



# Sphinx



## Word Similarities Game

By Teng Mao

Data Science Fellow from General Assembly

# Agenda

- Problem Statement
- Methodology
- Corpora
- Exploratory Data Analysis
- Model
- Bonus: Challenges
- Conclusion & Next Steps



# The Riddle of the Sphinx

"What is the  
creature that walks  
on four legs in the  
morning, two legs at  
noon and three in  
the evening?"

**Man**

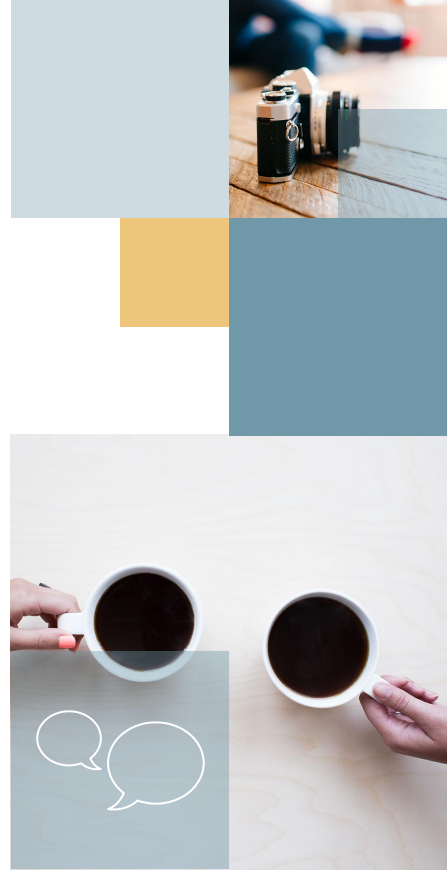
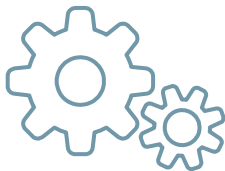
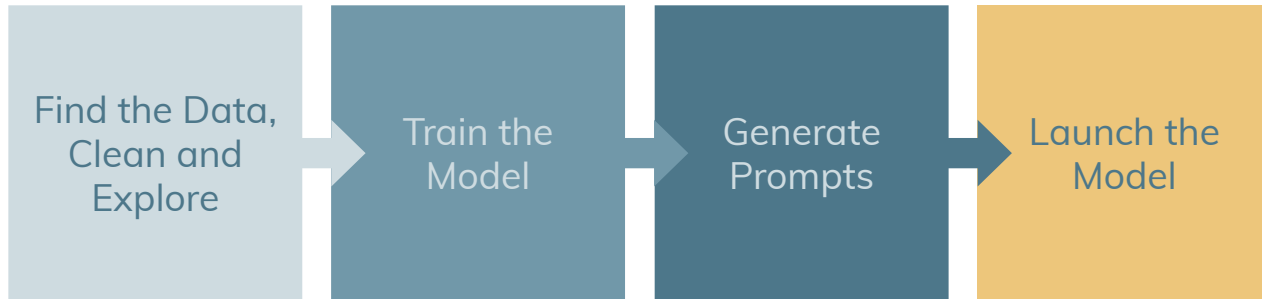




# Problem Statement:

How do you build a game using natural language processing with a neural net?

# Methodology





## A collage of images related to data science and Wikipedia. It includes a hand writing 'ideas' on a notepad, a Wikipedia logo, a globe icon, a bowl of green apples on a scale, and various geometric shapes in shades of green and blue.



300-dimensional  
vectors for 1  
million words  
trained on  
Wikipedia 2017,  
UMBC webbase  
corpus and  
statmt.org news  
dataset (16  
billion words).



