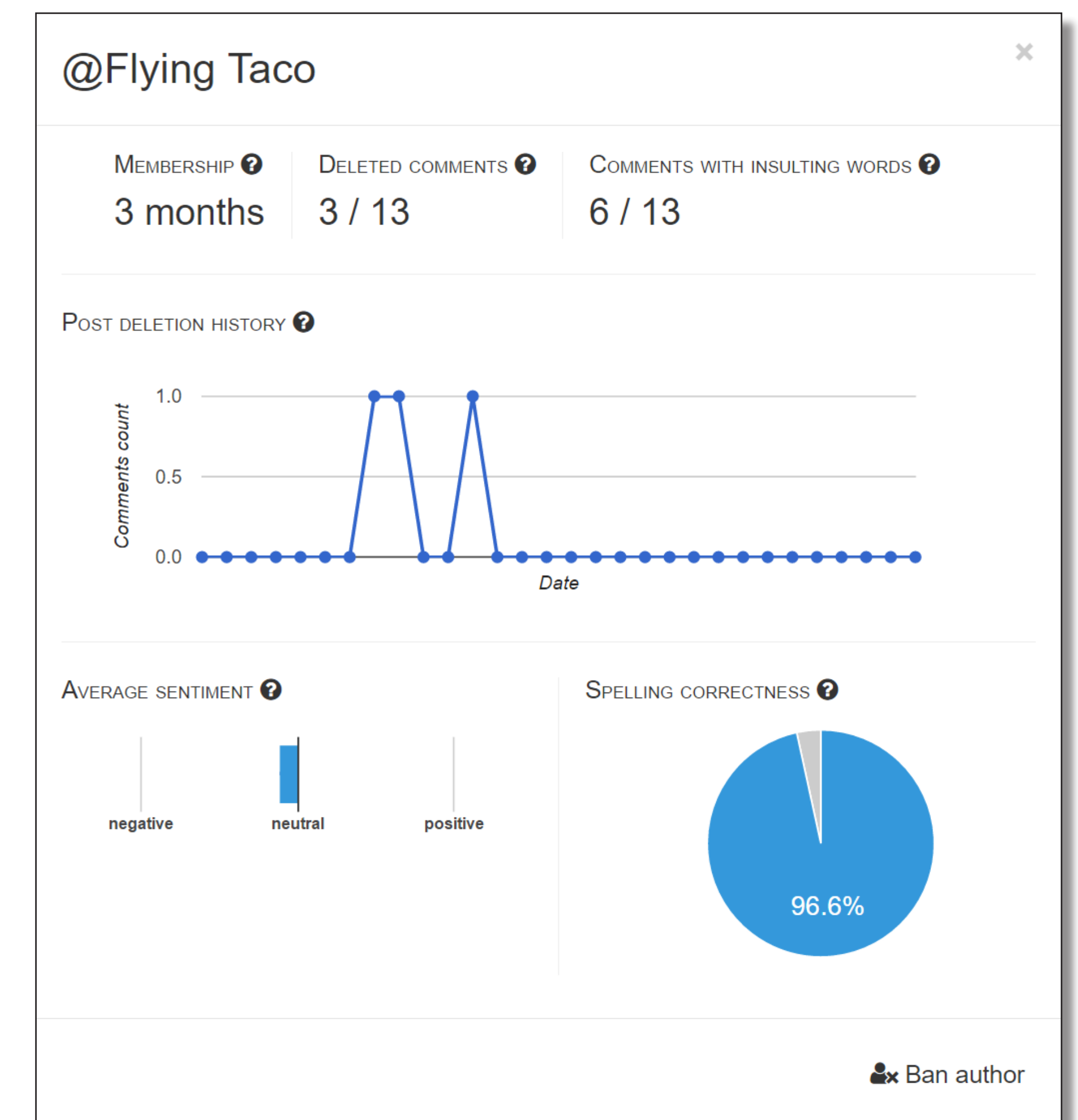


Figure 3. Screen with detailed author analysis results. These metrics allow the moderator to see whether author writes insulting comments or whether some of his comments were deleted by a moderator recently, etc.

### FEATURES

- duration of user's membership
- changes in user's behaviour
- user's reactiveness to another discussants
- another discussants reactiveness to user
- topic of discussion, where user writes