# Monitoring and Tagging Hate Speech in Social Media

Sophia Karakeva

Communications and Marketing Executive | Datascouting

Vice President | FIBEP

# Monitoring and Tagging Hate Speech in Social Media

## Some interesting facts about our very "hateful" social world

▸ Tech platforms now **control** the majority of online conversation

▸ Tech platforms have undertaken a shift towards **censorship** and **moderation**

▸ The American tradition of **free speech** on the internet is **no longer viable**

▸ Users, governments, and tech firms are all **behaving badly**

▸ **Free speech** has become a social, economic and political **weapon**

▸ **Racists, misogynists and oppressors** are allowed a voice

*Source: The Good Censor | Cultural context report, March 2018*

# Monitoring and Tagging Hate Speech in Social Media

## What is hate speech?



▶ No universal definition

▶ Freedom of expression vs hate speech

▶ Relying on hate speech policies

*"* *#HateSpeech is just a word for pussies who can't stand up for themselves so they need legislation on language*

*"* *#HateSpeech is free speech. Otherwise you're not allowed to hate nazis.*

*"* *Truth is now considered #HateSpeech. Thanks progressives.*

## What is the role of social media when it comes to hate speech?

▸ Code of conduct

▸ Still not there

▸ Cooperation but diversity



Countering illegal
hate speech online
#Noplace4hate

# Monitoring and Tagging Hate Speech in Social Media

What is the role of social media when it comes to hate speech?
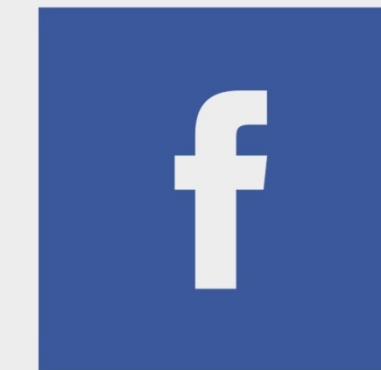
## Still not there

**YouTube**

" *Sometimes our systems get it wrong … Our system sometimes make mistakes in understanding context and nuances…*

**Johanna Wright**
*YouTube VP of Product Management (2017)*

**Twitter**

" *We see voices being silenced on Twitter every day. We've been working to counteract this for the past 2 years… We prioritized this in 2016. We updated our polices and increased the size of our teams. It wasn't enough.*
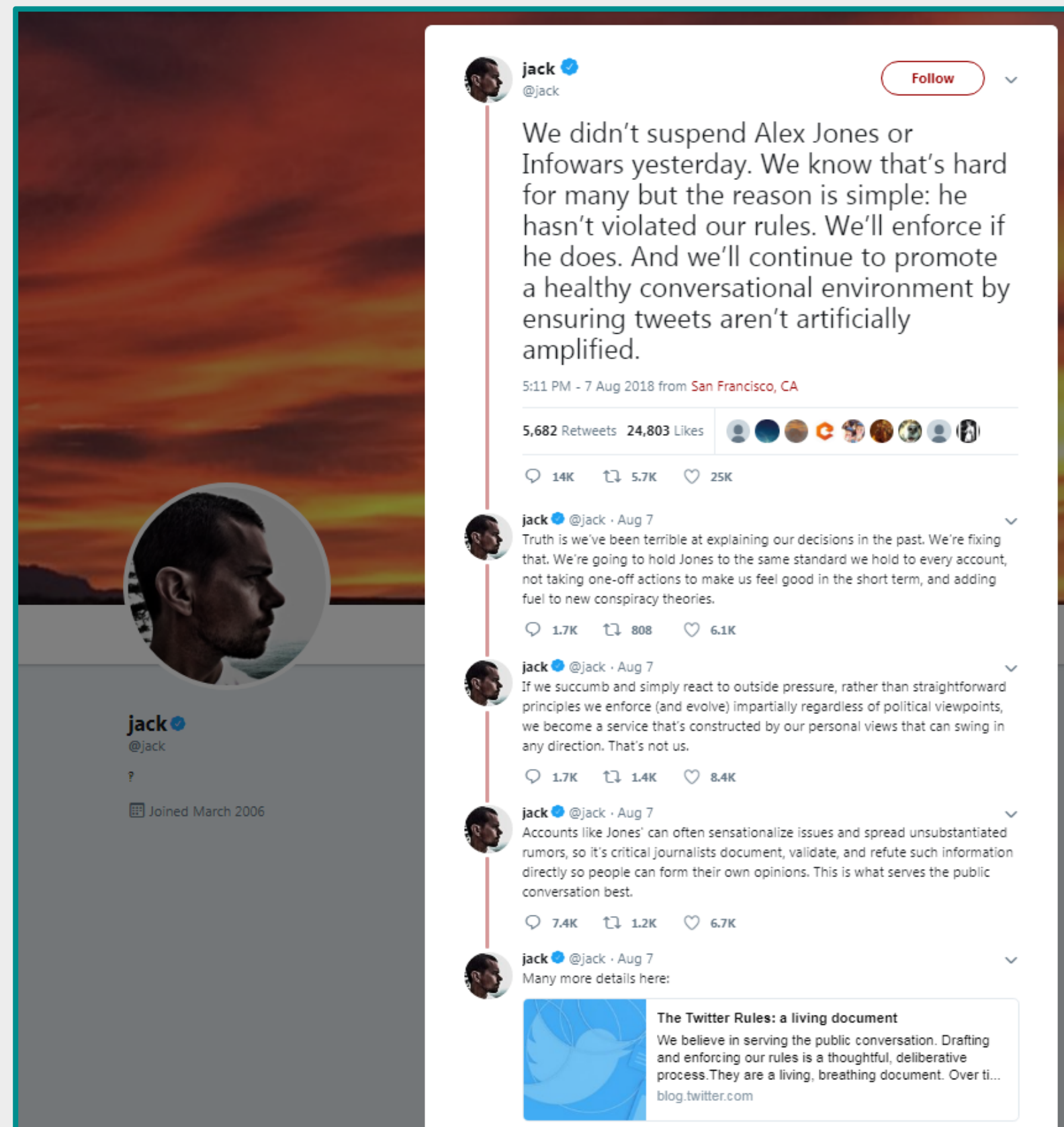
**Jack Dorsey**
*Twitter CEO (2017)*

**f**

" *… It's much easier to build an AI system that can detect a nipple that it is to determine what is linguistically hate speech.*