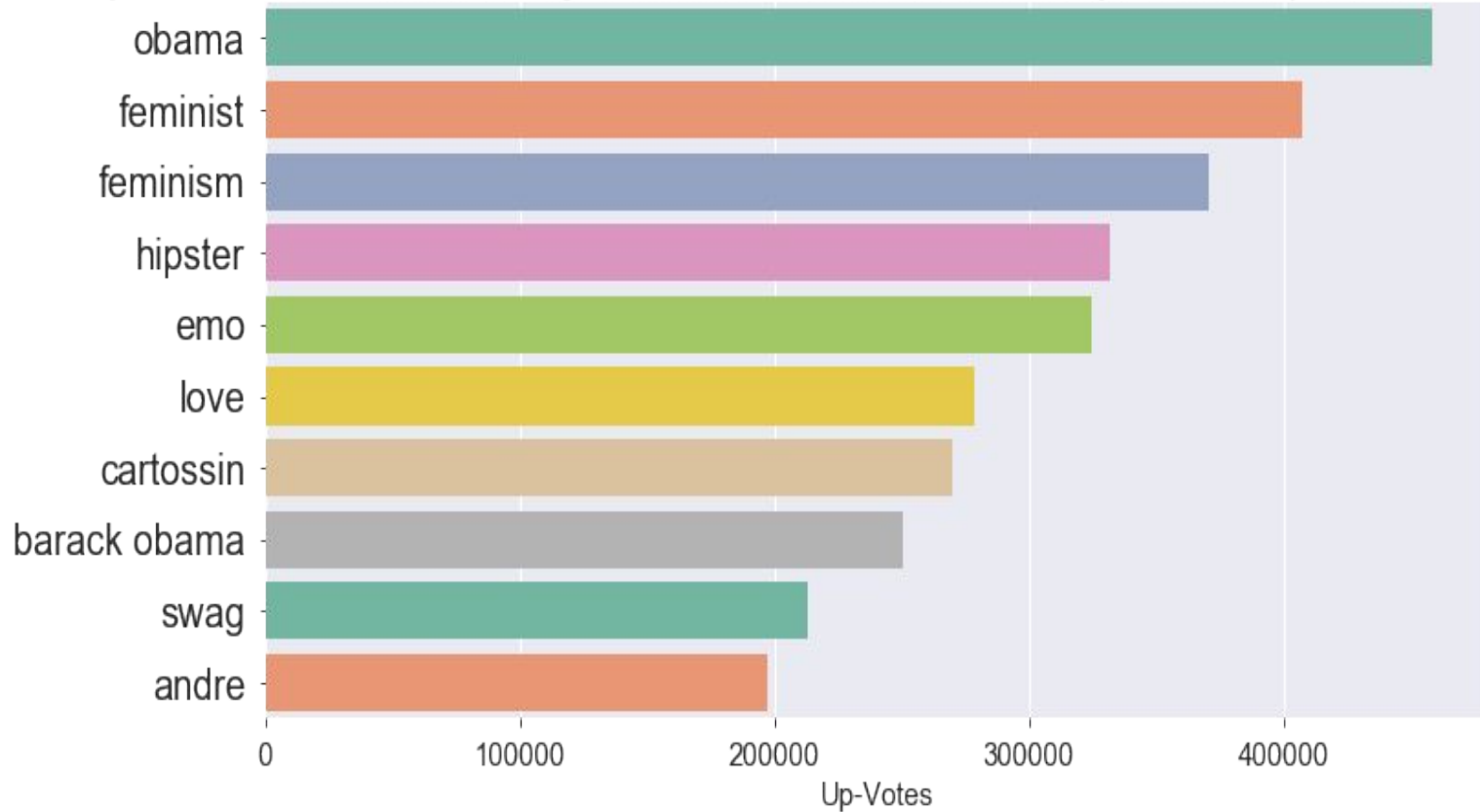
2.6 million words (1.6 million unique), authors, up-votes and down-votes, as well as definitions were scraped and put on kaggle for a competition.



## A kitchen scale with a metal bowl containing several green apples. The scale is white with a circular dial and a black needle. The bowl is silver and filled with about six green apples. The background is a white tiled wall.

34,886 movies were scraped from Wikipedia and put on kaggle for a competition. Data includes release year, title, origin/ethnicity, director, cast, genre, and plot.

