# Multilingual narrative tracking in the news - real-time experiments

Unstructured Data Meetup Berlin

2024-05-07

Robert Caulk, PhD & Elin Törnquist, PhD

*CEO*                                   *Director of Transparency*

Team:

Timothy Pogue, Wagner Costa Santos, Emre Suzen

Emergent Methods

# Our background

**Engineers and Researchers** committed to FOSS

- **Applying** AI to real-time adaptive modeling challenges
- **Scaling** software in all directions
- **Enriching** data for other businesses
- **Performing** research

## AskNews

News context engineering
https://asknews.app

## flowdapt

Real-time cluster orchestration
https://flowdapt.ai

FOSS 🤗

## FreqAI

AI/ML for algo-trading
https://www.freqtrade.io/en/stable/freqai/

FOSS 🤗

## MELISSA

Large-scale deep-learning for supercomputers
https://melissa.gitlabpages.inria.fr/melissa/

FOSS 🤗

## DATASIEVE

Data pipelining
https://github.com/emergentmethods/datasieve

FOSS 🤗

## Manifest
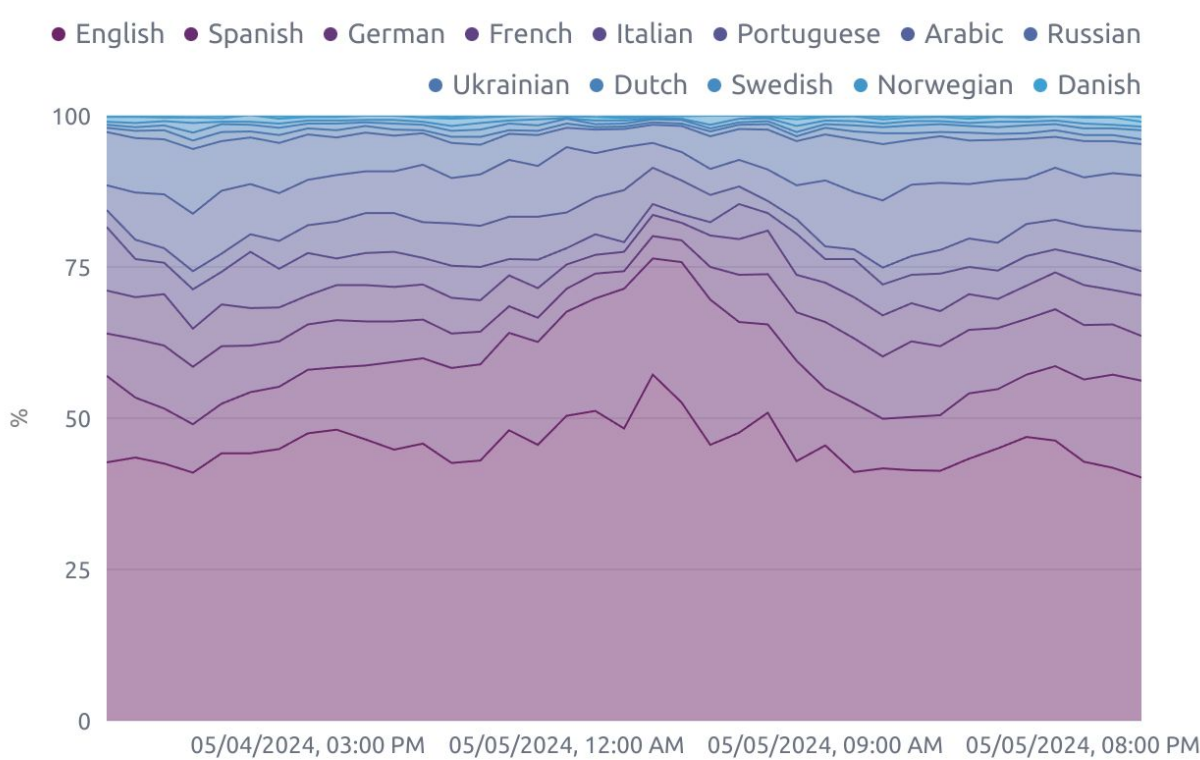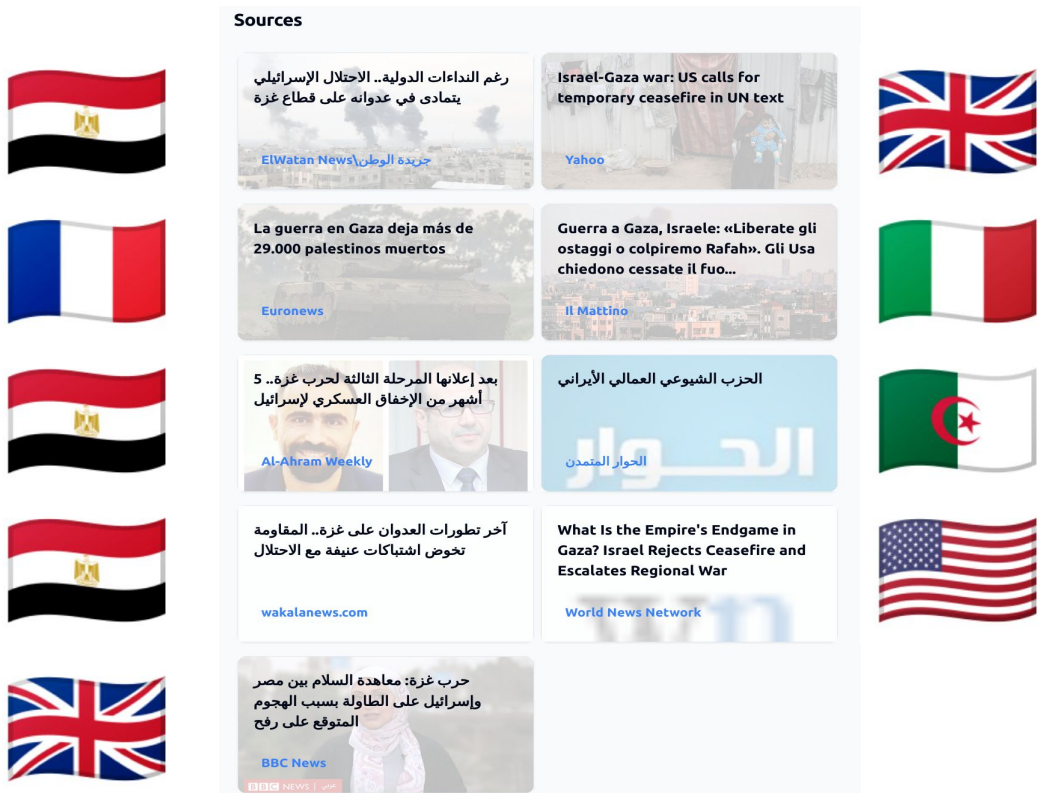
Modern configuration
https://github.com/emergentmethods/python-manifest

