# People Still Care About Facts: Twitter Users Engage More with Factual Discourse than Misinformation

**Luiz Giovanini**[*][1], *Shlok Gilda*[*][1], **Mirela Silva**[**][1], **Fabrício Ceschin**[**][2], **Prakash Shrestha**[1], **Christopher Brant**[1], **Juliana Fernandes**[1], **Catia S. Silva**[1], **André Grégio**[2], and **Daniela Oliveira**[1].

[1] University of Florida.

[2] Federal University of Paraná, Brazil.

[*] Authors have equal contribution.

[**] Authors have equal contribution.

# The Big Question

How do Twitter users engage with COVID-19 misinformation and factual-information?

# Research Questions

**RQ1:** Are COVID-19 misinformation tweets more engaging than COVID-19 factual tweets?

**RQ2:** Are general topic misinformation tweets more engaging than general topics factual tweets?

**RQ3:** Which features are most correlated with engagement in COVID-19 vs. general topics misinformation tweets?

**RQ4:** Which features are most correlated with engagement in COVID-19 vs. general topics factual tweets?

# Data Collection

2.1M tweets.

Four primary datasets:

- COVID-19 misleading claims.

- COVID-19 factual claims.

- Misleading claims on general topics.

- Factual claims on general topics.

# COVID-19 Twitter Data Sources

| Source | Description | False | True |
|---|---|---|---|
| Shahi et al. | Fact-checked Coronavirus-related tweets | 1345 | 41 |
| Schroeder et al. | Tweets linking COVID-19 with 5G conspiracy theories | ~58K | N/A |
| CoAID | News articles & social media posts with fake and factual claim labels | 484 | 8092 |
| Paka et al. (CTF) | Labeled and unlabeled tweets related to COVID-19 | ~17K | ~18K |
| Muric et al. | Tweets related to anti-vaccine narratives for COVID-19 | ~1.8M | N/A |

# General Topics Twitter Data Sources

| Source | Description | False | True |
|---|---|---|---|
| Mitra and Gilbert (CREDBANK) | Crowdsourced tweets related to real-world news events | N/A | ~1.94M |
| Russian Troll Tweets Kaggle | Tweets from malicious accounts connected to Russia's Internet Research Agency | 200K | N/A |
| Vo and Lee | Fact-checked tweets based on news articles from Snopes and Politifact | ~59K | ~14K |
| Jiang et al. | Tweets labeled across a spectrum of fact-check ratings | 1264 | 231 |

# Data Cleaning and Preparation

- **Discarded Tweets:** Removed duplicates, non-English, and text-less entries.

- **Collected Metadata:** Gathered details on tweets, authors, and engagement.

# Data Cleaning and Preparation

| | Factual | | Misinformation | |
|---|---|---|---|---|
| | COVID-Related | General Topics | COVID-Related | General Topics |
| $N$ | 9,111 (0.43%) ← | 1,243,913 (58.84%) | 828,501 (39.19%) | 32,243 (1.52%) |
| $n_{strata}$ | 4,814 | 4,448 | 4,533 | 4,147 |
| $\mu$ | 368.5 | 9,791.6 | 2,214.3 | 3,014.7 |
| $\sigma$ | 7,157.9 | 73,305.6 | 10,051.9 | 28,727.4 |
| Mean Rank | 2407.5 | 2244.5 | 2267.0 | 2074.0 |

