

NEW YORK TIMES Archives Data Analysis and Visualization using Prefect, DuckDB, and Streamlit

**By Katie Hutchinson
DS 3022: Data Engineering**



Project Summary

This project is an end-to-end data pipeline that analyzes 10 years of New York Times headlines from 2015–2025 in over 400,000 records to uncover shifts in public discourse.

By processing unstructured text, the system identifies how dominant cultural topics, like "COVID-19," "Elections," and "AI," have risen and fallen over the decade.



Technologies Used

- ★ **Prefect:** manages the workflow, ensuring resilient data processing and error handling.
- ★ **Pandas:** for efficient data manipulation.
- ★ **BERTopic:** clusters semantic meaning from headlines to discover topics
- ★ **DuckDB:** provides a high-performance, serverless SQL database optimized for analytical queries
- ★ **Streamlit and Plotly:** power an interactive frontend dashboard for real-time trend exploration.
- ★ **Conda/Mamba:** Manages the project's virtual environment and ensures reproducible dependency handling.



Project Steps

1. **Ingestion.py:** Obtain & Process Data from New York Times Archives from 2015 - 2025 using API Key
2. **Pipeline.py:** Transform raw text into structured insights
 - a. Load the raw CSV into a Pandas DataFrame, verifying file integrity.
 - b. Trained the BERTopic model on a randomized subset (20,000 articles).
 - c. The trained model then predicts and assigns a Topic ID to the remaining 380,000+ articles.
3. **App.py:** Visualize the data using Plotly
 - a. Using Streamlit, create the app with a drop-down menu and trend chart
 - b. Use Plotly to plot individual lines
 - c. View the key insights for each topic

Step 1: Obtain & Process Data from Archives

```
import requests
import pandas as pd
import time
from datetime import datetime
import os

# --- CONFIGURATION ---
API_KEY = "MKyesAaNHHea0smCGlNCQ9x45gvfxLpr"
BASE_URL = "https://api.nytimes.com/svc/archive/v1/{year}/{month}.json"

# Fetches the entire archive for a specific month
def fetch_month_archive(year, month):

    url = BASE_URL.format(year=year, month=month)
    params = {'api-key': API_KEY}

    print(f"Fetching Archive: {year}-{month:02d} ...")

    try:
        response = requests.get(url, params=params)

        if response.status_code == 429:
            print("Rate limit hit. Waiting 60 seconds...")
            time.sleep(60)
            response = requests.get(url, params=params)

        response.raise_for_status()
        data = response.json()

        docs = data['response']['docs']
        print(f" -> Found {len(docs)} articles.")
        return docs

    except Exception as e:
        print(f"Error fetching {year}-{month}: {e}")
        return []
```

First, I fetched all of the archives from the year 2020-2025 from the NYT Archives API. In order to comply with the request limit, I slept 60 seconds when the limit was hit.

Next, I cleaned the data by filtering for just news articles. Then I added each cleaned line with the following data (date, headline, etc.) to a list of dictionaries.

```
def clean_archive_data(docs):
    cleaned = []
    for doc in docs:
        if doc.get('type_of_material') not in ['News', 'Article', 'Review', 'Op-Ed']: # Filter for just "News" type material
            continue

        cleaned.append({
            'date': doc.get('pub_date'),
            'headline': doc.get('headline', {}).get('main', ''),
            'abstract': doc.get('abstract'),
            'keywords': [k['value'] for k in doc.get('keywords', [])],
            'section': doc.get('section_name')
        })

    return cleaned
```

Step 1: Main Method

```
if __name__ == "__main__":
    all_data = []

    years = range(2020, 2025)

    now = datetime.now()
    current_year = now.year
    current_month = now.month

    for year in years:
        for month in range(1, 13):
            if year > current_year or (year == current_year and month > current_month):
                break

            raw_docs = fetch_month_archive(year, month)
            clean_docs = clean_archive_data(raw_docs)
            all_data.extend(clean_docs)

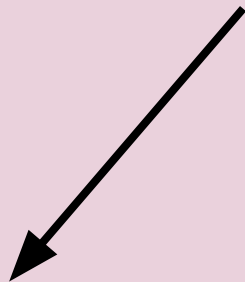
            print(f" -> Total collected so far: {len(all_data)}")

            time.sleep(5)

    df = pd.DataFrame(all_data)
    filename = "nyt_2020_2024.csv"
    df.to_csv(filename, index=False)

    print(f"Done! Saved {len(df)} articles to {filename}")
```

To finish the ingestion script, I looped through the years 2020-2025 to fetch and clean the archival data, sleeping 5 seconds between each month. I saved all this data to a csv file.



