# ARElight
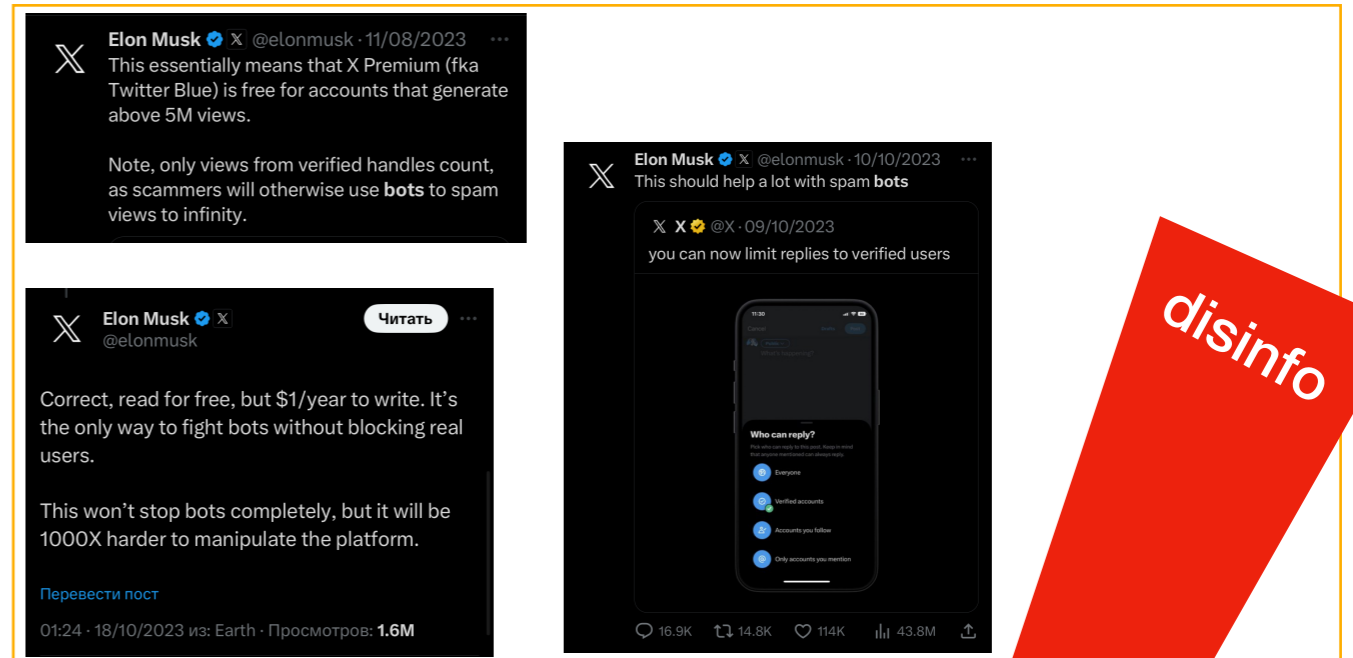## AI & graph powered text analysis tool

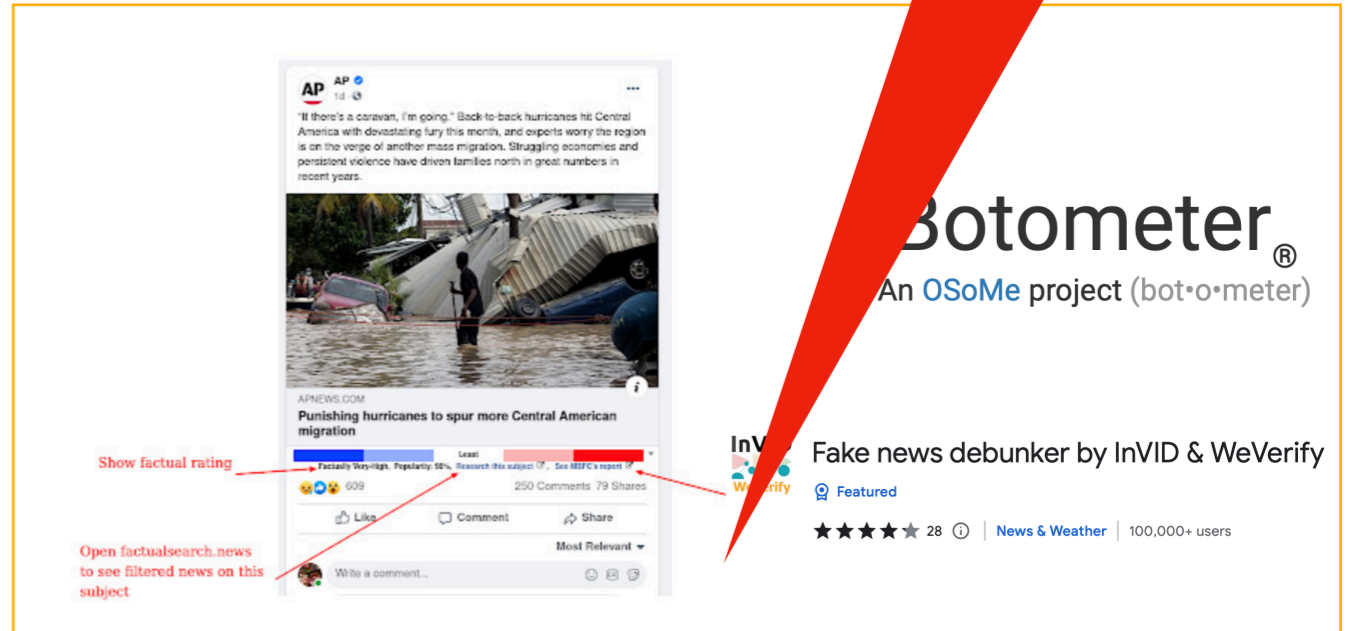Nikolay R., **Maxim K.**

# Structure

- The problem

- Description of ARElight

- Demo

- Conclusions

# The problem

- When speaking about disinformation, there exist 2 walls of defence:

  - protection by social media company

  - protection by users themself



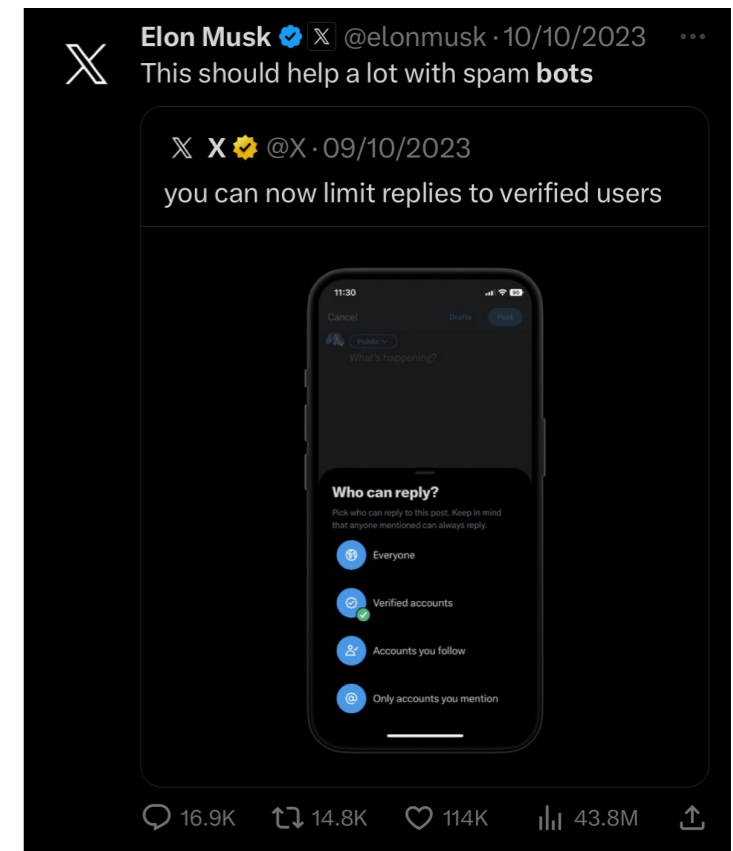measures by social media

disinfo

measures by users

# The problem

- Measures by SM examples:

  - Twitter/X efforts to decrease number of bots

  - Block/Shadowban of suspicious accounts

  - Tinder AI safety features



Elon Musk ✔ 𝕏 @elonmusk · 11/08/2023 ···
This essentially means that X Premium (fka Twitter Blue) is free for accounts that generate above 5M views.

Note, only views from verified handles count, as scammers will otherwise use **bots** to spam views to infinity.

Elon Musk ✔ 𝕏
@elonmusk                                    Читать

Correct, read for free, but $1/year to write. It's the only way to fight bots without blocking real users.

This won't stop bots completely, but it will be 1000X harder to manipulate the platform.