FOSS 🤗

# Why do we need to engineer news context?

# Motivations for engineering news context

- Enforcing journalistic standards 📰
  - Stating claims with supporting **evidence** and **attribution**
  - AP style-guidelines and formatting
- Enforcing source and language diversity for democratized representation 🇪🇺
  - Representing **diverse perspectives** on global issues
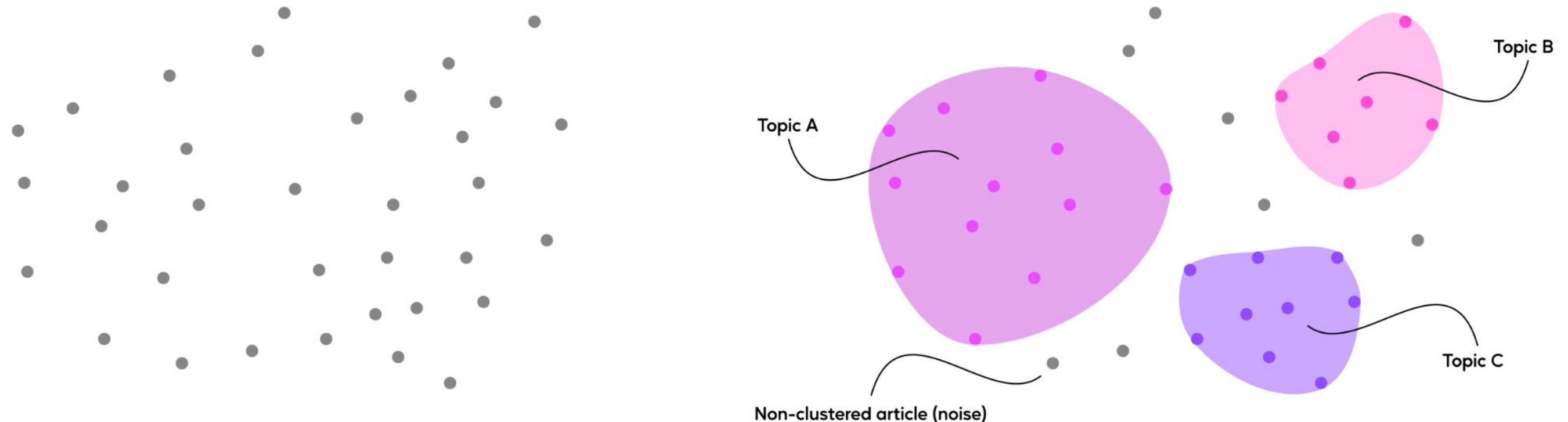
# Motivations for engineering news context

- Enforcing journalistic standards 📰
  - Stating claims with supporting **evidence** and **attribution**
  - AP style-guidelines and formatting

- Enforcing source and language diversity for democratized representation 🇪🇺
  - Representing **diverse perspectives** on global issues

- Avoiding stale/outdated reporting 🏇
  - Missing the latest news can cause **disinformation**, customer dissatisfaction, and **confusion**

- Minimizing hallucination 🤭
  - The cost of hallucination is too high, **a team of researchers is required**

- Scaling democratized news context to companies that are not interested in the logistics of tracking and diversifying 1 million articles per day 🚀

# Engineering the parameter space
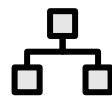
# Defining the objective

- A clean and well defined parameter space:

  - enables **clustering** of news topics across **diverse perspectives**

  - represents **entities**, especially those originating from small demographics

  - "normalizes" for **language differences**

Topic A

Topic B

Topic C

Non-clustered article (noise)

# Preparing the embedding

- LLM   **Enrich articles**

  - Translate

  - Summarize

  - Extract keywords, classification, and sentiment

Flexibility -> adaptability -> product opportunities



flowdapt workflow - **Summarize and embed**
1000 articles every 5 minutes

**Fetch news articles**
Structured news data

Title 1   Title 2   x1000   Title 1000
. . .

