



**Behavioral  
Data Science**

**Breaking free of the arms race**

Monitor, detect, assess and react  
to influence operations

# Data Science Institute



Dr Marian-Andrei RizoIU | Behavioral Data Science Lead  
Marian-Andrei.RizoIU@uts.edu.au  
<https://www.behavioral-ds.science>

# Red Queen effect



Content-based detectors are sensitive to adversarial training attacks – simply use the detector to train the attacker.

# Challenge: beyond content-based detectors

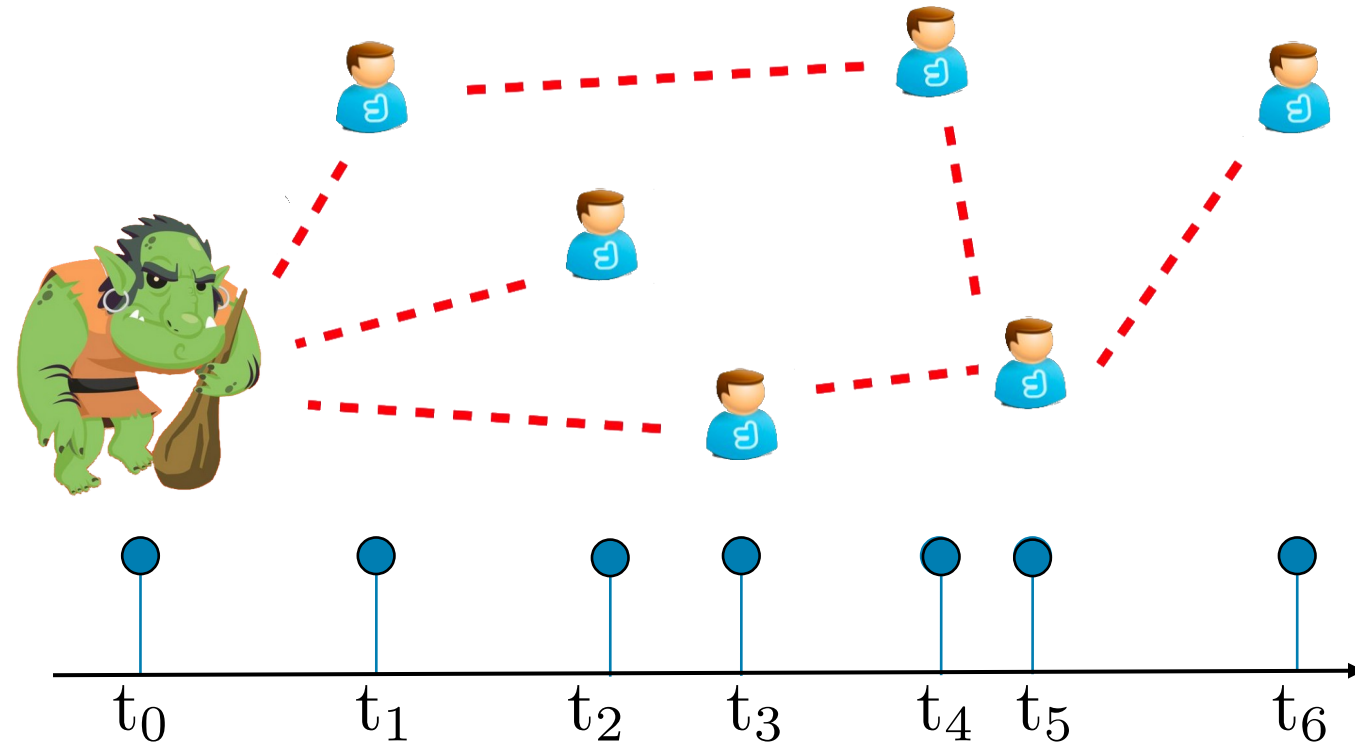
## Content- and user-based detection tools:

language nuances, language drift and adversarial attacks

IO are designed to elicit particular reactions from the target audience

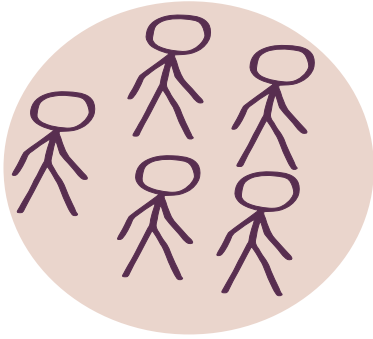
Distinguish users and content types based on the reaction of online social systems (**no content**)

Build early detection systems based on information spread patterns within the user population

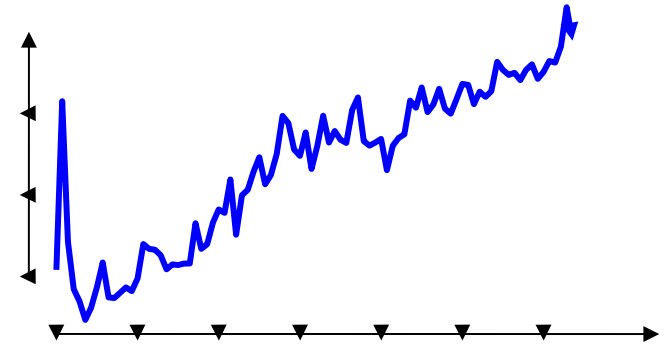


# The Behavioral Data Science

1.



