



Harnessing data to combat Propaganda and Disinformation

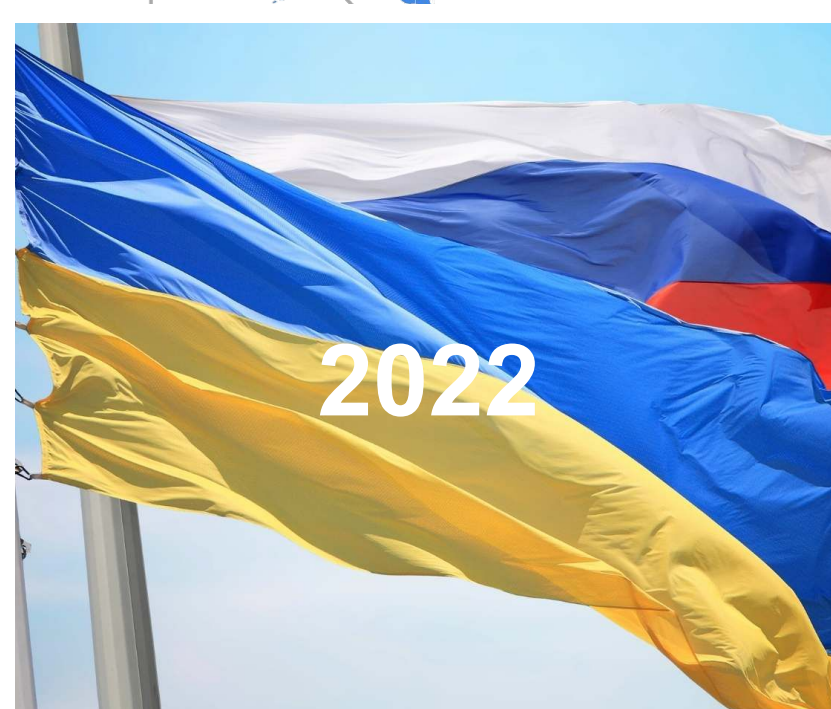
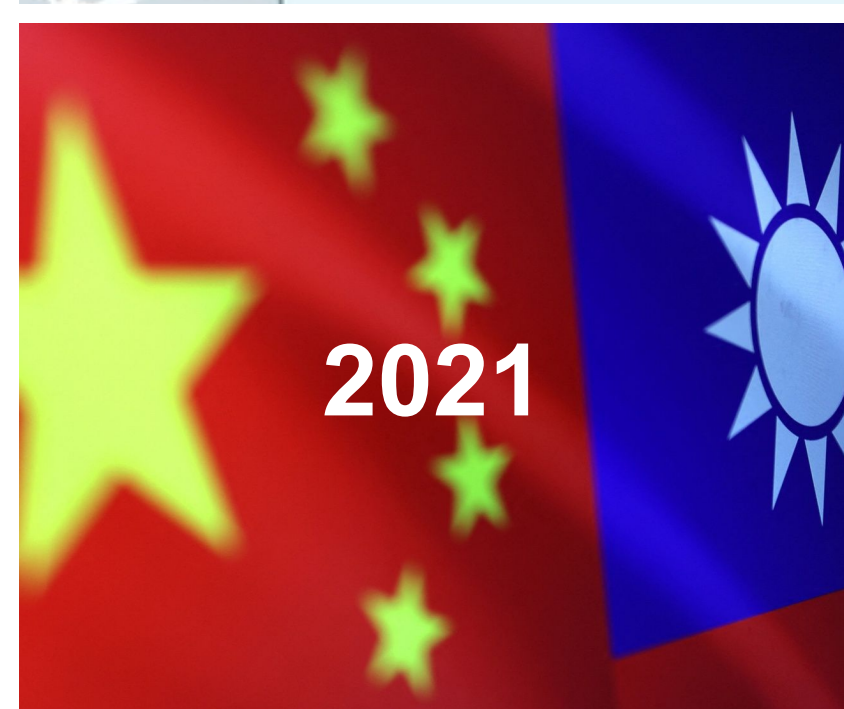
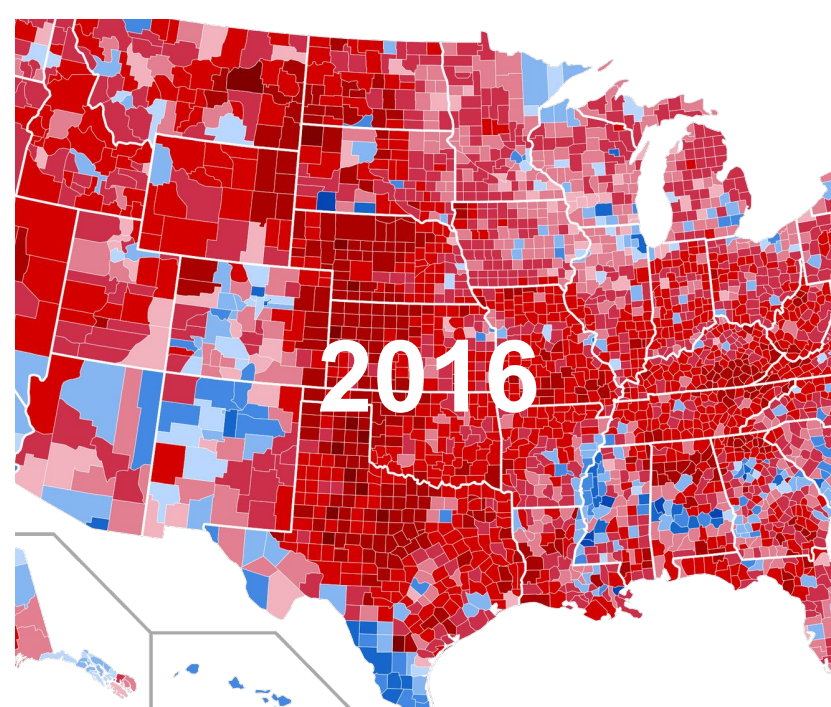