Descriptive statistics of our final four datasets based on the combined engagement metric.
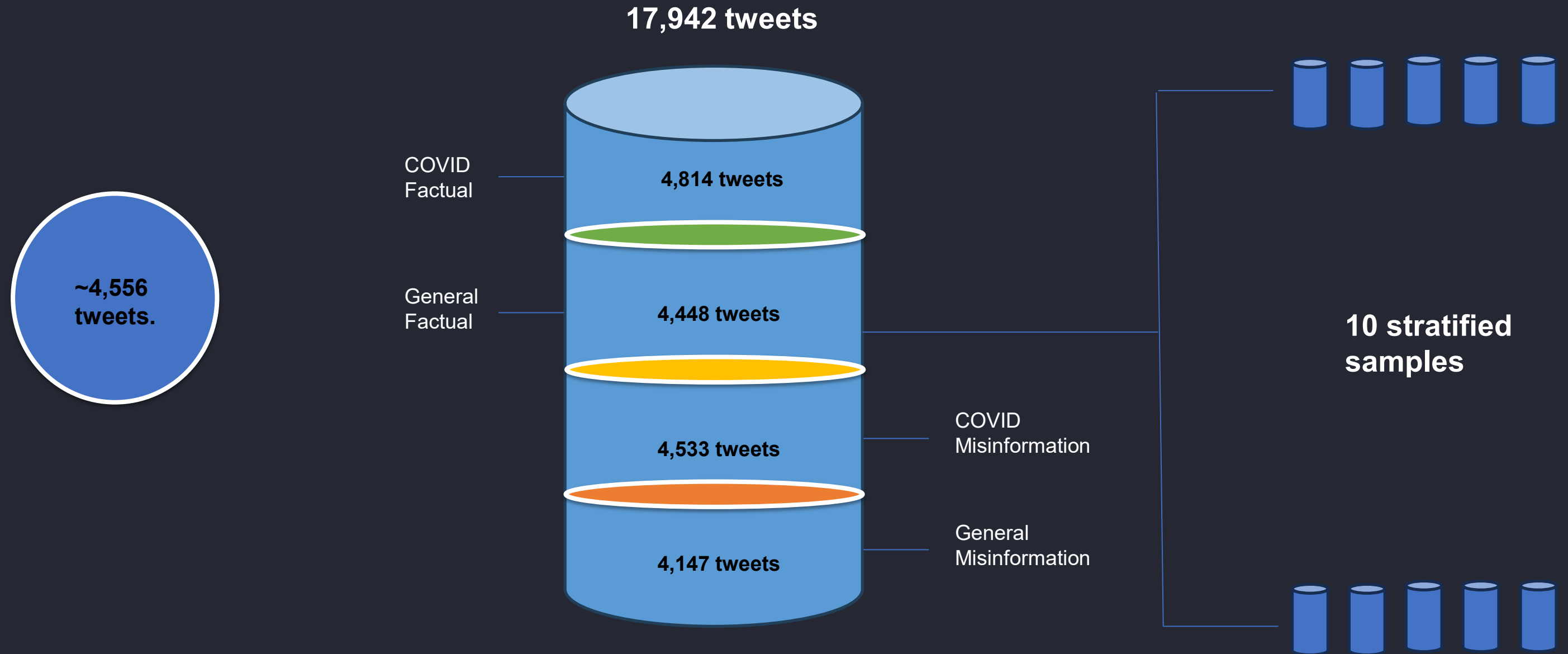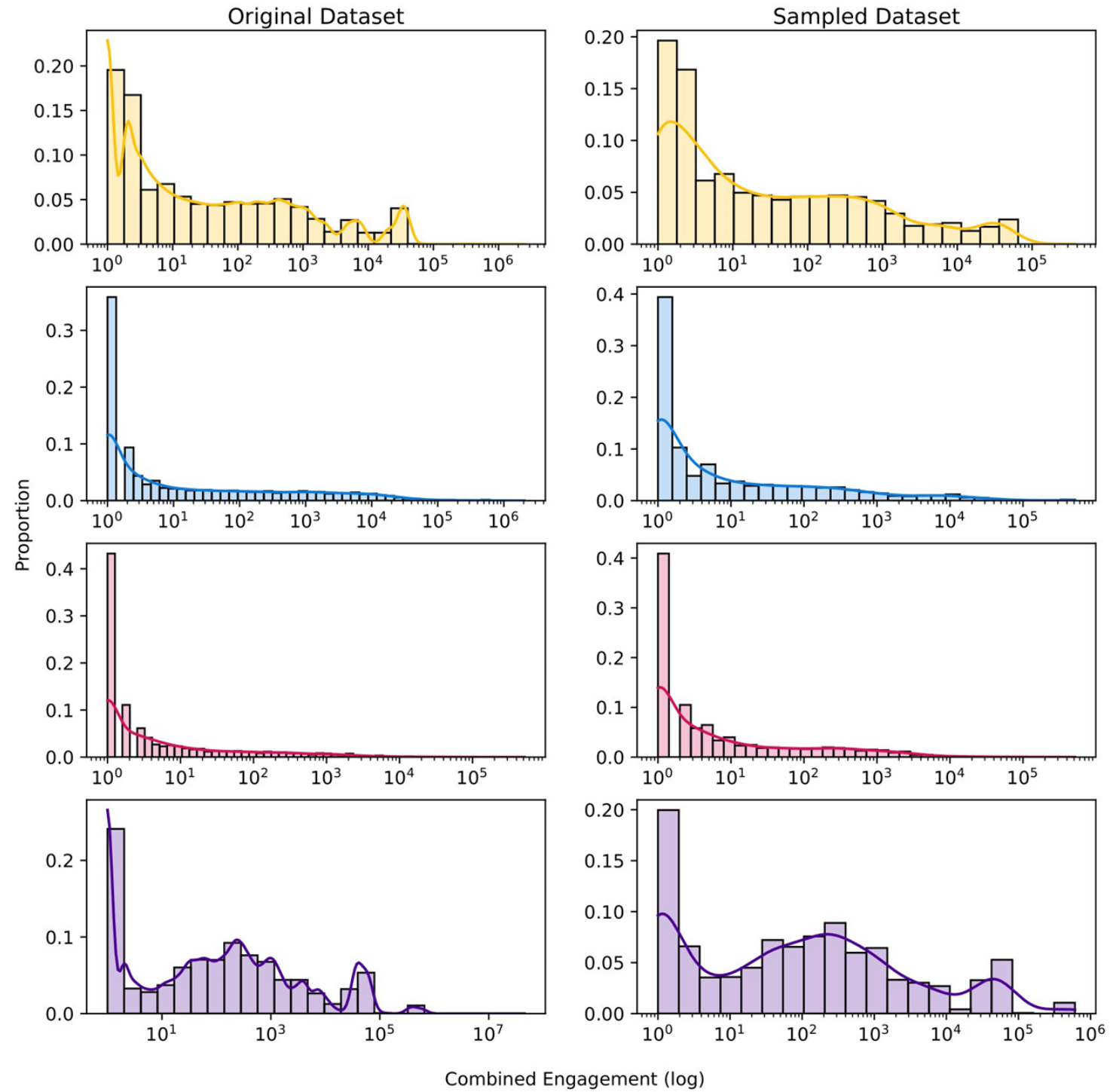
# Stratified Random Sampling



