# Detect-le-Defect

## A Python Project

Smart Censoring Application for Live-streaming platforms
To identify and classify Chat behavior
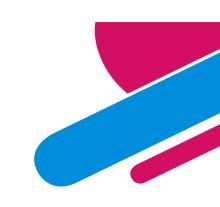Using Social Sentiment Analysis and Machine Learning

**DSTI - 21**

# Meet le Team



**Gustavo Chinchayan**
Operations



**Aleksandra Lazic**
Analytics



**Felipe Lopez**
Development



**Manuel Gawert**
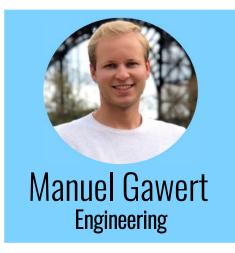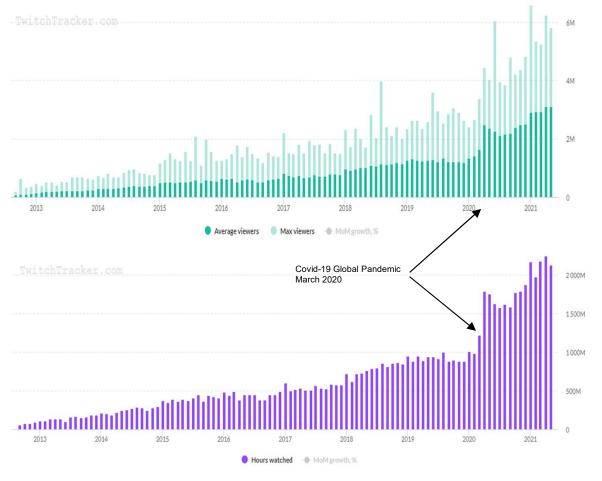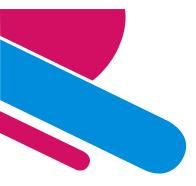Engineering
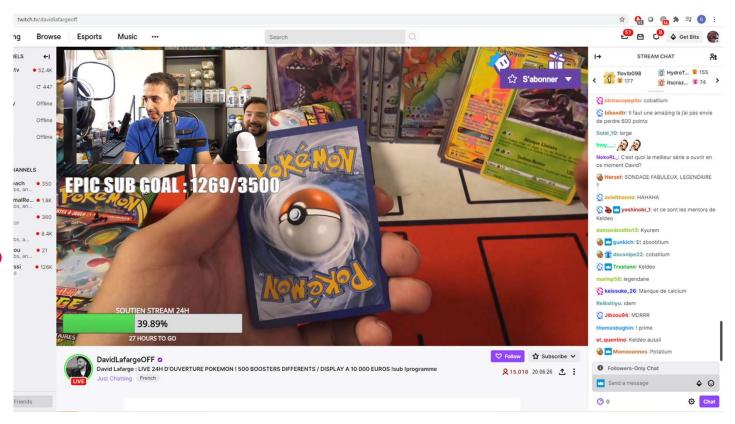
Covid-19 Global Pandemic
March 2020

# Twitch Platform

- Founded in 2011
- Acquired by Amazon for US$970 million
- Synergies with Amazon prime
- Optional subscription as a service
- Generates revenue via Twitch Partner/Affiliates and commercial deals
- Esports audience growth forecasting to be at 646 million by 2023

amazon

3

# Point of Interest

## Just Chatting Section



- ❑ Comprised of content creators/ influencers as Twitch Partners
- ❑ Community-driven based on their genre and taste
- ❑ Avg stream length 3-4 hours
- ❑ Engaged audience providing live feedback every second
- ❑ User feedback made up of strings, integers, and emoji (coded as text)

# Content Moderation

❑ **Mature content is easily accessible**

❑ **Hate Speech and Harassment is an issue**

❑ **Emojis/memes being used to perpetuate racism**

❑ **Limited control and analysis over audience feedback**



**Child Predators Use Twitch to Systematically Track Kids Livestreaming**
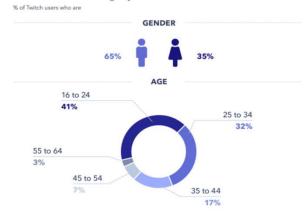
The gaming platform's design enables people to find and exploit kids in real-time

