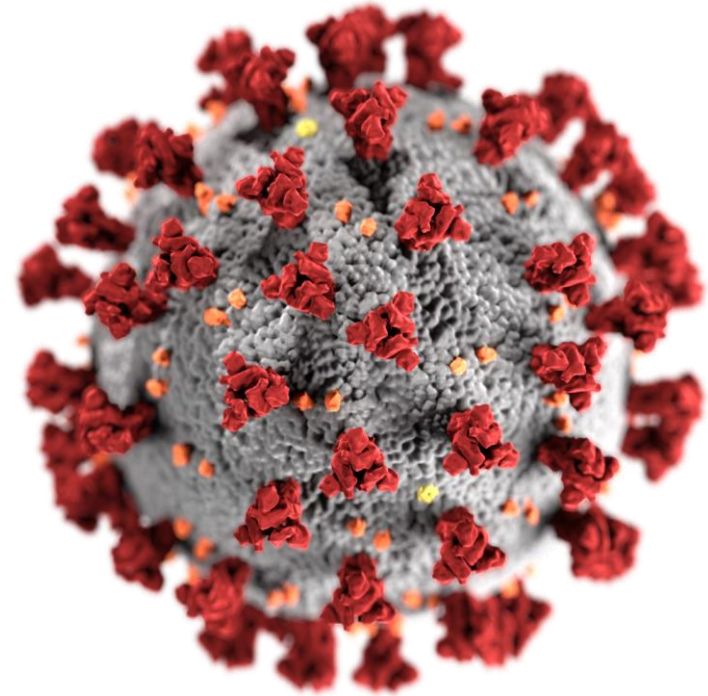# Social Networks: The Tracking of COVID-19 Misinformation Using NLP and Graph-based Approaches

Shogo Toyonaga
Aaryaman Kartha

# Motivation

- The Covid-19 pandemic has brought about millions of death globally
- Misinformation plays a big role through **vaccine hesitancy** and undermining of the virus's impact
- Social media has been used to spread **misinformation** at unprecedented rates, primarily through bots and influencers
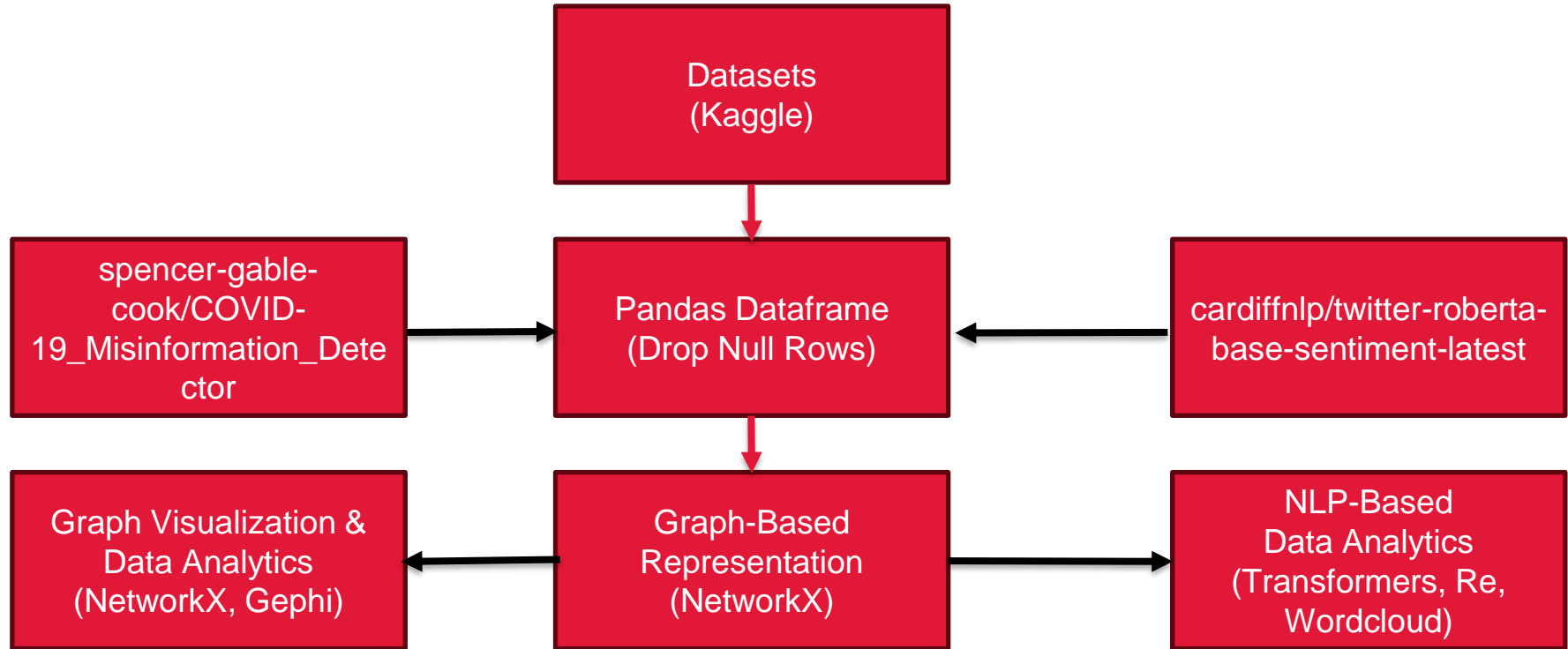
# Research Questions