Image courtesy of elgin.edu

Original Dataset — Sampled Dataset

Legend: Misinfo. COVID, Misinfo. General, Factual COVID, Factual General

Proportion

Combined Engagement (log)

# Feature Extraction

## Sociolinguistic Analysis

- Linguistic Inquiry and Word Count (LIWC) software.

- Emotional, cognitive, and structural components.

**1**     **2**     **3**

## Tweet Metadata

- Text-based features.

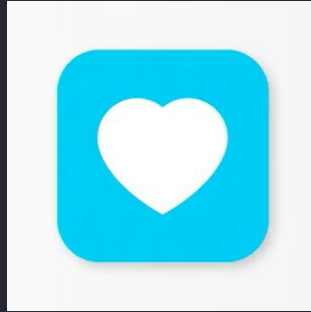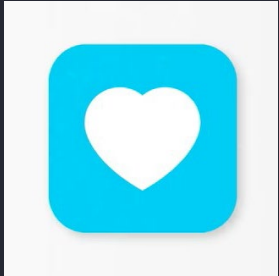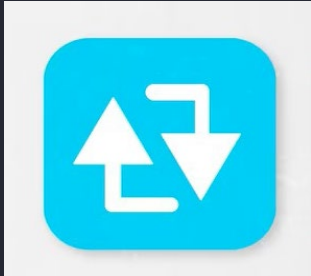- User-based features

- Network-based features.

## Sentiment Analysis

- NLTK VADER

# Tweet Metadata

| | |
|---|---|
| **# likes** | **# retweets** |
| **# combined engagement** | **# hashtags** |
| **# links** | **# emojis** |

# Tweet Metadata

## User-related features

- # of followers.
- # of friends.
- # of lists.
- # of favorited tweets.
- # of tweets made by the user.
- verified (binary).