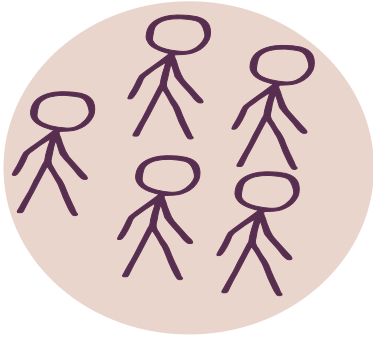
information diffusion  
epidemics spreading  
behavioral modeling



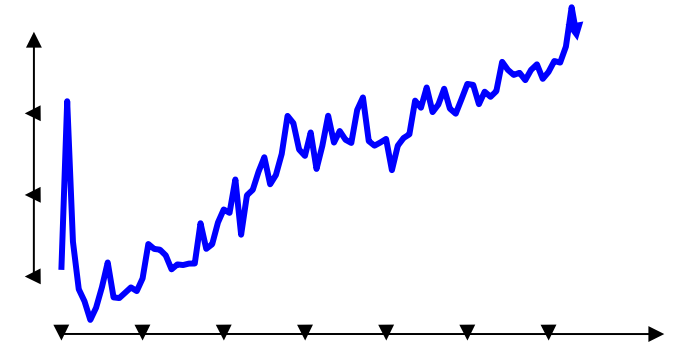


# The Behavioral Data Science

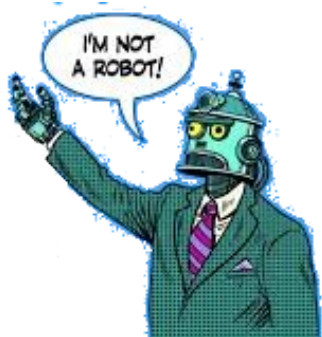
1.



information diffusion  
epidemics spreading  
behavioral modeling



2.



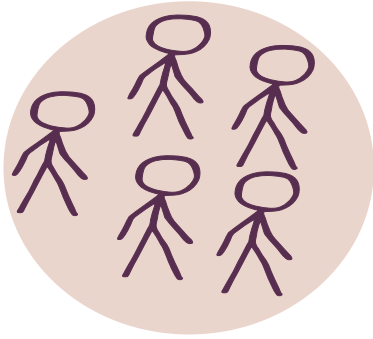
[Rizoiu et al ICWSM'18]



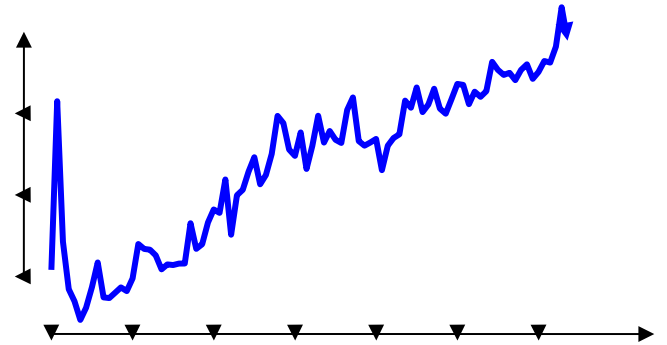
[Kim et al Journ.Comp.SocSci'19]

# The Behavioral Data Science

1.



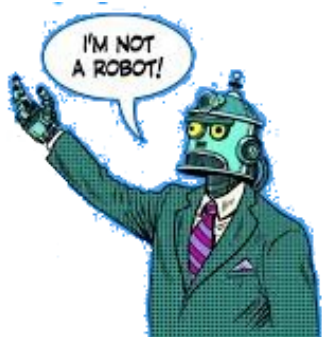
information diffusion  
epidemics spreading  
behavioral modeling



3.



2.

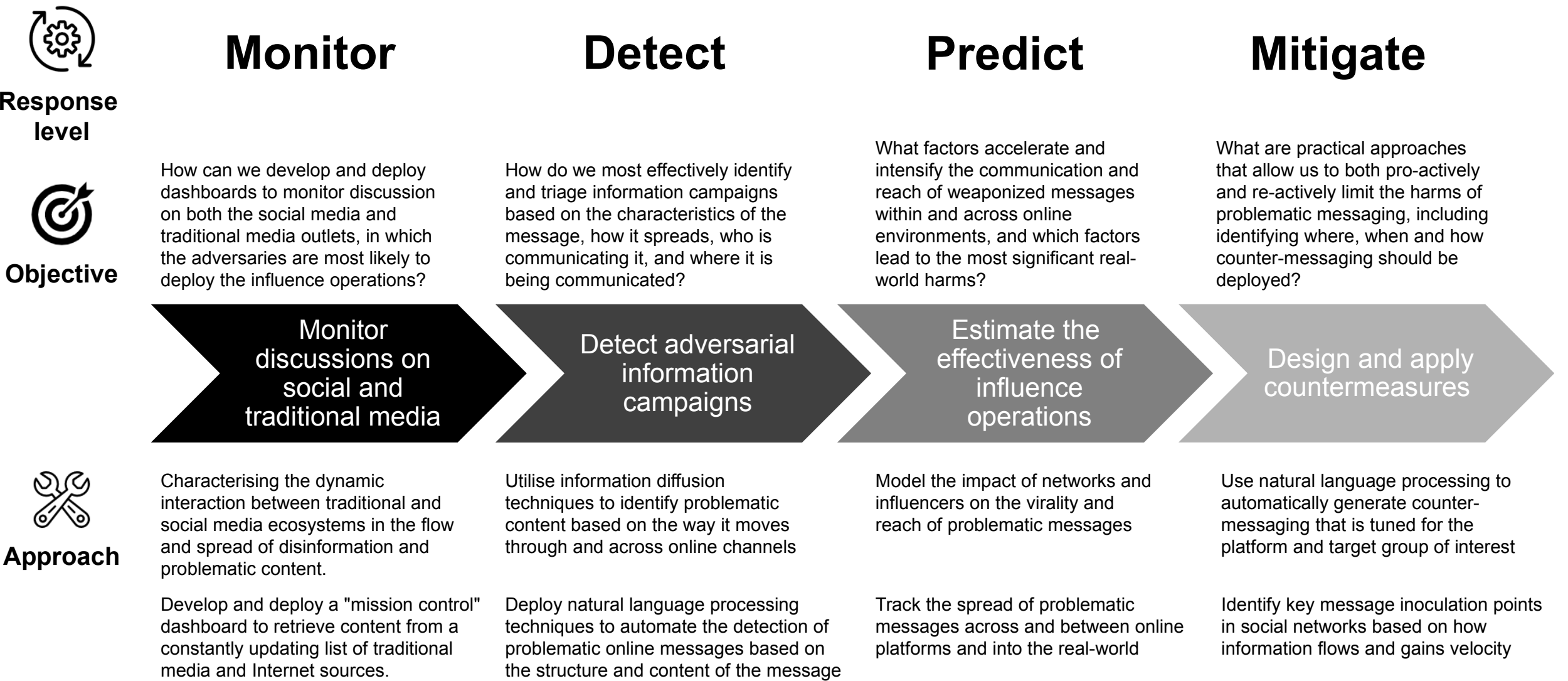


[Rizoiu et al ICWSM'18]