## Top 10 Urban Dictionary Words with the Most Up-Votes (Excluding Swears)



O. N.

D

O. N.

O. N.

O. N.

O. N.

O. N.

O. N.

O. N.

O. N.

O. N.

O. N.

O. N.

O. N.

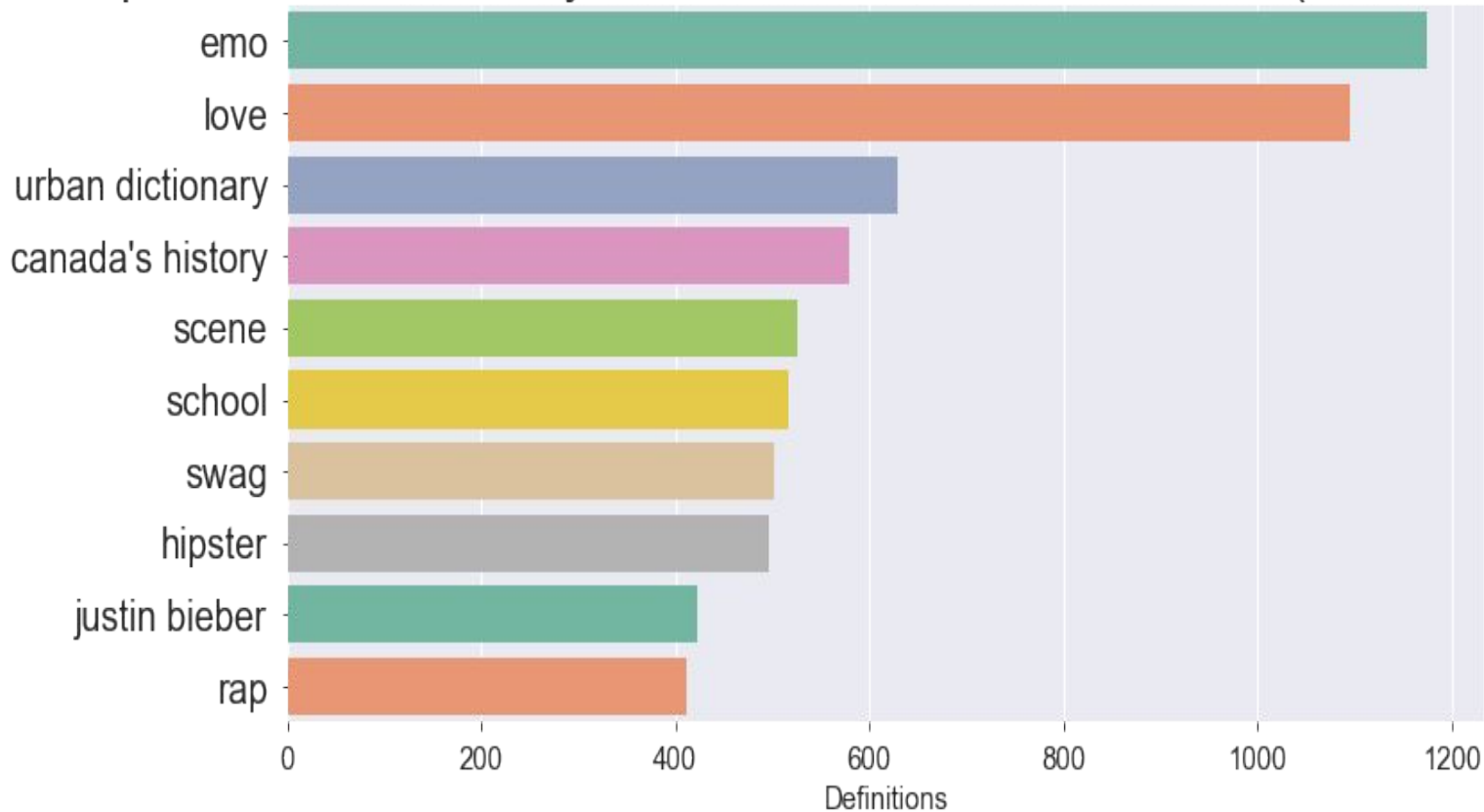
O. N.