1. What kinds of users participate in the dissemination of misinformation?

2. What Natural Language Processing (NLP) based patterns can be drawn from malicious users and their tweets/posts?

3. How is misinformation propagated in terms of a longitudinal time frame?

4. What are noticeable differences and similarities in how misinformation spreads between Reddit and Twitter?

YORK U

# Recap: What We Did Last Time

1. Twitter & Reddit Pipeline
2. Graph-based Analysis
   - Community Detection
   - PageRank
   - HITS
   - Null Model (Barabasi-Albert)
3. Data Visualization
   - Wordclouds
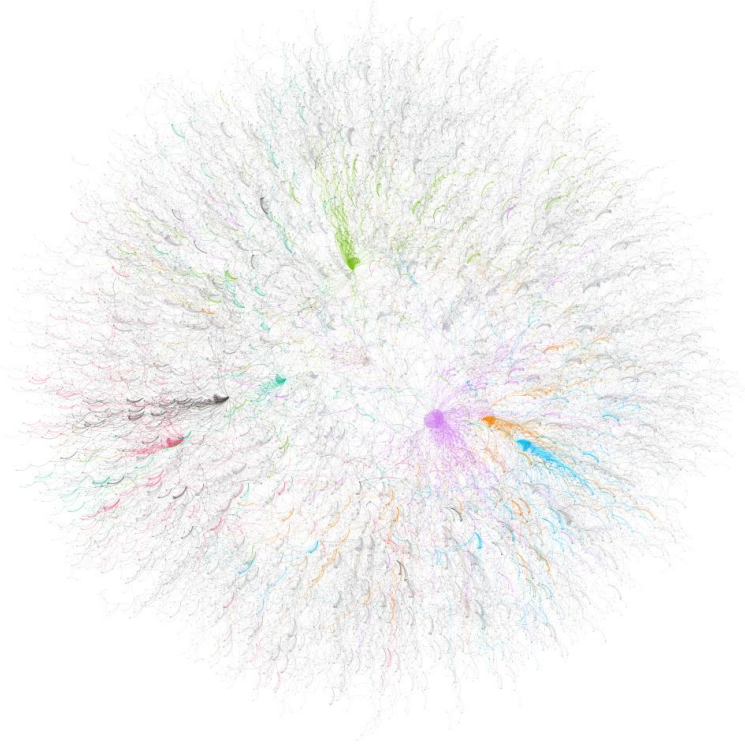   - Timeseries
   - Graphs (Gephi: Yifan Lu Layout)

YORK U

# Recap: Methodology (Twitter)

# Recap: Problem Definition

- Our aim is to create a **directed** social network graph **G = (V, E)** where:

  - **V** represents users who are participating in the spread of misinformation

  - **E** represents the relationships between users who share, view, or distribute misinformation

# Results and Analysis (Twitter) - 2



| Property | Value |
|---|---|
| Nodes | 35,147 |
| Edges | 35,146 |
| Degree Distribution | 2.99 |
| Triangles | 0 |
| Clustering Coefficient | 0 |
| Modularity | 0.980 |
| Communities | 112 |
| Diameter | 31 |

YORK U

# New Contribution: Interactive QA Dashboard



Goals:
1. Intuitive, Accessible, Efficient Visualizations to answer the Research Questions
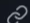2. Interactive, Scalable Solutions

# GS/EECS 6414 Twitter Dashboard

The goal of our project is to create a robust multifold analytics and visualization system for COVID misinformation in Twitter. Firstly, we will analyze several datasets which contain social media posts about COVID-19 in addition to user-level demographics (i.e., Number of Followers, Biography, etc.,). From this, we will apply NLP techniques which augment the available datasets for data mining and visualization purposes.