**Mark Zuckerberg**
*Facebook CEO (2018)*

# Monitoring and Tagging Hate Speech in Social Media

## What is the role of social media when it comes to hate speech?

### Collaboration

Alex Jones against CNN Senior Media Correspondent Brian Stelter: *"enemy . . . drunk on children's blood," "You will pay! You will fall!"*

Alex Jones against people of Muslim faith: *"And we are going to stop you again! Do you understand?!""*

Alex Jones against drag queens: *"We're going to destroy you."*

*Source: CNN*

# Monitoring and Tagging Hate Speech in Social Media

## The legal and discursive characteristics of hate speech

Create unmediated marketplaces of ideas **VS** Create well-ordered spaces for safety and civility

▶ Legal liability

▶ Danger of over-regulation

▶ Platforms vs editors / publishers



> *Platforms have to deny that they're media companies in order to retain their immunity from liability. But at the same time, they're exercising more influence as media companies...*
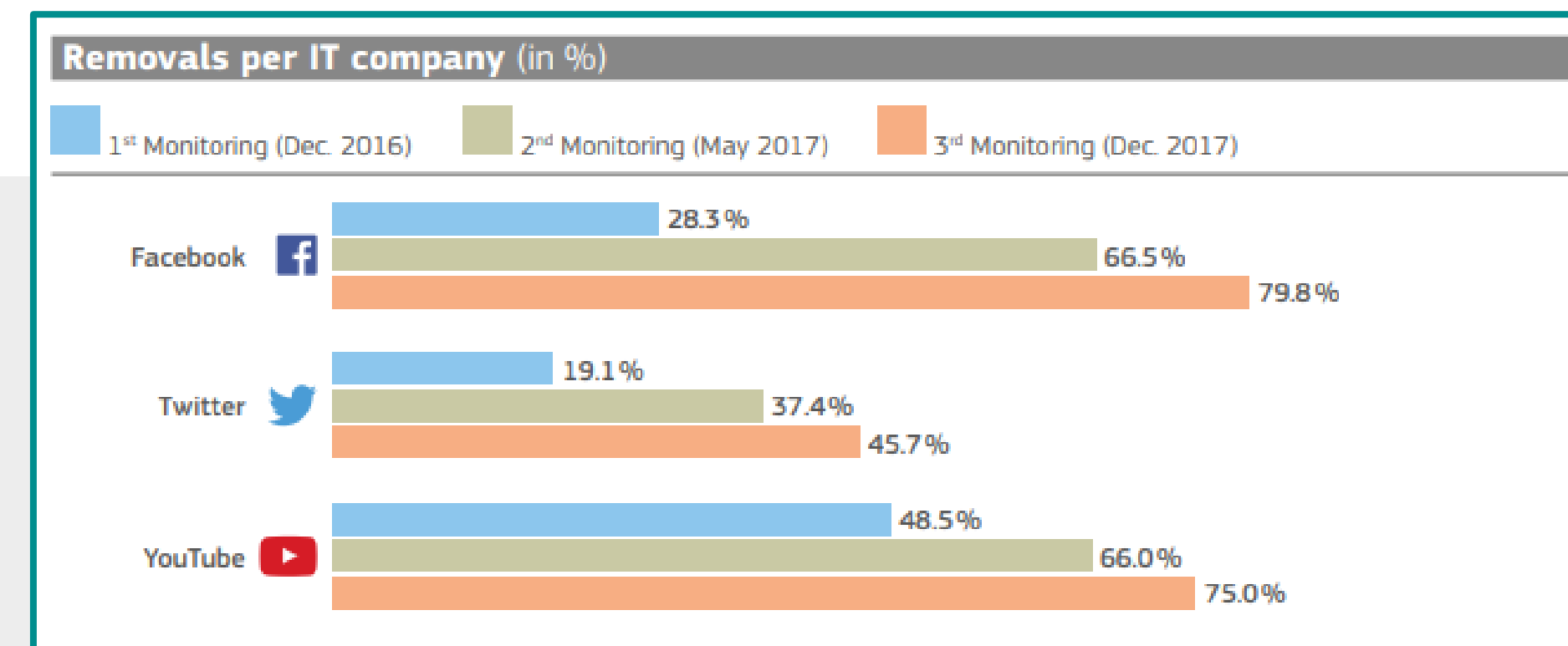
**Jeffrey Rosen**
*Professor of Law at the George Washington University and legal affairs editor of the The New Republic CEO (2018)*