- presence of profile image.
- use of default profile image.
- use of default profile.
- whether geolocation is enabled.
- user has an extended profile.
- user has a background tile.

# Sociolinguistic Analysis

## Language Metrics

- Total number of words.

- Average number of words per sentence.

- Number of words containing more than six letters.

- Number of words found in the LIWC dictionary

# Sociolinguistic Analysis

## Linguistic Indicators

- Function words.

- Grammar characteristics.

- Affective words.

- Social words.

- Cognitive process.

- Core needs.

- Informal speech.

# Sociolinguistic Analysis

## Summary Variables

- Analytical Thinking.

- Clout.

- Authenticity.

- Emotional Tone.

# Sociolinguistic Analysis

## Moral Frames

- Care.

- Fairness.

- Loyalty.

- Authority.

- Sanctity.

# Sentiment Analysis

# Correlation Analysis

## Pearson's correlation coefficient (*r*)

- Measure feature importance.
- *r* only captures linear relationships.

## Alternating Conditional Expectations

- Feature's fixed point of Maximal Correlation.

## Fisher *z*-transformation

- Reduce bias.
- Estimate population correlation.

# Statistical Analyses – RQ1 and RQ2

| Data | Measure | Measurement Statistics | |
|---|---|---|---|
| Combined Engagement (raw) | Shapiro-Wilk | Factual COVID-Related | $W = 0.7875$*** |
| | | Misinformation General Topics | $W = 0.8946$*** |
| | | Factual General Topics | $W = 0.9374$*** |
| | | Misinformation General Topics | $W = 0.7969$*** |
| Combined Engagement (log-norm) | Levene | Factual vs. Misinformation COVID-Related | $W = 378.89$*** |
| | | Factual vs. Misinformation General Topics | $W = 359.59$*** |
| | Two-Sample Kolmogorov-Smirnov | Factual vs. Misinformation COVID-Related | $K_2 = 0.2133$*** |
| | | Factual vs. Misinformation General Topics | $K_2 = 0.3459$*** |
| | Mann-Whitney U | Factual vs. Misinformation COVID-Related | $U = 7,662,279$***, $r = 0.35$ |
| | | Factual vs. Misinformation General Topics | $U = 5,725,193$***, $r = 0.31$ |

*** Significant at $p < .001$

Summary results for statistical tests conducted on engagement metrics and bot/user account labels.

# Findings – RQ1 and RQ2



**Factual tweets were more engaging than misinformation tweets, regardless of their topic.**

# Statistical Analyses – RQ3

| Feature Type | Feature | $r_z$ | (MC) $r_z$ |
|---|---|---|---|
| **Misinformation: COVID-Related** | | | |
| *LIWC* | Assent (Informal Speech) | 0.26 | 0.75 |
| | Colons (All Punctuation) | 0.34 | 0.75 |
| | Informal Speech | 0.19 | 0.69 |
| | Impersonal Pronouns | 0.06 | 0.64 |
| | Netspeak (Informal Speech) | 0.26 | 0.73 |
| | Quotation Marks (All Punctuation) | 0.10 | 0.50 |
| | Word Count | -0.10 | 0.51 |
| **Misinformation: General Topics** | | | |
| *User Metadata* | Followers Count | 0.28 | 0.73 |
| | Listed Count | 0.30 | 0.66 |
| | User Verified | 0.53 | 0.53 |

Summary of correlation analysis between the log normalized combined engagement metric and relevant features.