Step 2: Pipeline

```
@task
def get_data():
    if os.path.exists(CSV_FILE):
        return pd.read_csv(CSV_FILE)
    else:
        raise FileNotFoundError(f"Could not find {CSV_FILE}. Please run ingestion.py first")
```

First, I obtained the CSV data through pandas read function.

- ☆ Error checked the data for empty headlines & abstracts, replacing them with empty strings.
- ☆ Combine the headlines & abstracts to help train the model.
- ☆ Randomized 20,000 articles to train the model, requiring at least 150 articles per topic.
- ☆ Sort all the articles into the bucket topics in chunks of 5000.
- ☆ Transform: Attach each headline to their topic.
 - ☆ Convert the topics to a dictionary
 - ☆ Match each topic ID to the dictionary

```
def ingest_and_model(df):
    # error check for empty headlines / abstracts, fill with empty string
    df['headline'] = df['headline'].fillna('')
    df['abstract'] = df['abstract'].fillna('')

    # combine headline and abstract for more context to train model
    docs = (df['headline'] + " " + df['abstract']).tolist()

    # pick 20,000 random articles to train model
    from sklearn.utils import shuffle
    training_docs = shuffle(docs, random_state=42)[:20000]

    # a topic must have at least 150 articles
    topic_model = BERTopic(min_topic_size=300, verbose=True)
    topic_model.fit(training_docs)

    # transform in chunks of 5000
    chunk_size = 5000
    all_topics = []

    for i in range(0, len(docs), chunk_size):
        chunk = docs[i : i + chunk_size]
        # transform chunk
        chunk_topics, _ = topic_model.transform(chunk)
        all_topics.extend(chunk_topics)

    # attach headline to topic
    df['topic_id'] = all_topics

    # get the table of topics / topic ID and convert it to a dictionary
    topic_info = topic_model.get_topic_info()
    topic_map = topic_info.set_index('Topic')['Name'].to_dict()

    # match each topic id to the dictionary
    df['topic_name'] = df['topic_id'].map(topic_map)

    return df
```


Step 2: Pipeline continued

Connect to Duck DB and create a SQL table. This step transforms the 400,000+ processed records into a query-optimized SQL table

```
@task
def load_to_duckdb(df):
    conn = duckdb.connect(DB_FILE)
    conn.execute("CREATE OR REPLACE TABLE processed_articles AS SELECT * FROM df")
    conn.close()
    print("Loaded to DuckDB")
```

```
@flow(name="NYT-Pipeline")
def main_flow():
    raw_df = get_data()

    final_df = transform_and_model(raw_df)

    load_to_duckdb(final_df)

if __name__ == "__main__":
    main_flow()
```