By Cecilia D'Anastasio
Graphics and illustrations by Rachael Dottle
Published: September 21, 2022 | Updated: October 19, 2022

Twitch, the Amazon subsidiary where millions of people congregate every day to watch skilled gamers play franchises like Fortnite and Minecraft, is one of the most popular websites on the internet. But the factors that have contributed to its rapid growth, such as the ease with which anyone can open an account and begin broadcasting themselves live, have also enabled predators to target young users, according to an analysis from October 2020 through August 2022 by a researcher who studies livestreaming websites.

## The Twitch Demographic

% of Twitch users who are

**GENDER**

65%   35%

**AGE**

16 to 24 **41%**
25 to 34 **32%**
55 to 64 **3%**
45 to 54 **7%**
35 to 44 **17%**

**Question:** Which of the following best describes your gender?/How old are you?
**Source:** GlobalWebIndex Q2 2019 **Base:** 15,853 Twitch users aged 16-64 (excl. China)

5

# Solution Stages

Provide a content feedback SaaS based on live stream user experience

## Planning Stage

### Setting Baseline

Data extraction and cleaning the data Frame

## Analysis Stage

### Action Plan

Discover trends and patterns

## Modelling Stage

### Optimize
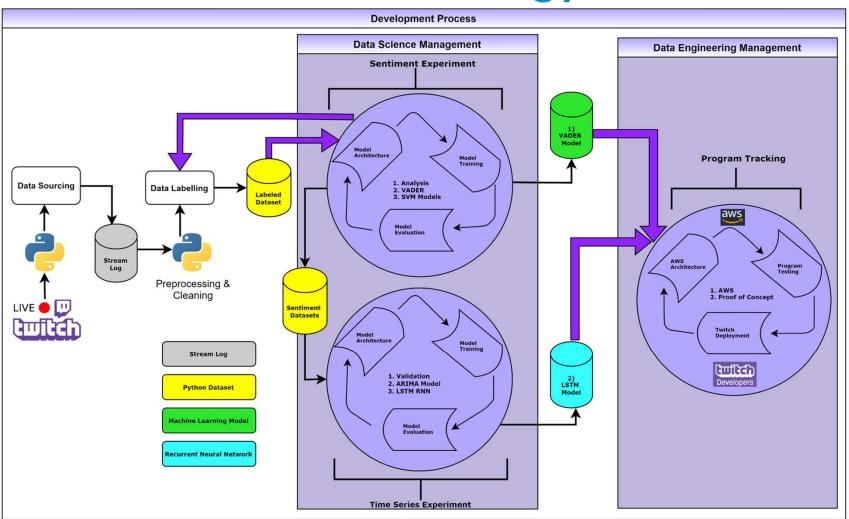
Machine Learning Model Tuning

## Final Stage

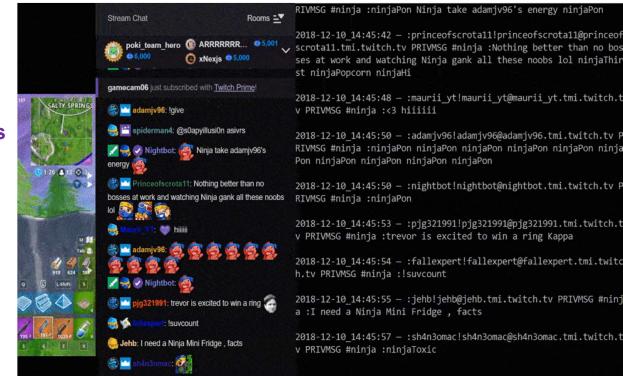### Implementation

Extrapolate and Integrate

# Methodology

# Data Extraction

**Web Scraping highly viewed streams**

**Via Twitch Internet Relay Chat (IRC)**

- Twitch Developer Account

- Using Python socket library

- Socket based in Time (seconds)