O. N.

O. N.

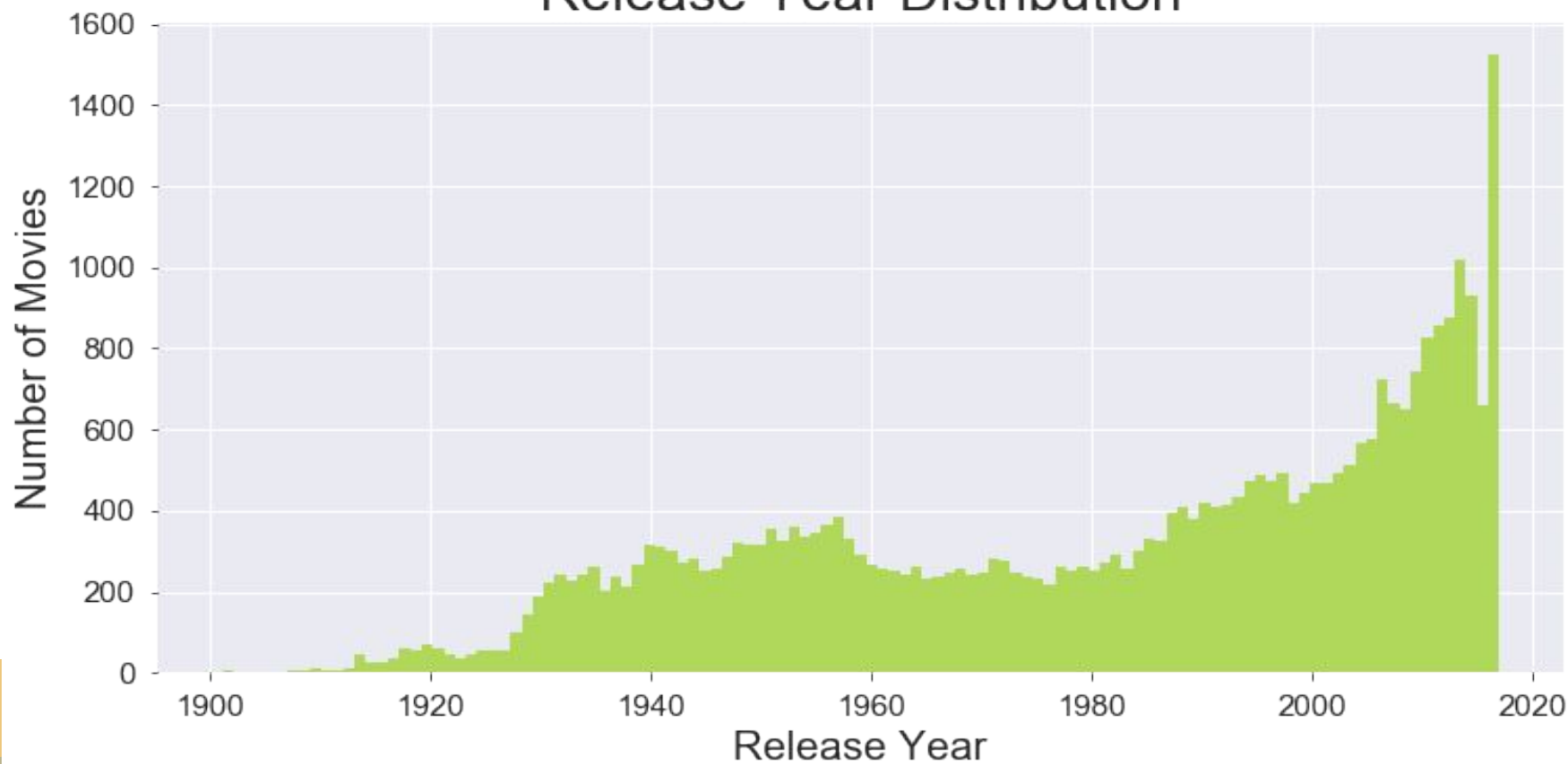
O. N.

