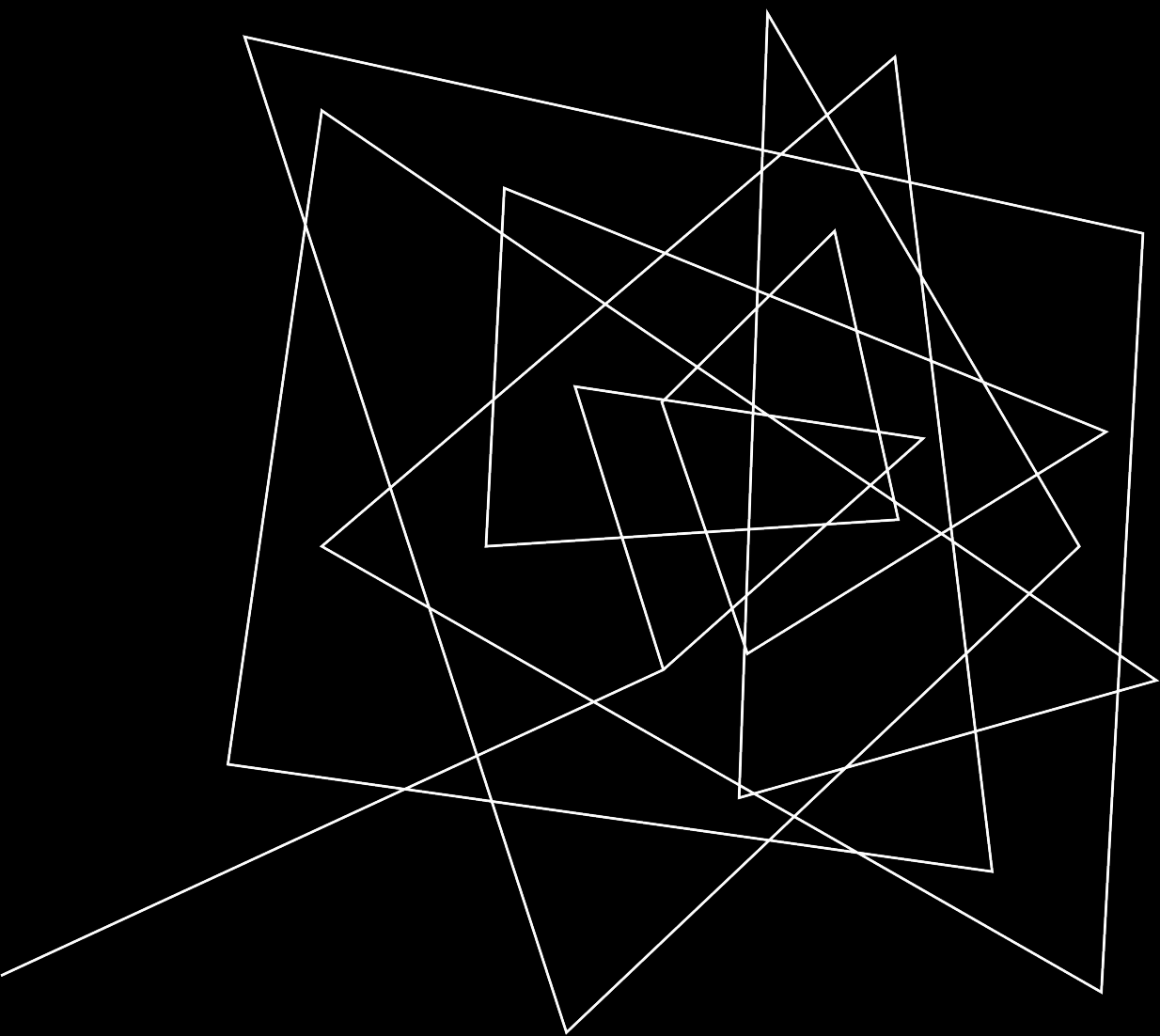# INTERNET POLARIZATION

Using Reddit Data to Examine Trends in Thread Comment Behavior
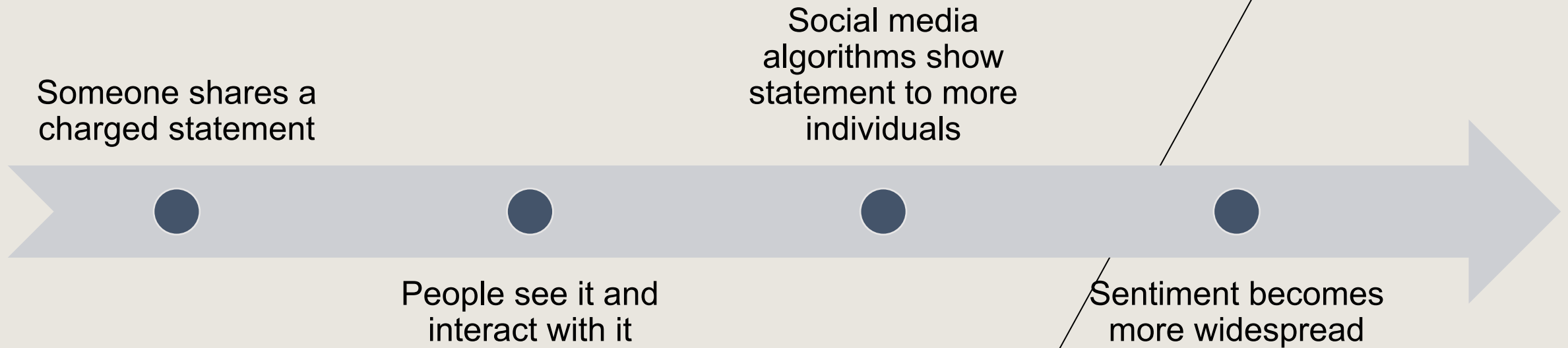
Neha Agarwal, David Amiel, Joel Cabrera