[Kim et al Journ.Comp.SocSci'19]

# Behavioral DS capabilities in Influence Operations space



# Our founders in the mis-, dis-, IO and IW spaces



**Australian Government**  
**Department of Defence**  
Defence Science and  
Technology Group

Real-time detection of  
disinformation campaigns



Information integrity initiative:  
fighting misinformation in Australia

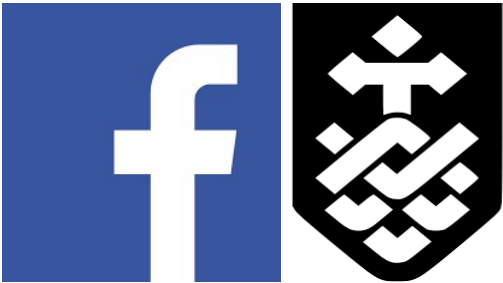


Effectiveness of Information  
Operations in the Pacific



**Australian Government**  
**Department of Defence**  
Defence Science and  
Technology Group

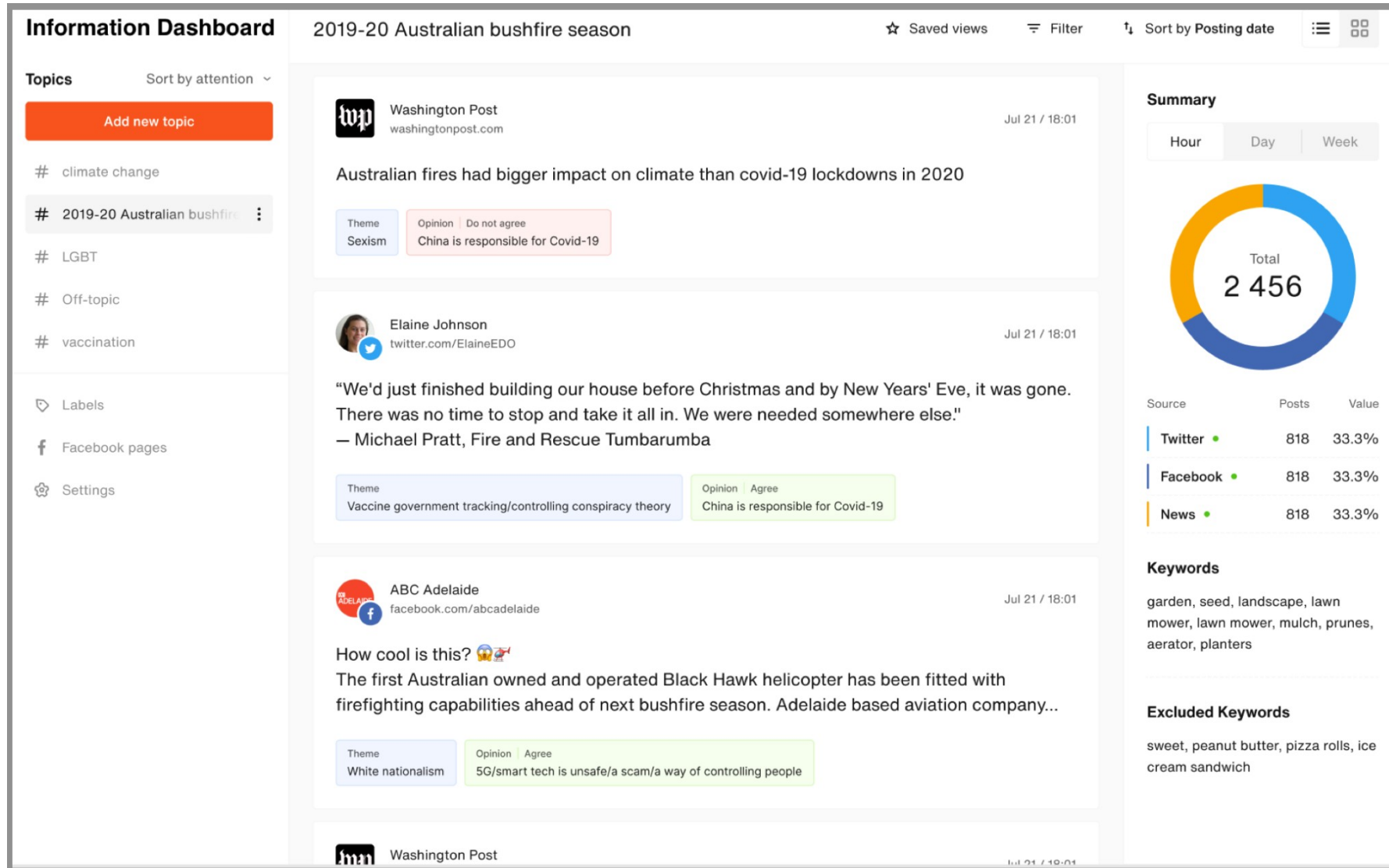
Information Warfare  
STaR Shot “Developing  
Situational Awareness”