**Enrich articles**

On-premise LLMs

- Fine-tuned LM (GLiNER-news)   **Enrich articles**

  - Generalist and lightweight entity extraction
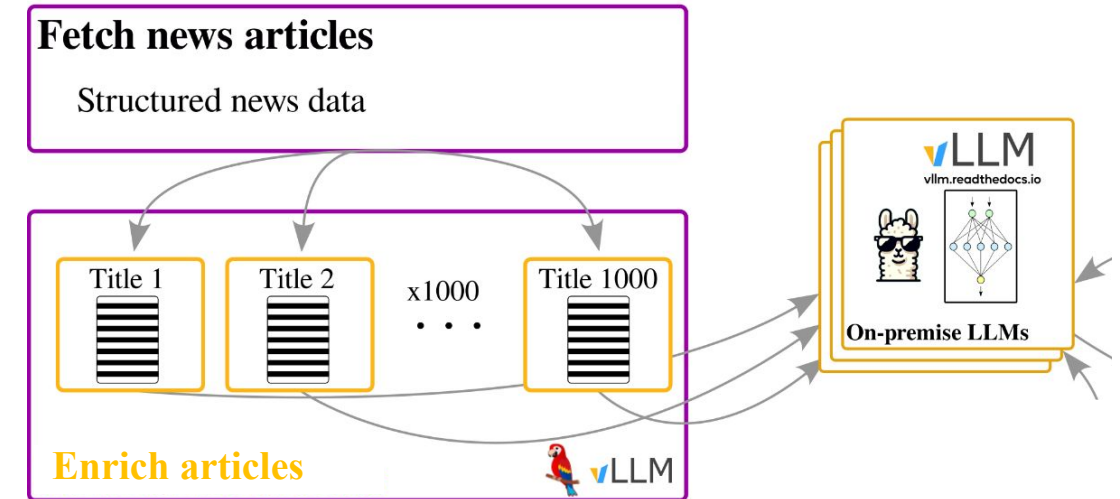
  - Based model by Zaratiana et al. (2023)

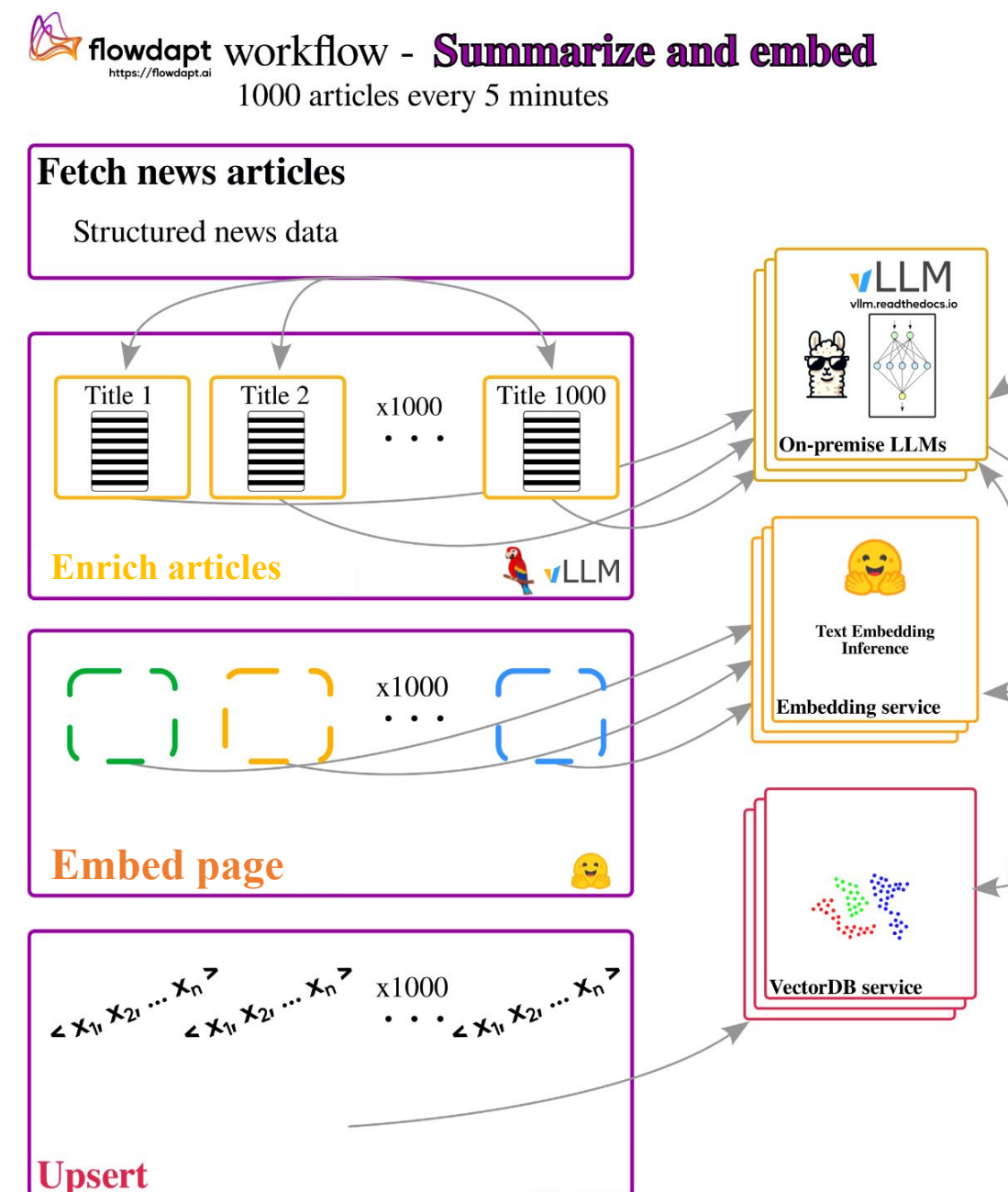Get your own lightsaber for Star Wars Day on May 4th!