# Monitoring and Tagging Hate Speech in Social Media

## Where does hate speech spread the most?
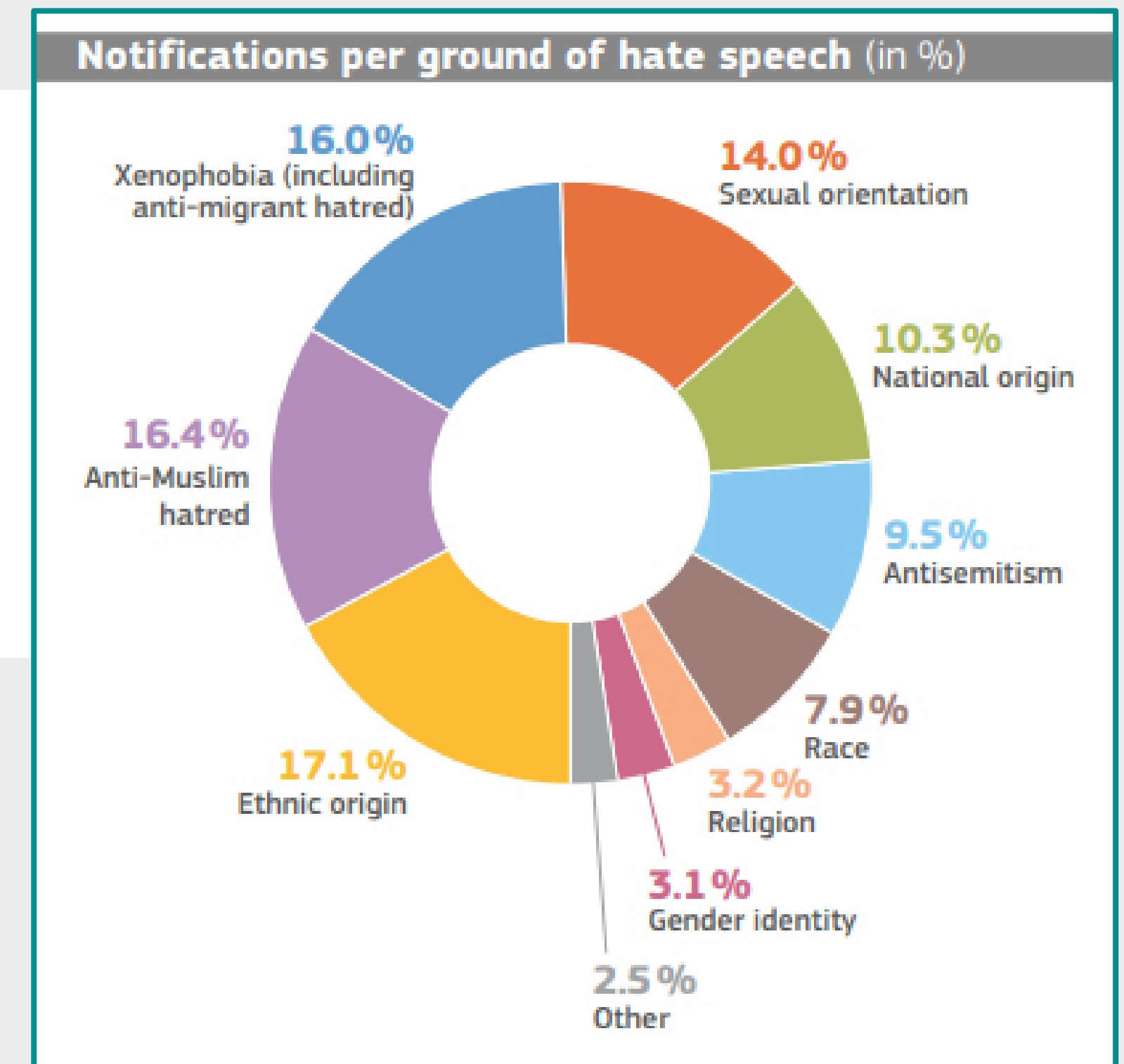


*Source:*
*The Flash Eurobarometer, September 2018*
*The 3rd monitoring of EU's Code of Conduct, January 2018*
*2017 Pew Research Center survey about online harassment*

▸ Social networks the most common online daily activity

▸ Malta has the highest online hate speech in the EU

▸ Facebook has the largest amount of notifications

▸ 70% of hate speech content was removed in 2017

# Monitoring and Tagging Hate Speech in Social Media

## Who is most targeted?

- Ethnic origin, anti-Muslim hatred, & xenophobia

- Gender matters to social media

- The developing world

European Commission

Code of Conduct on countering illegal hate speech online
**Results of the 3rd monitoring exercise**

Fact sheet | January 2018

**Notifications per ground of hate speech** (in %)

- **16.0%** Xenophobia (including anti-migrant hatred)
- **14.0%** Sexual orientation
- **10.3%** National origin
- **9.5%** Antisemitism
- **7.9%** Race
- **3.2%** Religion
- **3.1%** Gender identity
- **2.5%** Other
- **17.1%** Ethnic origin
- **16.4%** Anti-Muslim hatred

ifrro | ATHENS WORLD CONGRESS 2018

DATASCOUTING
Actionable Information

FIBEP

# Monitoring and Tagging Hate Speech in Social Media

---

## Who is most targeted?

## Gender matters to social media

---

" I think Twitter is the worst of the social media platforms, just because of the quickened and masked flow [of abuse] that happens. The content feels pretty similar across the platforms but the sheer volume of it on Twitter is what's different. "

Jessica Valenti, US journalist and writer



$#*%!

#ToxicTwitter

"It's like public lynching. It has made me frightened for my physical safety when I am out in the streets"

Amberin Zaman, Turkish journalist

# Monitoring and Tagging Hate Speech in Social Media

---

## Who is most targeted?

## Digital Colonization

---



> Facebook has been a useful instrument for those seeking to spread hate, in a context where for most users Facebook is the internet

United Nations,
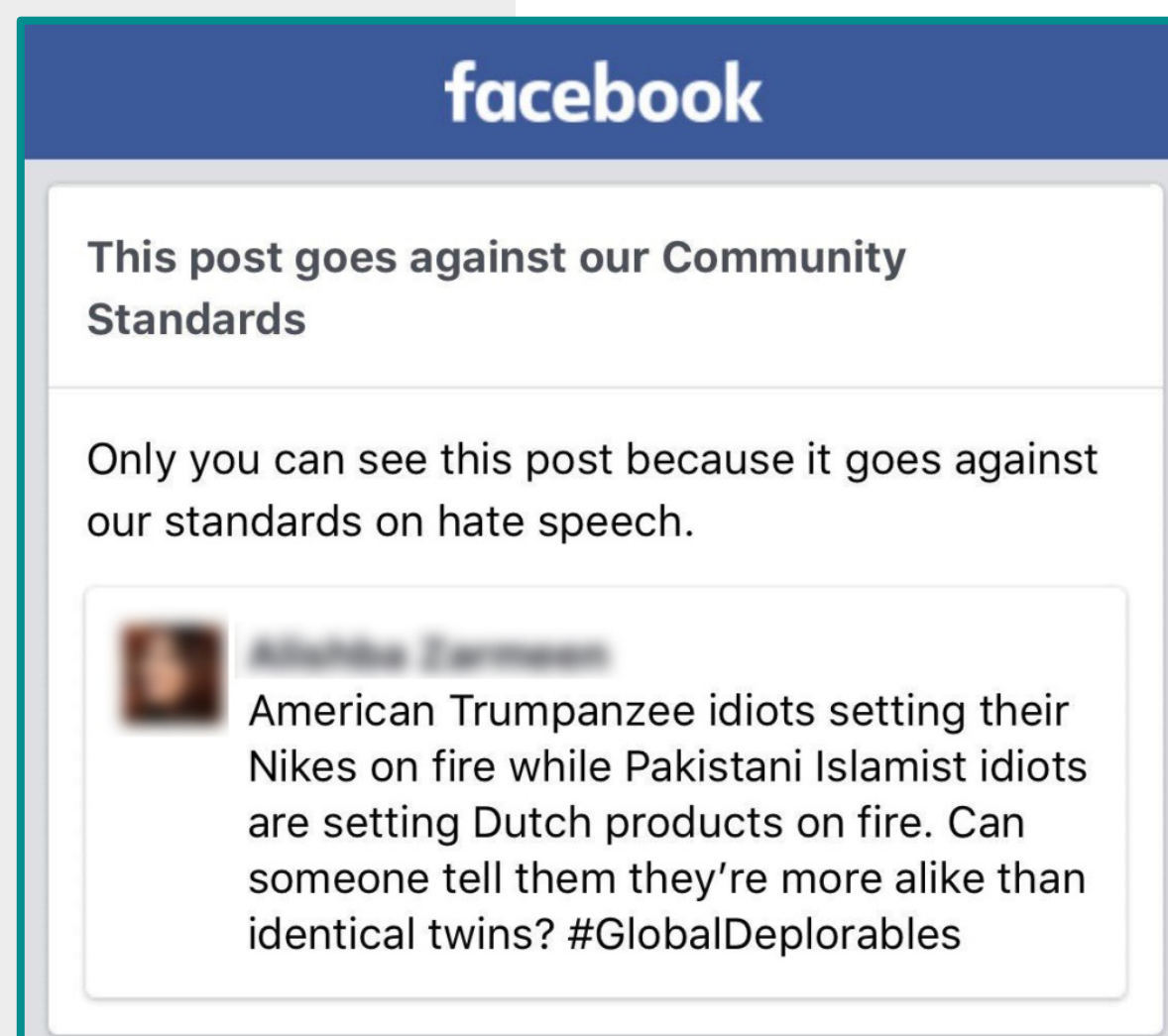Human Rights Council, Report of the Independent International Fact-Finding Mission on Myanmar, August 2018