Secondly, we will create directed social networks with malicious posts and tagged users to determine how misinformation travels. The resulting network characteristics will be reported in the methodology section. Using these characteristics, we will determine what null model most closely corresponds to our social networks.

Lastly, through the process of data analytics, graph visualization, and pattern recognition, we seek to answer the following research questions based on various subgraph motifs:

1. Which users participate the most in the dissemination of misinformation?
2. What NLP-based patterns can be drawn from users who spread misinformation?
3. How is misinformation propagated in terms of a longitudinal time frame?
4. What are notable similarities and differences in how misinformation spreads between Twitter and Reddit?

## Text-based Analytics & Visualizations
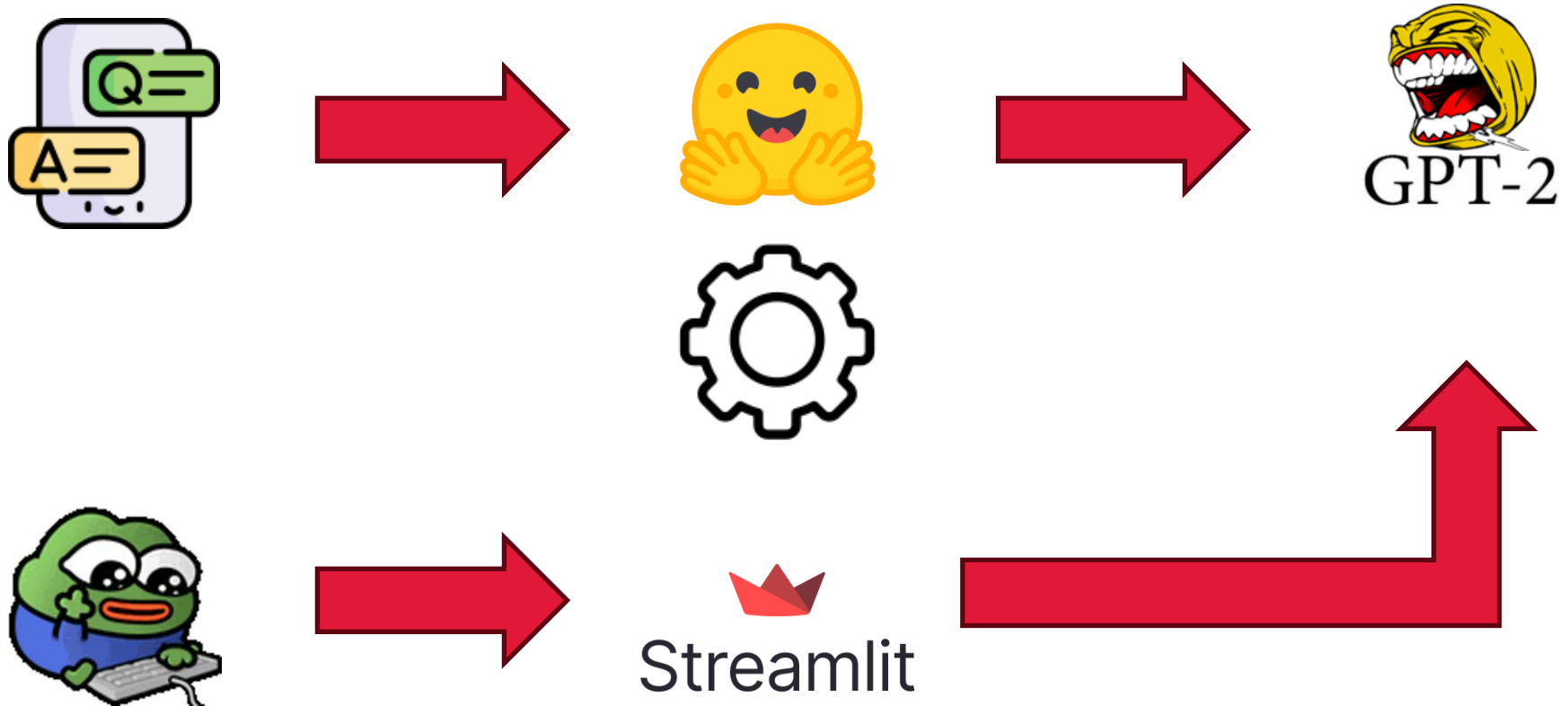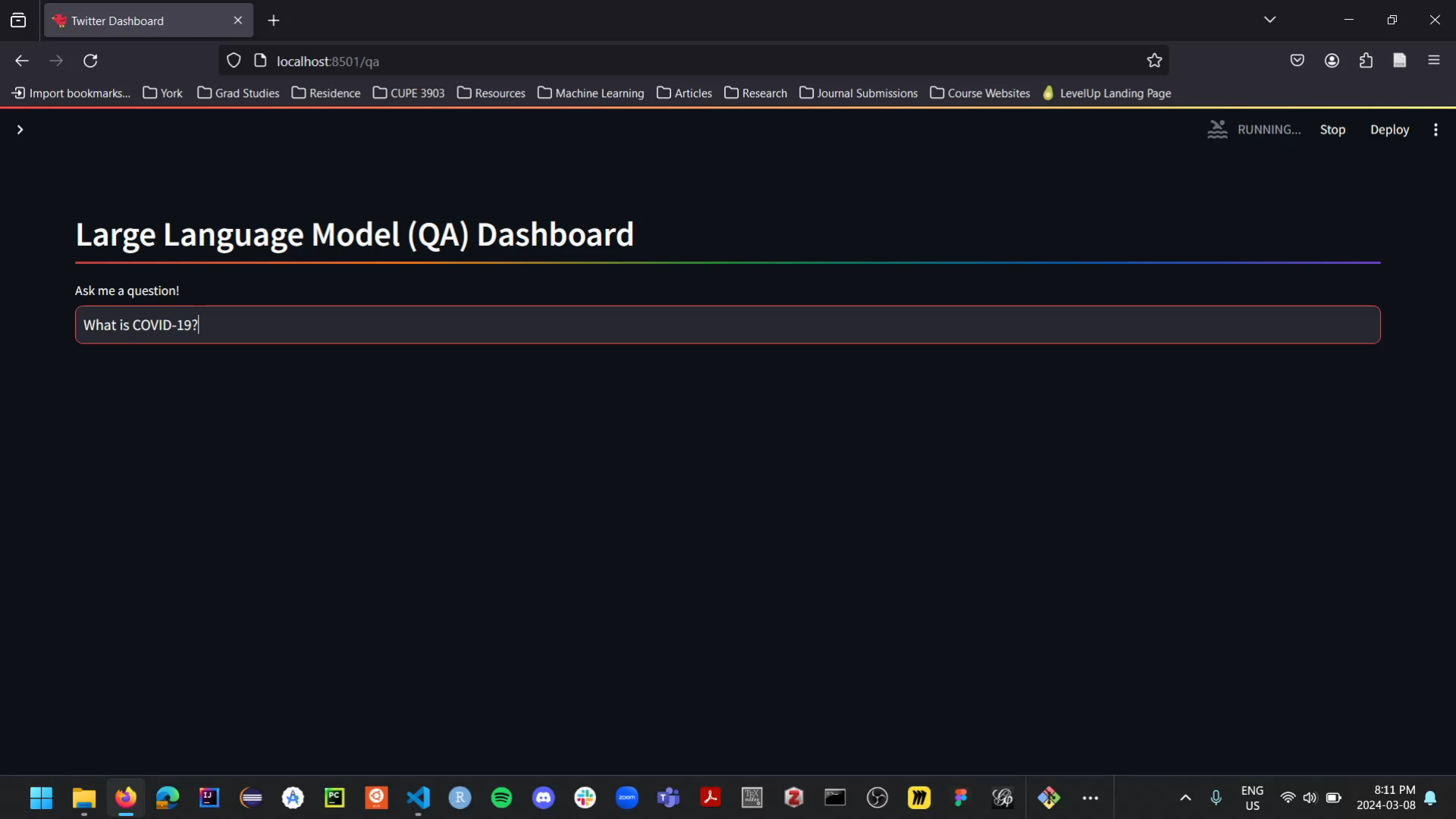
Please select a month