- Dataset downloaded as a Log

# Data Preprocessing & Cleaning

Original message

| | Date | Channel | Username | Message |
|---|---|---|---|---|
| 0 | 2021-10-18 17:18:04 | riotgames | xoocboots | EU i kcuf your women and slap your men LUL <3 ... |
| 1 | 2021-10-18 17:18:04 | riotgames | charaf54 | go home |
| 2 | 2021-10-18 17:18:10 | riotgames | secsso23 | 2-4 GG |
| 3 | 2021-10-18 17:18:10 | riotgames | gracejacky | tl had no dmg wtf |
| 4 | 2021-10-18 17:18:10 | riotgames | payab_005 | MAKE OR BREAK BatChest MAKE OR BREAK BatChest ... |

Stemming : Stemming is the process of reducing inflected (or sometimes derived) words to their word stem, base or root form and lowercasing

| | Date | Channel | Username | Message |
|---|---|---|---|---|
| 0 | 2021-10-18 17:18:04 | riotgames | xoocboots | eu i kcuf your women and slap your men LUL <3 ... |
| 1 | 2021-10-18 17:18:04 | riotgames | charaf54 | go home |
| 2 | 2021-10-18 17:18:10 | riotgames | secsso23 | 2-4 gg |
| 3 | 2021-10-18 17:18:10 | riotgames | gracejacky | tl had no dmg wtf |
| 4 | 2021-10-18 17:18:10 | riotgames | payab_005 | make or break batchest make or break batchest ... |

# Data Preprocessing & Cleaning

Removal of Stopwords (i, me, my, myself, we, our, ours, ourselves, you, you're, you've, you'll, you'd, your, yours)

| | Date | Channel | Username | Message |
|---|---|---|---|---|
| 0 | 2021-10-18 17:18:04 | riotgames | xoocboots | eu kcuf women slap men LUL <3 eu kcuf women sl... |
| 1 | 2021-10-18 17:18:04 | riotgames | charaf54 | go home |
| 2 | 2021-10-18 17:18:10 | riotgames | secsso23 | 2-4 gg |
| 3 | 2021-10-18 17:18:10 | riotgames | gracejacky | tl dmg wtf |
| 4 | 2021-10-18 17:18:10 | riotgames | payab_005 | make break batchest make break batchest make b... |

Remove duplicate words

| | Date | Channel | Username | Message |
|---|---|---|---|---|
| 0 | 2021-10-18 17:18:04 | riotgames | xoocboots | eu kcuf women slap men LUL <3 |
| 1 | 2021-10-18 17:18:04 | riotgames | charaf54 | go home |
| 2 | 2021-10-18 17:18:10 | riotgames | secsso23 | 2-4 gg |
| 3 | 2021-10-18 17:18:10 | riotgames | gracejacky | tl dmg wtf |
| 4 | 2021-10-18 17:18:10 | riotgames | payab_005 | make break batchest |

Remove words that only have one character

10

# Data Preprocessing & Cleaning

Remove the rare words

Remove all the links and page spam

if(('https' in message) or ('.com' in message))

We need to remove messages that contain only numbers.

Make an spelling check for all the words that are in our lexicon (Vader lexicon) + emotes/emojis
https://github.com/cjhutto/vaderSentiment/blob/master/vaderSentiment/vader_lexicon.txt

! After all this preprocessing we need to verify and remove if we have or not message

| | Date | Channel | Username | Message |
|---|---|---|---|---|
| 0 | 2021-10-18 17:18:04 | riotgames | xoocboots | slap LUL <3 |
| 1 | 2021-10-18 17:18:04 | riotgames | charaf54 | |
| 2 | 2021-10-18 17:18:10 | riotgames | secsso23 | gg |
| 3 | 2021-10-18 17:18:10 | riotgames | gracejacky | wtf |
| 4 | 2021-10-18 17:18:10 | riotgames | payab_005 | |