*Source: John Oliver, Last Week Tonight Show, 23 September 2018*

# Monitoring and Tagging Hate Speech in Social Media

## How are tech companies mismanaging hate speech?

Hate me, Hate me not

▸ Commercialized conversation

▸ Inconsistent interventions

▸ Lack of transparency

▸ Ineffective automation

▸ Underplaying the issue

facebook

**This post goes against our Community Standards**

Only you can see this post because it goes against our standards on hate speech.

Alishba Zarmeen
American Trumpanzee idiots setting their Nikes on fire while Pakistani Islamist idiots are setting Dutch products on fire. Can someone tell them they're more alike than identical twins? #GlobalDeplorables

Follow

This morning I woke up to a rape and death threat directed at my 5 year old daughter. That this is part of my work life is unacceptable.

1:04 PM - 27 Jul 2016

1,438    1,348

Replying to @maggieNYT

**Fuck you fuck you fuck you fuck you.** **Maggie** haberman is a lying bitch and an enabler of racism, misogyny, fascism and treason.

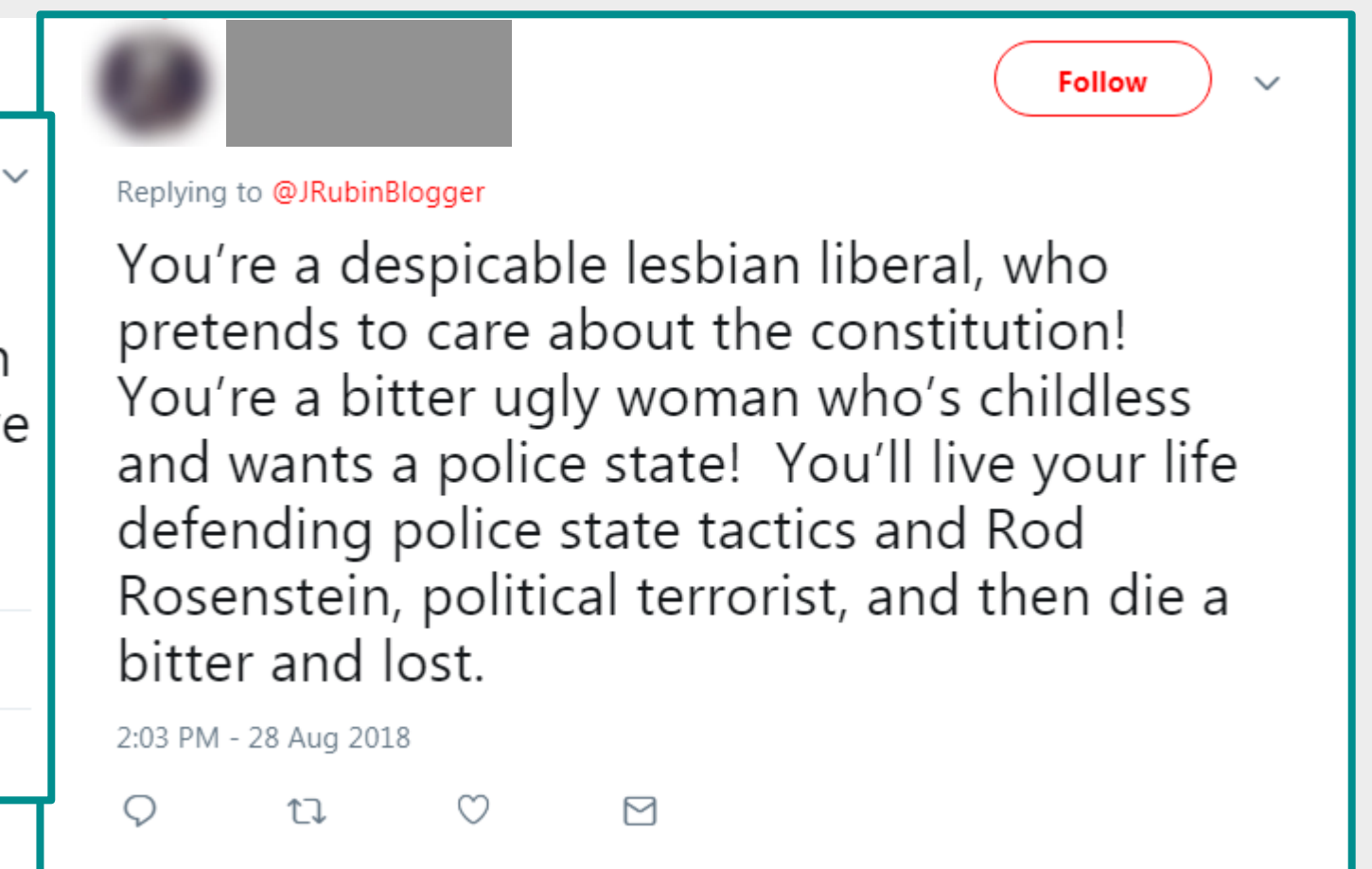ifrro | ATHENS WORLD CONGRESS 2018    DATASCOUTING Actionable Information    FIBEP