## Top 10 Urban Dictionary Words with the Most Definitions (Excluding Swears)

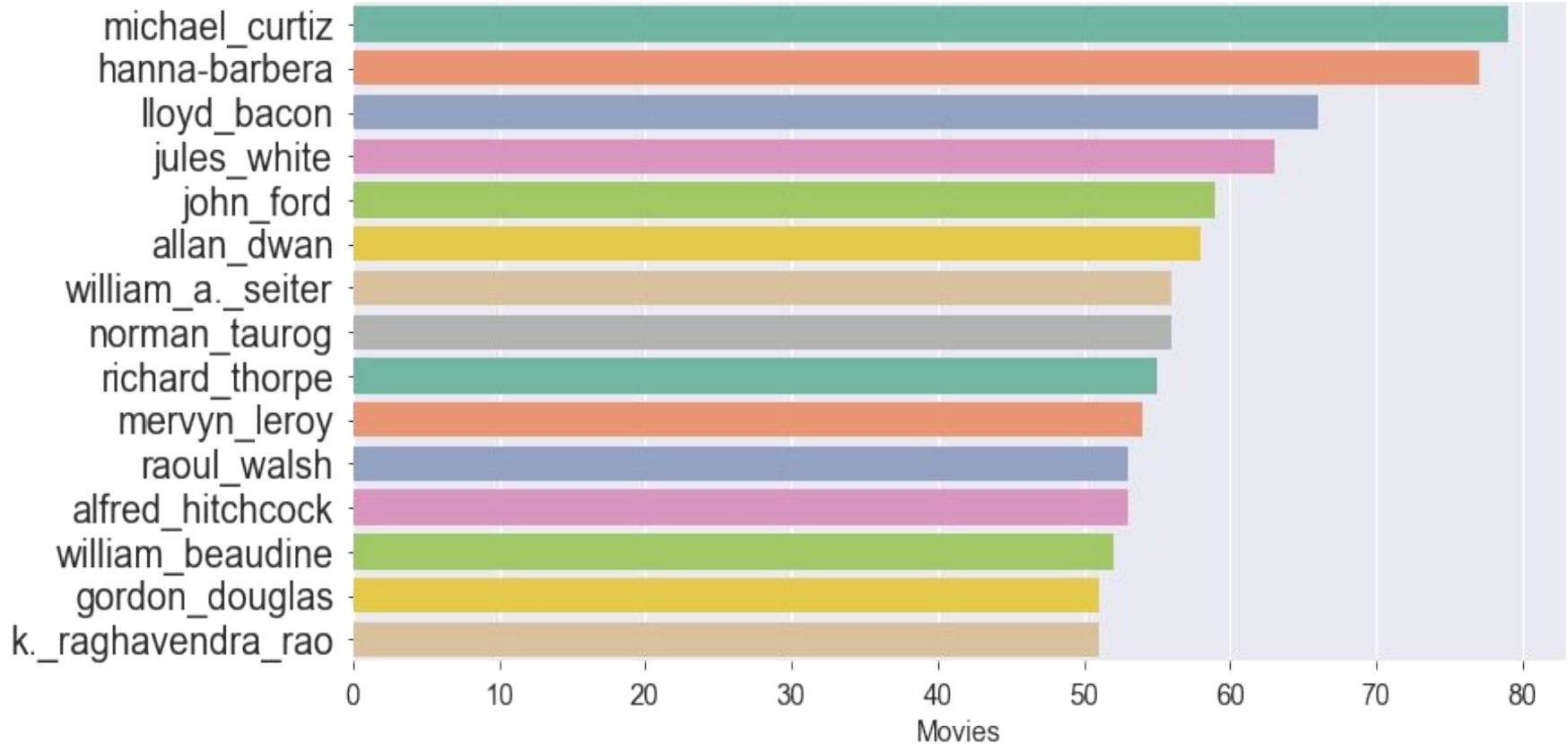