Cleanest message

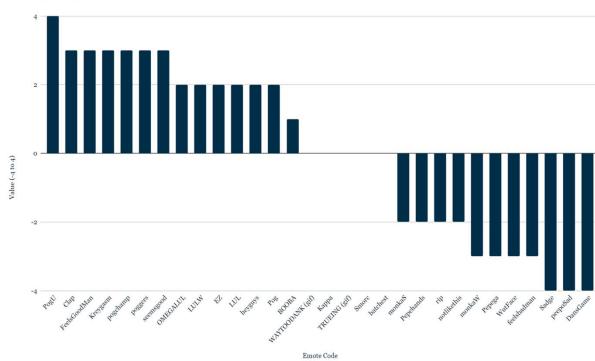| | Date | Channel | Username | Message |
|---|---|---|---|---|
| 0 | 2021-10-18 17:18:04 | riotgames | xoocboots | slap LUL <3 |
| 2 | 2021-10-18 17:18:10 | riotgames | secsso23 | gg |
| 3 | 2021-10-18 17:18:10 | riotgames | gracejacky | wtf |

# Sentiment Analysis

- Discovered 30 most common emotes based on Dataset
- Declare Sentiment and rating based on our judgement: LINK
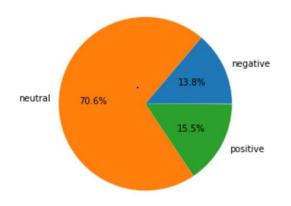- Declare Twitch emotes to Lexicon Dictionary



Value (-4 to 4) vs. Emote Code

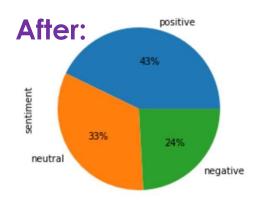# ValenceAware Dictionary for Sentiment Reasoning (VADER)

## Before:



VADER (Valence Aware Dictionary and Sentiment Reasoner) is a lexicon and rule-based sentiment analysis tool that is specifically attuned to sentiments expressed in social media.

It is used for sentiment analysis of text which has both the polarities i.e. positive/negative. VADER is used to quantify how much of positive or negative emotion the text has and also the intensity of emotion.
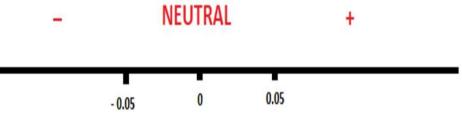
```
analyser.lexicon.update(new_emotes)
```

## After:



| | Date | Channel | Username | Message | Tag_owner | Cleanest_message | scores | comp_score | sentiment |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 2021-10-18 17:18:04 | riotgames | xoocboots | EU i kcuf your women and slap your men LUL <3 EU i kcuf your women and slap your men LUL <3 | 0 | slap LUL i | {'neg': 0.0, 'neu': 0.158, 'pos': 0.842, 'compound': 0.6523} | 0.6523 | positive |
| 1 | 2021-10-18 17:18:10 | riotgames | secsso23 | 2-4 GG | 0 | go | {'neg': 0.0, 'neu': 1.0, 'pos': 0.0, 'compound': 0.0} | 0.0000 | neutral |
| 2 | 2021-10-18 17:18:10 | riotgames | gracejacky | tl had no dmg wtf | 0 | wif | {'neg': 0.0, 'neu': 1.0, 'pos': 0.0, 'compound': 0.0} | 0.0000 | neutral |
| 3 | 2021-10-18 17:18:10 | riotgames | ritogames69420 | Win 4 borfday | 0 | win | {'neg': 0.0, 'neu': 0.0, 'pos': 1.0, 'compound': 0.5859} | 0.5859 | positive |
| 4 | 2021-10-18 17:18:10 | riotgames | tlev1n | LETS GO TL FeelsStrongMan Clap | 0 | Clap | {'neg': 0.0, 'neu': 0.0, 'pos': 1.0, 'compound': 0.4588} | 0.4588 | positive |