Перевести пост

01:24 · 18/10/2023 из: Earth · Просмотров: **1.6M**

Elon Musk ✔ 𝕏 @elonmusk · 10/10/2023 ···
This should help a lot with spam **bots**

𝕏 X 🏅 @X · 09/10/2023
you can now limit replies to verified users

💬 16.9K   ⇄ 14.8K   ♡ 114K   �󠀠 43.8M   ⬆
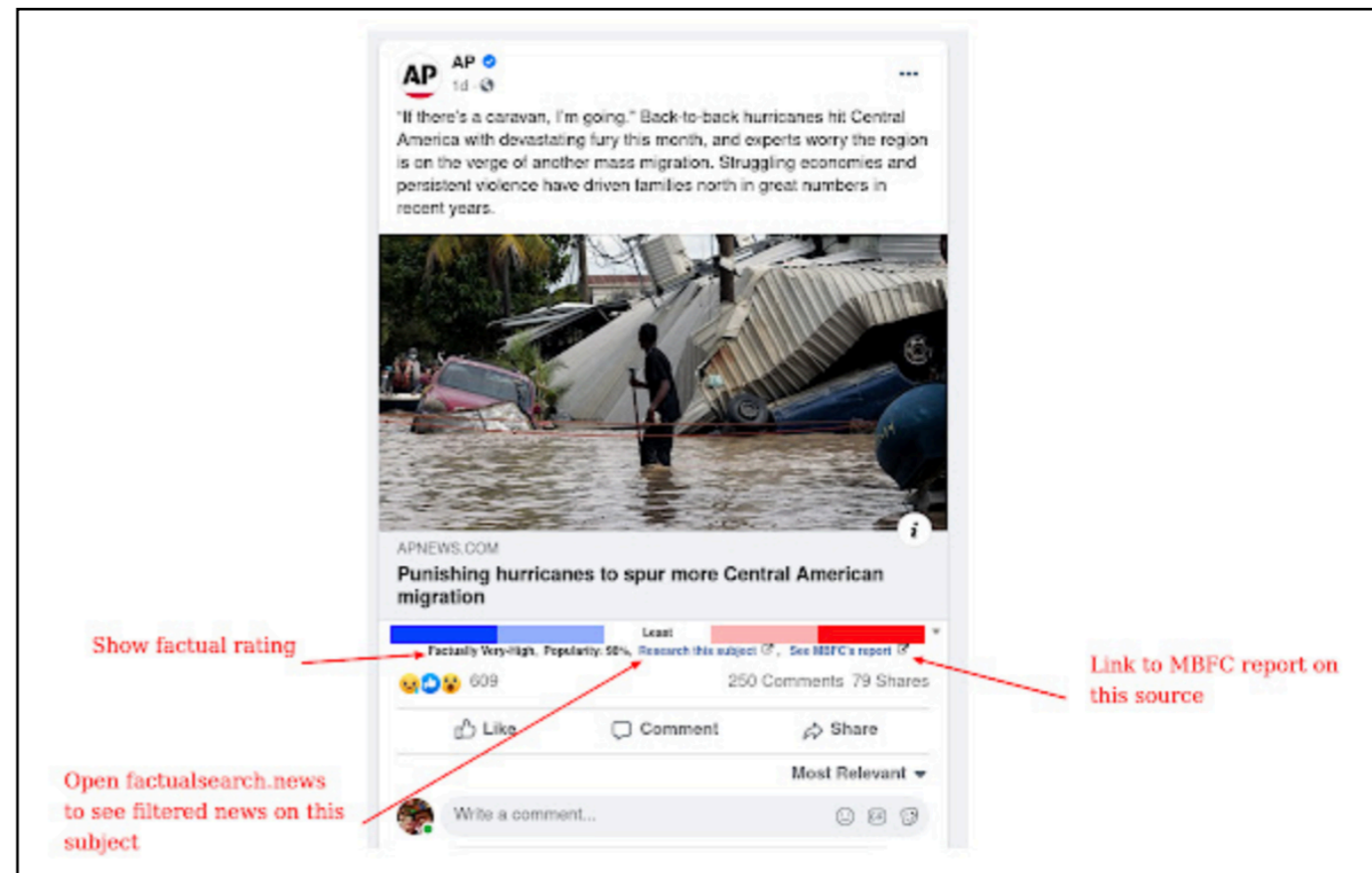


The Indian Express

## A look at shadowban policies on Instagram & Facebook and what you can do to appeal unfair bans

First, what is shadow banning? If you are a creator and have noticed a sudden, significant drop in engagement on your content then chances are...