Hate Speech propagation  
on Social Media

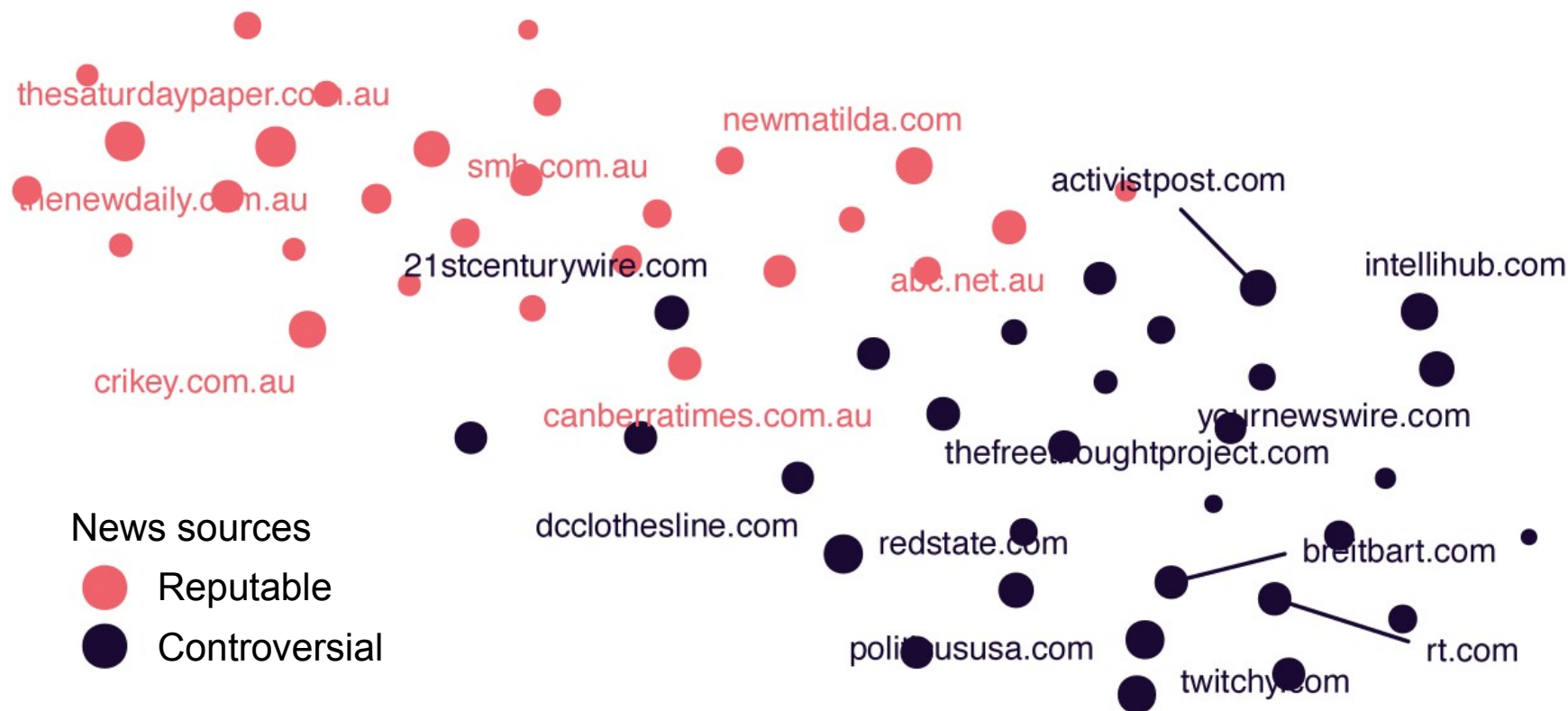


# Monitor: Monitoring discussion spaces (TRL: 3)



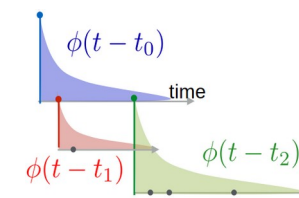
Graphical interface of the Information Dashboard

# Detect: separating controversial from reputable



Reputable and controversial sources are separable based solely on how their information spreads