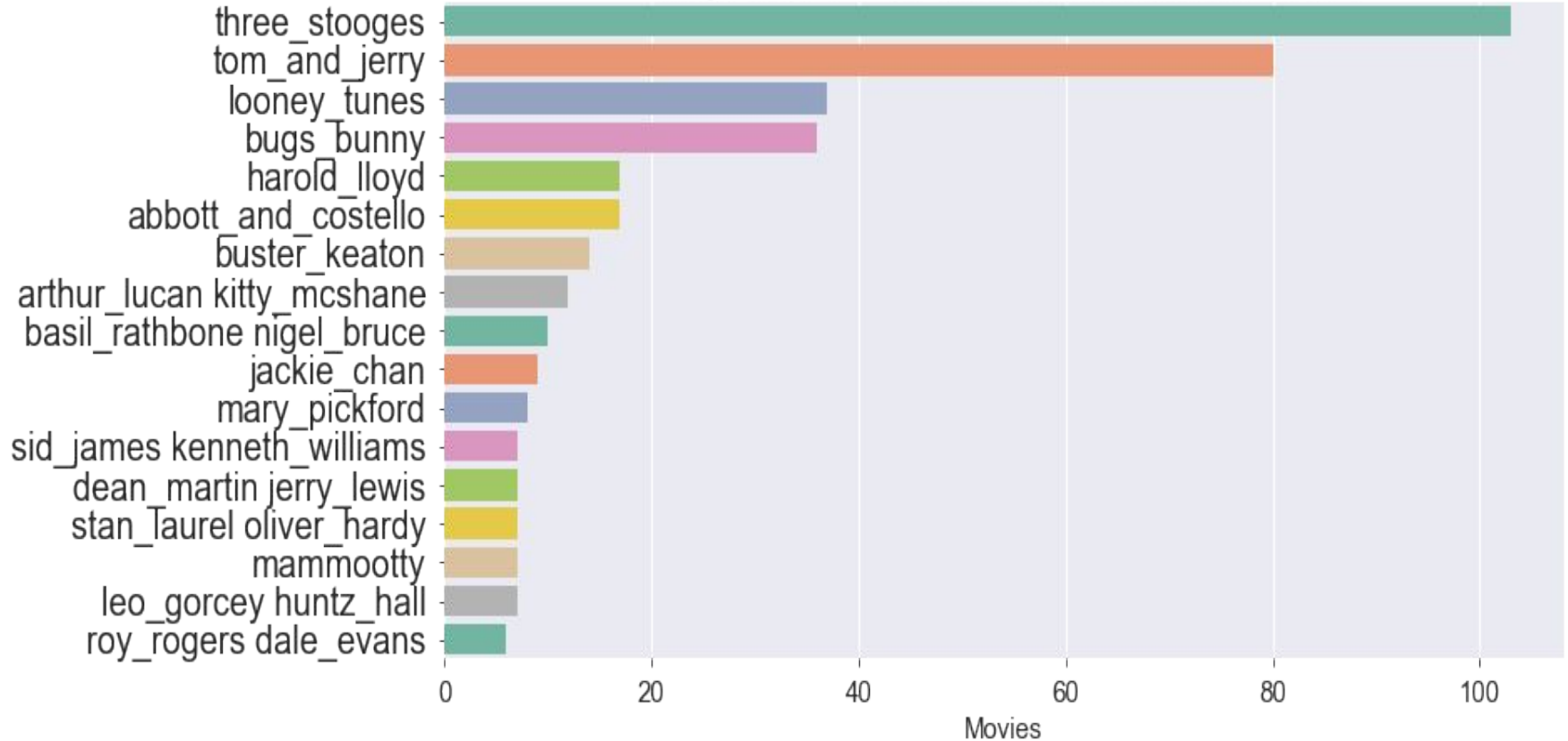
## Release Year Distribution



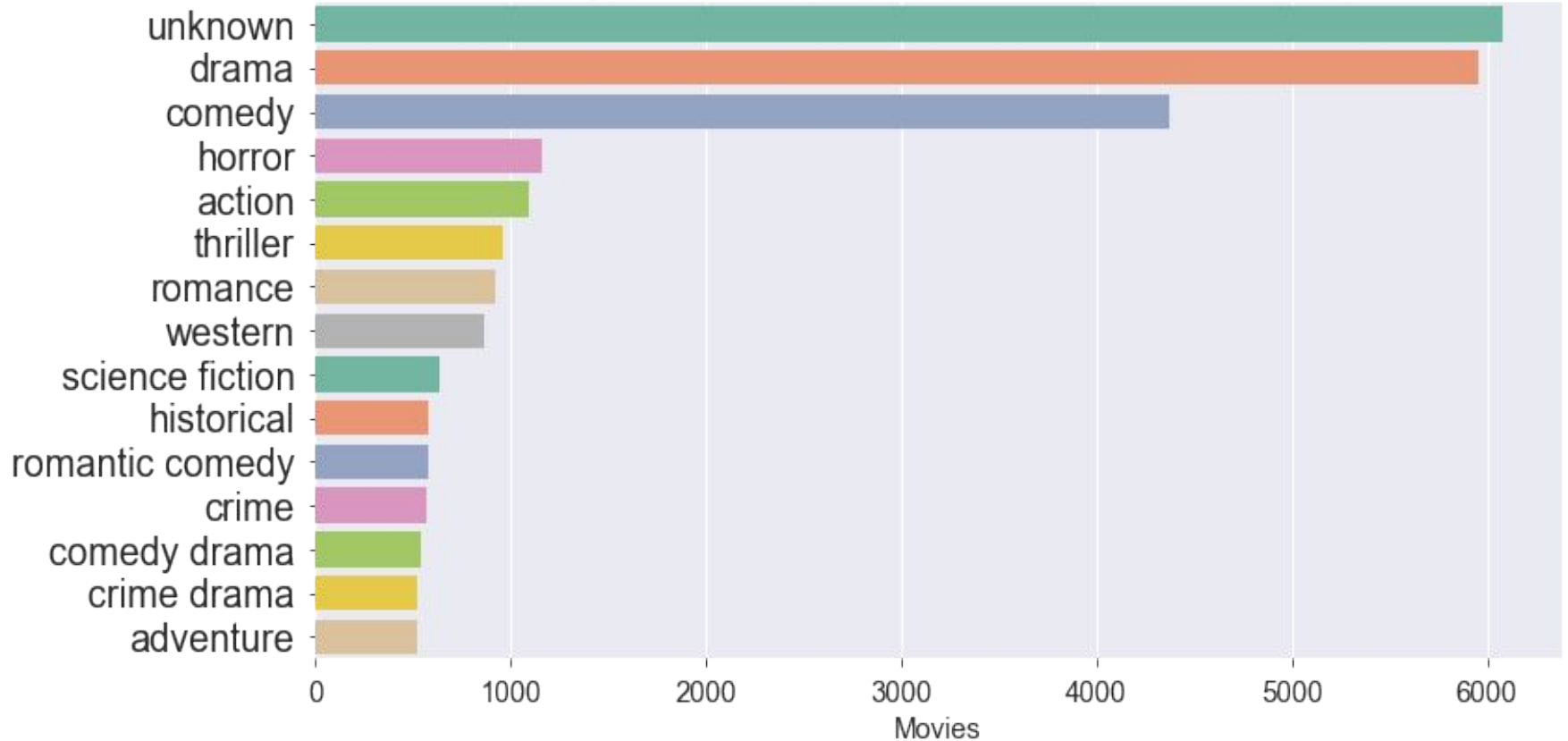
## Top 15 Directors with the Most Movies