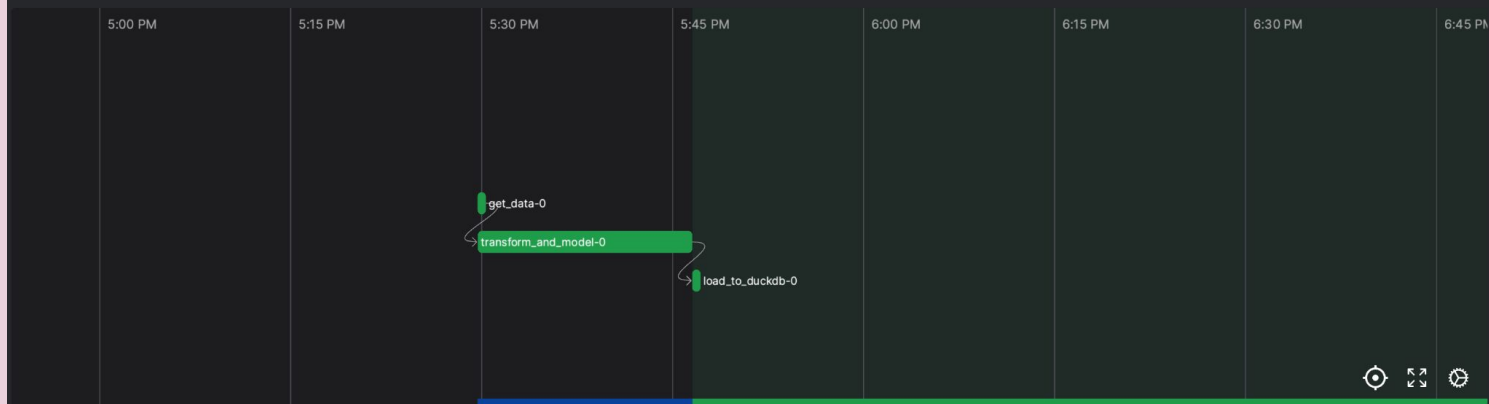
This function defines the pipeline's DAG. It is responsible for dependency management, ensuring that tasks execute in the strict order of Extraction → Transformation → Loading

A successful
Prefect flow
run!

Flow Runs / shiny-muskox

Completed 2025/12/10 05:29:42 PM 16m 54s 3 Task runs

Flow NYT-Pipeline



Logs Task Runs Subflow Runs Results Artifacts Details Parameters

Level: all

Oldest to newest

Dec 10th, 2025

INFO

Created task run 'get_data-0' for task 'get_data'

05:29:42 PM
prefect.flow_runs

INFO

Executing 'get_data-0' immediately...

05:29:42 PM
prefect.flow_runs

INFO

Finished in state Completed()

05:29:44 PM
get_data-0
prefect.task_runs

```
D SELECT topic_name, count(*) FROM processed_articles WHERE topic_id != -1 GROUP BY 1 ORDER BY 2 DESC LIMIT 10;
```

topic_name varchar	count_star() int64
0_trump_president_donald_the	34054
3_restaurant_food_chef_and	13904
2_novel_books_book_and	12293
1_police_shooting_man_was	12110
5_coronavirus_vaccine_covid_virus	11441
6_china_korea_north_chinese	10383
4_art_museum_artist_of	10252
7_met_couple_bride_groom	8976
8_fashion_designer_week_paris	7744
9_ukraine_russia_russian_putin	7623
10 rows	2 columns

The results of
the mapping
produced in
DuckDB!

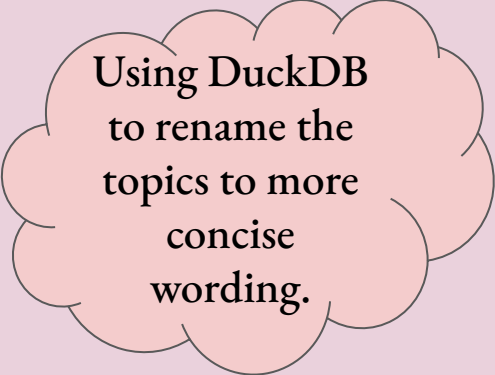
```
D SELECT * FROM processed_articles WHERE topic_name = '0_trump_president_donald_the' LIMIT 5;
```