Another space where AI will change the game is in-app safety and efficiency. Tinder already leverages machine learning for various safety features, but the opportunity to evolve these features through the use of AI can potentially make their impact even greater. For

# The problem

- Measures by users examples:

  - Tools that spot fake news & bots

  - Turing tests for chats

  - Bot-helpers to manage group chats

  - Manual human verification

# The problem

- For manual verification you need to read a lot of text - in fact, you need to scan previous activity of account

 Media outlet

 User

- For many accounts **you just can not do it manually**

# ARElight

SuspiciousCat
@suscat2023

Alice was his friend, so she will support Bob

12:16 AM · 2021-04-27 · Twitter for iPhone

**3,911** Retweets  **50** Quote Tweets  **27.9k** Likes

SuspiciousCat
@suscat2023

Chris thinks Bob deceived him

12:16 AM · 2021-04-27 · Twitter for iPhone

**3,911** Retweets  **50** Quote Tweets  **27.9k** Likes

SuspiciousCat
@suscat2023

Bob didn't like living in the USA
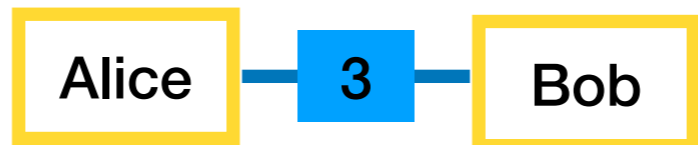
12:16 AM · 2021-04-27 · Twitter for iPhone

**3,911** Retweets  **50** Quote Tweets  **27.9k** Likes

# ARElight



- ARElight use AI to identify instances in text such as: Person, Organization, Location, Date

- Within small amount of text, identify relations and their sentiment (positive 😃, negative 😡, neutral 😐)