## Top 15 Casts with the Most Movies



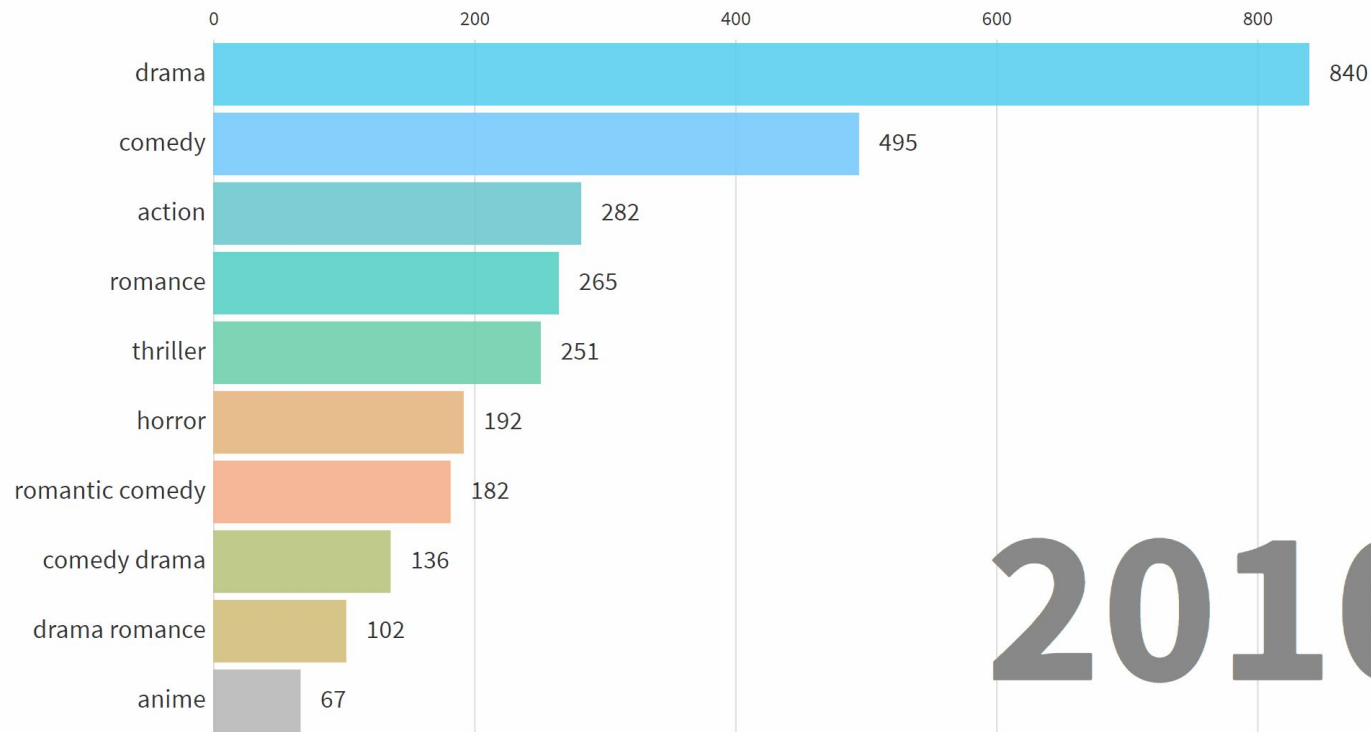
## Top 15 Genres with the Most Movies



# Top 10 Movie Genres by Decade (1900-2017)

from Wikipedia Movies Plots

🔄 Replay



2010

# Word2Vec Model

## horror

'horrors', 'vacui',  
'Horror', 'terror',  
'horror-', 'dread',  
'shock', 'disgust',  
'revulsion', 'carnage'

## Patriots

'Pats',  
'England\_Patriots',  
'Chargers', 'Colts',  
'Ravens', 'Jaguars',  
'Raiders', 'Seahawks',  
'Bengals', 'Belichick'

## turnt

'flipped\_up', 'f\*cked',  
'krunked', 'geeked',  
'fubared', 't\_u', 'ttgt',  
'troud', 'schywasted',  
'crunked'

## horror

'thriller', 'fantasy',  
'zombi\_kampung\_pisang',  
'jangan\_pandang\_belakang',  
'the\_creeping\_flesh', 'Suspense',  
'diamond\_is\_unbreakable\_chapte  
r\_i', 'susan\_denberg', 'the\_skull',  
'freddie\_francis'