date varchar	headline varchar	...	topic_id int64	topic_name varchar
2015-01-01T01:36:1...	Much of David Duke...	...	0	0_trump_president_...
2015-01-01T22:15:4...	Teach Congress a L...	...	0	0_trump_president_...
2015-01-01T22:43:2...	Jeb Bush Resigns F...	...	0	0_trump_president_...
2015-01-02T12:35:1...	Jeb Bush, and Two	0	0_trump_president_...
2015-01-02T16:38:1...	'To Make Men Free,...	...	0	0_trump_president_...
5 rows	9 columns (4 shown)			

```

D UPDATE processed_articles SET topic_name = 'Climate/Energy' WHERE topic_id = 19;
D UPDATE processed_articles SET topic_name = 'General Housing' WHERE topic_id = 20;
D UPDATE processed_articles SET topic_name = 'Israel/Palestine' WHERE topic_id = 21;
D UPDATE processed_articles SET topic_name = 'Macroeconomics' WHERE topic_id = 22;
D UPDATE processed_articles SET topic_name = 'Pop Music/Albums' WHERE topic_id = 23;
D UPDATE processed_articles SET topic_name = 'Puzzles/Games' WHERE topic_id = 24;
D UPDATE processed_articles SET topic_name = 'Dance/Ballet' WHERE topic_id = 25;
D UPDATE processed_articles SET topic_name = 'Classical Music' WHERE topic_id = 26;
D UPDATE processed_articles SET topic_name = 'Olympics' WHERE topic_id = 27;
D UPDATE processed_articles SET topic_name = 'Legislation/Tax' WHERE topic_id = 28;
D UPDATE processed_articles SET topic_name = 'NBA/Basketball' WHERE topic_id = 29;
D UPDATE processed_articles SET topic_name = 'Social Issues' WHERE topic_id = 30;
D UPDATE processed_articles SET topic_name = 'NFL/Football' WHERE topic_id = 31;
D UPDATE processed_articles SET topic_name = 'Brexit/UK Politics' WHERE topic_id = 32;
D UPDATE processed_articles SET topic_name = 'Autos/EVs' WHERE topic_id = 33;
D UPDATE processed_articles SET topic_name = 'Tennis' WHERE topic_id = 34;
D UPDATE processed_articles SET topic_name = 'NYT Internal Notes' WHERE topic_id = 35;
D UPDATE processed_articles SET topic_name = 'Soccer/World Cup' WHERE topic_id = 36;
D UPDATE processed_articles SET topic_name = 'NYC Local Politics' WHERE topic_id = 37;
D UPDATE processed_articles SET topic_name = 'Corrections' WHERE topic_id = 38;
D UPDATE processed_articles SET topic_name = 'Health/Disease' WHERE topic_id = 39;
D UPDATE processed_articles SET topic_name = 'Awards Shows' WHERE topic_id = 40;
D UPDATE processed_articles SET topic_name = 'Spelling Bee' WHERE topic_id = 41;
D UPDATE processed_articles SET topic_name = 'Social Media' WHERE topic_id = 42;
D UPDATE processed_articles SET topic_name = 'Metropolitan Diary' WHERE topic_id = 43;
D UPDATE processed_articles SET topic_name = 'Vatican/Pope' WHERE topic_id = 44;
D UPDATE processed_articles SET topic_name = 'College Football' WHERE topic_id = 45;
D UPDATE processed_articles SET topic_name = 'NYC Transit' WHERE topic_id = 46;
D UPDATE processed_articles SET topic_name = 'Real Estate Listings' WHERE topic_id = 47;
D UPDATE processed_articles SET topic_name = 'Golf' WHERE topic_id = 48;
D UPDATE processed_articles SET topic_name = 'France/Macron' WHERE topic_id = 49;
D UPDATE processed_articles SET topic_name = 'Parenting' WHERE topic_id = 50;
D UPDATE processed_articles SET topic_name = 'US-Mexico Border' WHERE topic_id = 51;
D UPDATE processed_articles SET topic_name = 'European Migration' WHERE topic_id = 52;

```



Using DuckDB
to rename the
topics to more
concise
wording.

