

A Distributed Architecture for Toxic Chat Detection

Wonhee Jung (wonheej2@illinois.edu) *University of Illinois Urbana-Champaign*, Kevin Mackie (kevindm2@illinois.edu) *University of Illinois Urbana-Champaign*, Cindy Tseng (cindyst2@illinois.edu) *University of Illinois Urbana-Champaign* and Conrad Harley (harley3@illinois.edu) *University of Illinois Urbana-Champaign*

Abstract—This paper presents the architecture, development, and use of a distributed system for toxic chat filtering. This system differentiates itself in that a) its filtering is based on machine learning and deeper contextual analysis, and b) it is deployed as a scalable and easily integrated web framework that can be adapted to any source of text for online interaction of any size. The platform is based on Docker and Kubernetes for easy deployment and dependency management, and uses state-of-the-art distributed systems technology to allow for fast scale-out to large systems. The system is presented in the context of a web chat application and Twitch chat bot as motivating examples.

Index Terms—toxic comment, web, chat, detoxifier, detection, classifier, distributed, architecture

I. INTRODUCTION

ONLINE platforms allow people to express their opinions freely, and stimulate collaboration across the globe. Unfortunately, online interaction may often come with loosened inhibitions in making profane, bigoted, or offensive remarks. We refer to such unwelcome remarks as “toxic chat”. Online systems may or may not have their own embedded profanity filtering, and those that do typically use pre-registered terms and simple pattern matching. This approach lacks the deeper contextual understanding needed to identify sentences that are toxic but that may not contain banned terms. Thus we propose a new toxic chat filtering system that differentiates itself in that a) its filtering is based on machine learning and deeper contextual analysis, and b) it is deployed as a scalable and easily integrated web framework that can be adapted to any source of text for online interaction of any size. The platform is based on Docker and Kubernetes for easy deployment and dependency management and to allow for fast scale-out to large systems. It uses state-of-the-art distributed systems technology for processing and storage, to allow for rapid scaling to any size while maintaining a shared file space (HDFS) between each Kubernetes Zone. This paper presents the architecture, development, and use of this system in the context of a web chat application and Twitch chatBot as motivating examples.

A. What we are going to make

We will create a prototype of PaaS/SaaS service that provides the following specific capabilities:

- machine-learning-based toxic chat identification and filtering engine
- integrated web chat application or chatbot that uses the engine to analyze a real-time stream of text

II. DATASETS

There is a dearth of labelled datasets for training classifiers to detect toxic comments. An online search and literature review were conducted on IEEE Xplore, Scopus, and Science Direct. We concluded that the best dataset available is Toxic Comment Classification Challenge dataset released by Jigsaw and Google on Kaggle in 2018, see [1]. Note that three additional datasets were identified - from Reddit [2], [3], Wikipedia [2], and Twitter [3], [4]. However, the datasets were not appropriate for use with our classifier.

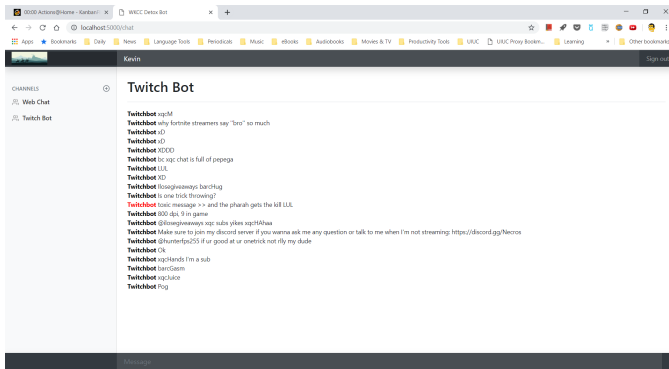
III. TECHNOLOGIES AND TOOLS

The following technologies and solutions were integrated to provide a general framework that can scale to high volume/traffic in the future.

- Flask, websockets, javascript, bootstrap - web application framework and client-server sockets
- Docker, Kubernetes - easy deploy and scale out
- CI/CD pipeline - to automate the build, integration, and deployment
- AWS, GCP, Heroku, etc - to deploy the solution into a mainstream PaaS infrastructure
- HDFS or similar - to store big data and share it between systems
- Scikit-learn or Apache Spark + MLlib - machine learning for the detox engine
- RESTful APIs - to help other applications integrate detox engine in the system

A. Web application

The user interface is implemented as a web application reachable via public URL. Python flask was chosen for the web application development framework, due to its simplicity, modularity, and compatibility with deployment to Google Cloud Platform. Client-side programming is done in Javascript with Bootstrap for the GUI framework and websockets for client-server socket connectivity for real-time chat.



Detoxifier Web Application

In the initial prototype, the web application connects the user to both a webchat and to the “Twitch Bot” that monitors a specific Twitch TV chat channel. Text typed into the web chat or received from Twitch is passed to the classifier. Toxic messages are marked in red with a prefix indicating that the message is toxic.

In future versions of the application a load balancer can be used to distribute traffic among several instances. Also, a complete application will allow plugs ins for different chat sources (Twitter, web forums, Reddit, Wikipedia), with the user able to specify the chat source and channel (currently this is hard-coded).

B. Deployment and scaling

- Docker, Kubernetes - easy deploy and scale out

C. Automated build, integration, deployment

- CI/CD pipeline - to automate the build, integration, and deployment

D. IaaS/PaaS infrastructure

- AWS, GCP, Heroku, etc - to deploy the solution into a mainstream PaaS infrastructure

E. Shared storage

- HDFS or similar - to store big data and share it between systems

F. Machine Learning Framework

- Scikit-learn or Apache Spark + MLlib - machine learning for the detox engine

G. Application Programming Interfaces

- RESTful APIs - to help other applications integrate detox engine in the system

IV. METHOD/DESIGN

A. Sub-section

V. PRELIMINARY EVALUATION/RESULTS

A. Sub-section

VI. DISCUSSION

A. Sub-section

VII. RELATED WORK (OPTIONAL FOR THIS MILESTONE)

A. Sub-section

VIII. FUTURE WORK

A. Sub-section

IX. DIVISION OF WORK (MAY OVERLAP)

A. Sub-section

X. CLOUD SYSTEM INTEGRATION AND DEVELOPMENT

ACKNOWLEDGMENT

The authors would like to thank ...

REFERENCES

- [1] Jigsaw, “Toxic comment classification challenge,” <https://www.kaggle.com/c/jigsaw-toxic-comment-classification-challenge/data..>
- [2] R. K. Éloi Brassard-Gourdeau, “Impact of sentiment detection to recognize toxic and subversive online comments,” <http://arxiv.org/pdf/1812.01704v1>, 2018.
- [3] R. G. Chandra Khatri Behnam Hedayatnia, “Detecting offensive content in open-domain conversations using two stage semi-supervision,” <http://arxiv.org/pdf/1811.12900v1>, no. NIPS CONVAI Workshop 2018, 2018.
- [4] R. K. Betty van Aken Julian Risch, “Challenges for toxic comment classification: An in-depth error analysis,” <http://arxiv.org/pdf/1809.07572v1>, nos. ALW2: 2nd Workshop on Abusive Language Online to be held at EMNLP 2018 (Brussels, Belgium), October 31st, 2018, 2018.

Wonhee Jung (wonheej2@illinois.edu)

Kevin Mackie (kevindm2@illinois.edu)

Cindy Tseng (cindyst2@illinois.edu)

Conrad Harley (harley3@illinois.edu)