# Reddit Auto-Moderation by Evaluating Community Opinion

Ayush Baid, Ankur Bhardwaj, Tarun Pasumarthi
CSE6240 Spring 2020
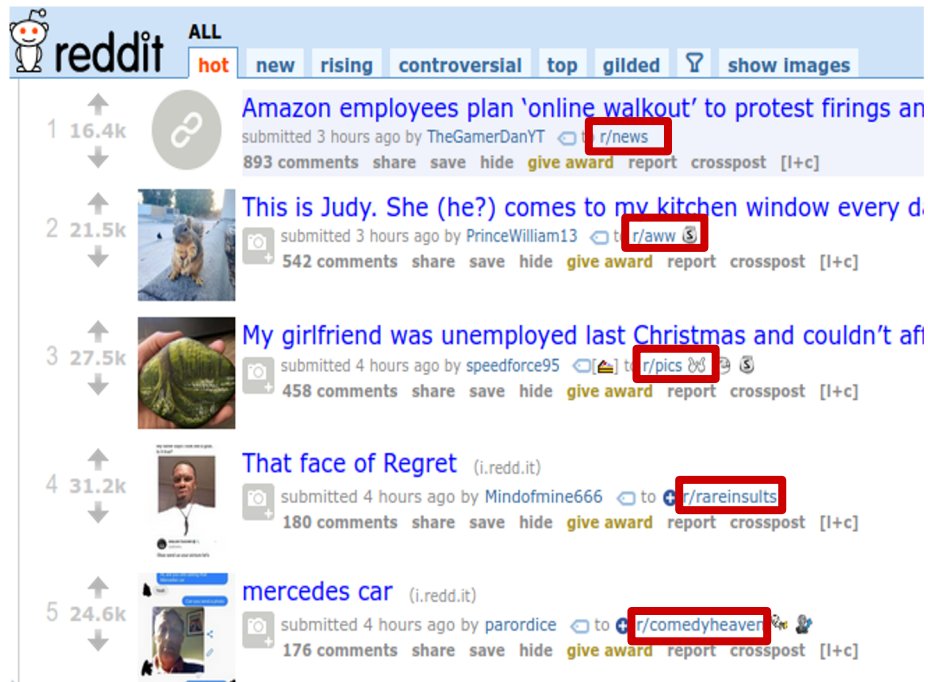
# What is Reddit?

## Reddit

- #5 most visited in the US
- #13 most visited in the world

## Subreddits

- Groups for specific topic discussion
- Moderators decide rules for content being posted on their subreddit

# What is the problem?

## Controversial content

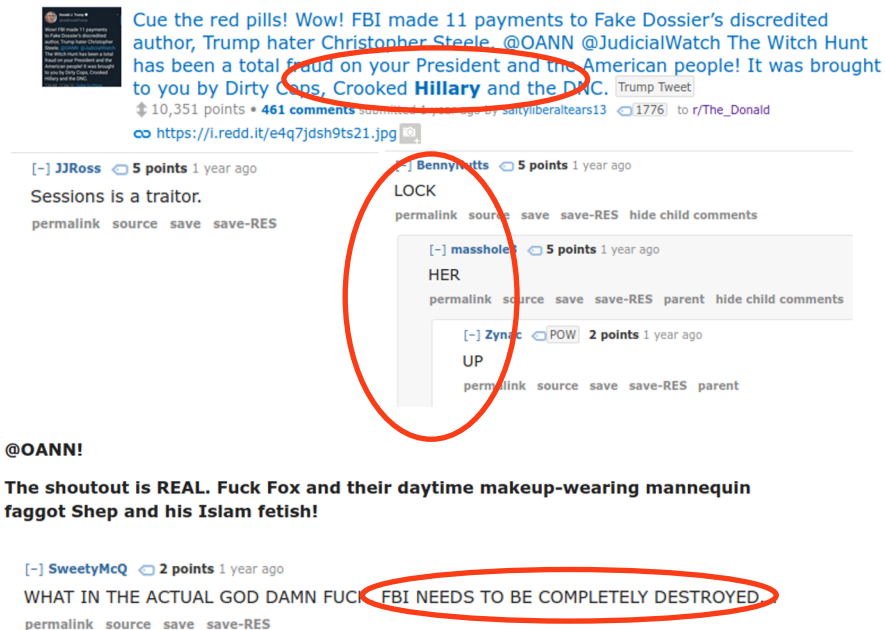https://en.wikipedia.org/wiki/Controversial_Reddit_communities

## 'The_donald' subreddit on Hillary?

# Motivation?

**Why would Reddit care?**

- **Revenue:**
  - https://qz.com/1246087/opinion-reddits-advertising-strategies-still-hide-hate-speech/
  - Reddit 'whitelists' subreddits which are safe and advertiser-friendly
- **User-base:**
  - Users don't want to see hate-speech and bad content

**Problem statement?**

- **A framework which can:**
  - Flag hate-speech and polarized content
  - Provide insights and reasons for moderators to ban such subreddits

# Existing work on Auto-Moderation?

Jhaver et al., '19: **Human-Machine collaboration**, the case of Reddit Automoderation

- automated systems filter out individual comments by searching hand-designed **regex patterns**

Jiang et al., '20: **Modbot**: Automated Comment Moderation

- A combination of **manually crafted rules** and basic NLP to predict 0/1 classification

**Limitations?**

- We want to **analyze millions of comments**

- Explore advanced NLP techniques – **"Summarization", "Contextual text completion"** etc.