HATRED, ANGER, AND VIOLENCE CAN DESTROY US: THE POLITICS OF **POLARIZATION** IS DANGEROUS.
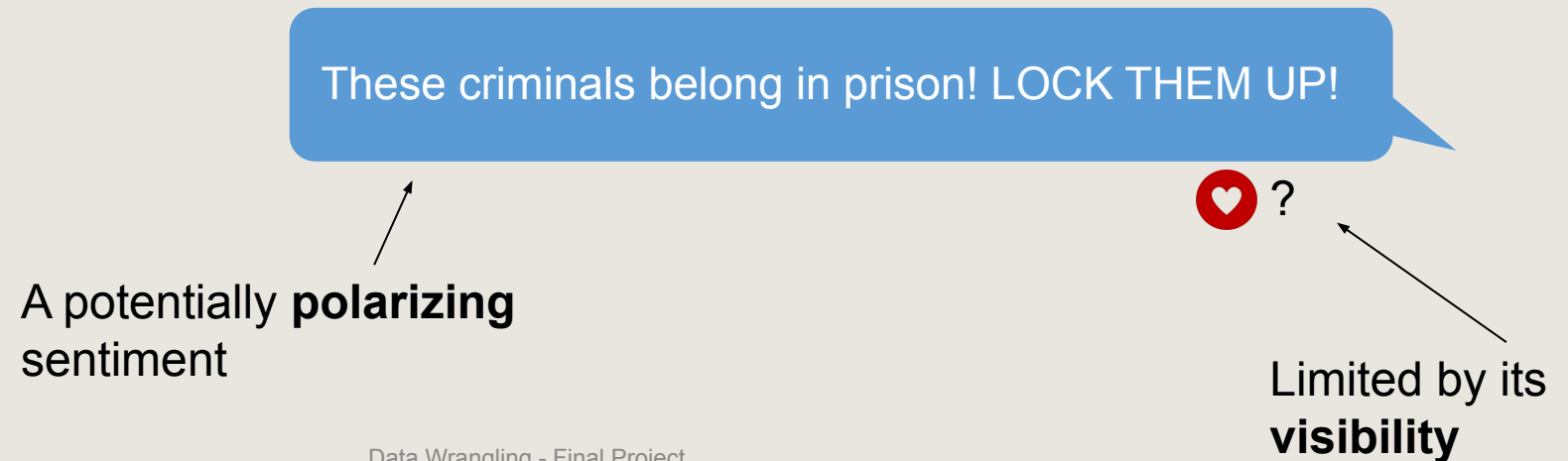
RAHUL GANDHI

# AN EXAMPLE…

How many likes do you think a comment like this would receive?

These criminals belong in prison! LOCK THEM UP!

♥ ?

A potentially **polarizing** sentiment

Limited by its **visibility**

# AN EXAMPLE…

How many likes do you think a comment like this would receive?

Politicians outlaw ice cream in all restaurants nationwide, citing concerns about rising costs. "This will be a major blow to Big Dairy," the president shared in a press conference last week.

❤️ -1,382,438

These criminals belong in prison! LOCK THEM UP!

❤️ LOTS

# AN EXAMPLE…

How many likes do you think a comment like this would receive?

Last, week thousands of migrants from Ukraine entered the United States seeking refuge from ongoing turmoil in their homes. US services welcomed them to safety. #StandWithUkraine

♥ 3,583,294

These criminals belong in prison! LOCK THEM UP!

♥ -LOTS

# HOW CAN WE UNDERSTAND WHAT'S GOING ON HERE?

**Sentiment** of 'parent' post

Politicians outlaw ice cream in all restaurants nationwide, citing concerns about rising costs. "This will be a major blow to Big Dairy," the president shared in a press conference last week.

♥ - 1,382,438

**Sentiment** of 'child' post

These criminals belong in prison! LOCK THEM UP!

♥ LOTS

**Popularity** of 'parent' post

**Popularity** of 'child' post

Where better to look than **Reddit**?

# OUR RESEARCH QUESTIONS

- What relationships exist among the content (defined by sentiment and topic) and popularity of a post and the content and popularity of its replies?
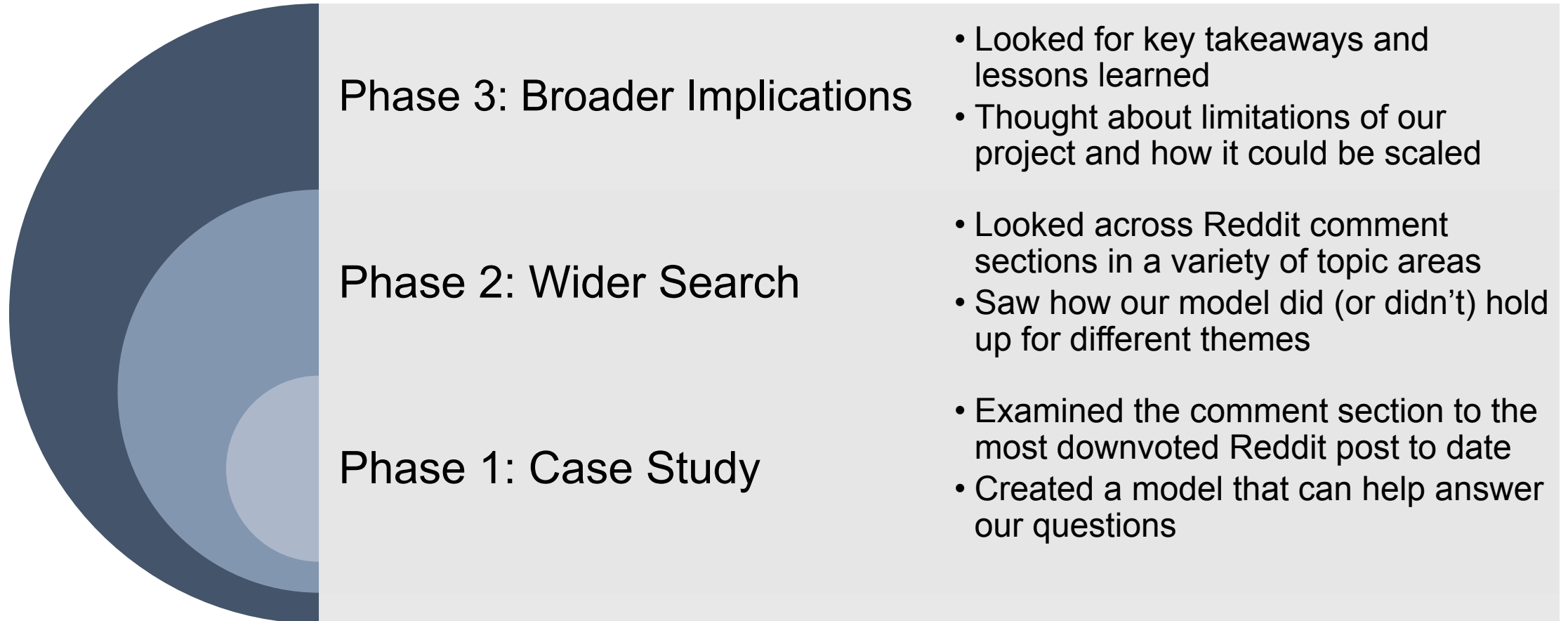
  **We think** posts with negative sentiment in response to unpopular posts will receive more likes, and vice versa.
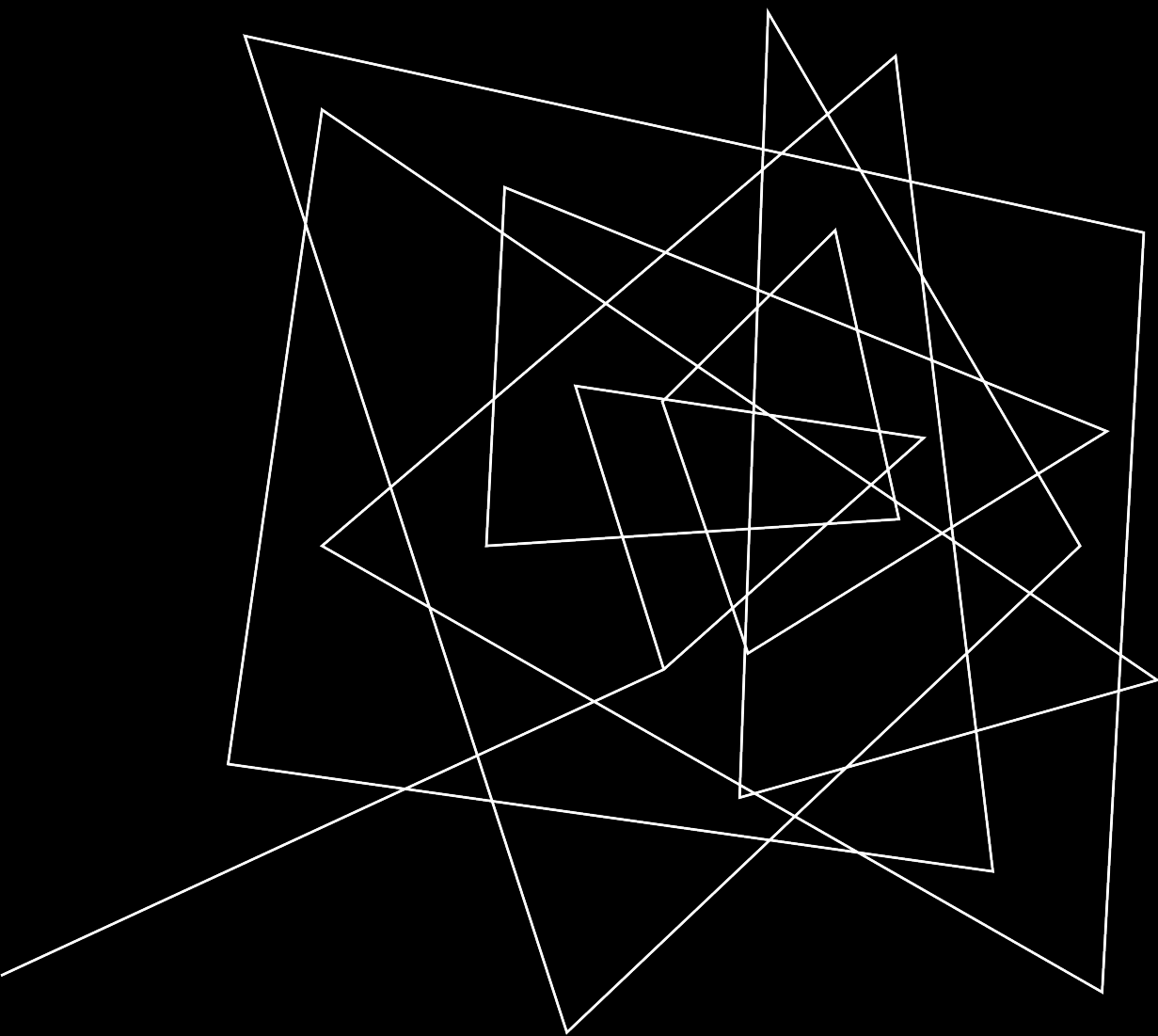
- Are these relationships universal, or do they look different in different areas (platform, topic)?