# Statistical Analyses –RQ4

| Feature Type | Feature | $r_z$ | (MC) $r_z$ |
|---|---|---|---|
| **Factual: COVID-Related** | | | |
| *LIWC* | Affective Processes | 0.53 | 0.71 |
| | All Punctuation | -0.05 | 0.58 |
| | Assent (Informal Language) | 0.65 | 0.74 |
| | Clout | 0.36 | 0.56 |
| | Colon (Punctuation) | 0.34 | 0.54 |
| | Dictionary Words | 0.13 | 0.56 |
| | Past Focus | 0.49 | 0.66 |
| | Informal Speech | 0.62 | 0.72 |
| | Insight (Cognitive Processes) | 0.32 | 0.68 |
| | Male Referents (Social Words) | 0.77 | 0.88 |
| | Netspeak (Informal Language) | 0.66 | 0.77 |
| | Positive Emotion (Affect Words) | 0.52 | 0.78 |
| | Person Pronouns (Linguistic Dimensions) | 0.31 | 0.56 |
| | Question Marks (All Punctuation) | -0.31 | 0.53 |
| | Reward (Drives) | 0.33 | 0.67 |
| | Sad (Affect Words) | 0.48 | 0.65 |
| | 3rd Person Singular (Function Words) | 0.81 | 0.91 |
| | Words > 6 Letters | -0.26 | 0.59 |
| | Social Words | 0.41 | 0.63 |
| | Time (Relativity) | 0.21 | 0.51 |
| *Sentiment* | VADER Compound | 0.19 | 0.66 |
| **Factual: General Topics** | | | |
| *LIWC* | Assent (Informal Speech) | 0.36 | 0.68 |
| | Colons (All Punctuation) | 0.20 | 0.52 |
| | Informal Speech | 0.29 | 0.62 |
| | Netspeak (Informal Speech) | 0.32 | 0.63 |
| | Prepositions (Function Words) | 0.02 | 0.54 |

# Findings – RQ3

**RQ3:** **Which features are most correlated with engagement in COVID-19 vs. general topics misinformation tweets?**

COVID-19:

- Grammar (e.g., use of informal speech).

Factual:

- User metadata (e.g., verified user).

# Findings – RQ4

**RQ4:** **Which features are most correlated with engagement in COVID-19 vs. general topics factual tweets?**

COVID-19:

- Grammar (e.g., use of netspeak).

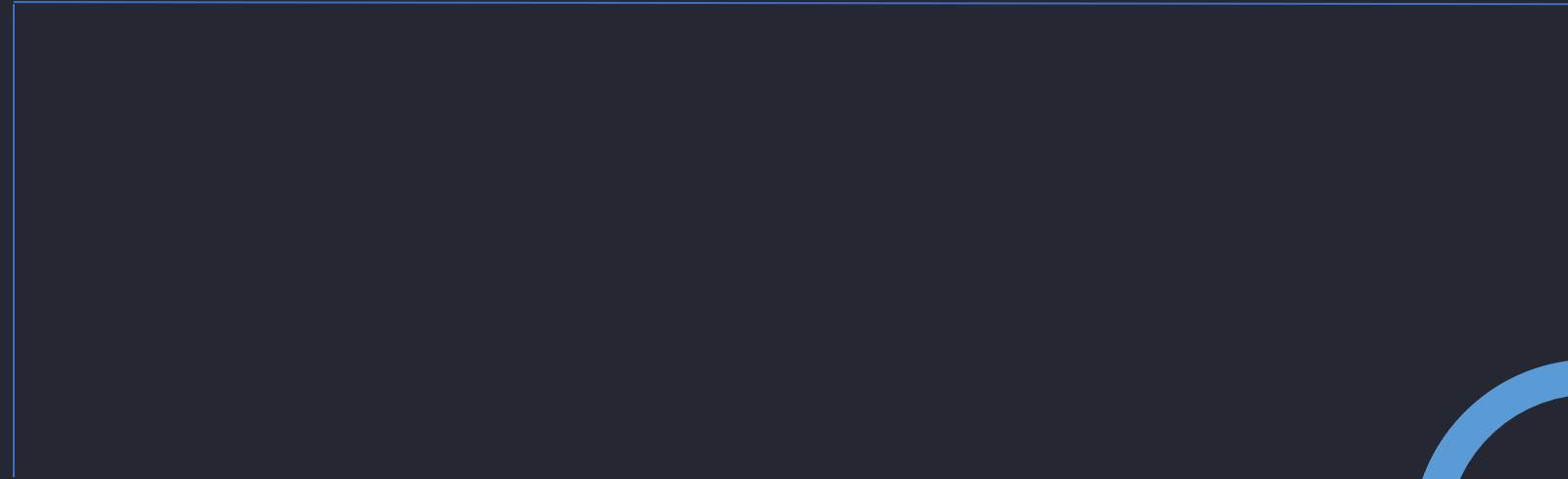- Emotion (both positive and negative).