February

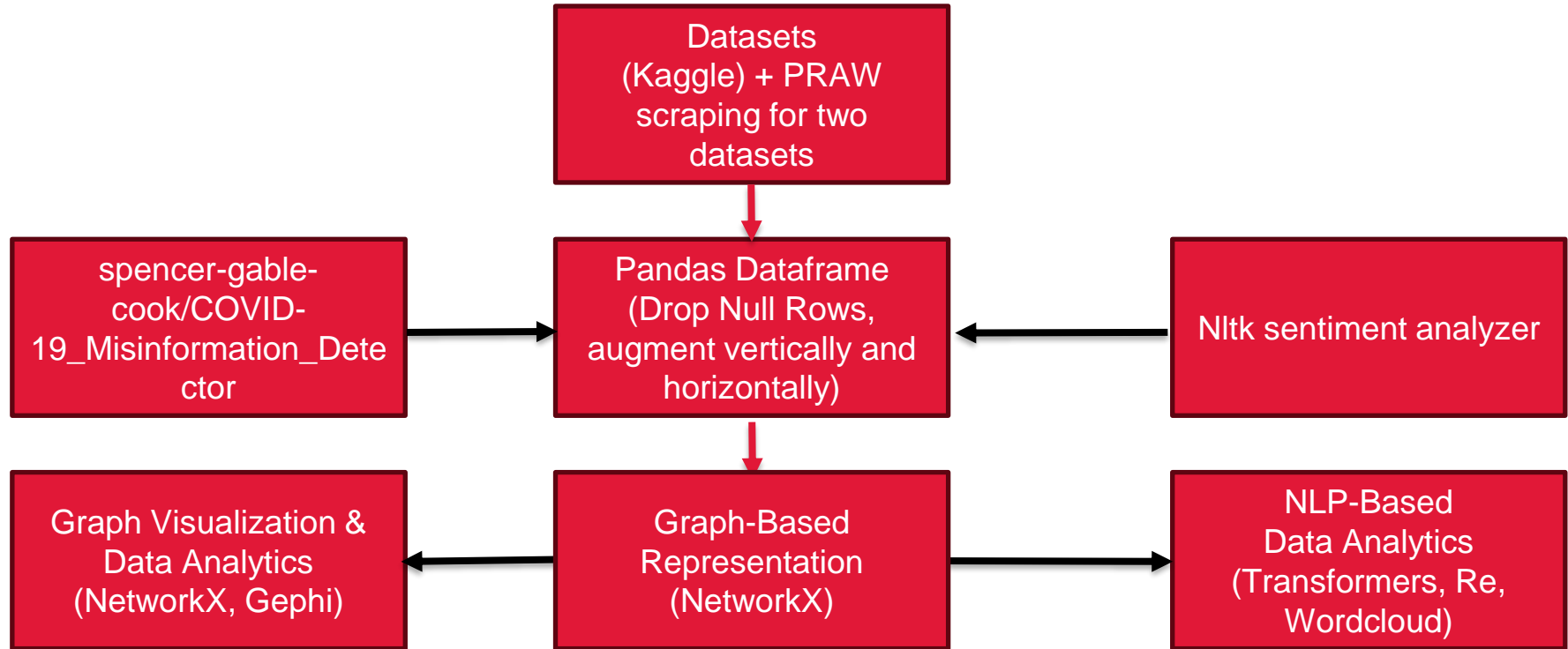# New Contribution: QA Large Language Model