  **We think** they are not universal. Topics like politics, climate change, or abortion might have more polarized comment sections.

- Do these relationships tell us anything about online polarization of opinion? What next steps might these findings suggest?

# WHAT WE DID

**Phase 3: Broader Implications**

- Looked for key takeaways and lessons learned
- Thought about limitations of our project and how it could be scaled

**Phase 2: Wider Search**

- Looked across Reddit comment sections in a variety of topic areas
- Saw how our model did (or didn't) hold up for different themes

**Phase 1: Case Study**

- Examined the comment section to the most downvoted Reddit post to date
- Created a model that can help answer our questions

**PHASE 1:** CASE STUDY

*EA's Blunder*

# /R/StarWarsBattlefront

About 4 years ago…

**Seriously? I paid 80$ to have Vader locked?**
This is a joke. I'll be contacting EA support for a refund... I can't even playing f***ing Darth Vader?!?!? Disgusting. This age of "micro-transactions" has gone WAY too far. Leave it to EA though to stretch the boundaries.

♥ **202,000**

The intent is to provide players with a sense of pride and accomplishment for unlocking different heroes.
As for cost, we selected initial values based upon data from the Open Beta and other adjustments made to milestone rewards before launch. Among other things, we're looking at average per-player credit earn rates on a daily basis, and we'll be making constant adjustments to ensure that players have challenges that are compelling, rewarding, and of course attainable via gameplay.
We appreciate the candid feedback, and the passion the community has put forth around the current topics here on Reddit, our forums and across numerous social media outlets.

♥ **-667,700**

# SCRAPING REDDIT DATA

```
> glimpse(commentData)
Rows: 474
Columns: 10
$ url        <chr> "https://www.reddit.com/r/StarWarsBattlefront/comments/7cff0b/seriously_i_paid_80_to_have_vader_lo…
$ author     <chr> "EACommunityTeam", "bookem_danno", "kibber", "artycharred", "[deleted]", "[deleted]", "[deleted]",…
$ date       <chr> "2017-11-12", "2017-11-12", "2017-11-12", "2017-11-12", "2017-11-12", "2017-11-13", "2017-11-13", …
$ timestamp  <dbl> 1510513883, 1510524156, 1510523516, 1510514150, 1510524964, 1510534011, 1510534762, 1510535423, 15…
$ score      <dbl> -667702, 23119, 1991, 14405, 4561, 1066, 1367, 456, 234, 228, 165, 192, 256, 263, 89, 72, 41, 37, …
$ upvotes    <dbl> -667702, 23119, 1991, 14405, 4561, 1066, 1367, 456, 234, 228, 165, 192, 256, 263, 89, 72, 41, 37, …
$ downvotes  <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,…
$ golds      <dbl> 124, 2, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 2, 1, 0, 1, 1, 1, 0, 1, 1, 1, …
$ comment    <chr> "The intent is to provide players with a sense of pride and accomplishment for unlocking different…
$ comment_id <chr> "1", "1_1", "1_2", "1_3", "1_3_1", "1_3_1_1", "1_3_1_1_1", "1_3_1_1_1_1", "1_3_1_1_1_1_1", "1_3_1_…
```