Detect controversial news without content analysis



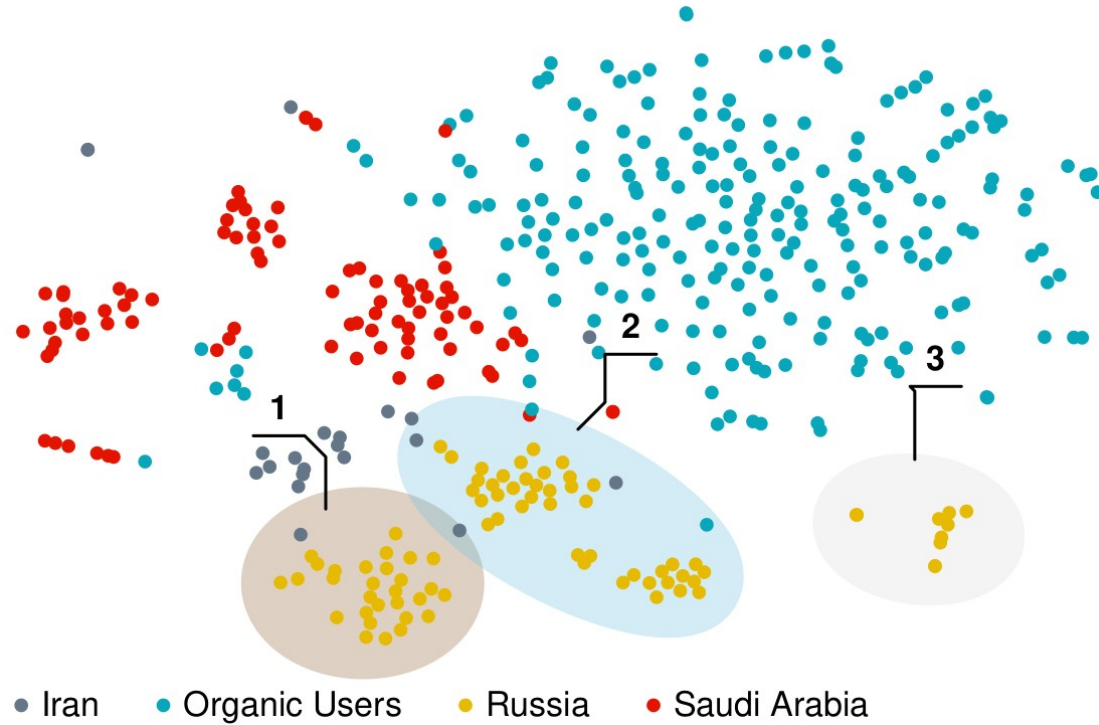
evently

## The technical detail:

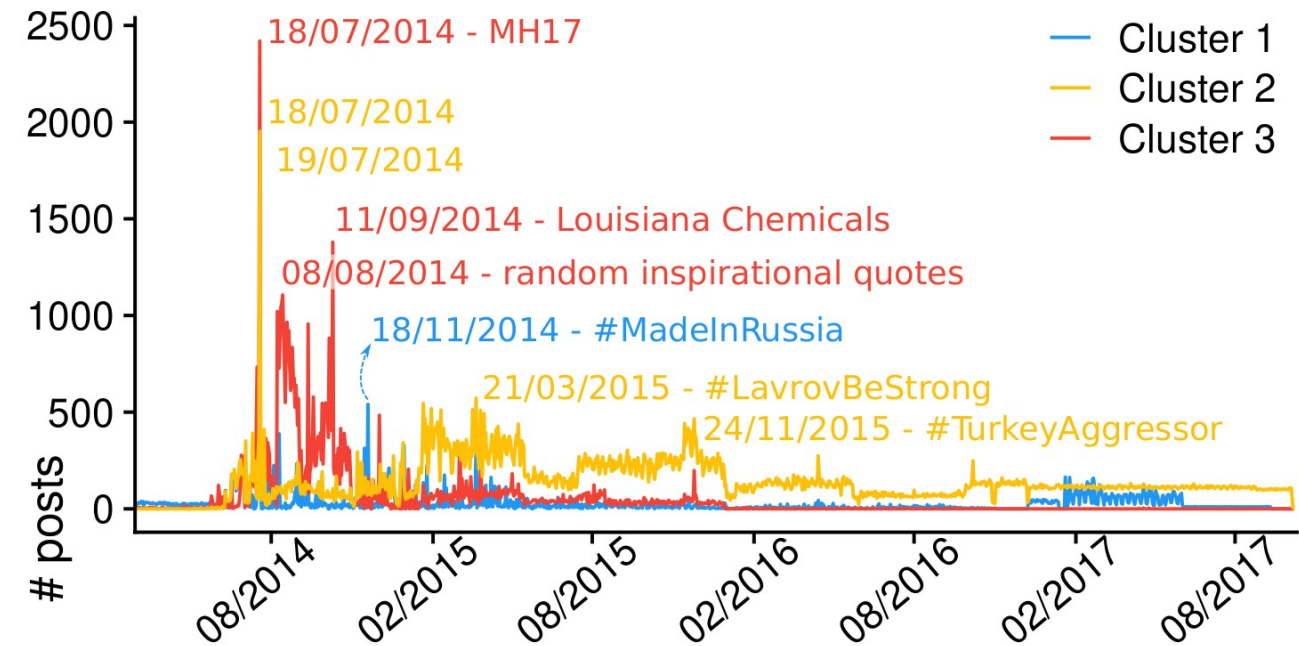
Mathematical generative modelling; Hawkes processes; joint modelling

[https://www.behavioral-ds.science/theme1\\_content/evently/](https://www.behavioral-ds.science/theme1_content/evently/)

# Detect: identify agent types and coordinated behavior



IC-TH clusters IO agents from specific countries based solely on the timing of the cascades in which they participate; it identifies even individual “troll farms”.



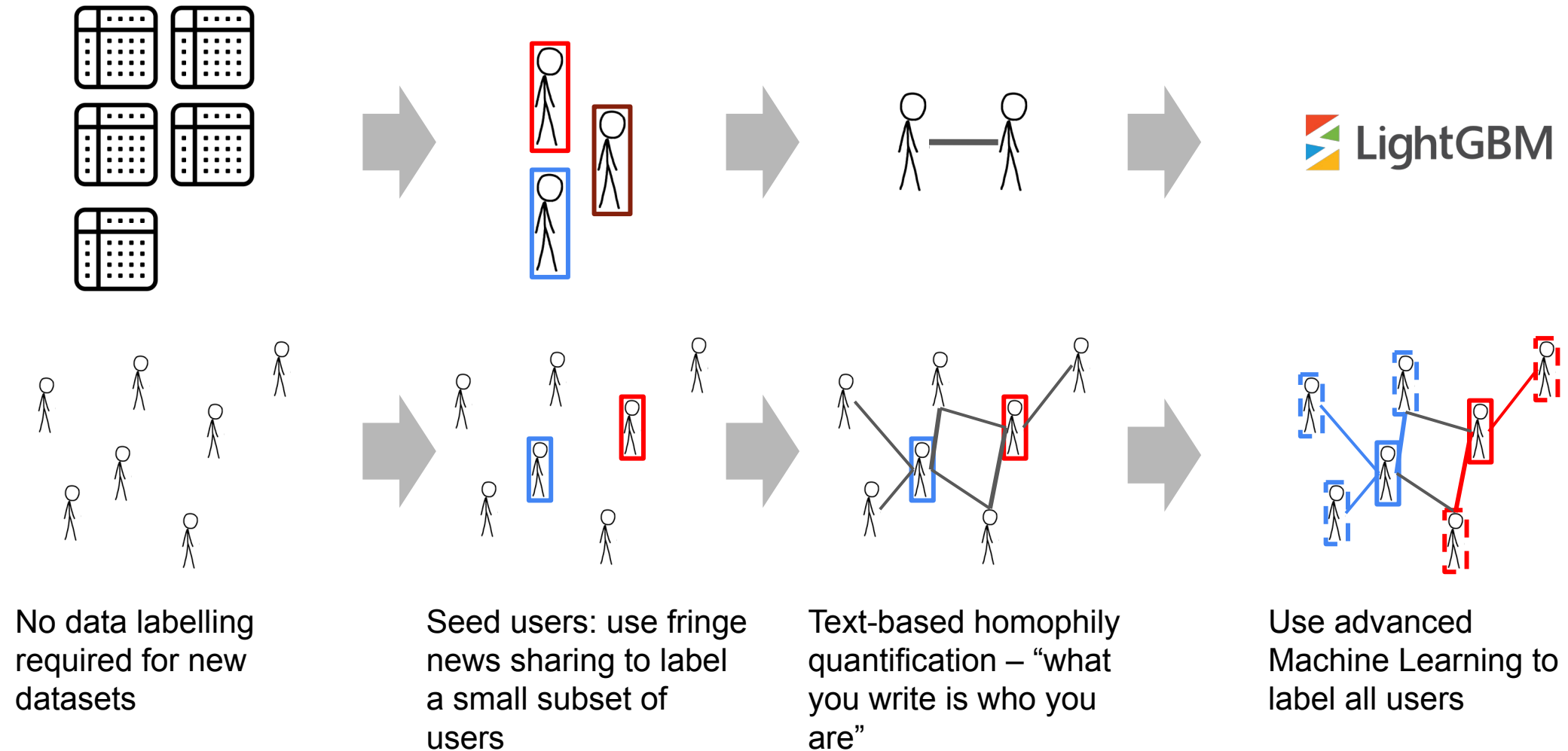
Qualitative investigations uncovers strategies of Russian trolls farms:  
C1: Russian news with patriotic framing;  
C2: Regional and conservative news;  
C3: tweet in English, *#music*, *#usa*, relationship advice