Nesin Veli,

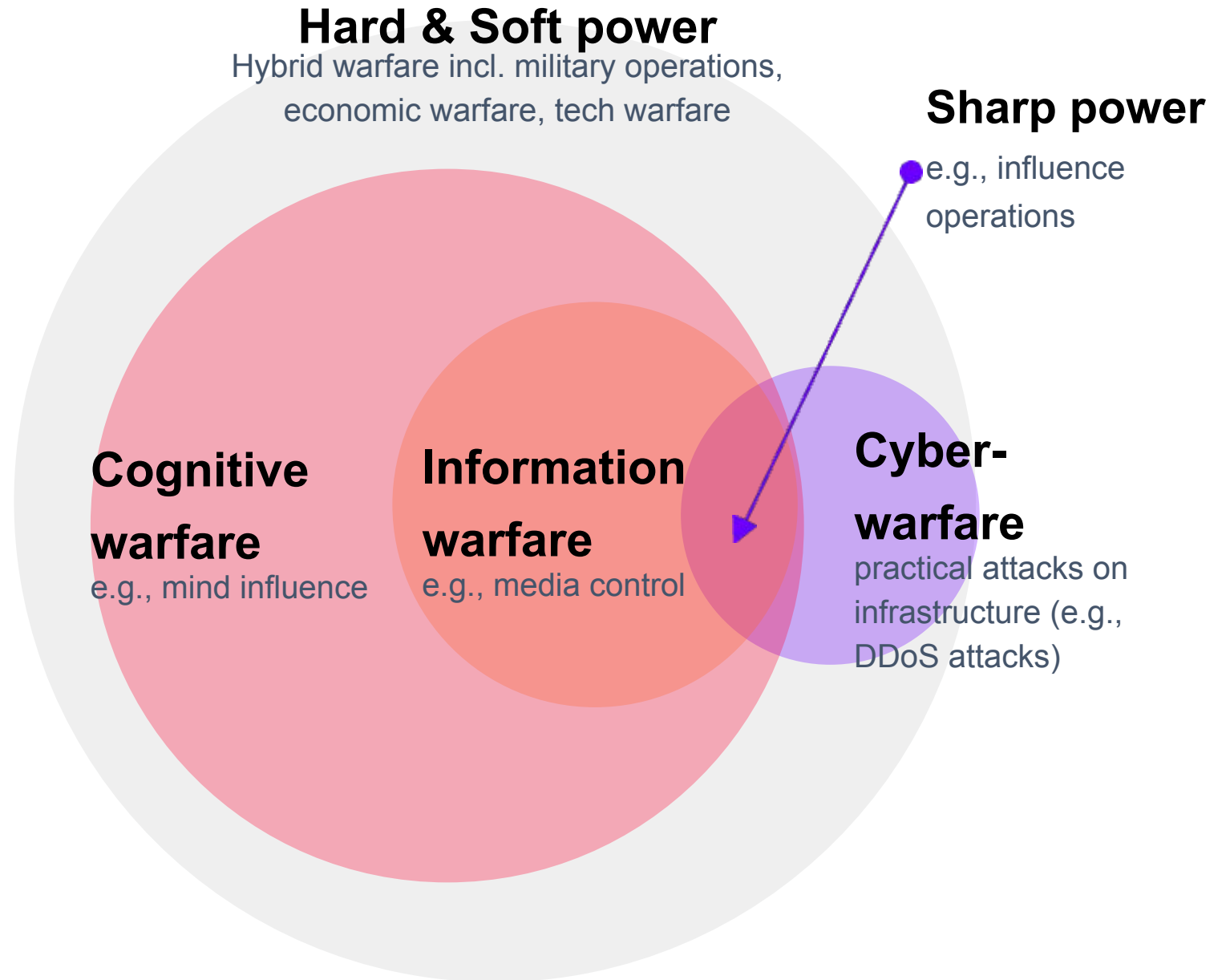
Products & Services Manager @ Identrics

Web Data Extraction Summit | October 2023

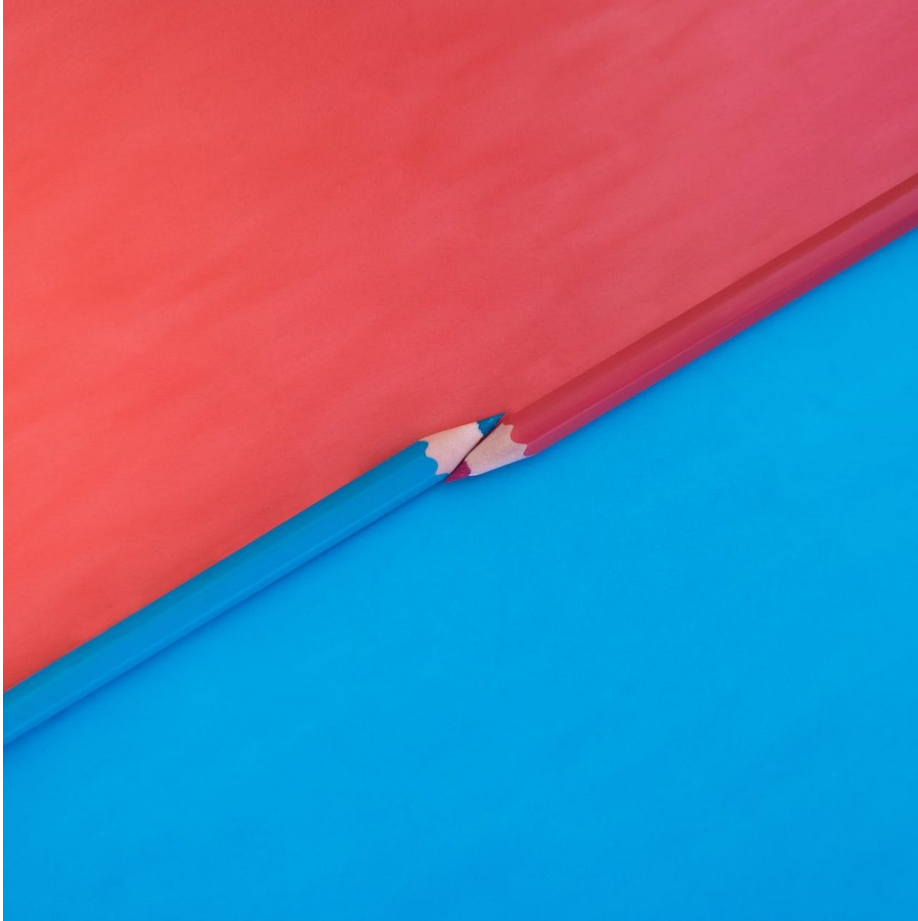
Cognitive warfare



Cognitive warfare



Cognitive warfare



Division and distrust

societal split

government distrust

enable saviour figures

Cognitive warfare



Hearts and minds

it is not just about facts and lies

it is about angst

Cognitive warfare



Actors

rise of non-state actors

Cognitive warfare