product   event   date

# Embedding and storing the page

- ## Build and embed page  **Build/Embed page**

  - ### Model choice affects everything:

    - Retrieval speed/quality

    - Storage costs

    - Clustering compute costs

  - ### Doc structure vs expected Query structure

- ## VectorDB  **Upsert**

  - ### The beating heart of the architecture:

    - Robustness

    - Parallelizability

    - Metadata filtering flexibility

    - Quantization

    - Performance



flowdapt workflow - **Summarize and embed**
1000 articles every 5 minutes

# Tracking narratives in our parameter space

# What is a news narrative?

- 📜 A series of related news reports

- 👬👩‍🤝‍👨 Multiple points of view

- ❌ Errors - accidental and purposeful

# Finding related news reports

## Parameter space

Topic identification through HDBSCAN

Topic B

Topic A

Topic C

Non-clustered article (noise)

Time window 0

Characterized by our embedded enrichments

Clusters of semantic similarity

Multiple countries
Multiple languages
Multiple sources
**Competing perspectives**

# Identifying the series of events

## Parameter space

Topic identification through HDBSCAN

Topic A

Topic B

Topic C

Non-clustered article (noise)

Time window 0

## Temporal clustering

Clutering at time interval Δt

Time window 0

Time window 1

## Cluster connection methods

- Adaptively (re)train one binary classifier per cluster
- Track medoid/centroid drift with hyper spheres
- Use overlap clustering techniques

# Long range niche tracking

May 1, 2024

May 6, 2024

**Madonna's Copacabana Beach Concert Set to Be a Historic Event with Over 1.5 Million Fans**

South America

Updated 01/05/2024, 18:48 (6 days ago)



Image prompt

Toggle citations

**Madonna's Rio Concert Shatters Records with 1.6 Million Fans, Ignites Local Economy**

South America

Updated 06/05/2024, 02:49 (1 day ago)



Image prompt

Toggle citations

# Long range niche tracking

# Context polishing (enforcing diversity)

- Pruning a single cluster
  - diversity enforcement
- Confirming continuity
- Reranking the cluster

Prompt engineering

- Document formatting
- Citation control
- Journalistic guidelines

Topic A

</doc>🇩🇪</doc>

</doc>🇫🇷</doc>

</doc>🇺🇸</doc>

</doc>🇺🇦</doc>

</doc>🇷🇺</doc>

# Report cluster alignment (quantifying diversity)

Topic A

- Identify alignment and contradictions
- Compute confidence levels

</doc>🇩🇪</doc>

</doc>🇫🇷</doc>

</doc>🇺🇸</doc>

</doc>🇺🇦</doc>

</doc>🇷🇺</doc>

# Tracking the death of Alexei Navalny

🇺🇸 🇫🇷 🇷🇺 comparison

Percentage of total country news coverage devoted to Navalny narrative



■ 0.14% of Russian News Coverage   ■ 0.53% of French News Coverage   ■ 0.29% of U.S. News Coverage

🇷🇺   🇫🇷   🇺🇸

Emergent Methods
www.emergentmethods.ai

# comparison

Emergent Methods
www.emergentmethods.ai

## Russian Opposition Leader Alexei Navalny Dies in Prison

- Russian News Coverage
- French News Coverage
- U.S. News Coverage

**% of Total News from Country**

Timeline values:
- Feb. 16, 12 p.m.: 0.13, 0.65, 0.35, 0.95
- Feb. 16, 4 p.m.: 0.17, 0.46, 0.54
- Feb. 16, 8 p.m.: 0.61, 0.27
- Feb. 16, 12 a.m.: 0.58, 1.93
- Feb. 17, 4 a.m.: 0.39
- Feb. 17, 8 a.m.: 0.54, 0.52
- Feb. 17, 12 p.m.: 0.46, 0.51
- Feb. 17, 4 p.m.: 0.16, 0.22, 0.25
- Feb. 17, 8 p.m.: 0.2, 0.12
- Feb. 17, 12 a.m.: 0.11, 0.13
- Feb. 18, 8 a.m.: 0.08, 0.11, 0.48
- Feb. 18, 12 p.m.: 0.07, 0.27, 0.51
- Feb. 18, 4 p.m.: 0.08, 0.15, 0.37
- Feb. 18, 8 p.m.: 0.27, 0.41
- Feb. 18, 12 a.m.: 0.17, 1
- Feb. 19, 4 a.m.: 0.14, 0.63
- Feb. 19, 8 a.m.: 0.19, 1.14
- Feb. 19, 12 p.m.: 0.08, 0.21, 0.36
- Feb. 19, 4 p.m.: 0.05, 0.27, 0.39
- Feb. 19, 8 p.m.: 0.11, 0.75
- Feb. 20, 4 a.m.: 0.22, 0.23, 0.81
- Feb. 20, 8 a.m.: 0.43, 0.25
- Feb. 20, 12 p.m.: 1.21, 0.29, 0.36