# ARElight



- Small texts united into a single network, that represent the narrative of text

# ARElight

- Demo: https://guardeec.github.io/arelight_demo/template.html



(c) visualisation model selector

(a) dataset selector

(b) visualisation options

(d) force layout visualization model

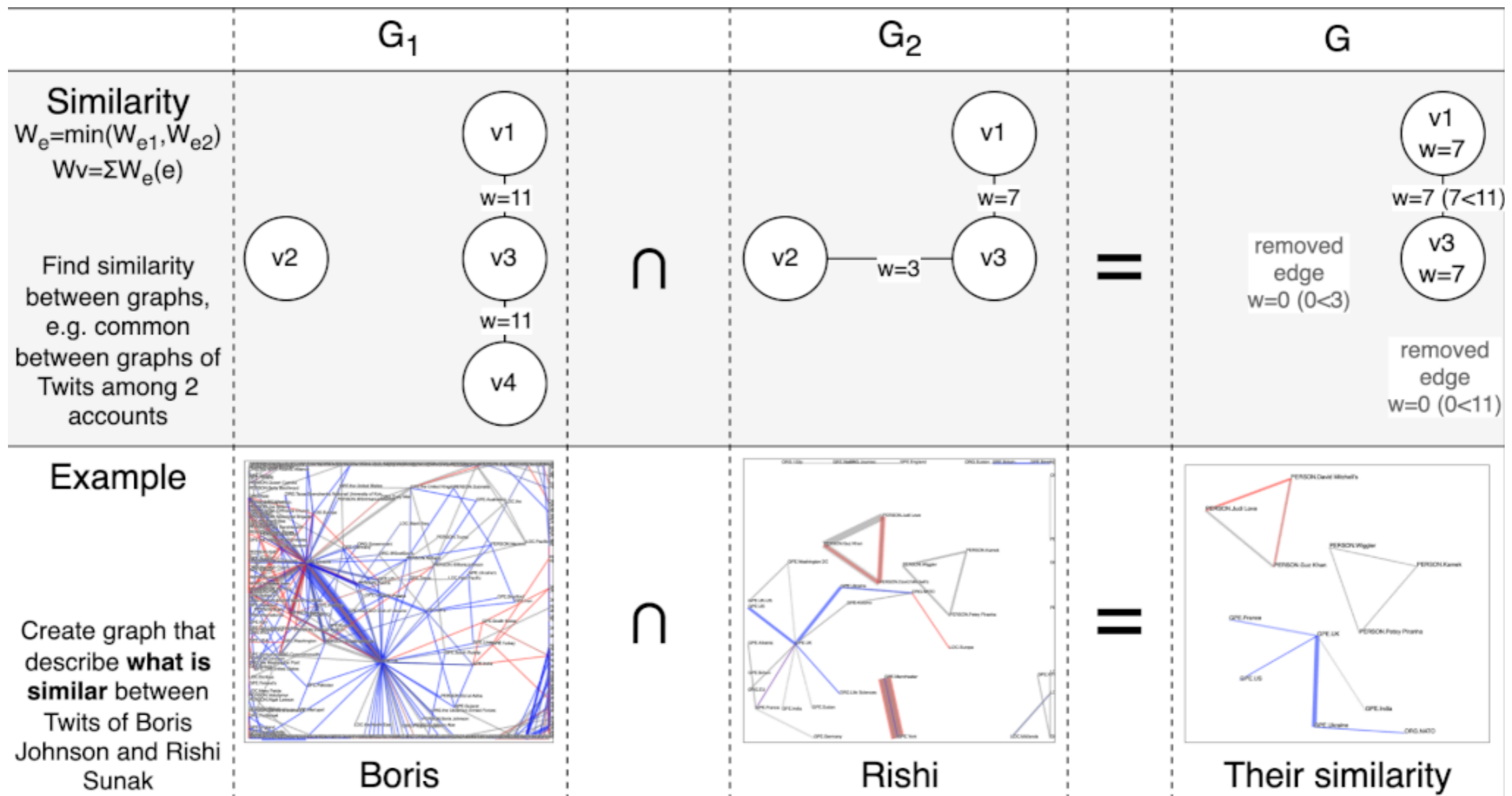(d) radial layout visualization model

# ARElight

- Operations - UNION
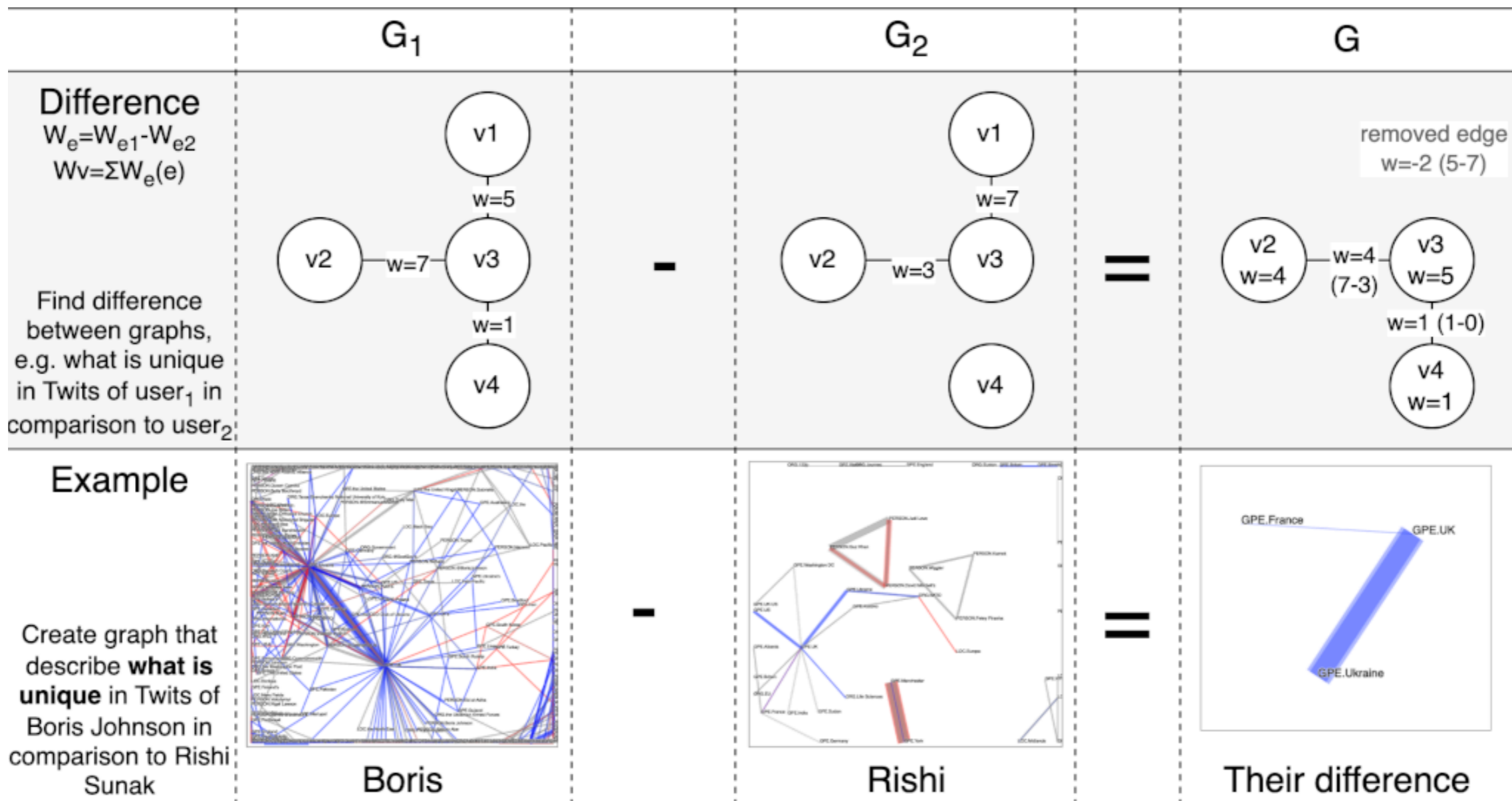
- Merge multiple texts' graphs together