**commentData** <- tibble(get_thread_content("https://www.reddit.com/r/StarWarsBattlefront/
comments/7cff0b/seriously_i_paid_80_to_have_vader_locked/dppum98/")$comments)
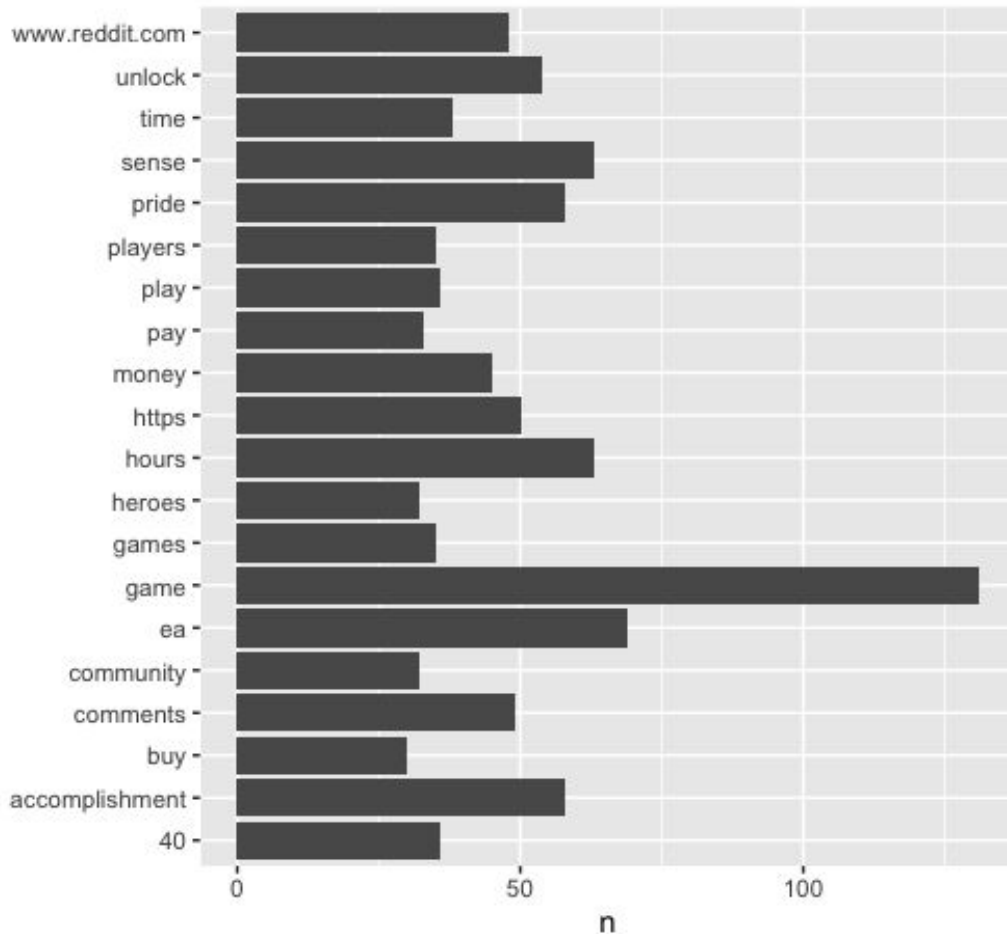
...

**overallSentiment** = positive - negative,
**parent_ID** = substr(comment_id, 1, regexpr("\\_[^\\_]*$", comment_id)-1))
**parentSentiment** = workingData[match(parent_ID, workingData$comment_id), 14]
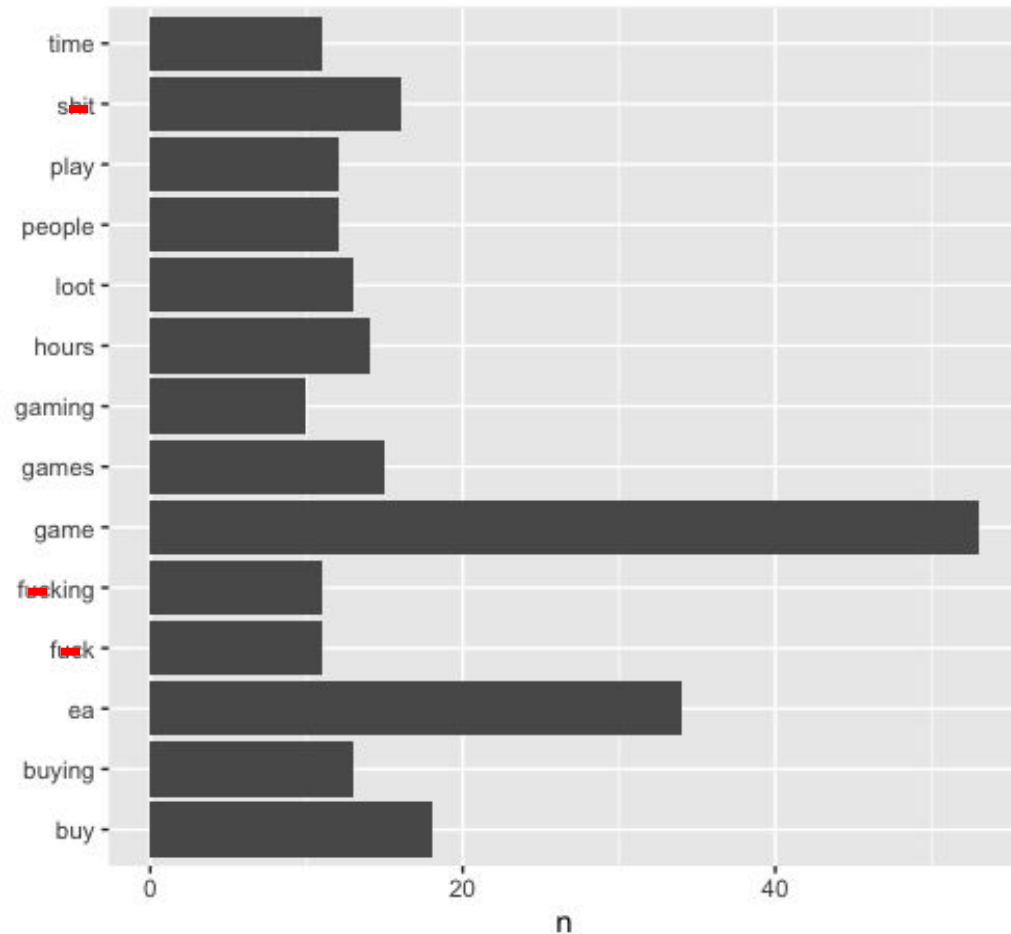**parentPopularity** = workingData[match(parent_ID, workingData$comment_id), 13]