## The technical detail:

Interval-censored Transformer Hawkes; Twitter Moderation Research Consortium dataset; partial data loss

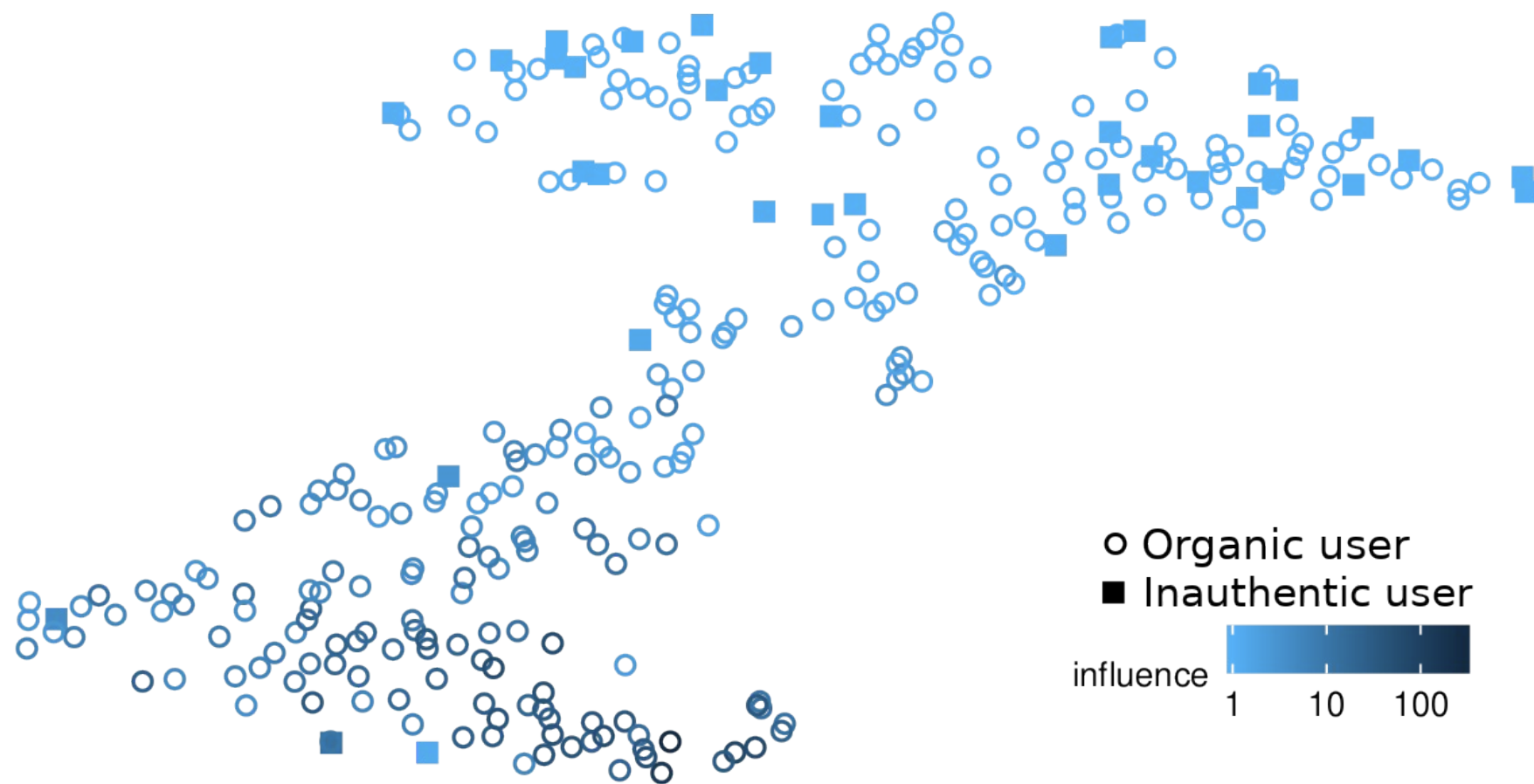
# React: far-right and extremist ideology detection



## The technical detail:

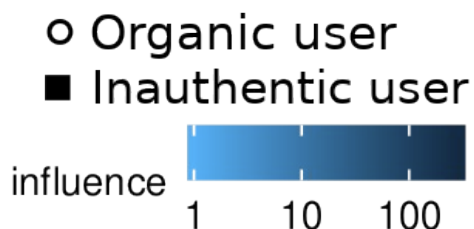
Ideology proxies; homophily lenses (text, follower, URLs); automatic user labelling

# React: Identify influential inauthentic users (bots)



Identify users engaged in influence operations

Estimate their impact on the wider community



**birdspotter**

[https://www.behavioral-ds.science/theme2\\_content/birdspotter/](https://www.behavioral-ds.science/theme2_content/birdspotter/)