13

# VADER Analysis



NEUTRAL

$-$        NEUTRAL        $+$

-0.05     0     0.05

positive sentiment: compound score >= 0.05
neutral sentiment: (compound score > -0.05) and (compound score < 0.05)
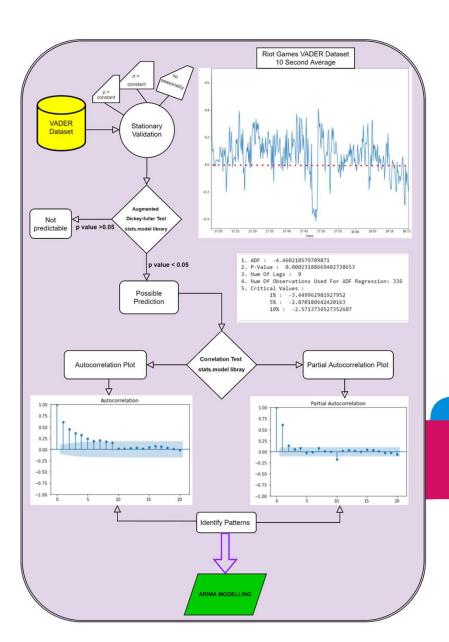negative sentiment: compound score <= -0.05

# SVM Model

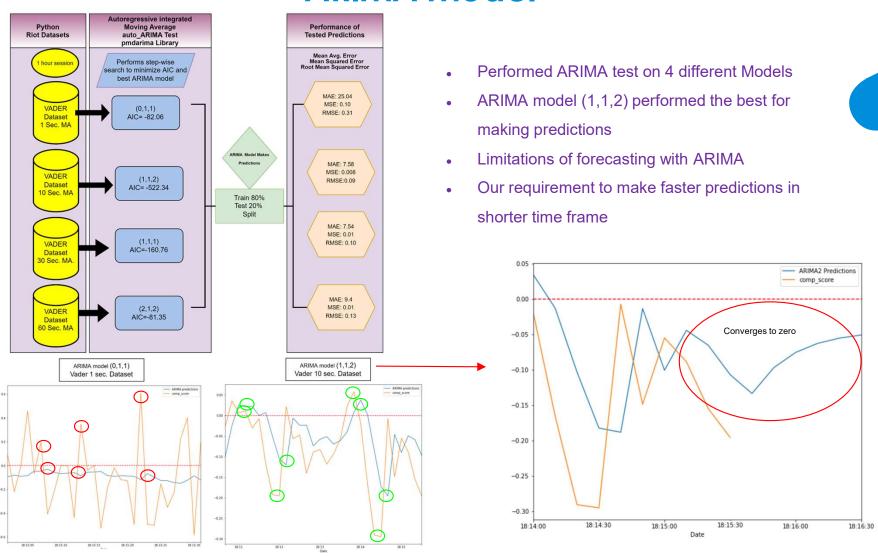| model_name | Mean Accuracy | Standard deviation |
|---|---|---|
| LinearSVC | 0.927473 | 0.026591 |
| LogisticRegression | 0.916484 | 0.016258 |
| MultinomialNB | 0.922344 | 0.012259 |
| RandomForestClassifier | 0.671062 | 0.021732 |

# Time Series: Validation

- Stationary Test
- Augmented Dicker Fuller Test
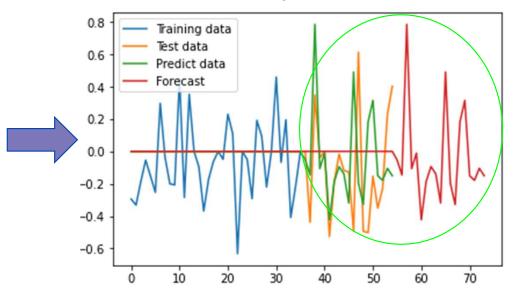- Correlation Test (ACF & PACF)
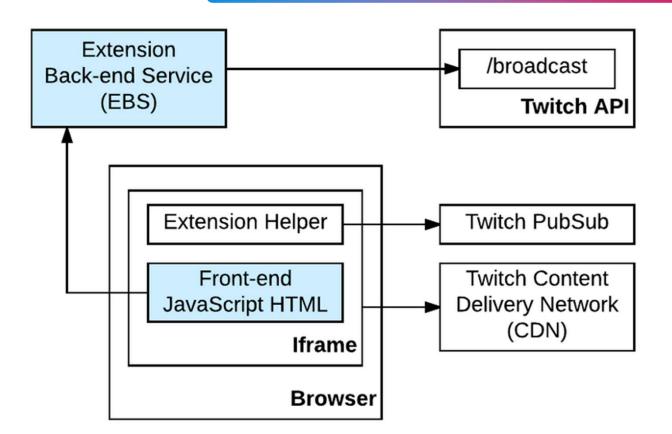
# ARIMA Model