# EXPLORING THE COMMENT DATA

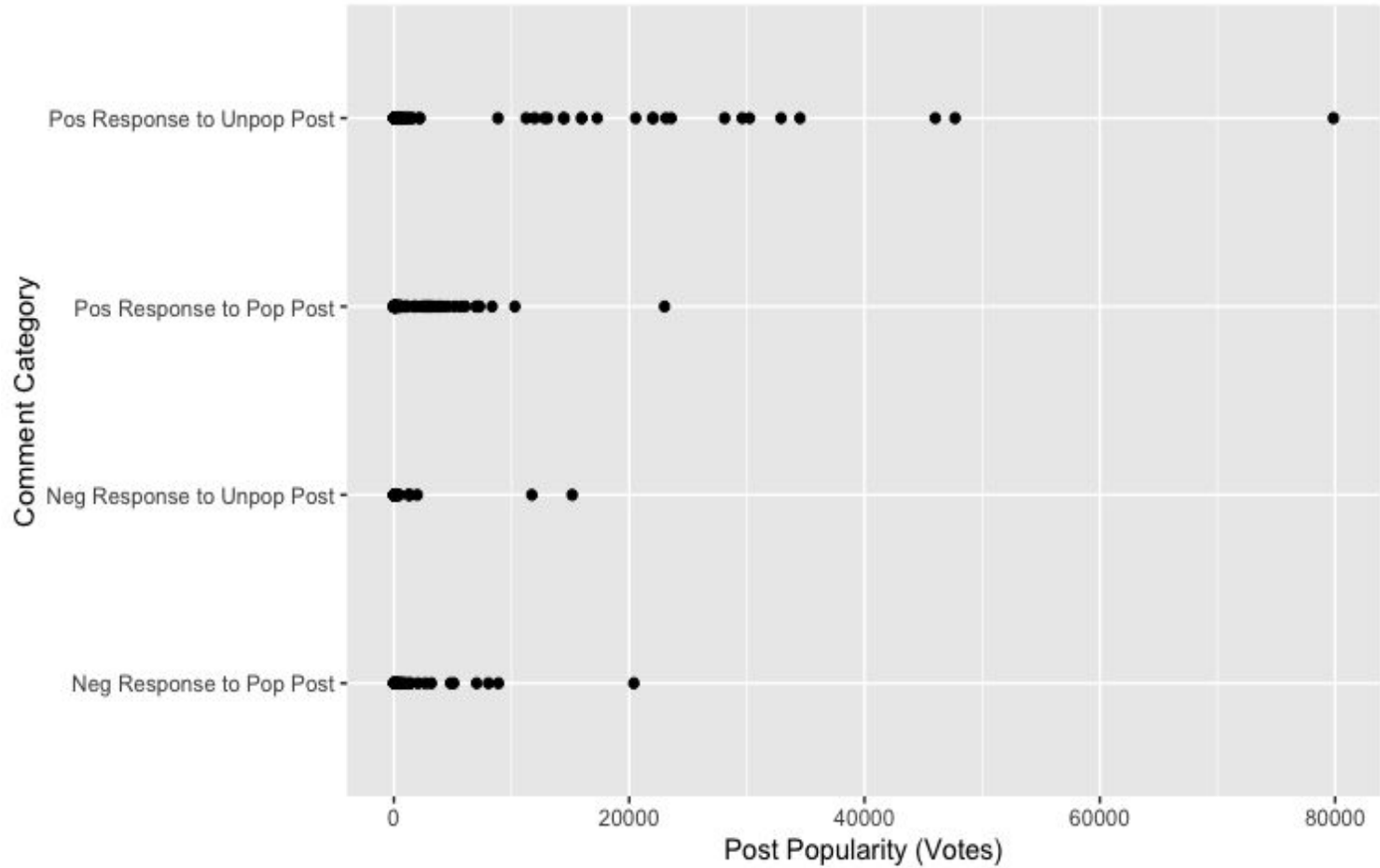### Most Common Words in **Positive** Comments

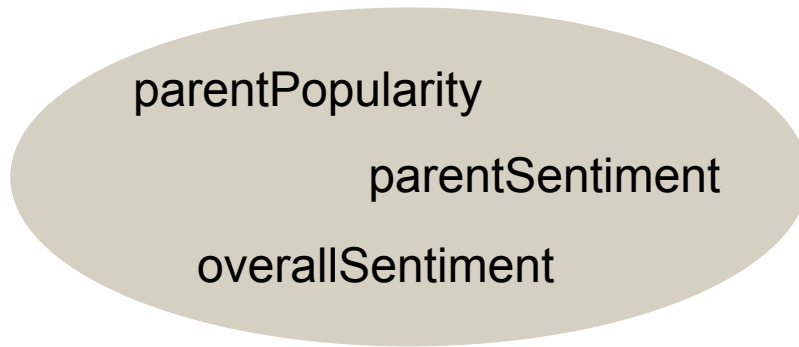### Most Common Words in **Negative** Comments



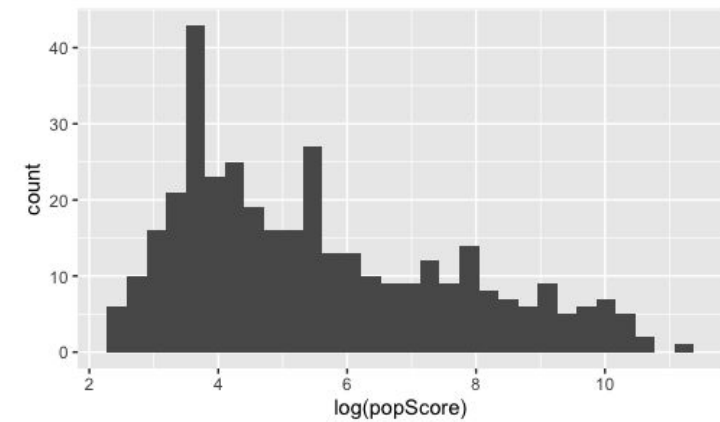*Knowing this, we removed words like "game" with stop words*