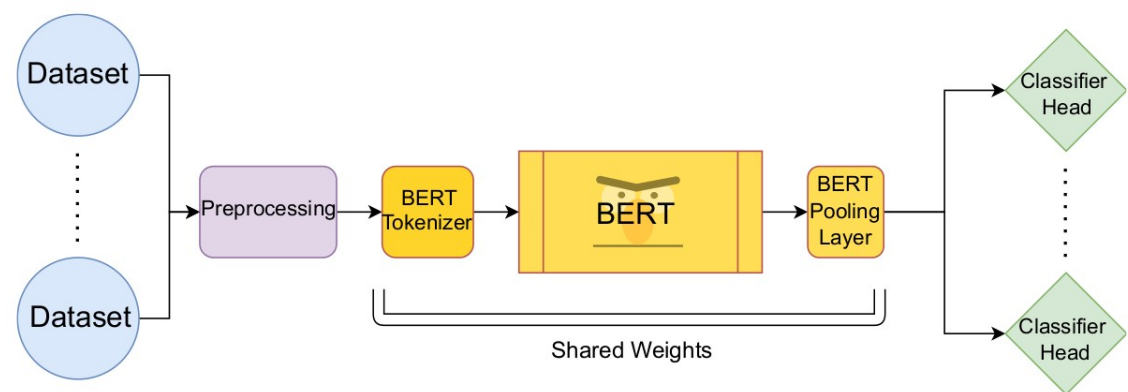


**The technical detail:**

Influence estimation using stochastic modelling; content-free analysis




# React: Detecting Hate Speech in Unseen Domains



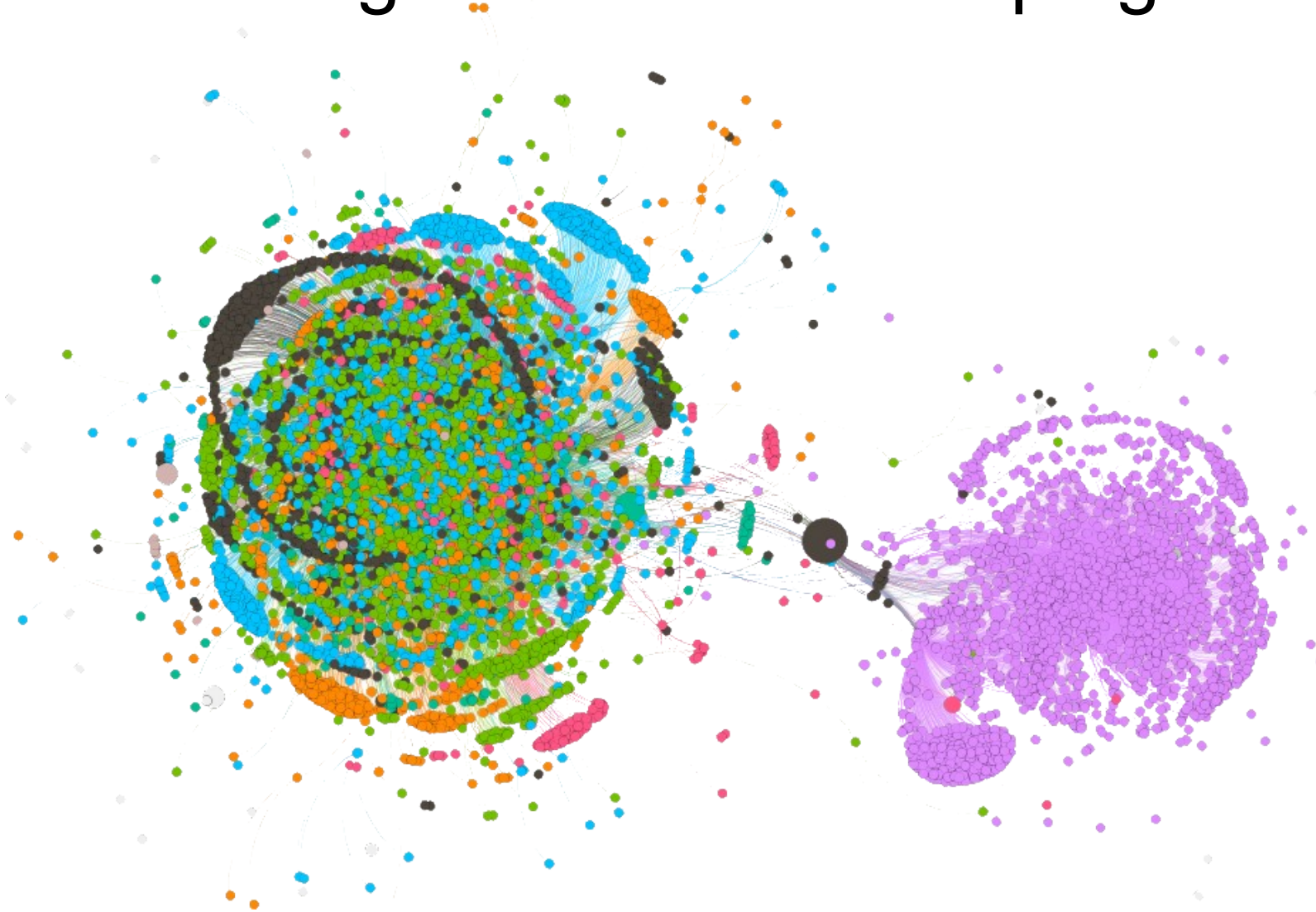
Novel learning paradigm to leverage many disparate datasets to learn a single hate speech representation

Improved performances over the state-of-the-art, generalizable to novel datasets.