# Monitoring and Tagging Hate Speech in Social Media

Hate me, hate me not

**Commercialized conversation**

▶ **Hate me not**: hate speech is not allowed on our platforms

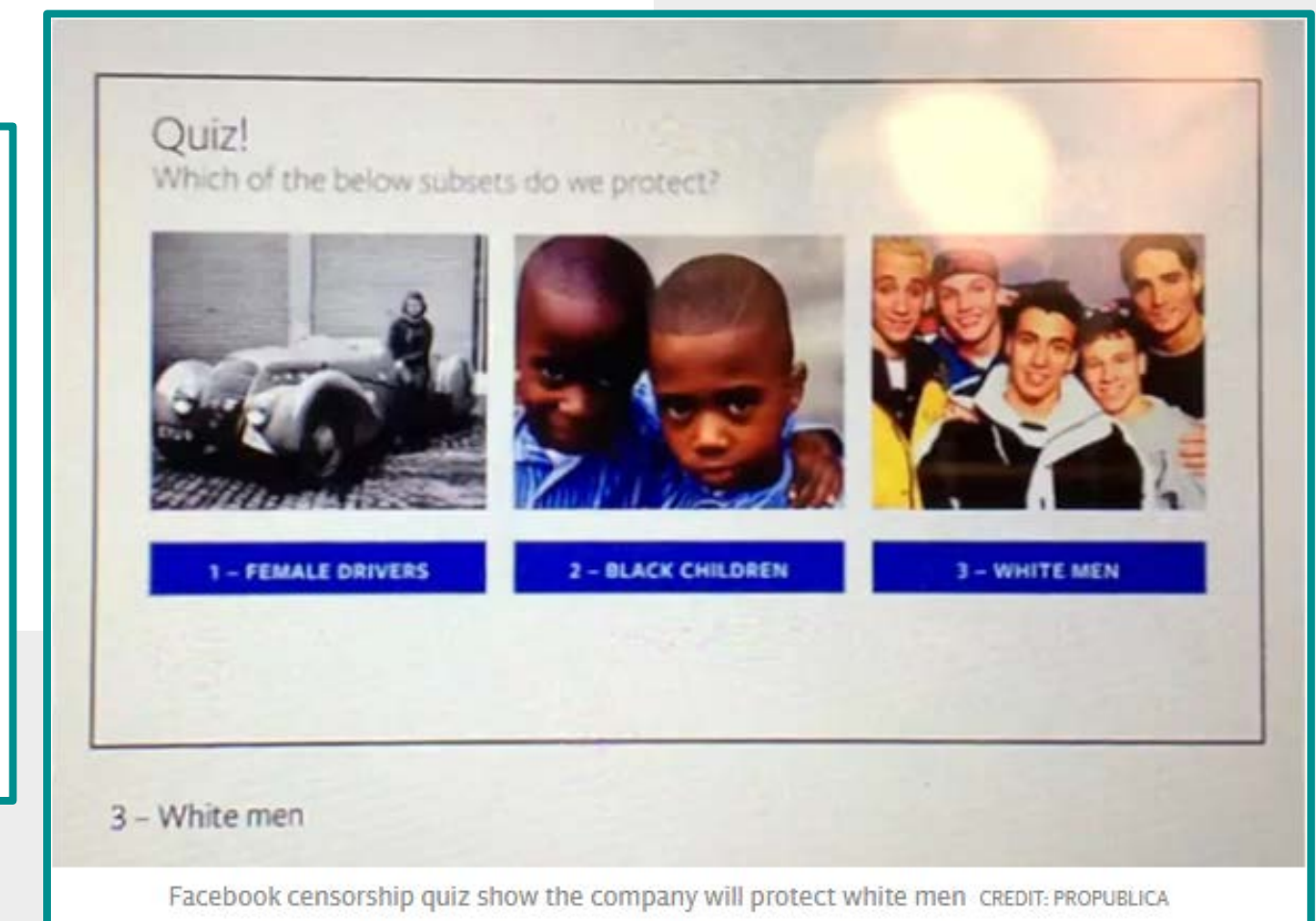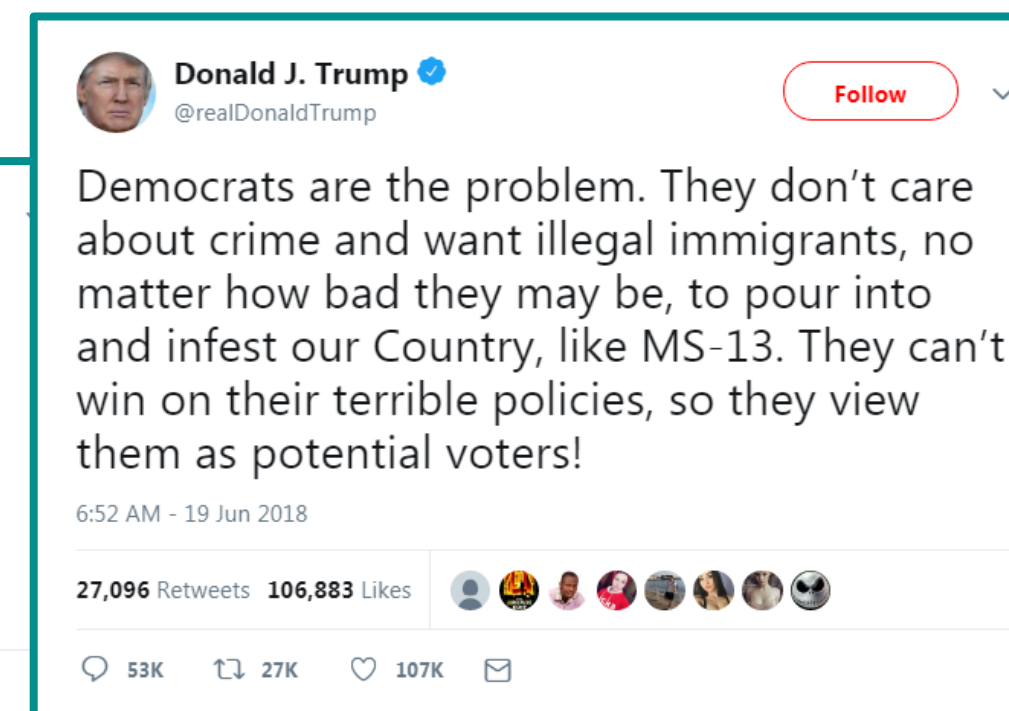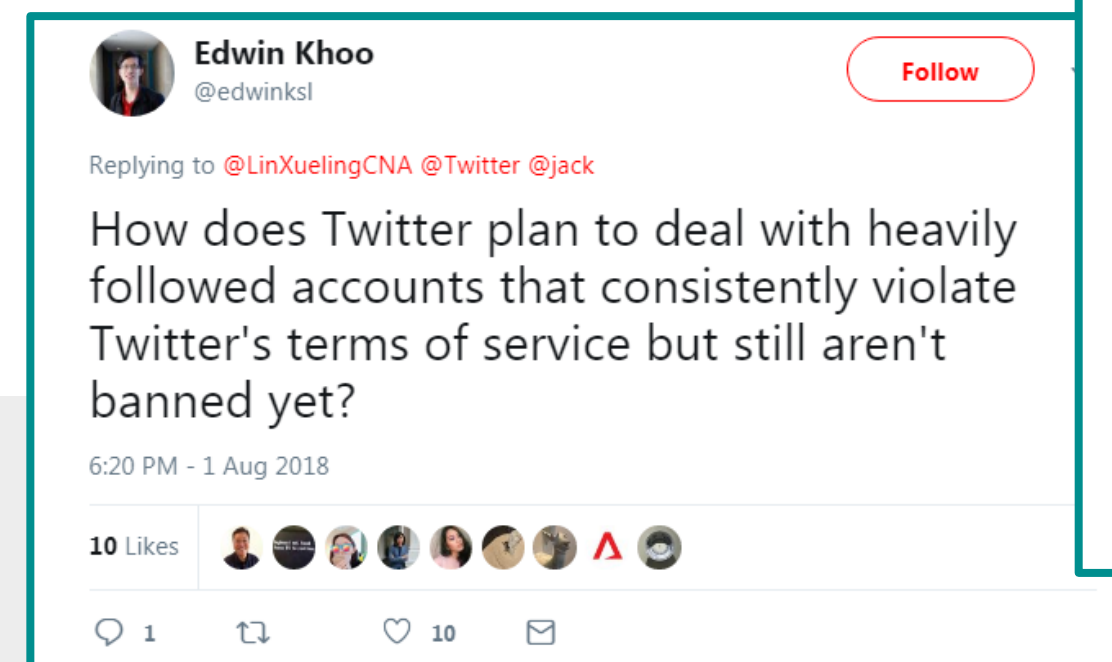▶ **Hate me**: hate speech is often the most engaged content