- Writer's confidence.

Factual:

- Grammar (e.g., use of colons or prepositions).

# Discussion

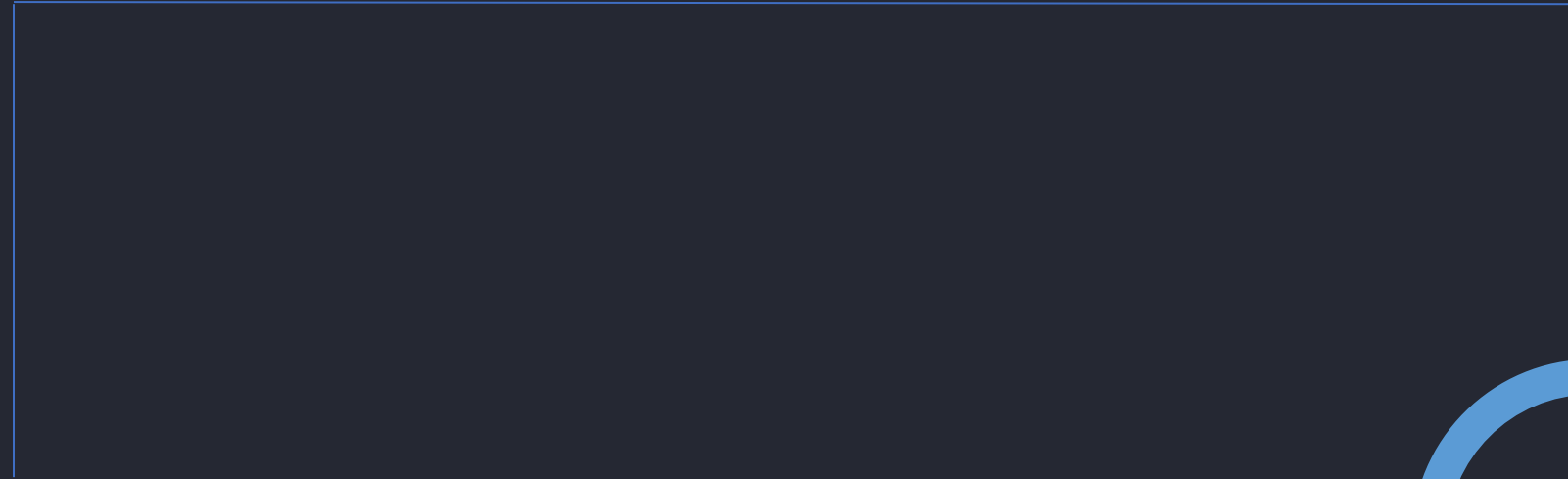- General topic misinformation ➔ User metadata

# Discussion

- Semantic content of tweet not relevant *(except for factual COVID-19 tweets)*.

- Factual COVID ➔ Tweet syntax.

- General topic ➔ Tweet syntax.

# Discussion

- Factual COVID-19 ➔ Sentiment.
  ➔ Cognitive processing keywords.