Comfort info bubbles

algorithm-pushed news calcified info bubbles

targets are more easily agitated

Cognitive warfare



Media literacy

ideal and too slow

Cognitive warfare



Tech resilience

automation enabling human decision

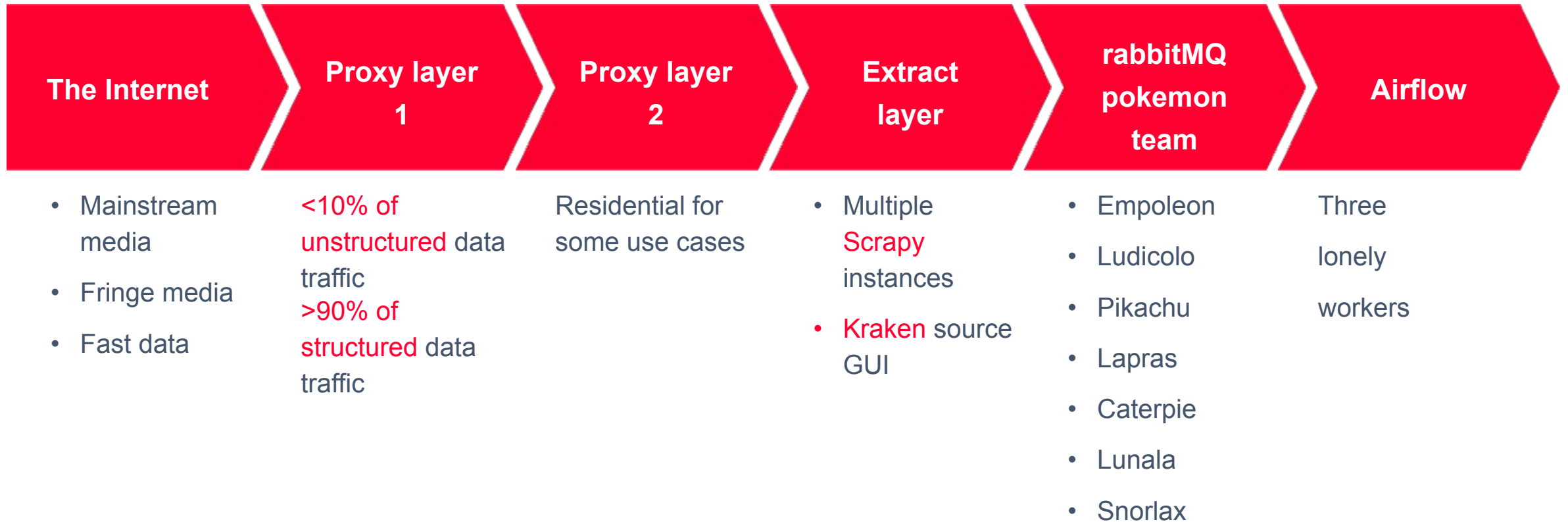
monitoring and alerting systems



Data extraction

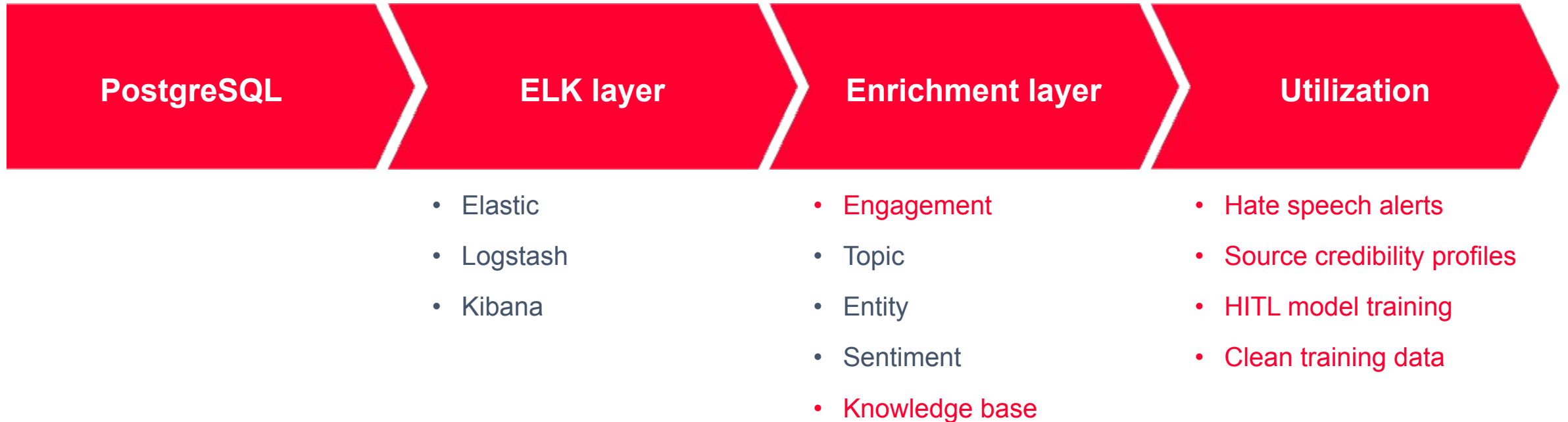
Extraction ecosystem

Encompass all aggregated data including disinfo monitoring



Extraction ecosystem cont.