Replying to @redsteeze @nypost

Twitter helps you hate people you don't even know; Facebook helps you hate people you've known all your life.

5:59 AM - 21 Aug 2018

1 Retweet 1 Like

Replying to @JRubinBlogger

You're a despicable lesbian liberal, who pretends to care about the constitution! You're a bitter ugly woman who's childless and wants a police state! You'll live your life defending police state tactics and Rod Rosenstein, political terrorist, and then die a bitter and lost.

2:03 PM - 28 Aug 2018

# Monitoring and Tagging Hate Speech in Social Media

---

Hate me, hate me not

**Inconsistent interventions & lack of transparency**

---

▸ **Hate me not**: we take strong measures to take hate speech down

▸ **Hate me**: not everything has the same value





> Edwin Khoo
> @edwinksl
> Replying to @LinXuelingCNA @Twitter @jack
> How does Twitter plan to deal with heavily followed accounts that consistently violate Twitter's terms of service but still aren't banned yet?
> 6:20 PM - 1 Aug 2018
> 10 Likes



> Donald J. Trump
> @realDonaldTrump
> Democrats are the problem. They don't care about crime and want illegal immigrants, no matter how bad they may be, to pour into and infest our Country, like MS-13. They can't win on their terrible policies, so they view them as potential voters!
> 6:52 AM - 19 Jun 2018
> 27,096 Retweets  106,883 Likes
> 53K   27K   107K



Quiz!
Which of the below subsets do we protect?

1 – FEMALE DRIVERS   2 – BLACK CHILDREN   3 – WHITE MEN

3 – White men

Facebook censorship quiz show the company will protect white men  CREDIT: PROPUBLICA