### Event annotations

**Death of Alexei Navalny in Prison Shakes Russian Opposition and Sparks Global Outrage**

**Biden and Harris Accuse Putin of Responsibility for Navalny's Death Amid Global Condemnation**

**Global Condemnation and Protests Follow Alexei Navalny's Death in Russian Custody**

**Search for Alexei Navalny's Body Continues Amidst Global Outcry and Protests in Russia**

**Global Outcry as Kremlin Withholds Navalny's Body Amid Suspicions of Foul Play**

**Russia Detains Over 400 Amid Mourning for Opposition Leader Navalny**

**Detentions Surge Across Russia Amid Mourning for Navalny, Family Accuses Authorities of Hiding Body**

**Arrests Surge in Russia as Global Leaders Denounce Navalnys Death, Biden Blames Putin**

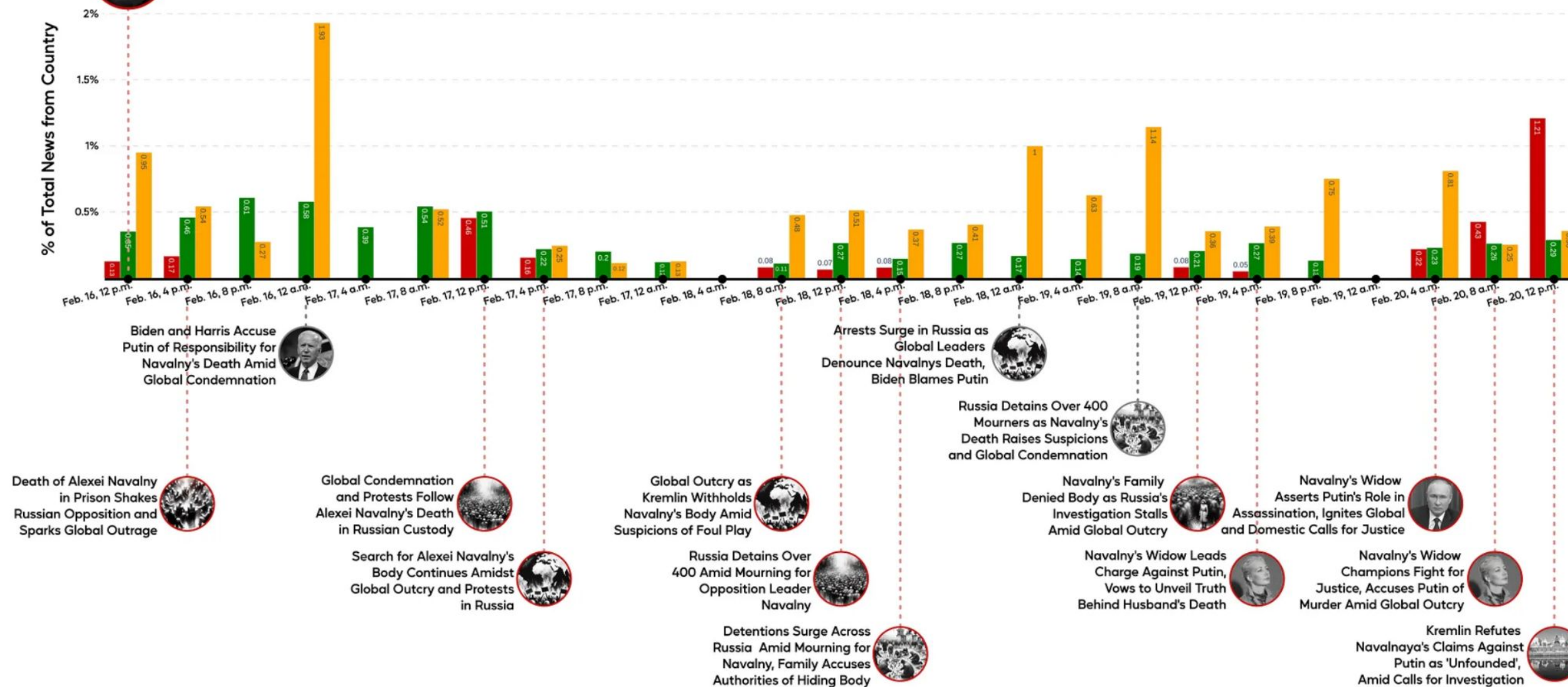**Russia Detains Over 400 Mourners as Navalny's Death Raises Suspicions and Global Condemnation**