# SAMPLED COMMENTS

# CREATING THE MODEL

parentPopularity

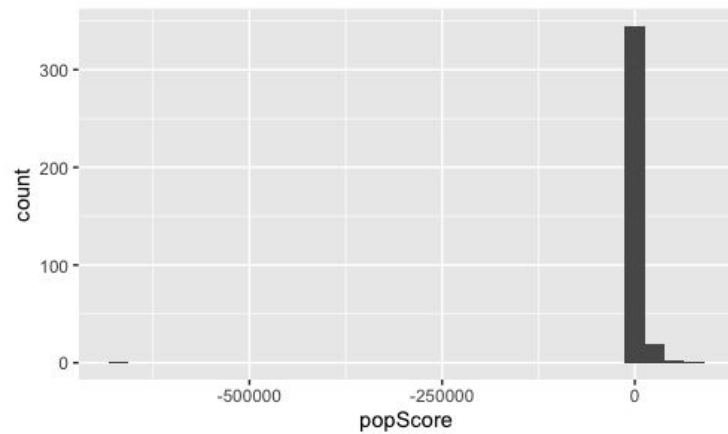parentSentiment     ⟶     Comment's Popularity

overallSentiment

```
Coefficients:
                              Estimate Std. Error t value Pr(>|t|)
(Intercept)                    2.28324    0.89197   2.560 0.011609 *
log(abs(parentPopularity))     0.41798    0.11254   3.714 0.000301 ***
parentSentiment               -0.18718    0.05761  -3.249 0.001470 **
overallSentiment               0.06315    0.09033   0.699 0.485725
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

**Note**: Data used for the model was one where the sentiment of the comment and the popularity of its parent were aligned. This was, in part, because using the entire pool of comments yielded a useless model, though this choice is also theoretically justified.

# CASE STUDY TAKEAWAYS

Popularity of parent post is a big indicator of a post's potential to become popular.
*High coefficient for parent popularity in linear model*
Makes sense to us – commenting on an unpopular post doesn't give you traction
- harder to 'go viral' from a lower-visibility starting point

Unpopular posts lend themselves to more popular comment sections.
*Negative coefficient for parent sentiment in linear model*
Makes sense to us – more opportunities for a 'clap back,' opportunity to defend your belief
- interesting relationship with social media algorithms & negativity bias

The sentiment of a post seems to be less important than the sentiment of the post it responds to.
*Minute coefficient for post sentiment in linear model, little significance*