- Performed ARIMA test on 4 different Models
- ARIMA model (1,1,2) performed the best for making predictions
- Limitations of forecasting with ARIMA
- Our requirement to make faster predictions in shorter time frame

# LSTM RNN

Take a minute sample

Create the model



VADER
Riot Games Dataset

80% Train
20% Test
Split

```
input_size = 1
hidden_size = 3
num_layers = 1

num_classes = 1

lstm = LSTM(num_classes, input_size, hidden_size, num_layers)
```

Training and test

# LSTM RNN

Forecast

| | |
|---|---|
| 2021-10-18 18:15:31 | -0.053558 |
| 2021-10-18 18:15:32 | -0.146200 |
| 2021-10-18 18:15:33 | 0.785068 |
| 2021-10-18 18:15:34 | -0.106826 |
| 2021-10-18 18:15:35 | -0.010469 |
| 2021-10-18 18:15:36 | -0.422834 |
| 2021-10-18 18:15:37 | -0.186522 |
| 2021-10-18 18:15:38 | -0.094237 |
| 2021-10-18 18:15:39 | -0.139344 |
| 2021-10-18 18:15:40 | -0.319031 |
| 2021-10-18 18:15:41 | 0.490356 |
| 2021-10-18 18:15:42 | -0.195136 |
| 2021-10-18 18:15:43 | -0.328921 |
| 2021-10-18 18:15:44 | 0.182459 |
| 2021-10-18 18:15:45 | 0.314759 |
| 2021-10-18 18:15:46 | -0.150587 |
| 2021-10-18 18:15:47 | -0.179027 |
| 2021-10-18 18:15:48 | -0.104804 |
| 2021-10-18 18:15:49 | -0.150647 |