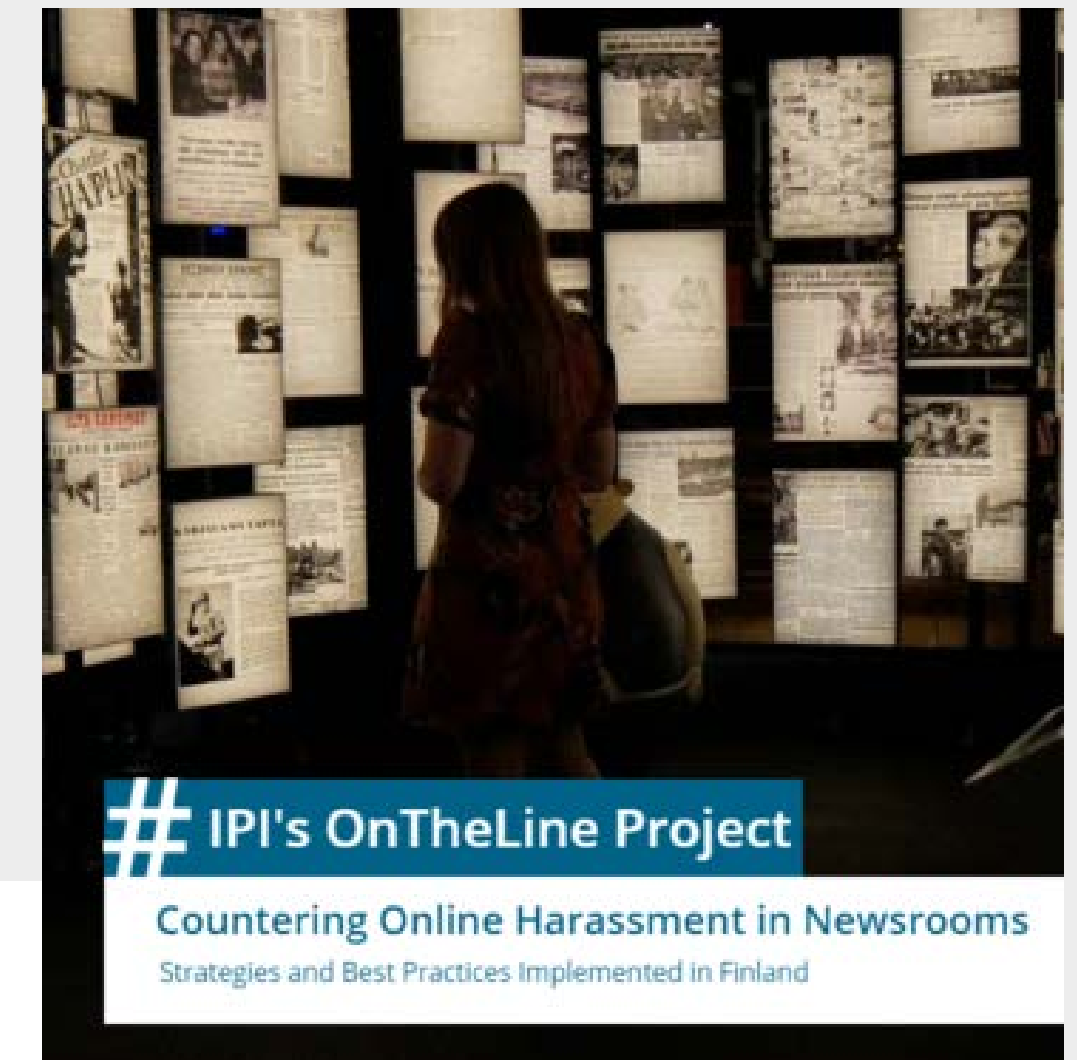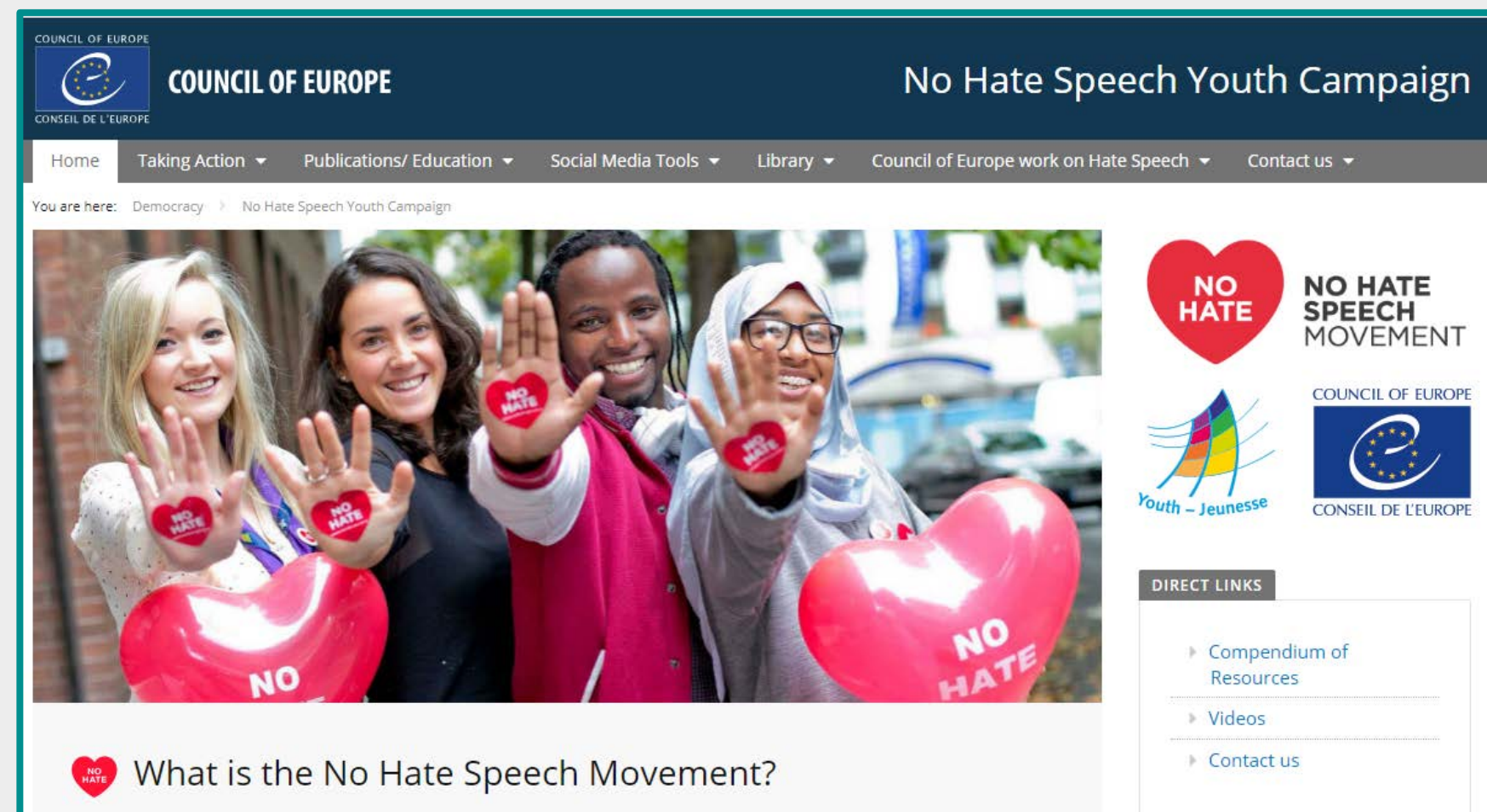# Monitoring and Tagging Hate Speech in Social Media

Hate me, hate me not

**Ineffective automation**

▶ **Hate me**: hate speech is highly reactive with high error rate

▶ **Hate me not**: automation and AI tools



Most common words found in Hate Comments

Most common words found in Non-hate Comments

This data representation shows common terms used in hateful comments on Reddit, Twitter, and other social media sites compared to non-hateful comments. (courtesy of the Anti-Defamation League)

# Monitoring and Tagging Hate Speech in Social Media

Hate me, hate me not

**Underplaying hate speech**





▸ **Hate me**: hate speech becomes the "new normal"