Step 3: Building the App using StreamLit

```
@st.cache_data # data only loaded once
def load_data():

    try:
        conn = duckdb.connect(DB_FILE, read_only=True)
        # count the number of articles per topic, grouped by month
        query = f"""
        SELECT
            topic_name,
            date_trunc('month', CAST(date AS TIMESTAMP)) as month,
            count(*) as count
        FROM processed_articles
        WHERE topic_id != -1
        GROUP BY 1, 2
        ORDER BY 2
        """
        df = conn.execute(query).df()
        conn.close()
        return df
    except Exception as e:
        st.error(f"Error loading DuckDB data. Have you run pipeline.py yet? Error: {e}")
        return pd.DataFrame()
```

Streamlit is an open-source Python framework that allows data scientists and machine learning engineers to create interactive web applications and dashboards using only Python

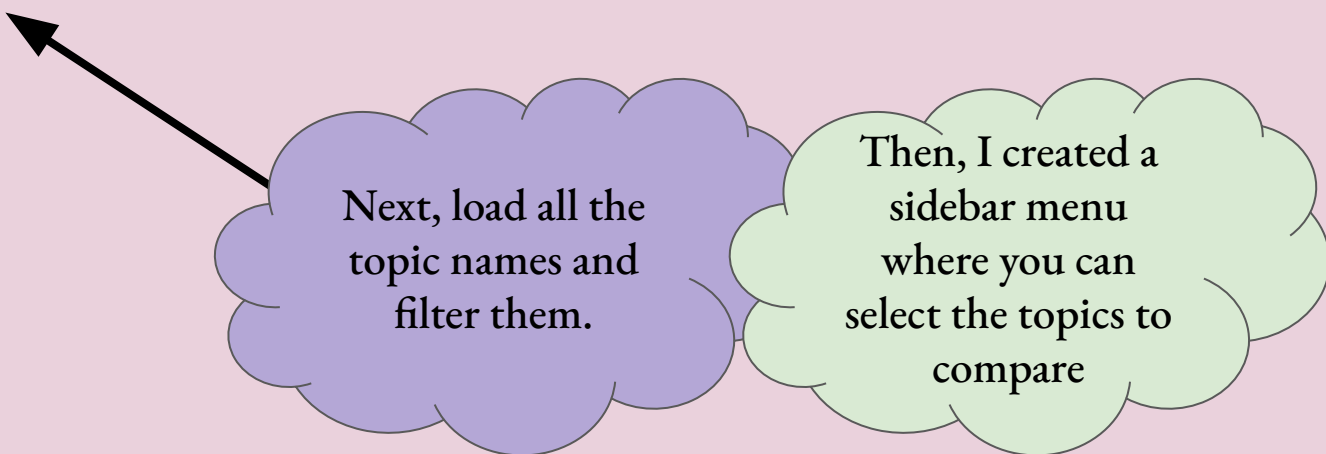
First step: load all of the data from the DuckDB file


```
if not df.empty:
    st.sidebar.header("Filter Trends")

    df['topic_name'] = df['topic_name'].fillna('Unclassified').astype(str)
    df['display_topic'] = df['topic_name']
    all_display_topics = sorted(df['display_topic'].unique())

    initial_filter = [t for t in all_display_topics if t not in ['Noise Random Stuff', 'Abstract', 'Main']] # filter

# allow user to select multiple topics
selected_topics = st.sidebar.multiselect(
    "Select Topics to Compare",
    all_display_topics,
    default=initial_filter[:0])
```



Next, load all the
topic names and
filter them.