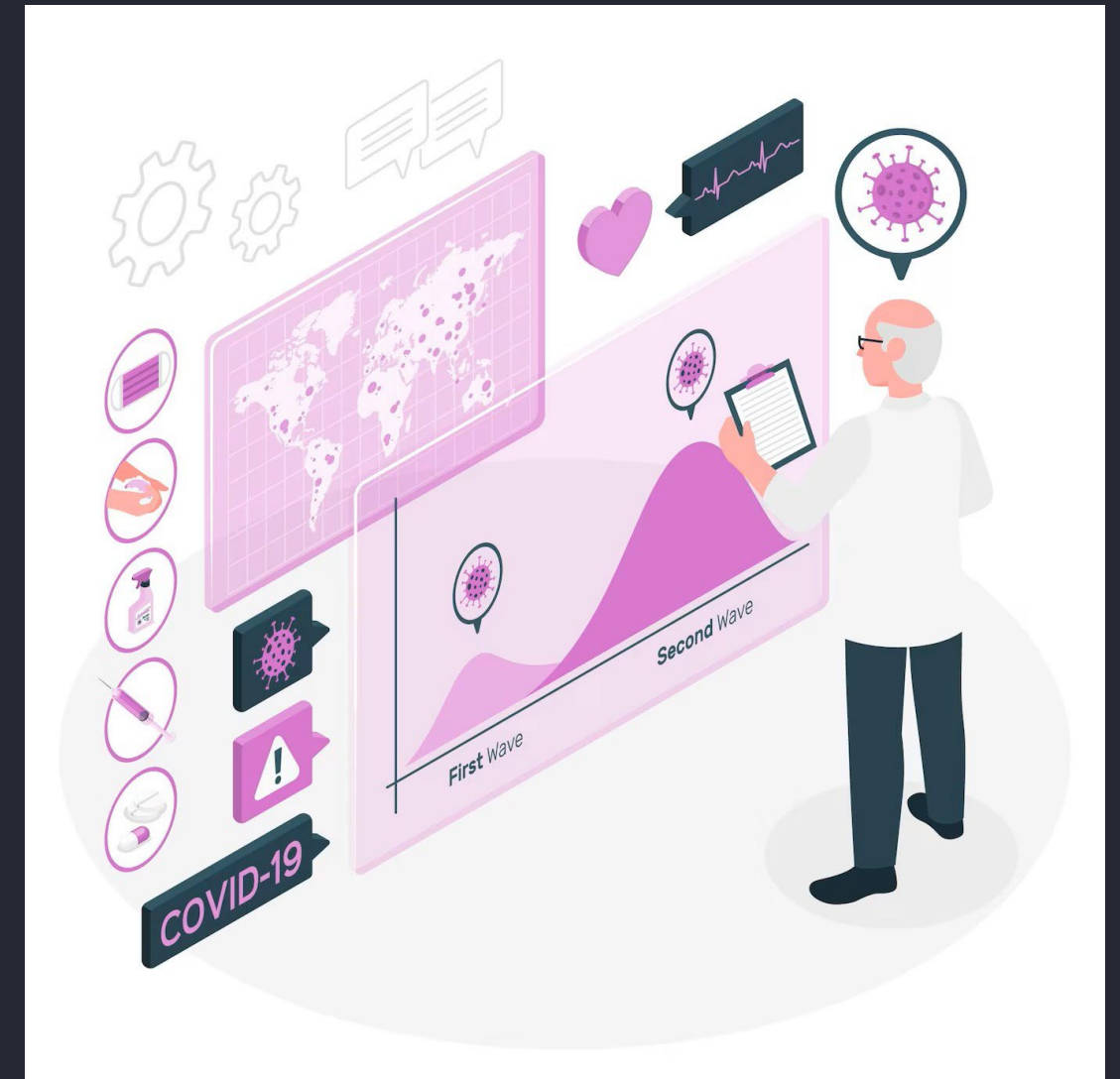- Misinformation COVID-19 ➔ Clear, straightforward language.

# Discussion

- Factual tweets > Misinformation tweets
  *(in terms of engagement)*

# Study Limitations & Future Work

## Dataset Imbalance

- Highly imbalanced dataset.

- Findings may not apply to other contexts.

- Explore temporal trends in tweet data.

# Study Limitations & Future Work

## Feature Engineering

- Did not measure presence of images.

- Utilize automated feature extractors.