More important seems to be the (mis)match in sentiment from parent and response
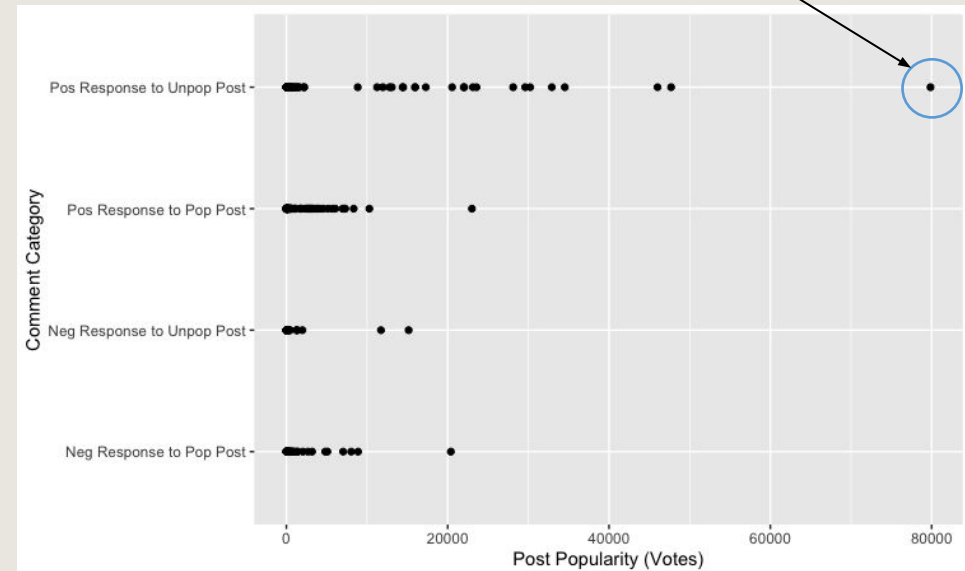
# CASE STUDY LIMITATIONS

Length of posts versus words contained in lexicons
  roughly 10% of words are classified
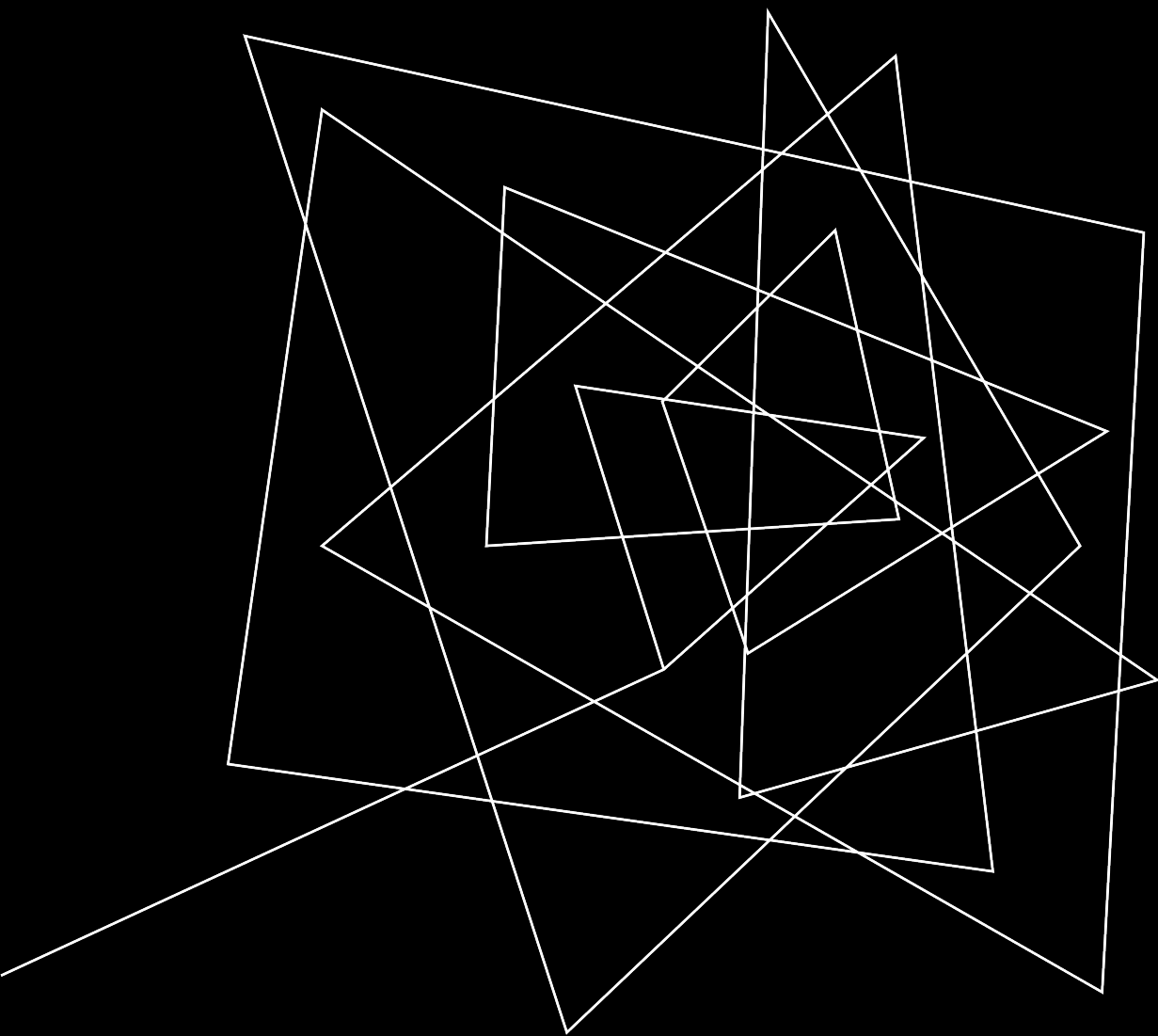  most comments are already short

  *look at sentiment score per word?*

Lack of posts with unpopular parents
  functionally, only EA's post filled this role
  difficult to distinguish between negative
  parent and depth in thread

SARCASM!

*"I wonder if Burger King wants to sell me a sense of pride and accomplishment by making me work 10 hours for my f\*\*king fries."*
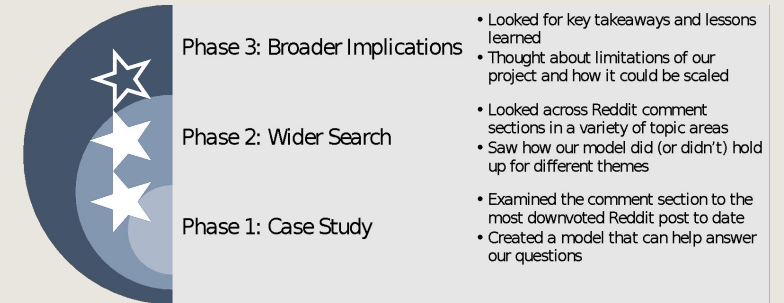
**PHASE 2:** WIDER SEARCH

*Model Cross-Topic Generalizability*

# SELECTING OTHER THREADS



Phase 3: Broader Implications
- Looked for key takeaways and lessons learned
- Thought about limitations of our project and how it could be scaled

Phase 2: Wider Search
- Looked across Reddit comment sections in a variety of topic areas
- Saw how our model did (or didn't) hold up for different themes

Phase 1: Case Study
- Examined the comment section to the most downvoted Reddit post to date
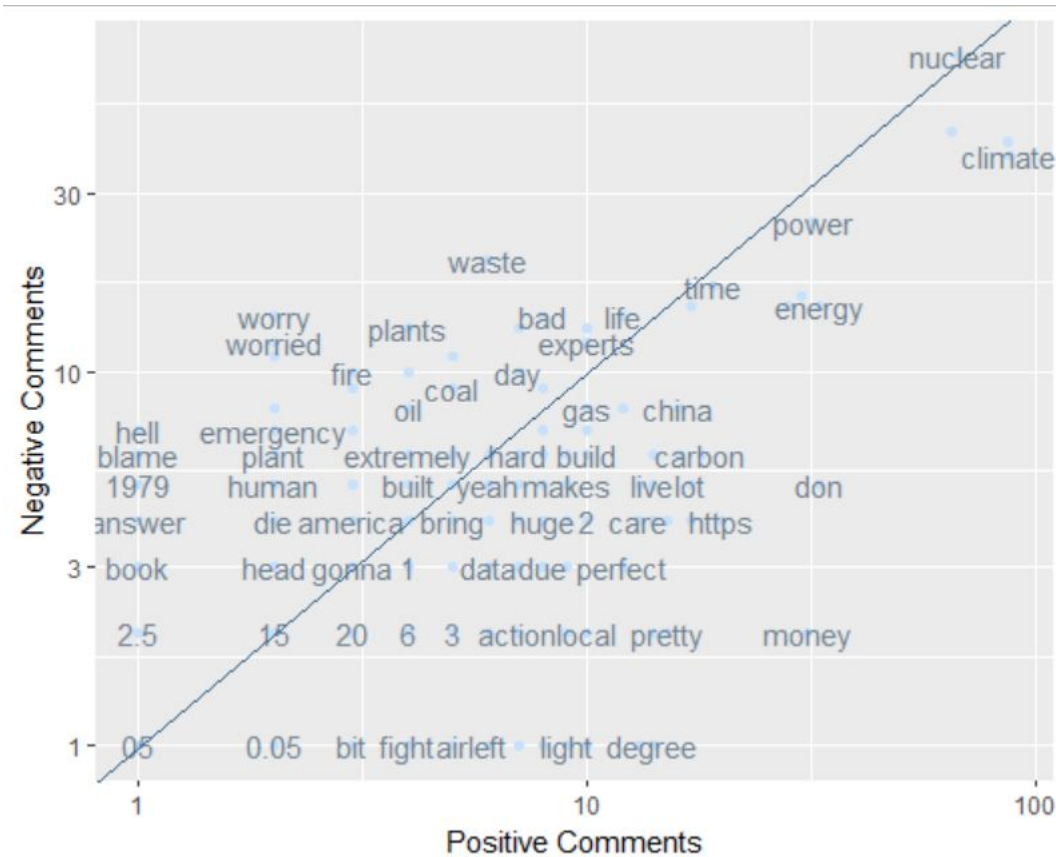- Created a model that can help answer our questions

Topics Selected Using Implications from Case Study:
- Looking for longer posts, more 'charged' language
- Looking for fewer instances of sarcasm
- More mix of up and down votes

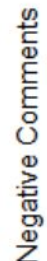Climate Change                    Politics

# CLIMATE CHANGE