Then, I created a
sidebar menu
where you can
select the topics to
compare

For the chart, for each selected topic, Plotly draws a line with the topic frequency by counting the number of articles under that topic for each month

```
if selected_topics:
    filtered_df = df[df['display_topic'].isin(selected_topics)] # filter based on user se

    # Create interactive line chart with Plotly
    fig = px.line(
        filtered_df,
        x='month',
        y='count',
        color='display_topic',
        title="Topic Frequency Over Time (2015-2024)",
        labels={'count': 'Number of Articles', 'month': 'Date', 'display_topic': 'Topic'},
        template="plotly_white",
        height=600
    )

    # Customize chart appearance
    fig.update_xaxes(dtick="M12", tickformat="%Y") # Show labels annually
    fig.update_traces(mode="lines", hovertemplate="%{y} articles<br>%{x}|%b %Y}")
    fig.update_layout(hovermode="x unified")

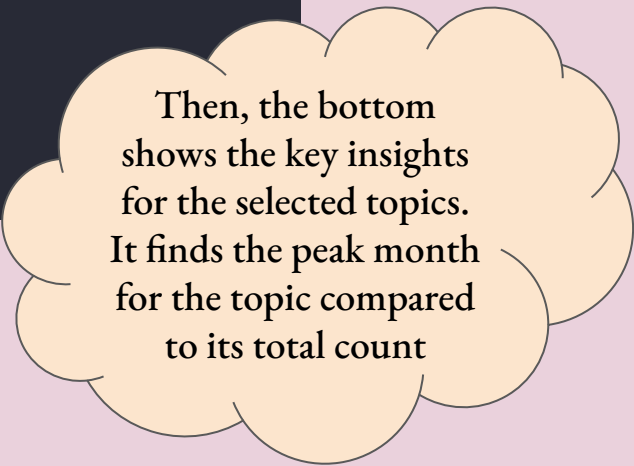
    st.plotly_chart(fig, use_container_width=True)
```

```
# INSIGHTS
st.subheader("Key Insights")
cols = st.columns(3)

for idx, topic in enumerate(selected_topics):
    topic_data = filtered_df[filtered_df['display_topic'] == topic]

    # Find the peak month and total count for the topic
    total_articles = topic_data['count'].sum()
    peak_row = topic_data.sort_values('count', ascending=False).iloc[0]
    peak_month_str = peak_row['month'].strftime('%b %Y')
    peak_count = peak_row['count']

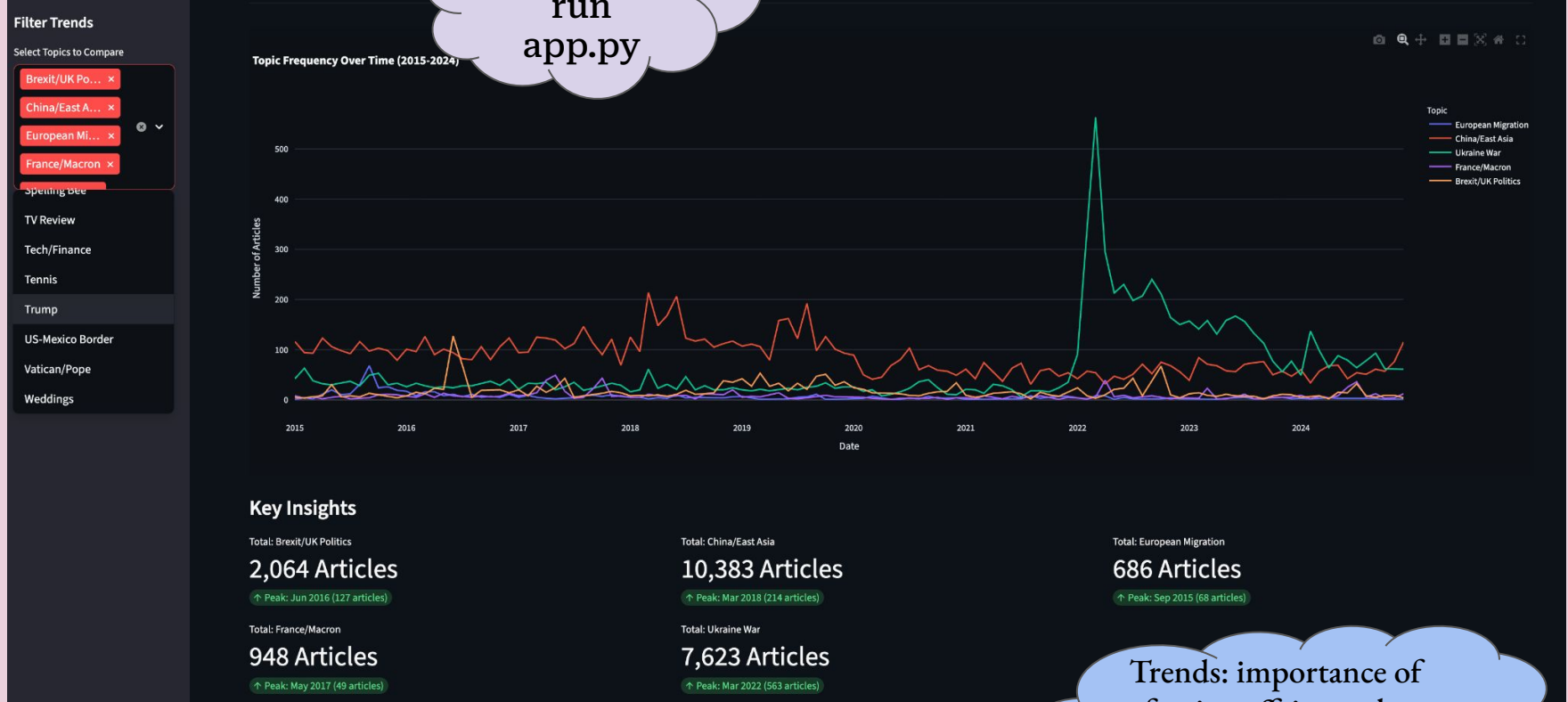
    with cols[idx % 3]:
        st.metric(
            label=f"Total: {topic}",
            value=f"{total_articles:,} Articles",
            delta=f"Peak: {peak_month_str} ({peak_count:,} articles)"
        )
```