Encompass all aggregated data including disinfo monitoring



Data per use case

Fast data

- new social media
- mainstream news user comments

"Fringe" media

- mimic mainstream design
- infect the mainstream
- replace the mainstream

Mainstream

- prone to misinformation
- prone to narrative hijack
- prone to algorithm mimicry

Narrative tracking

Origin to end point

Insights

Actionable alerting

Control vs censorship

Disinfo playbook

Propaganda, disinfo, hate speech

Propaganda

- **Broad, general** messages
- **No fact checking**

Disinformation

- **Manipulated content**
- **Manipulated context**

Malinformation

- **Direct cyber attacks**
- **Information leaks**

Misinformation

- **Unaware actors**
- **Main way of dissemination**

Hate speech

- **Marginalizing groups**
- **Media narrative hijack**
- **Media regulatory fines**

Natural **emotional** language



Loaded language

Outrage as Donald Trump suggests injecting disinfectant to kill virus.



Exaggeration or minimization

Coronavirus '**risk to the American people remains very low**', Trump said.



Name calling

WHO: Coronavirus emergency is '**Public Enemy Number 1**'.



Doubt

Can the same be said for the Obama Administration?



Repetition

I still have a **dream**. It is a dream deeply rooted in the American **dream**. I have a **dream** that one day . . .



Slogans

"**BUILD THE WALL!**" Trump tweeted.



Appeal to fear

A dark, impenetrable and "**irreversible**" winter of persecution of the faithful by their own shepherds will fall.



Bandwagon

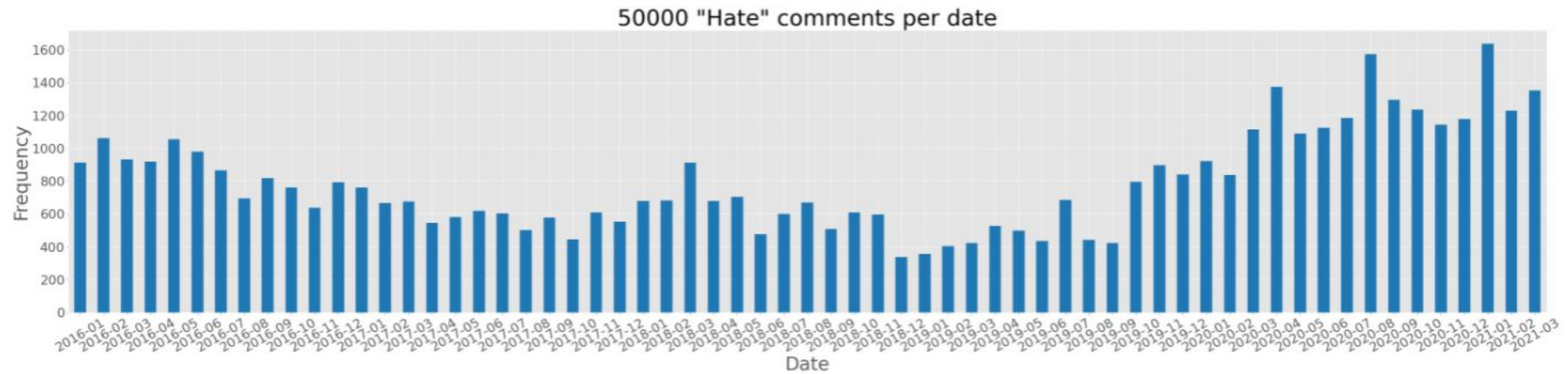
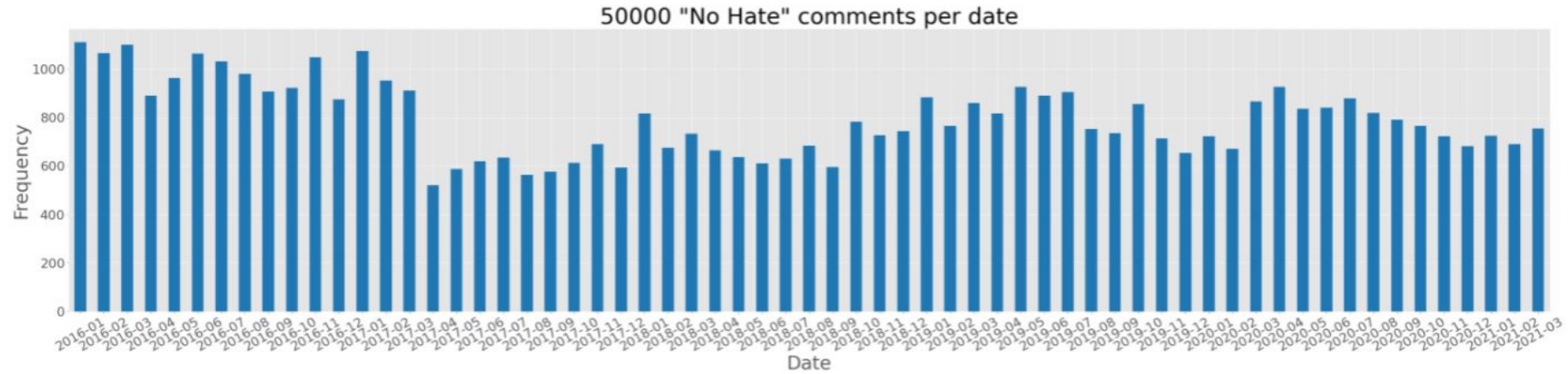
He tweeted, "**EU no longer considers # Hamas a terrorist group. Time for US to do same.**"