## Netflix and chill

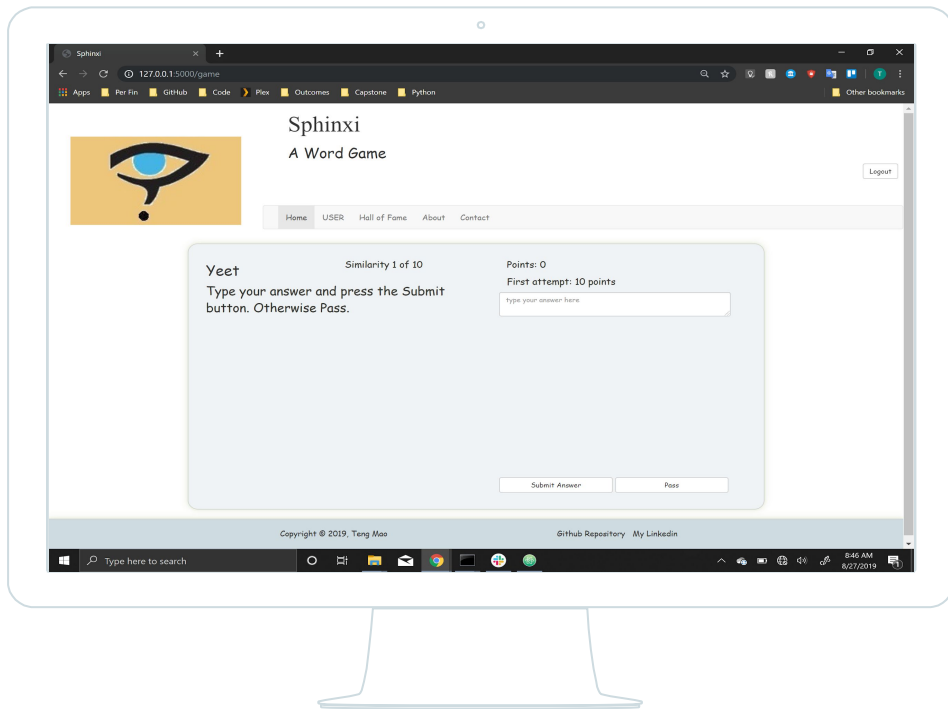
'netflix\_and\_cuddle', 'netflix\_an\_chill',  
'hulu\_and\_hangout',  
'imax\_and\_climax',  
'youtube\_and\_blow', 'hulu\_and\_chill',  
'netflix\_and\_chill\_alone',  
'hulu\_and\_hang', 'p\*\*nhub\_and\_chill',  
'netflix\_and\_jill'





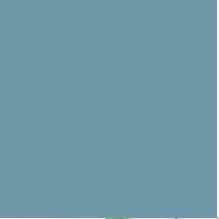
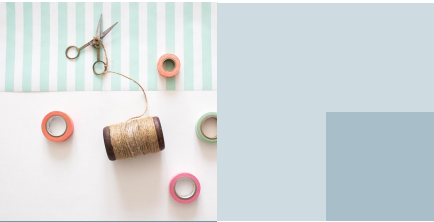
# Sphinx on the Desktop

Using Flask,  
Sphinx is able to  
be launched in a  
browser.



# Challenges

- Deciding on a capstone
- Size of the data
- Offensive content
- Scope of the project
- Deciphering someone else's code
- Momentum



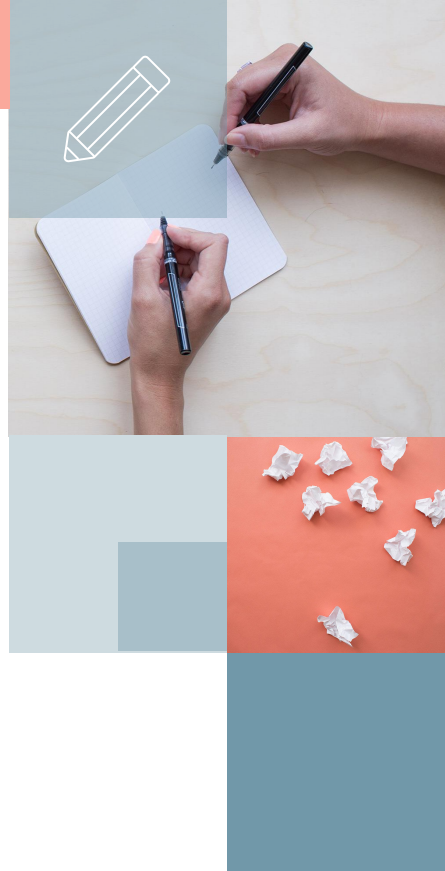


## Conclusions

- Natural Language Processing - powerful, but flawed

## Next Steps

- Deploy model to AWS
- Finish the game in Flask
- Generate more prompts
- Create analogies option
- Gather more data



Thanks!

# Any questions?

You can find me at @ [github.com/TengCXXI/](https://github.com/TengCXXI/) and  
[linkedin.com/in/tengmao](https://www.linkedin.com/in/tengmao)