publicdomainvectors.org

# Study Limitations & Future Work

## Classification Models

- Relied only on pairwise correlation.
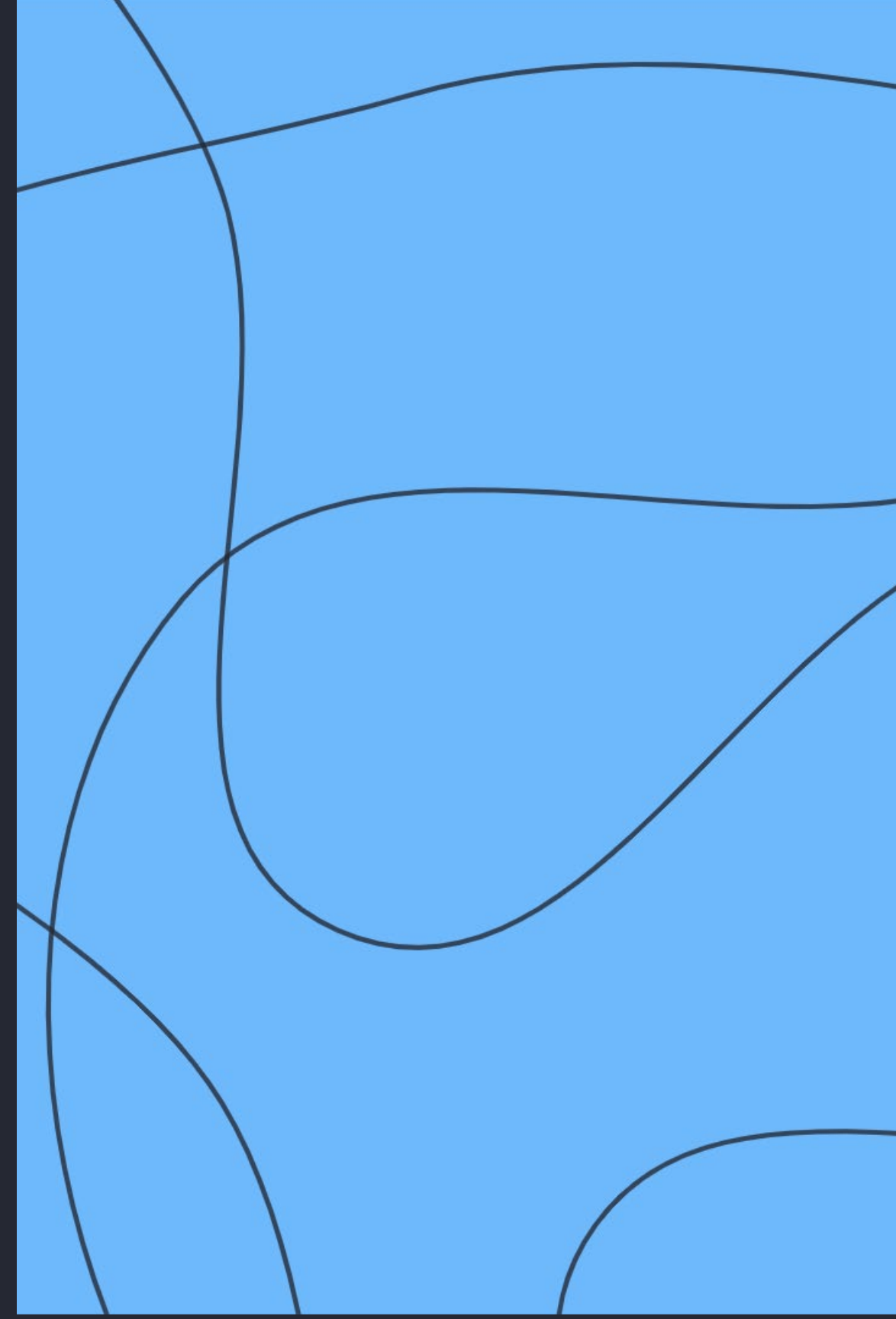
- Study multivariate analyses.

# CONCLUSION

- Dataset of 2.1M COVID-19 and non-COVID related tweets.

- Misinformation tweets less engaging than factual tweets.

- Tweet features correlating with engagement vary based on veracity.

# Acknowledgements

Images courtesy of freepik.com.

# People Still Care About Facts: Twitter Users Engage More with Factual Discourse than Misinformation

Luiz Giovanini: lfrancogiovanini@ufl.edu

Shlok Gilda: shlokgilda@ufl.edu

Mirela Silva: msilva1@ufl.edu

Fabrício Ceschin: fjocescin@inf.ufpr.br

Daniela Oliveria: daniela@ece.ufl.edu

## Questions?

We welcome questions and further discussion.