# Large Language Model (QA) Dashboard

Ask me a question!
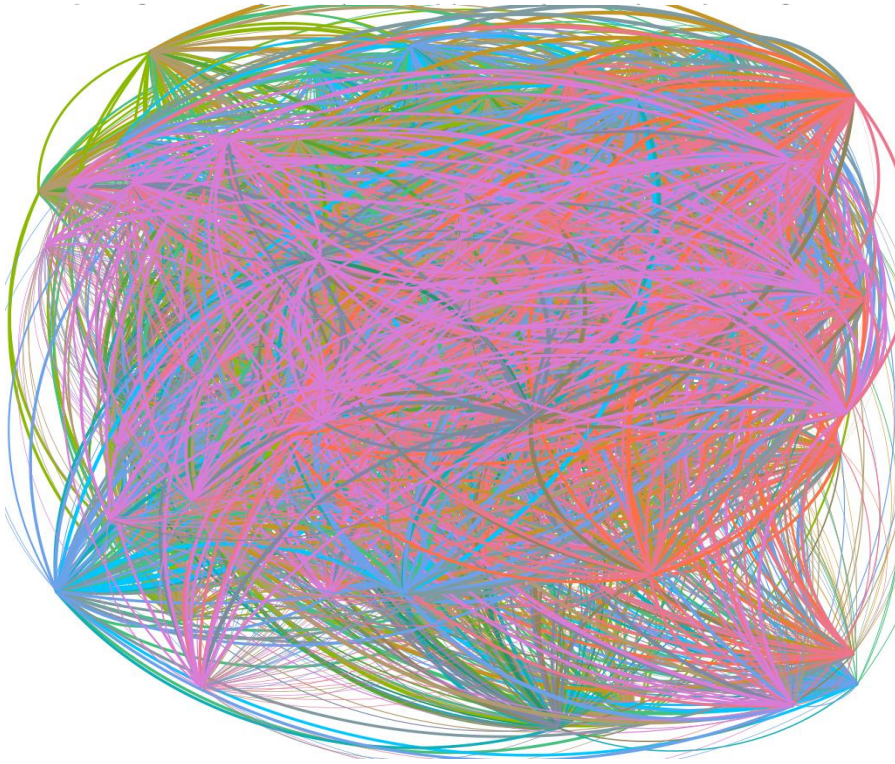
What is COVID-19?

# Methodology (Reddit)

# Results and Analysis (Reddit)



| PageRank | Top 10 PageRank Users in the Reddit Graph |
|---|---|
| | User throwaway742858 has score contribution: 0.02105 |
| | User BipolarRN has score contribution: 0.02105 |
| | User zombiephish has score contribution: 0.02073 |
| | User jdfertig has score contribution: 0.02035 |
| | User ghafgarionbaconsmith has score contribution: 0.01824 |
| | User TheDarkKitten95 has score contribution: 0.01764 |
| | User gebba54 has score contribution: 0.01739 |
| | User littlestircrazy has score contribution: 0.01721 |
| | User captaindata1701 has score contribution: 0.01721 |
| | User volcomp has score contribution: 0.0172 |
| Hubs | Top 10 Authority Users in the Reddit Graph |
| | User BipolarRN has score contribution: 0.021 |
| | User throwaway742858 has score contribution: 0.0207 |
| | User zombiephish has score contribution: 0.02031 |
| | User jdfertig has score contribution: 0.0203 |
| | User ghafgarionbaconsmith has score contribution: 0.01827 |
| | User Kerboq has score contribution: 0.01743 |
| | User gebba54 has score contribution: 0.01733 |
| | User volcomp has score contribution: 0.01677 |
| | User TheDarkKitten95 has score contribution: 0.01676 |
| | User virgilash has score contribution: 0.01638 |
| Authorities | Top 10 Hub Users in the Reddit Graph |
| | User BipolarRN has score contribution: 0.021 |
| | User throwaway742858 has score contribution: 0.0207 |
| | User zombiephish has score contribution: 0.02031 |
| | User jdfertig has score contribution: 0.0203 |
| | User ghafgarionbaconsmith has score contribution: 01827 |
| | User Kerboq has score contribution: 0.01743 |
| | User gebba54 has score contribution: 0.01733 |
| | User volcomp has score contribution: 0.01677 |
| | User TheDarkKitten95 has score contribution: 0.01676 |
| | User virgilash has score contribution: 0.01638 |

YORK U

# Conclusion and Analysis (1)

# Conclusion and Analysis (2)

- Differences between **Reddit** and **Twitter**:
  - Reddit seemed to have much more neutrality and positivity in their user-content in comparison with Twitter
  - Node importance (i.e., PageRank, HITS) in Twitter shows great **disparity** between user influence while Reddit is much more **equitable**.
  - Reddit has more moderation with stricter rules, coupled with the ability to quarantine communities
- Limitations:
  - Quarantined subreddits have **small network diameter** and high user **overlap** between each other.
  - Twitter Dataset is a **premature snapshot** of the misinformation network; more longitudinal evaluations are required.

YORK U