**Navalny's Family Denied Body as Russia's Investigation Stalls Amid Global Outcry**

**Navalny's Widow Leads Charge Against Putin, Vows to Unveil Truth Behind Husband's Death**

**Navalny's Widow Asserts Putin's Role in Assassination, Ignites Global and Domestic Calls for Justice**

**Navalny's Widow Champions Fight for Justice, Accuses Putin of Murder Amid Global Outcry**

**Kremlin Refutes Navalnaya's Claims Against Putin as 'Unfounded', Amid Calls for Investigation**

Russian, French, and U.S. news coverage of important events following the death of Alexei Navalny. Timestamps are in UTC.
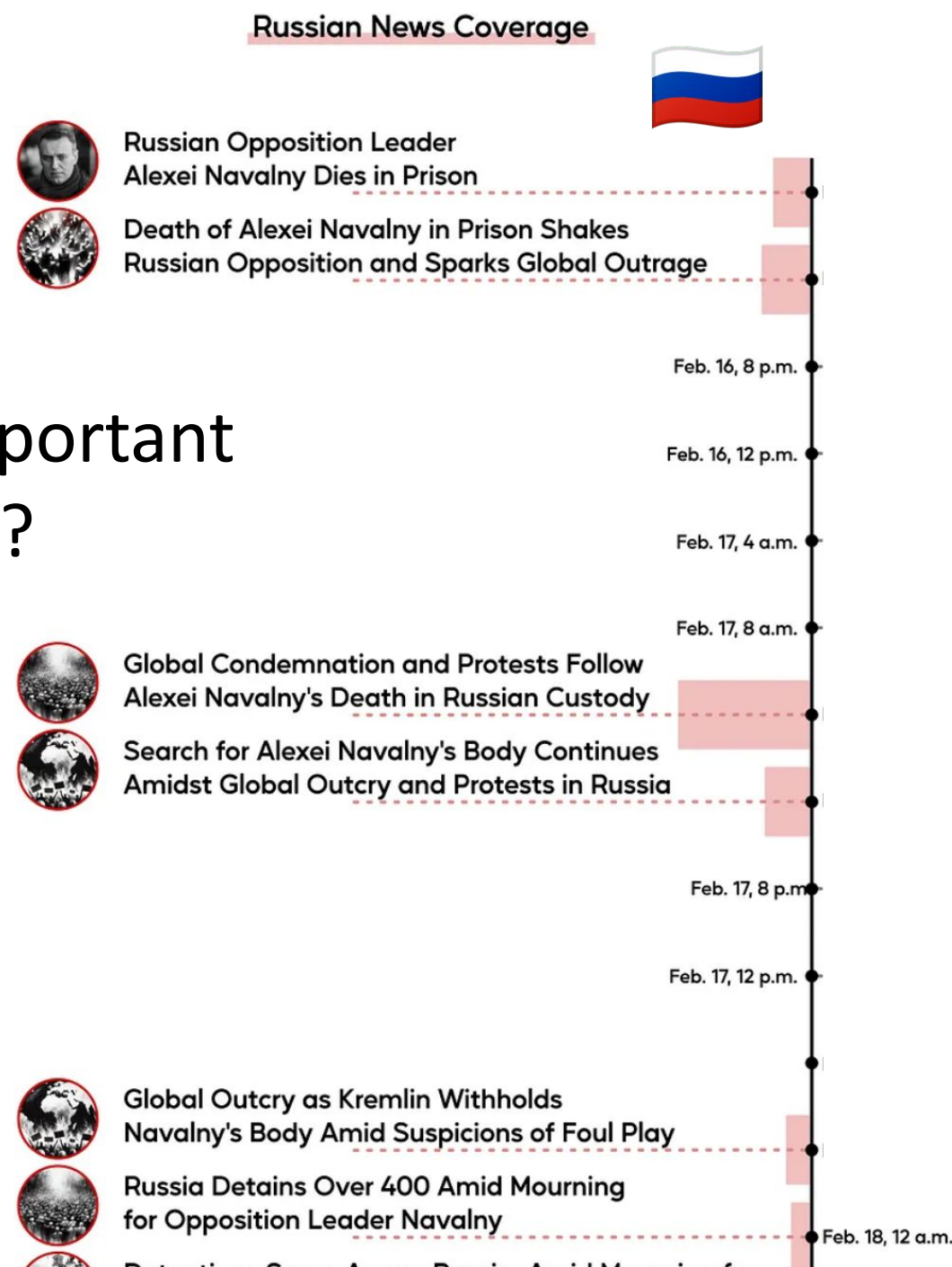
# Uncovering non-reporting

- Which aspects of the narrative were least reported by Russian sources?
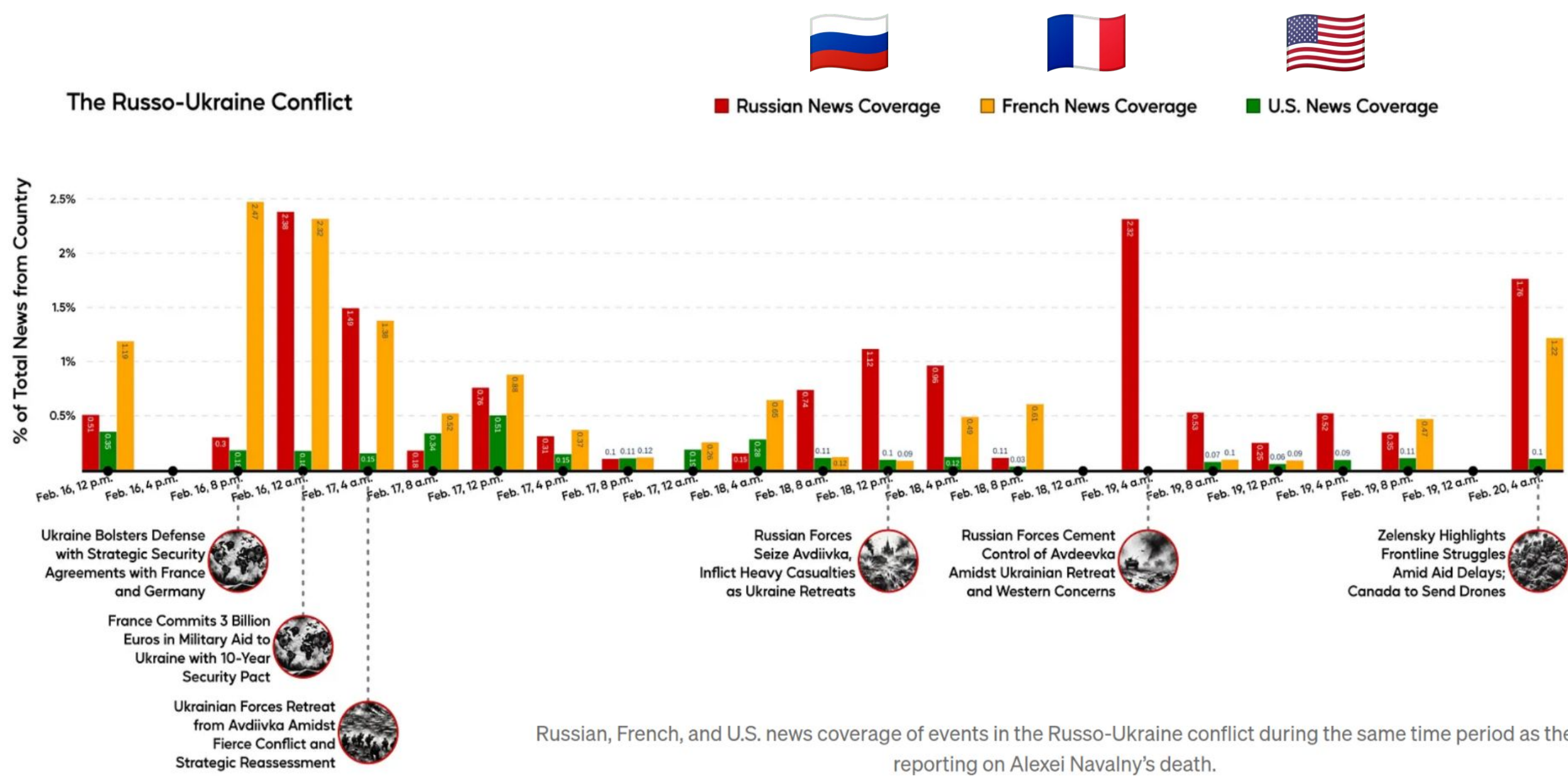
No Russian News Coverage