		Testing Dataset									# Wins
		DAVIDSON	WASEEM	REDDIT	GAB	FOX	STORMFRONT	MANDL	HATEVAL	PUBFIGS-L	
MTL	MTL-NCH	<b>0.6822</b>	0.3801	<b>0.8456</b>	<b>0.8738</b>	<b>0.6150</b>	<b>0.6826</b>	0.5312	<b>0.6449</b>	0.6175	6
	MTL-MV	0.6455	0.4048	0.8263	0.8660	0.6030	0.6771	0.4834	0.6315	<b>0.6231</b>	1
BERT baseline trained on:	DAVIDSON		0.5556	0.5914	0.6731	0.4932	0.4597	0.5690	0.5414	0.5469	0
	WASEEM	0.6136		0.6000	0.6427	0.5519	0.5356	0.5099	0.5784	0.5611	0
	REDDIT	0.6135	0.4957		0.8083	0.5229	0.5559	0.4900	0.5741	0.5402	0
	GAB	0.5720	0.4595	0.8375		0.5075	0.5645	0.4277	0.5664	0.5185	0
	FOX	0.4285	0.4249	0.4234	0.4651		0.3865	0.4159	0.4490	0.3926	0
	STORMFRONT	0.4533	0.5467	0.5822	0.6487	0.5740		0.5104	0.5664	0.5659	0
	MANDL	0.3336	0.4822	0.4066	0.4582	0.4010	0.3518		0.4546	0.3633	0
	HATEVAL	0.5849	0.5824	0.5700	0.5796	0.5532	0.5466	0.5348		0.5432	0
	PUBFIGS-L	0.6351	<b>0.6048</b>	0.5970	0.6600	0.5546	0.5249	<b>0.5963</b>	0.5858		2

 **The technical detail:**  
Transfer learning; language models fine-tuning;

# Detecting coordinated campaigns



Clear structure with two clusters:  
disinformation (right) and debunking (left)

**Disinformation cluster:** tightly connected,  
coordinated and timed retweeting

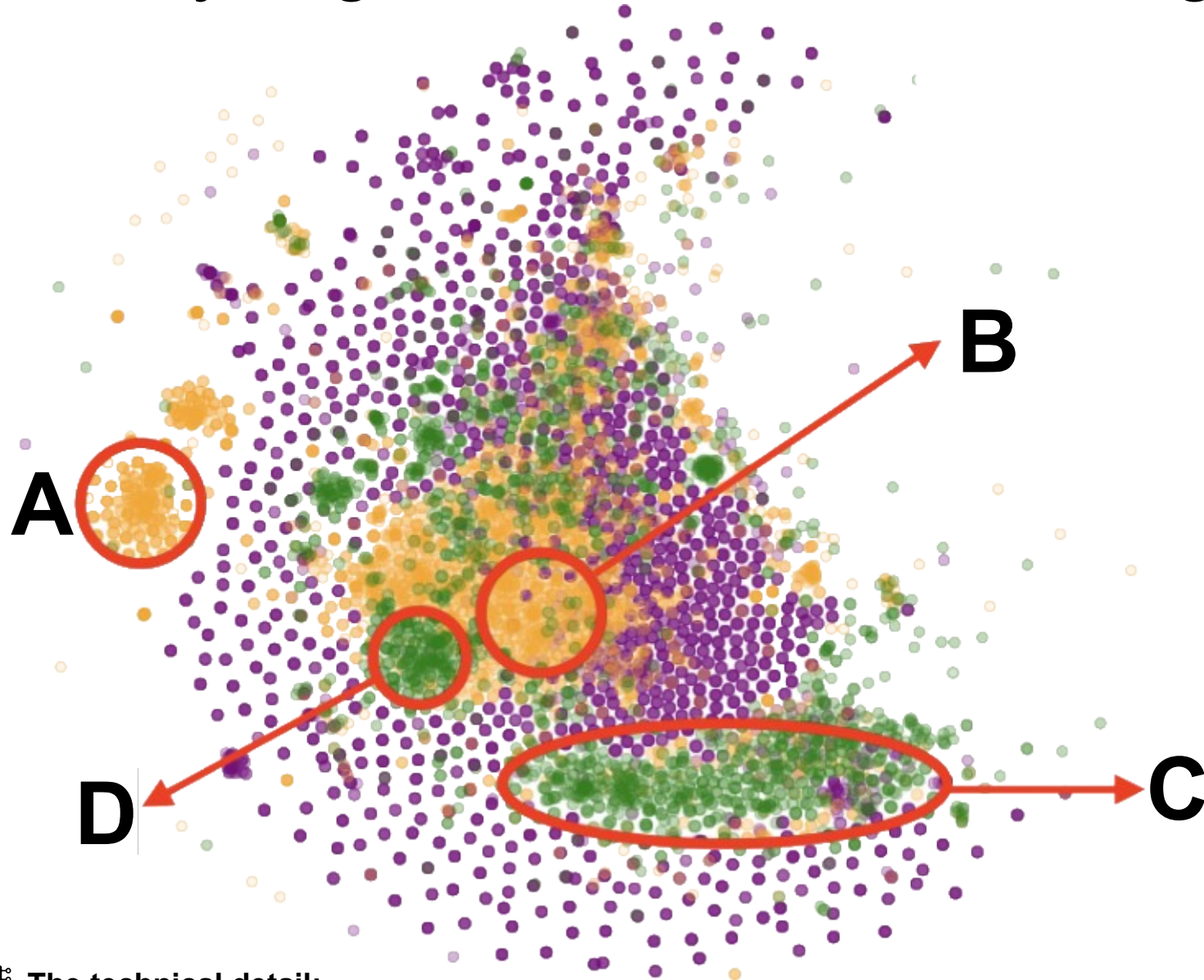
**Debunking cluster:** organic retweeting,  
reactionary, loosely connected, multiple  
communities



## The technical detail:

Map information networks from social media; content, interactions, structure and diffusions analyse; social network analysis

# Analysing coordinated troll strategies



(yellow) right trolls: focused MAGA  
(magenta) left trolls: surround discussion  
(green) news trolls: selective highlighting

**A** – (right trolls) Hillary cannot be trusted  
*#ThingsMoreTrustedThanHillary*

**B** – (right trolls) Mimic black Trump supporters  
*#Blacks4Trump*

**C** – (news trolls) News about violence and civil unrest  
*#news*

**D** – (news trolls) Federal politics, policy and regulation  
*#politics*



The technical detail:

Semantic edit distance; dimensionality reduction; Twitter trolls