Hate speech alerting system

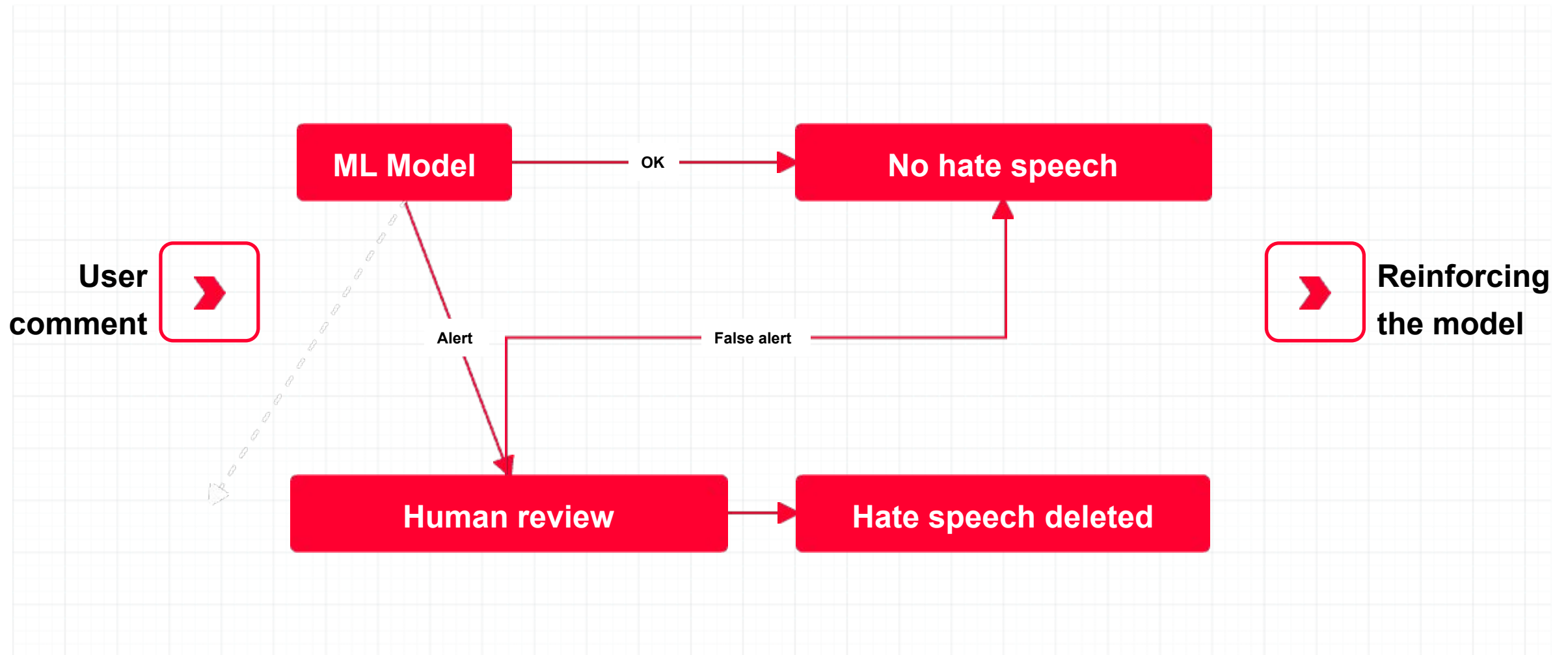
Hate speech raw dataset

Moderator team annotation	Comment count
No hate speech	4,000,000
Hate speech	150,000
Adverts and spam	90,000
Different language	12,000
Political campaign	300
Not GDPR compliant	120

No hate vs hate speech histogram



Hate speech model reinforcement



Neural network results

Corpus: 419306 train + 52414 dev + 52413 test

Results:

- F-score (micro) 0.7993
- F-score (macro) 0.7951
- Accuracy 0.7993

By class:

	precision	recall	f1-score	support
ok	0.8150	0.8339	0.8243	29596
hate	0.7779	0.7544	0.7660	22817
micro avg	0.7993	0.7993	0.7993	52413
macro avg	0.7964	0.7942	0.7951	52413
weighted avg	0.7988	0.7993	0.7989	52413
samples avg	0.7993	0.7993	0.7993	52413

Source credibility profiles

Outlet and author credibility



Original content

document clustering tracks copied content across sources



No info (clickbait) titles

recognition model



Topic dispersion

topic annotation



Staff transparency

author names, about us, editorial



Ownership transparency

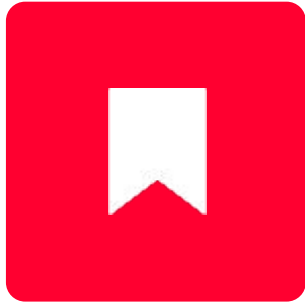
owner, publisher, PEP connections



Author profile

all of the above

Campaign pattern data



Topic dispersion

topic annotation



Komsomolskaya Pravda

on the eve of the Crimea annexation
the outlet starts publishing
exclusively sports content



Red flags

outlet and author profiles
along with historical patterns
are used to train recognition models

The **five** **questions** of disinfo campaigns

- Is there an active campaign?
- What is the origin point and current spread?
- Who are the actors?
- What are the goals?
- Is the pattern new or known?