Then, the bottom shows the key insights for the selected topics. It finds the peak month for the topic compared to its total count

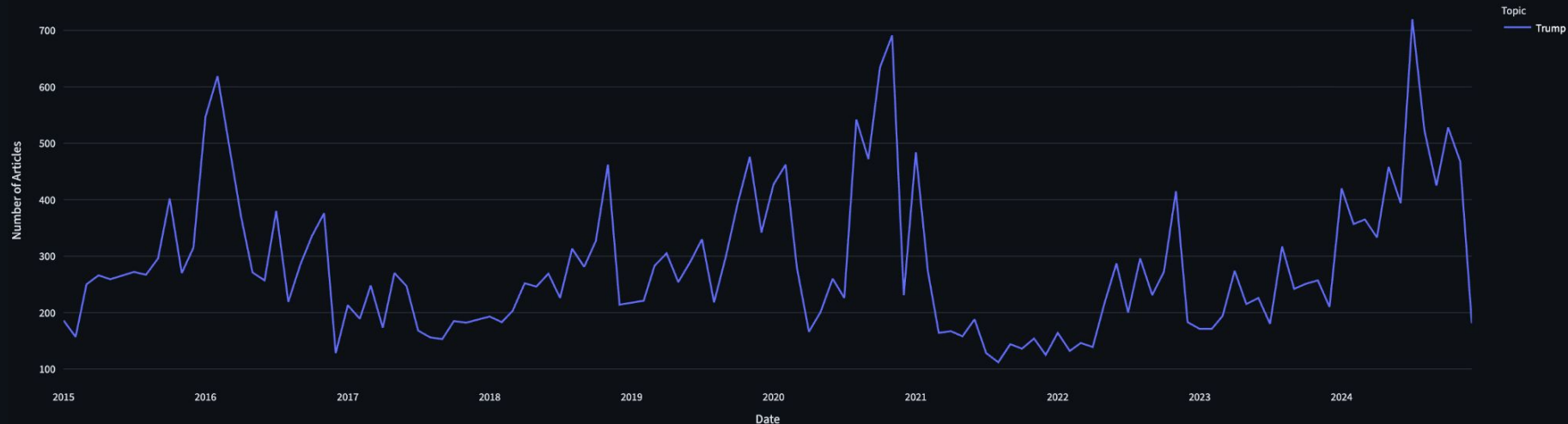
Running the App:

Command:
streamlit
run
app.py



Trends: importance of foreign affairs to the NYT and American citizens

Topic Frequency Over Time (2015-2024)



Key Insights

Total: Trump

34,054 Articles

↑ Peak: Jul 2024 (720 articles)

Trends: election cycles,
legislation, major events (e.g.
attempted assassination)