```
Coefficients:
                    Estimate Std. Error t value Pr(>|t|)
(Intercept)         7.395639   1.554480   4.758 3.31e-06 ***
overallSentiment    0.105425   0.539530   0.195    0.845
parentSentiment    -0.396221   0.483903  -0.819    0.414
parentPopularity    0.034231   0.008534   4.011 7.98e-05 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

# POLITICS



```
Coefficients:
                        Estimate Std. Error t value Pr(>|t|)
(Intercept)             1.371698   0.249946   5.488 8.66e-08 ***
overallSentiment        0.030817   0.046947   0.656    0.512
parentSentiment        -0.001759   0.048007  -0.037    0.971
log(parentPopularity)   0.273813   0.045292   6.045 4.42e-09 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```
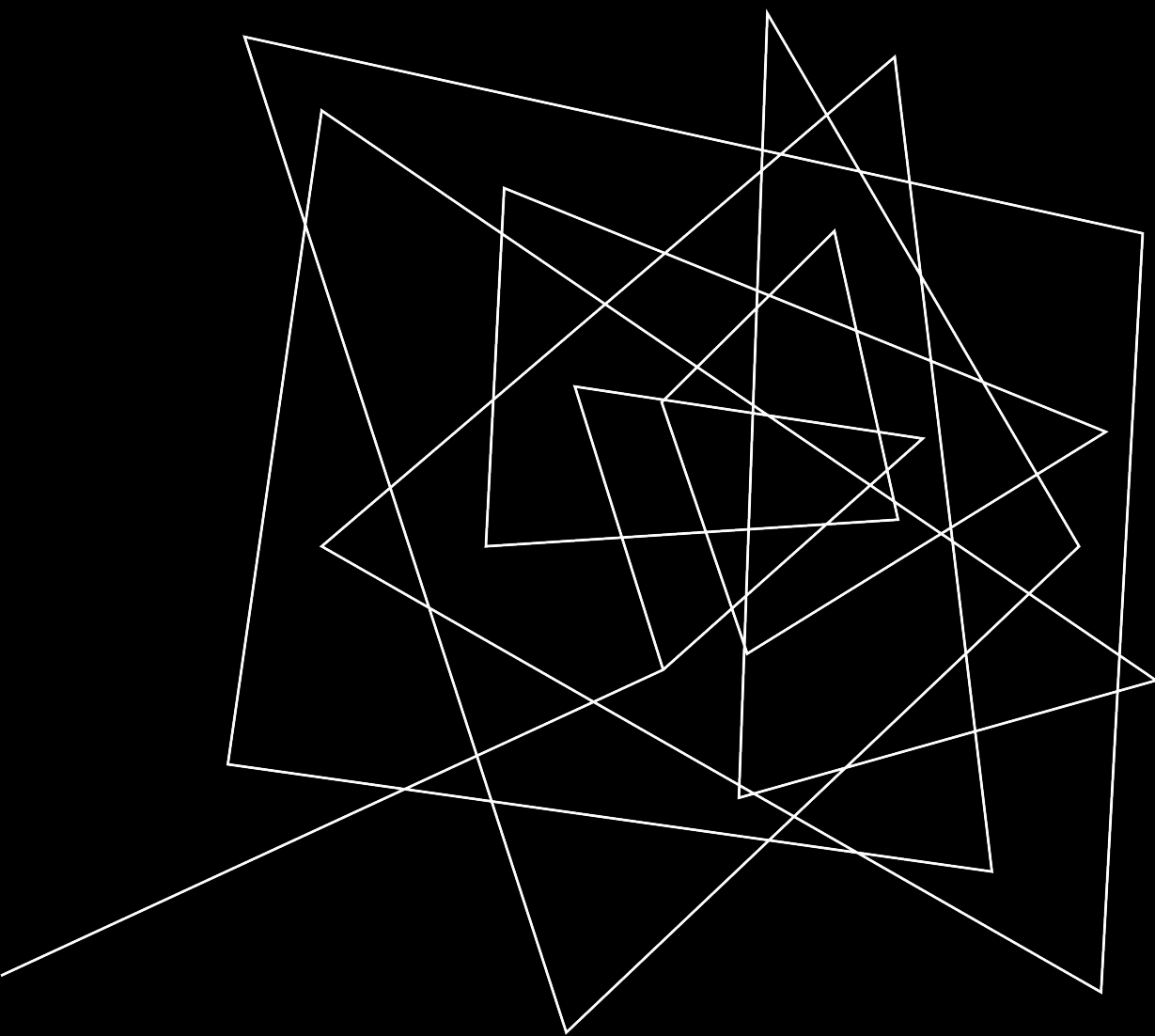
# COMPARISON ACROSS TOPICS

Regression Coefficients:

|  | **Parent Popularity** | **Parent Sentiment** | **Post Sentiment** |
|---|---|---|---|
| Star Wars | 0.41* | -0.19* | 0.06 |
| Climate Change | 0.03* | -0.40 | 0.11 |
| Politics | 0.27* | -0.001 | 0.03 |

* Statistically Significant at at least the 0.05-level

**PHASE 3:** BROADER IMPLICATIONS

*What Now?*

# WHAT WE'VE FOUND, OVERALL

1. Serious variability in indicators of post popularity
   *Makes sense - social media is complex; if we knew the answer, we'd be rich*
   *However, seems to be some relationships with how 'serious' and* **polarized** *the topic is*
   *We're looking at a small piece of the puzzle (parent and child popularity and sentiment)*

2. Limited capabilities of sentiment analysis
   *The process explored in class works well with much larger inputs*
   *Our lexicons don't account for context: sarcasm*
   *Two factors really compounded this problem: internet comments are short & few words can be characterized*
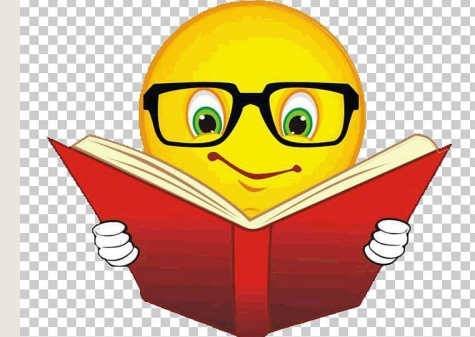
3. Interesting qualitative data
   *This process can be applied and reapplied to other threads*
   *Interesting visualizations to explore by thread and characteristic, more info in our report*

# IMPLICATIONS FOR FUTURE WORK

1. More thought is needed for small-batch sentiment analysis and how to account for word context (even beyond looking at n-grams)

2. Reddit is just one of the major platforms that contributes to polarization.
   a. Would be interesting to look across networks (Twitter, Facebook, Instagram) for similar analyses.
   b. Do the same topics look different on different platforms?
   c. Unlike Reddit, other platforms do not have "downvotes," how to account for this?

3. Application of more robust statistics. Simple regression models capture some, but not all, variation in comment popularity, and to varying degrees of significance.
   a. Perhaps future work could start with more available data and look to identify major contributors to post popularity (through PCA, perhaps).

4. Beyond just comment popularity, what other metrics can be used to better capture "polarization," and how can data science (with more robust application) help us understand this?
   *(In early iterations of this project, we attempted to assign each comment a "polarization quotient" to capture a mix of qualitative and quantitative properties.)*

# REFERENCES

- Ivan-Rivera. "Ivan-Rivera/Redditextractor: A Minimalistic R Wrapper for the Reddit API." *GitHub*, https://github.com/ivan-rivera/RedditExtractor.
- "R/Askanamerican - Are You Concerned about Climate Change?" *Reddit*, https://www.reddit.com/r/AskAnAmerican/comments/rtm63v/are_you_concerned_about_climate_change/.
- "R/Politics - Republican Who Refuses to Bend the Knee to Trump Surges in Ohio Senate Race." *Reddit*, https://www.reddit.com/r/politics/comments/uekpka/republican_who_refuses_to_bend_the_knee_to_trump/.
- "R/Starwarsbattlefront - Comment by U/EACommunityTeam on 'Seriously? I Paid 80$ to Have Vader Locked?".'" *Reddit*, https://www.reddit.com/r/StarWarsBattlefront/comments/7cff0b/seriously_i_paid_80_to_have_vader_locked/dppum98/.