LLMs

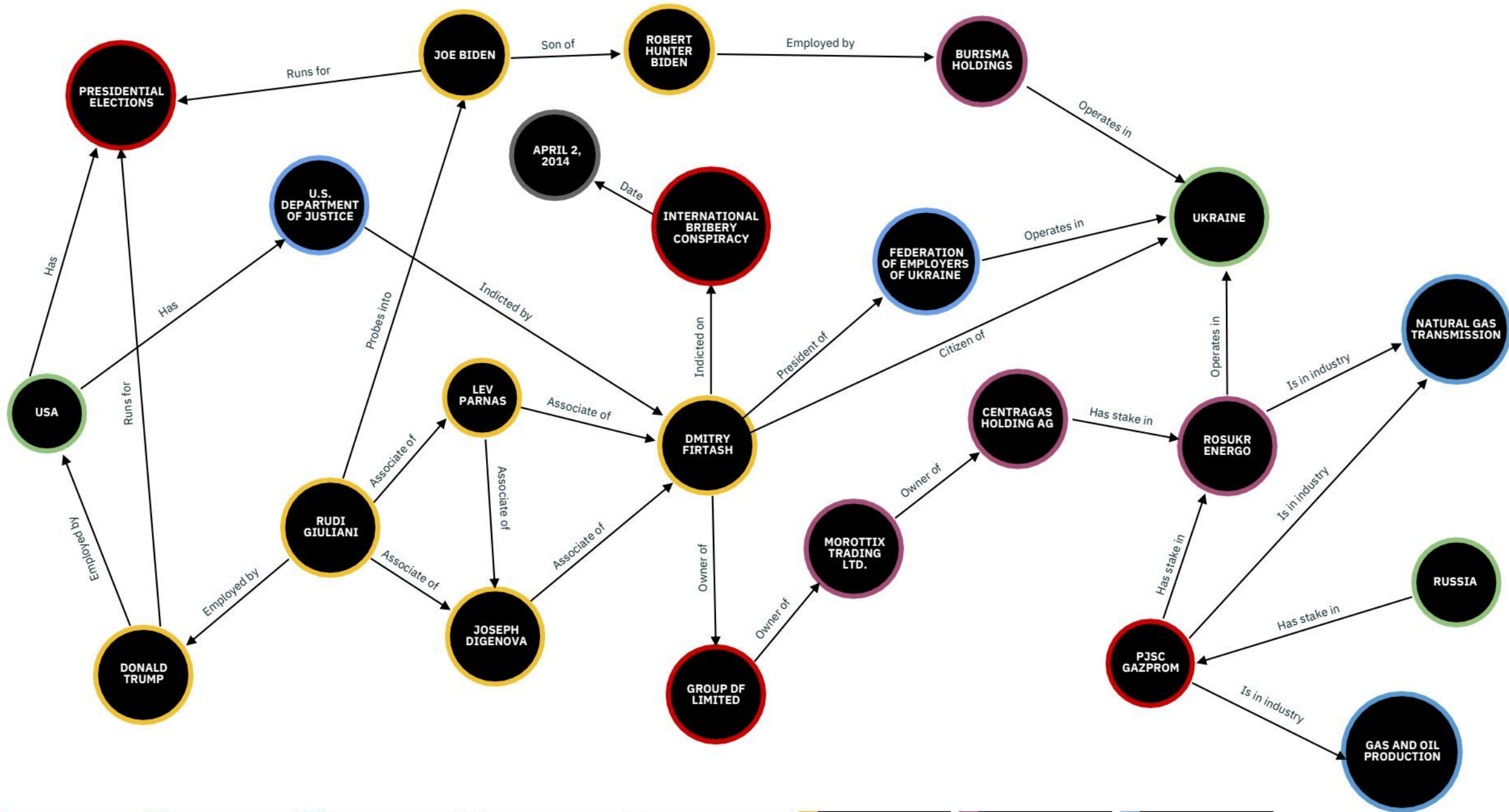
LLM and disinformation

Threats

- Synthetic disinfo data fed back to the model ☆
- Personalised disinfo ☆
- Enable new bad actors ☆
- Boost bot/cyborg networks ☆
- Deep fakes and composite pictures ☆

Opportunities

- Unstructured to structured data extraction ☆
- Source credibility profiles ☆
- Campaign pattern tracking ☆
- Supervised fact-checking ☆
- Cross border narrative tracking ☆



NEWS/EVENTS

COUNTRIES

INDUSTRIES

TIMESTAMPS

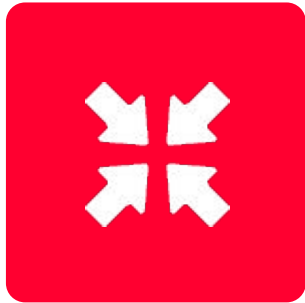
RELATIONSHIPS

PEOPLE

COMPANIES

INSTITUTIONS

Custom models



Focused