# ARElight

- Operations - INTERSECTION
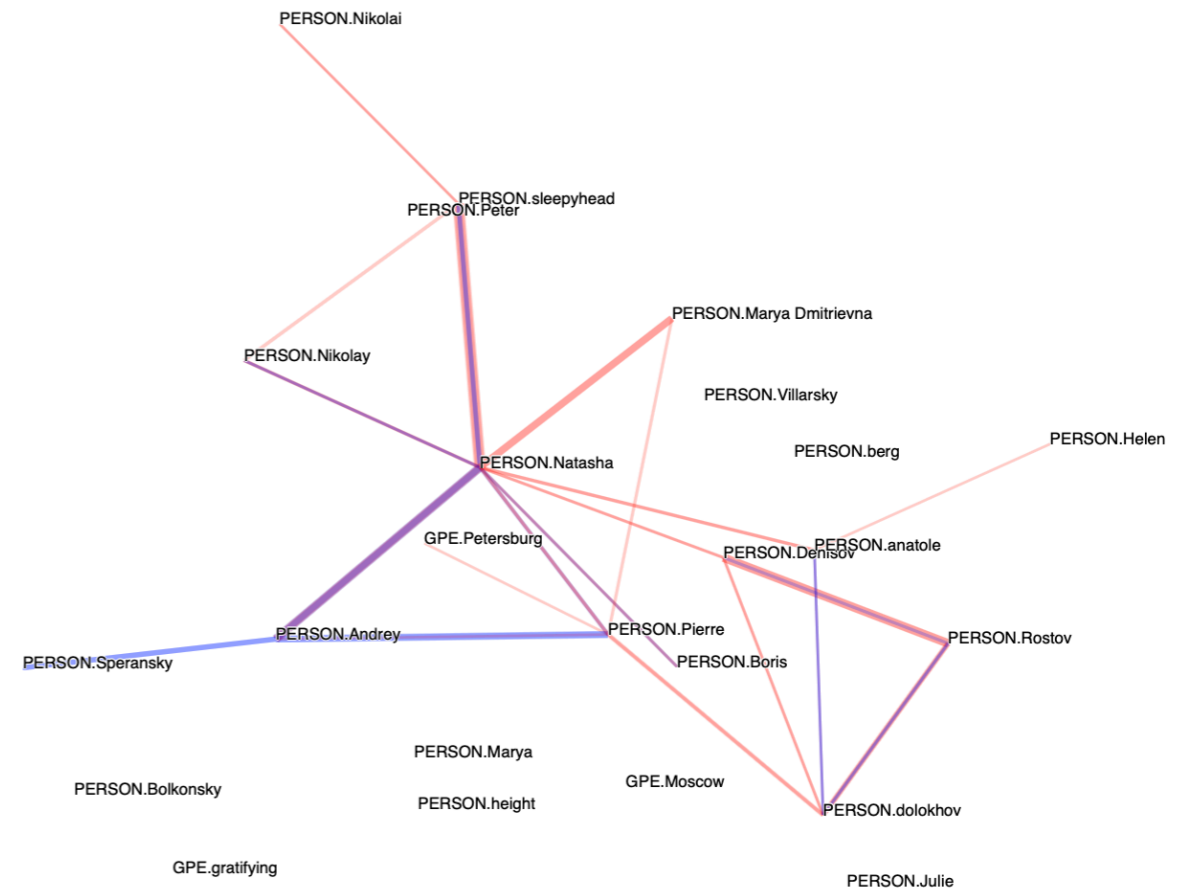
- What is similar between 2 texts?

# ARElight

- Operations - DIFFERENCE (not commutative)