Feb. 16, 12 a.m.

Feb. 16, 4 p.m.

Global Leaders Accuse Putin of Responsibility for Navalny's Death, Kremlin Denies Allegations

Biden and Harris Accuse Putin of Responsibility for Navalny's Death Amid Global Condemnation

Global Outcry as Russian Opposition Leader Alexei Navalny Dies in Prison

Global Outcry and Protests Demand Justice for Alexei Navalny as Leaders Accuse Putin

Feb. 17, 12 a.m.

Feb. 17, 4 p.m.

Mystery Surrounds Navalny's Death as Family Seeks Body Amidst Global Outcry

Global Leaders Demand Accountability as Russia Detains Mourners

Feb. 18, 4 a.m.

Feb. 18, 8 a.m.

Feb. 18, 12 a.m.

# Uncovering non-reporting

- Which details were important to the Russian sources?

**Russian News Coverage**

**Russian Opposition Leader Alexei Navalny Dies in Prison**

**Death of Alexei Navalny in Prison Shakes Russian Opposition and Sparks Global Outrage**

Feb. 16, 8 p.m.

Feb. 16, 12 p.m.

Feb. 17, 4 a.m.

Feb. 17, 8 a.m.

**Global Condemnation and Protests Follow Alexei Navalny's Death in Russian Custody**

**Search for Alexei Navalny's Body Continues Amidst Global Outcry and Protests in Russia**

Feb. 17, 8 p.m.

Feb. 17, 12 p.m.

**Global Outcry as Kremlin Withholds Navalny's Body Amid Suspicions of Foul Play**

**Russia Detains Over 400 Amid Mourning for Opposition Leader Navalny**

Feb. 18, 12 a.m.

# Coverage of the Russo-Ukraine conflict

The Russo-Ukraine Conflict

Russian News Coverage | French News Coverage | U.S. News Coverage

Russian, French, and U.S. news coverage of events in the Russo-Ukraine conflict during the same time period as the reporting on Alexei Navalny's death.