Disinfo specific

Bad actor modus operandi



Clean

Remove known hate speech

No synthetic disinfo contamination

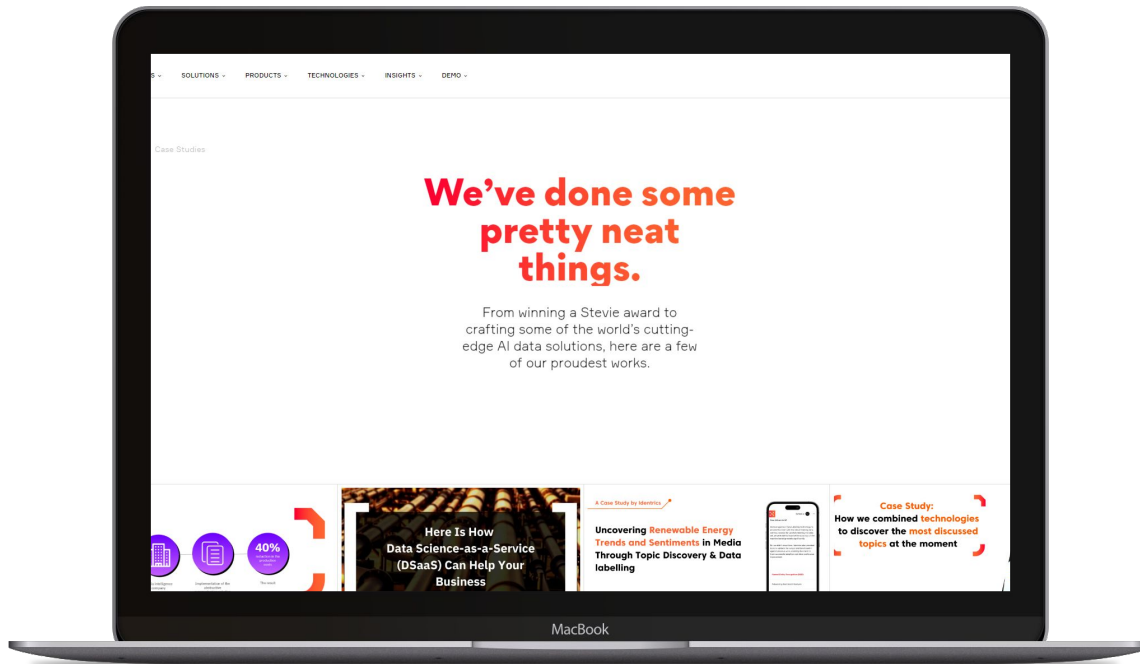


Use case specific

Hate speech per outlet

Manageable corpora

Useful links



- Identrics case studies
- Identrics blog
- Bellingcat Online Investigation Toolkit
- The Code of Practice on Disinformation

Citations

“How China's Cognitive Warfare Works: A Frontline Perspective of Taiwan's Anti-Disinformation Wars”



“SemEval-2020 Task 11: Detection of Propaganda Techniques in News Articles”





THANK
YOU

nesin.veli@identrics.ai

www.identrics.ai