- What is unique in one text in comparison to another?
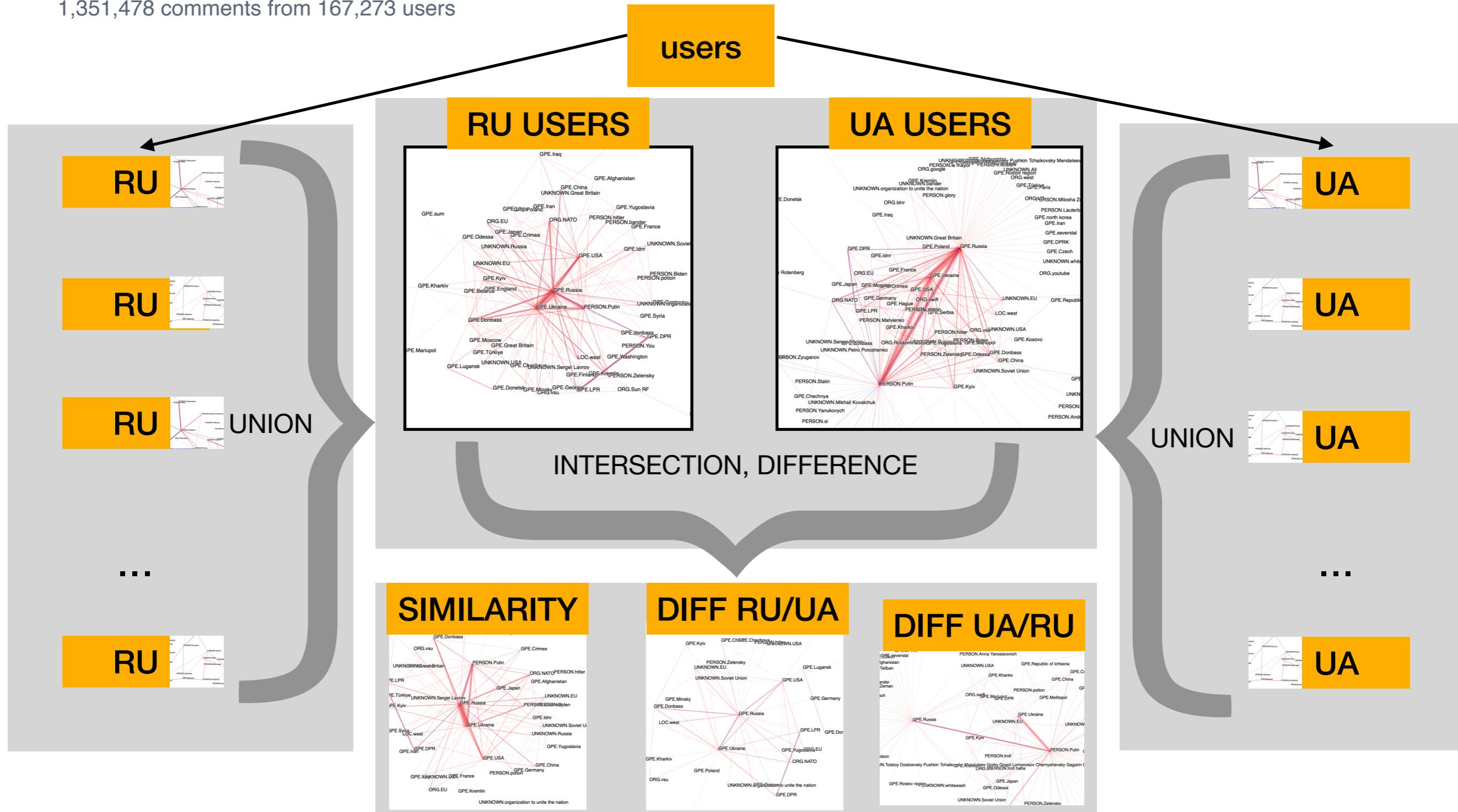
# Demo

- Examples:

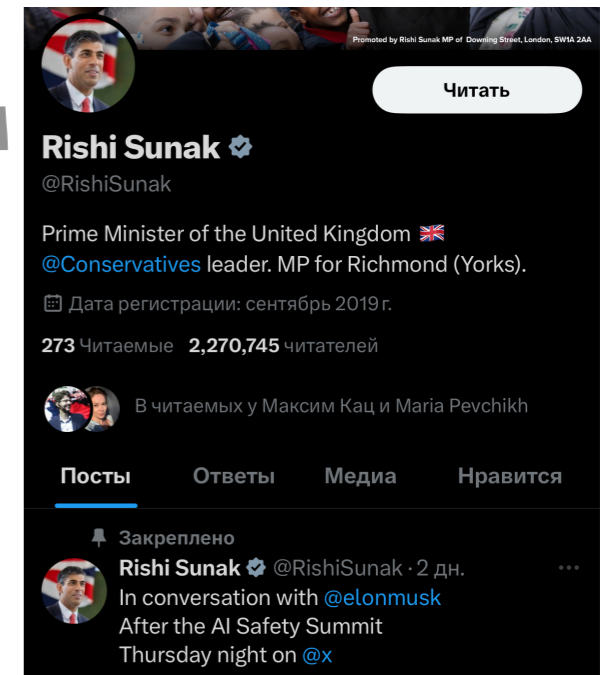- War & Peace book

# Demo
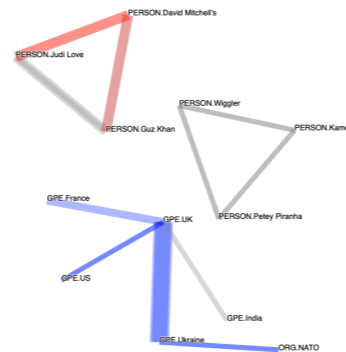
- Examples: RU/UA comments

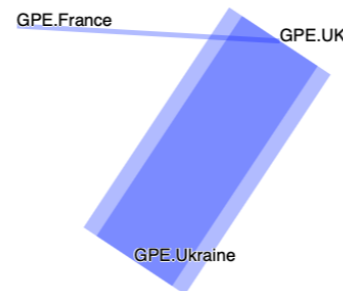1,351,478 comments from 167,273 users

# Demo

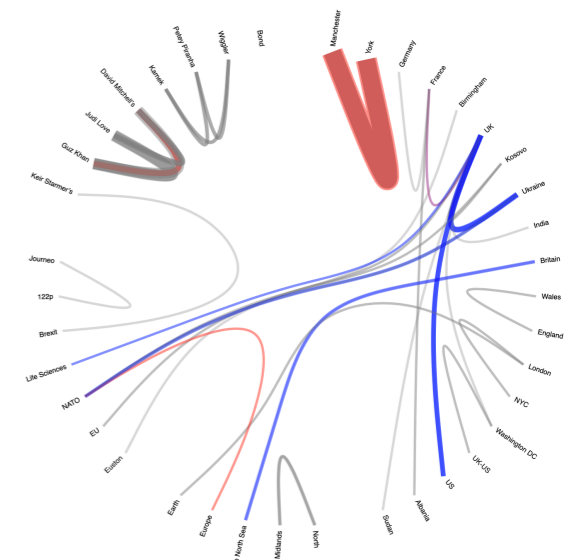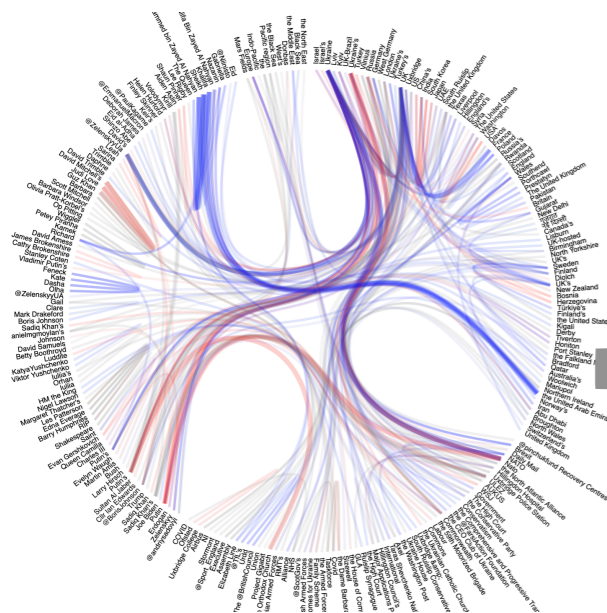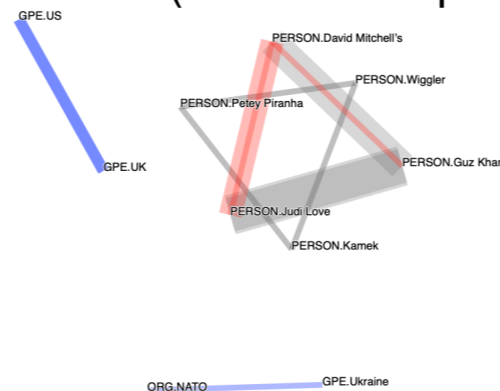- Examples: Boris Johnson VS Rishi Sunak

1000 last Twits



INTERSECTION (What is similar?)

DIFFERENCE (What is unique in Boris?)

DIFFERENCE (What is unique in Rishi?)

# Overview

- ARElight **helps to extract/analyse narrative** and represent it in understandable manner for:

  ➡users - for accounts verification and media checks

  ➡researchers - extract/compare narratives of accounts/disinformation

  ➡STEM - present texts in form of graphs for ML/automation/etc.

- ARElight is a tool in several versions:

  ➡Python Library & Demo [complete] - for STEM

  ➡Web Service [planned] - for researchers

  ➡Twitter/X browser extension [planned] - for users

- We are ready to help if you would like to use it in your research.


- Python Library: https://github.com/nicolay-r/ARElight/

- Demo: https://guardeec.github.io/arelight_demo/template.html

Nikolay R. (rusnicolay@gmail.com), **Maxim K.** (maksim.kalameyets@newcastle.ac.uk)