- Analyze large amounts of text data: Language Models to the rescue!

# Text Completion Baselines?



**Fixed-window Language Model**

**RNN Language Model**

# Our Approach

Train a model on each subreddit
- Use pre-trained deep networks

Use model to generate opinions
- On manually curated topics

Provide short summaries for moderators
- To take actions using insights

# Our Approach

**Novelty?**

- Using cutting-edge NLP techniques for transfer learning (ULMFit, GPT2).
- Design a system which combines machine intelligence with human expertise

**Why it should perform better?**

- ULMFit/GPT2 has millions of parameters and pretrained on huge datasets
  - Good at context, semantic similarity, natural language understanding
- A human moderator can get insights from generated summaries and take appropriate actions
  - Instead of 0/1 classification's black box behaviour

# Demo

# Data

- Used PushShift's API

- 14 subreddits

- Similar number of comments per month from Sep-Dec 2019 to cover a variety of discussions

- 100k+ comments per subreddit for training, 60k+ comments for validation



Average sentiment for each subreddit

# Experiments and Evaluations?

**Experiment?**
- Train language model
- Manually complete 350 sentences

**Evaluation?**
- Is it actionable for a moderator?
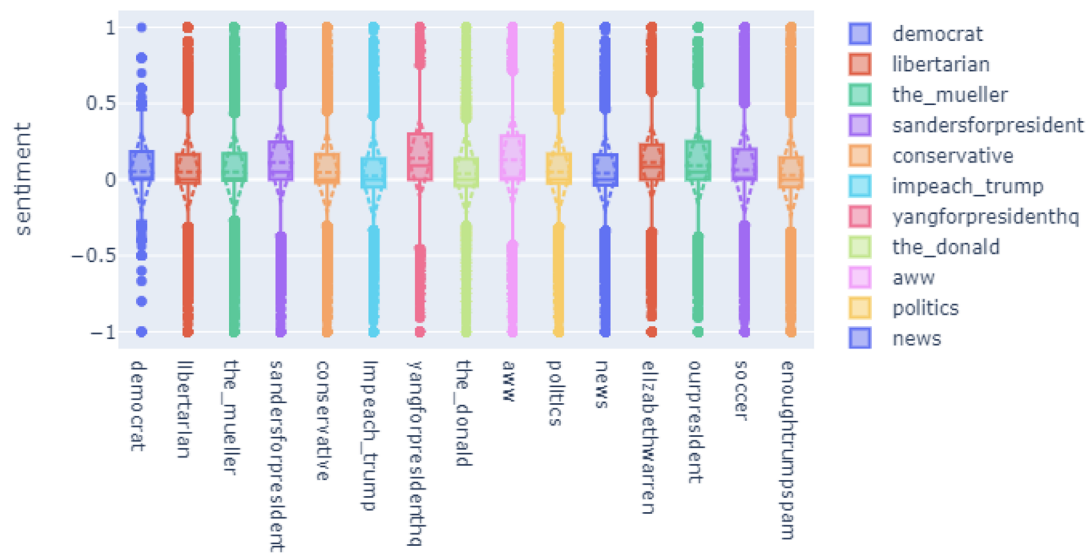- Is it relevant to the discussion on the subreddit?

**Example Subreddit: Impeach_trump (anti trump leaning)**
- Trump is a traitor (usability=    )
- Trump is a senator (usability=    ; incorrect sentence)
- Trump is a man (usability=    ; not providing any insight)
- Trump is Pelosi should be made do this (usability=    ; incoherent)

# Results - Usability Metric

| Subreddit | Fixed Window (baseline) | LSTM based (baseline) | ULMFit | GPT2 |
|---|---|---|---|---|
| r/politics | 1.34% | 13.08% | 8.17% | 55.35% |
| r/the_donald | 0.68% | 10.82% | 18.31% | 51.68% |
| r/OurPresident | 0.73% | 8.70% | 6.48% | 33.49% |

# Results – Samples

## r/politics

| Input | RNN (Baseline) | ULMFit | GPT2 |
|-------|----------------|--------|------|
| Trump is a | traitor | Russian puppet | Evil puppet |
| Biden should | Not be the nominee | Sanders | be in the Senate if he had a shot at the nomination |
| UBI | Is a good thing | Not a job | is not an automatic negative incentive |

## r/the_donald

| Input | RNN | ULMFit | GPT2 |
|-------|-----|--------|------|
| Trump is a | cuck | Handling ISIS | Smarter than the Dems |
| Biden should | Be the same | Be president | be on the cutting edge of all of this |
| UBI | *null* | is socialism | ! I need an ANTIFA card! |

# Conclusion and Future Work

**Conclusion?**

- Transfer learning showing a lot of potential:
    - Deep models can be applied to small and unlabeled datasets (100k comments)
- Enables human moderators to go through a small set of results and take action which are explainable

**Future work?**

- Integrate the NLP phrase completions with sentence-emotion mapping
    - Additional input to the moderator
- Multi-word embeddings instead of a word-by-word model for stronger contexts

# Thank you!

We'll be live soon…