VADER
Riot Games Dataset



19

# Twitch Extension

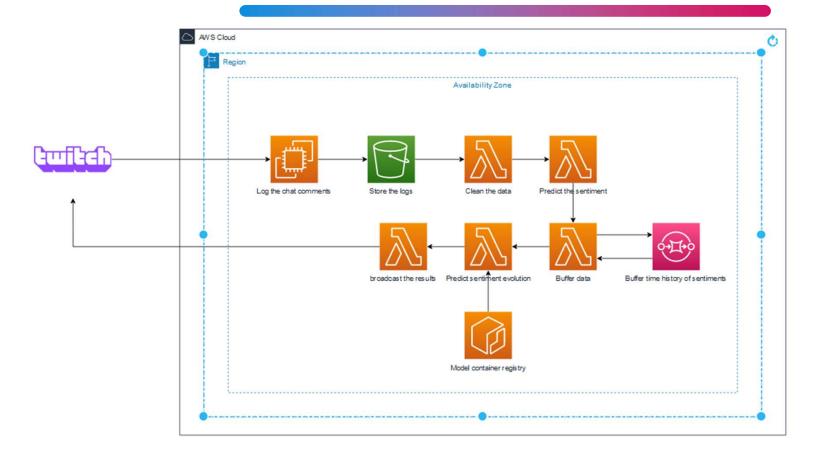# EBS Design Choices

## Hosting Solution

- Architecture is cost driven => Serverless Solution
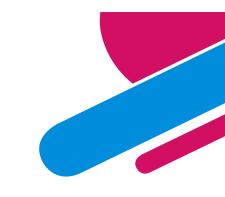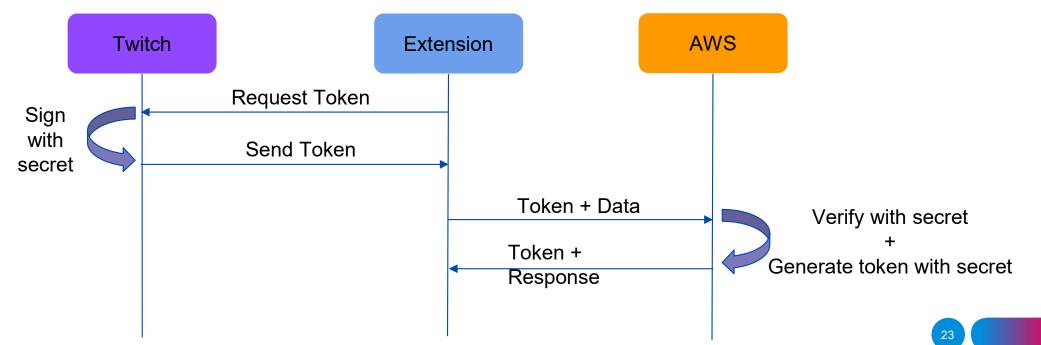- Chosen host : AWS

## Stream Data Collection

- Communication with Twitch via PubSub / Websocket
- EC2 endpoint : t3.micro more cost effective due to number of connections

# AWS Architecture



AWS Cloud

Region

Availability Zone

- Log the chat comments
- Store the logs
- Clean the data
- Predict the sentiment
- broadcast the results
- Predict sentiment evolution
- Buffer data
- Buffer time history of sentiments
- Model container registry

22

# Data Transfer with JWT

# MVP

# DEMO

```
ssh -i "test2.pem" ubuntu@ec2-13-38-99-32.eu-west-3.compute.amazonaws.com
```

```
ubuntu@ip-172-31-17-216:~$ cd twitch-import
```

Export the environment variables (credentials)

```
ubuntu@ip-172-31-17-216:~/twitch-import$ nano launch_listeners.py
```

```
# List of channels to connect to
channels_to_listen_to = ['BlackDesertGame']
```

```
# Scrape live chat data into raw log files. (Duration is seconds)
bot.listen(LOGDIR, channels_to_listen_to, duration = 1000)
```
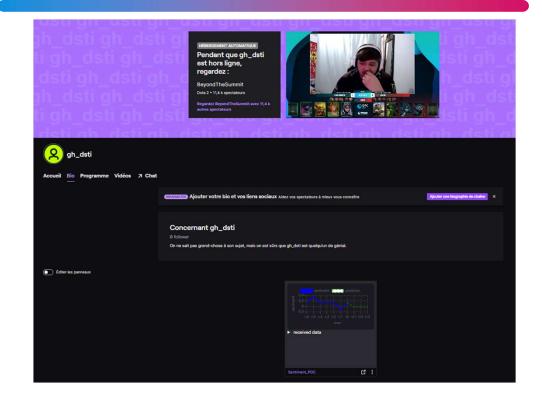
```
ubuntu@ip-172-31-17-216:~/twitch-import$ ./launch.sh
```
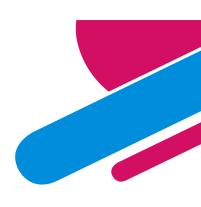
```
Open channel BeyondTheSummit
Uploading file BeyondTheSummit.log.2021-12-11_09-28-36 ...
Removing local file BeyondTheSummit.log.2021-12-11_09-28-36 ...
```

▶ received data

# DEMO



https://vimeo.com/703474180

# Future Improvements

- Continually update Lexicon for new Emotes

- Enhance the cloud infrastructure

- Find the best values for the LSTM models.

- Upgrading the architecture for scaling to 10s of simultaneous users

- Add additional features and metrics to the app

# Revenue Scheme

**Twitch represents an important example of a major service looking for new opportunities to monetize its enthusiastic streaming community.**

- Subscription model
- 4 tier model for pricing
- Based on user traffic and added features