# Using NLP to Sort between r/Games and r/IndieGaming Subreddits
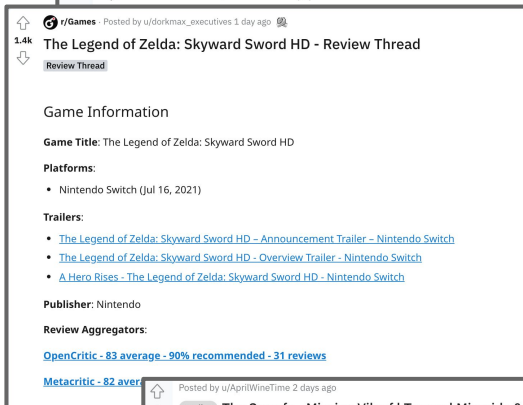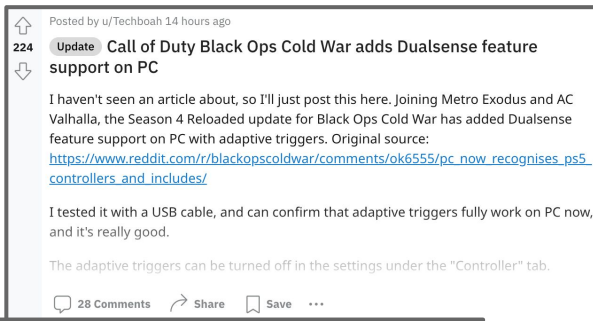
Julia Tsaltas
July 16th, 2021

# Outline

- Gaming discussions

- Subreddit data

- Model setup

- CountVectorizer insights

# Gaming Discussion

r/Games

AAA = Games built by a large studio with a big budget

**AAA Discussion**

**Reviews**

**Trailers**

3
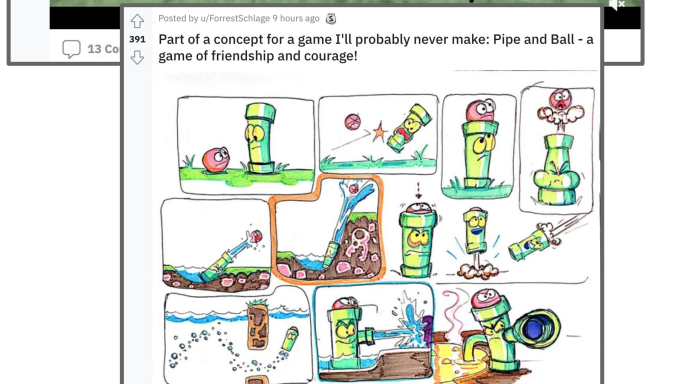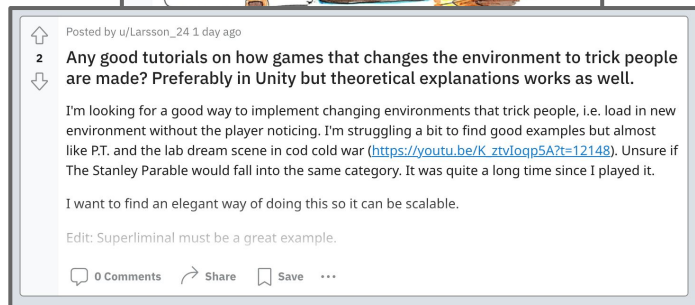
# Gaming Discussion

r/IndieGaming

Indie = Game developed by a small group not supported by a large developer
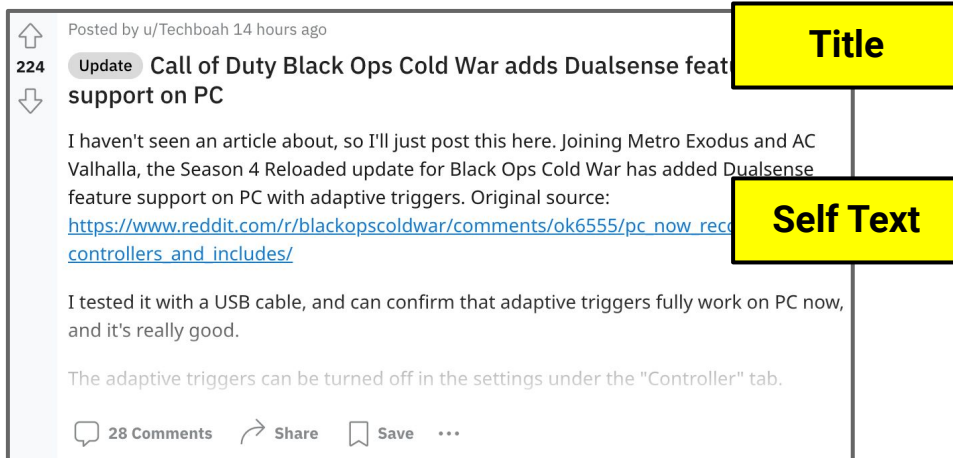
**Progress**

**Concept Art**

**Advice**

4

# Subreddit Data

- 1126 posts / subreddit

- Title + selftext

- 20,000+ unique terms

- Pipeline =

  CountVectorizer + Model



Title

Self Text

# Model Results

| Best Accuracy | Best Variance/Bias | Worst Performer |
|---|---|---|
| **Random Forest** | **Logistic Regression** | **KNN** |
| **99.9%** - Train<br>**86.0%** - Test | **81.2%** - Train<br>**80.8%** - Test | **71.5%** - Train<br>**64.7%** - Test |
| max_depth = None<br>min_sample_split = 4<br>N_estimators = 150 | C = 0.001<br>penalty = Ridge | k = 15 |

# CountVectorizer Insights

- ↓max_features ↑variance ↑score

- ↑max_df ↓variance ↓score

- ngram range ➡ best 2 or 3

# Summary

- Models do capture differences between gaming forums

- All models performed better than baseline

- Overfitting present on all models