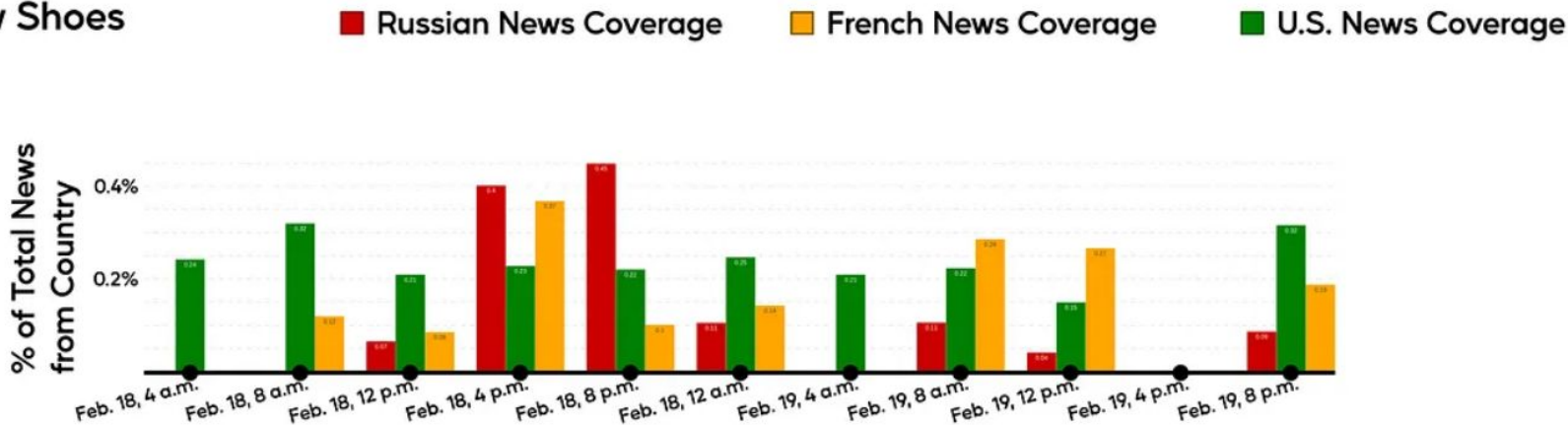# Coverage of the Russo-Ukraine conflict



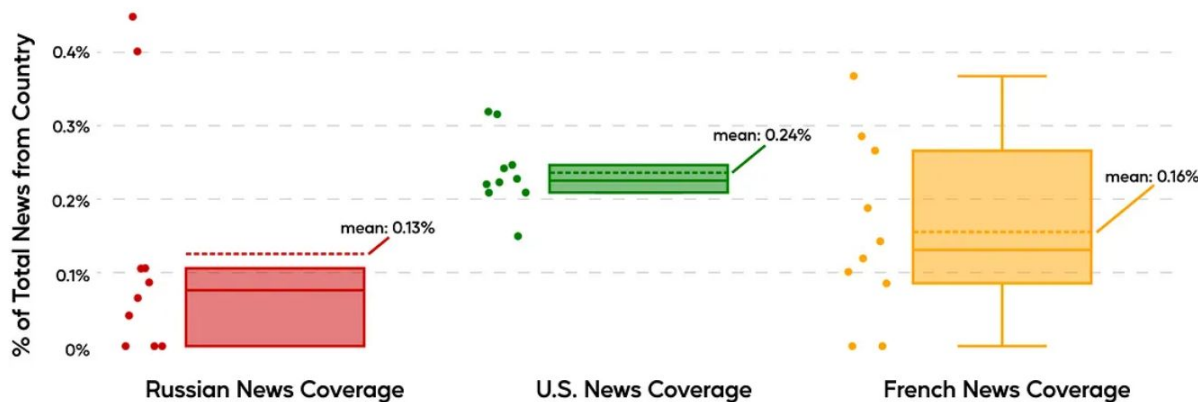The Russian coverage of of the Russo-Ukraine conflict is the same as from French news outlets, while U.S. news coverage is significantly below Russian.

# Coverage of Trump's new shoes

Trump's New Shoes

■ Russian News Coverage  ■ French News Coverage  ■ U.S. News Coverage

News coverage of U.S. ex-president Donald Trump's release of the "Never Surrender High-Tops" shoes.



mean: 0.13%  mean: 0.24%  mean: 0.16%

Russian News Coverage    U.S. News Coverage    French News Coverage

Russian coverage of the release of U.S. ex-president Donald Trump's new shoes matches that of France and U.S.

# Compared to other topics

Equivalent volume of Russian news coverage for:

- Trump's new shoes
- Death of Navalvy's death



Russian news coverage of the death of Alexei Navalny, the Russo-Ukraine conflict, and Trumps new shoes during the topics' overlapping time period.

Data and analysis available in Google Colab

# Blog article available for more details



Identifying Media Bias with AI

## Identifying media bias with AI

Russian news coverage of the death of Alexei Navalny

Emergent Methods · Following
4 min read · Feb 23, 2024

👏 69    💬                          🔖  ▶  ⬆  ⋯

Blog article on [Medium](#)

www.emergentmethods.ai

# Outsource your news context to AskNews

Blog article on [Medium](Medium)

```python
from asknews import AskNewsSDK

ask = AskNewsSDK()

query = "The effect of President Xi's visit to France on the Russo-Ukraine conflict"

news_context = ask.news.search_news(query).as_string

# Now run your Chat completion as usual:
system = {
    "role": "system",
    "content": f"A chat between a curious user and an artificial intelligence Assistant. The Assistant has access to the following news articles that may be useful for answering the User's questions: {news_articles}"
}
response = oai.chat.completions.create(
    model="gpt-3.5-turbo",
    messages=[system, user]
)

print(response.choices[0].message.content)
```

Emergent Methods
www.emergentmethods.ai

# Outsource your news context to AskNews

Blog article on [Medium](Medium)

```python
from asknews import AskNewsSDK

query = "The effect of President Xi's visit to France on the Russo-Ukraine conflict"

# Grab a prompt-optimized string ready to go for your LLM:
news_context = an_client.news.search_news(
    query=query, # any natural language query
    n_articles=10, # control the number of articles to include in the context
    return_type="string",  # you can also ask for "dicts" if you want more metadata and formatted documents
    method="nl"  # use "nl" for natural language for your search, or "kw" for keyword search
).as_string
```

# Outsource your news context to AskNews

- 🌍 Global coverage - 300k articles per day
- 🏎️ Low-latency (100 ms) - aimed at tight spots in your LLM stack
- 🧺 Stories + clustering, Chat, Finance, Citation control
- 📈 Hot topic following/tracking/filtering
- 📊 Usage tracked - pay for what you use
- 👥 Reddit perspective - include social context
- FREE Free news api [https://my.asknews.app/plans](https://my.asknews.app/plans)

Blog article on [Medium](Medium)

**AskNews**

[https://docs.asknews.app](https://docs.asknews.app)

# Thanks for your attention!



Engineering news context
https://asknews.app

Robert Caulk, PhD & Elin Törnquist, PhD
*CEO*                    *Director of Transparency*



https://emergentmethods.ai