▸ **Hate me not**: education, awareness campaigns, media ethics, regulation

# Monitoring and Tagging Hate Speech in Social Media

## Typology and coding

- ▸ Scale vs human annotation
- ▸ Typology and coding
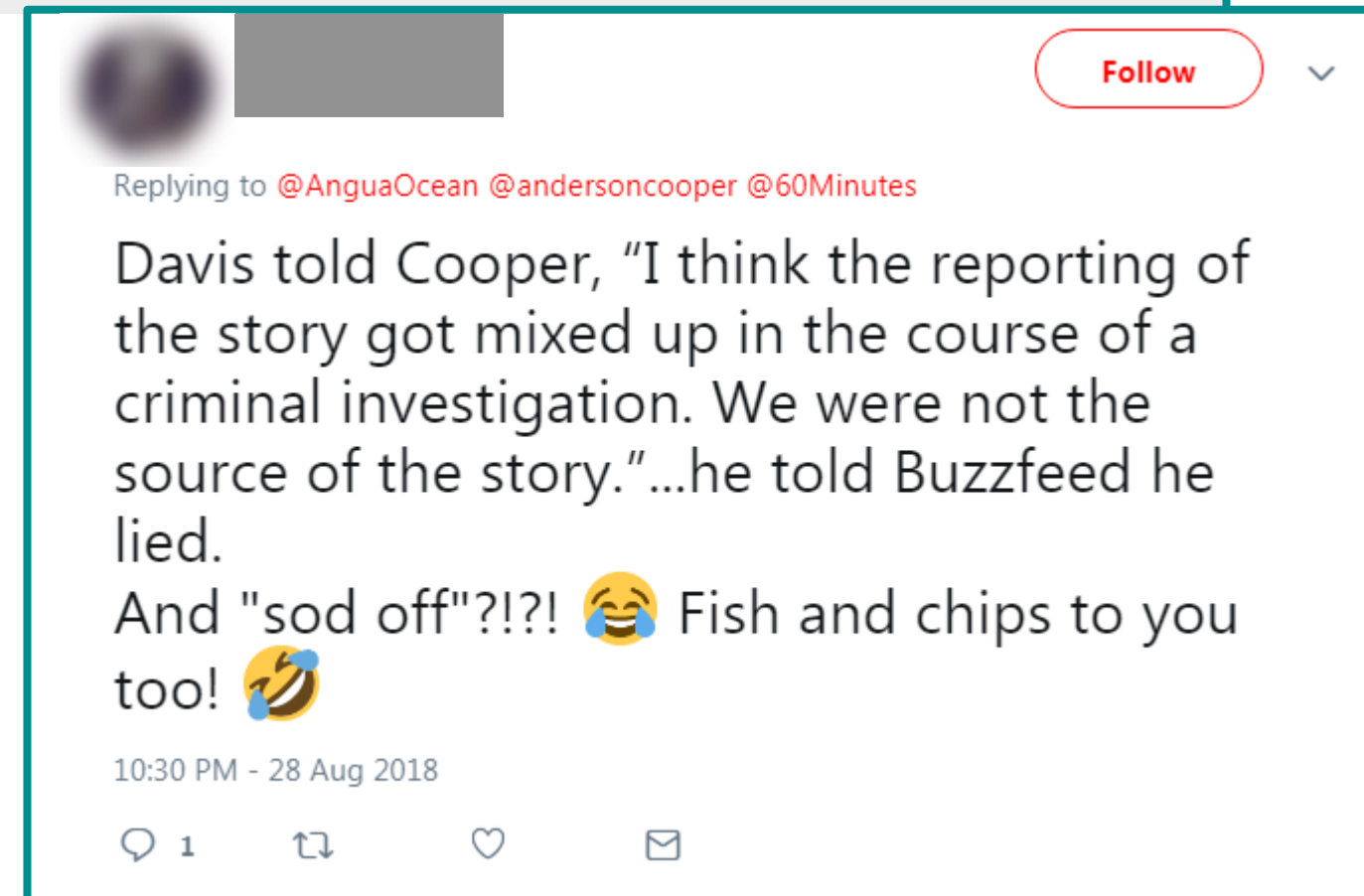- ▸ Main challenges

### Hate Speech

Highlights

- Q: What is our stance on white supremacy, white nationalism and white separatism?
  - We don't allow praise, support and representation of **white supremacy** as an ideology. Eg. "White supremacy is the right thing"; "I am a white supremacist"; "Join the next White Supremacy rally!"
  - We allow praise, support and representation of white nationalism as an ideology. Eg. "White nationalism is the only way"; "I am a proud white nationalist"
  - We allow praise, support and representation of white separatism as an ideology. Eg. "White separatism is the perfect solution to America's problems"; "I am a white separatist". By the same token, we allow to call for the creation of white ethno-states (Eg. "The US should be a white-only nation")

# Monitoring and Tagging Hate Speech in Social Media

## Typology and coding

### Coding

- A global set of policies to achieve consistency
- YES or NO **OR** multiple choice questions
- Test cases as words take on new meaning
- Offensive speech vs hateful territory

**Follow**

Replying to @AnguaOcean @andersoncooper @60Minutes

Davis told Cooper, "I think the reporting of the story got mixed up in the course of a criminal investigation. We were not the source of the story."...he told Buzzfeed he lied.
And "sod off"?!?! 😂 Fish and chips to you too! 🤣

10:30 PM - 28 Aug 2018



**Migrants:** people who cross an international border with intent to establish residency in a new country, regardless of whether their motivation is economic or political.

# Monitoring and Tagging Hate Speech in Social Media

Typology and coding

## Main challenges

- ▸ Language diversity

- ▸ Ineffective automation

- ▸ Incubating fake news and misinformation

- ▸ Commercialized conversation

- ▸ Slow corrections – unfair censorship

- ▸ Global inconsistency

- ▸ Geopolitical situations

- ▸ Historical and cultural implications

# Monitoring and Tagging Hate Speech in Social Media

## What's being done to fight it?

- ▸ Expanded hate speech rules
- ▸ More human annotation
- ▸ Sophisticated technology
- ▸ Academic & scientific research

**Rights, Equality and Citizenship Programme**

European Commission

**DACHS**

**A Data-Driven Approach to Countering Hate Speech**

DATASCOUTING
actionable information

European Journalism Centre
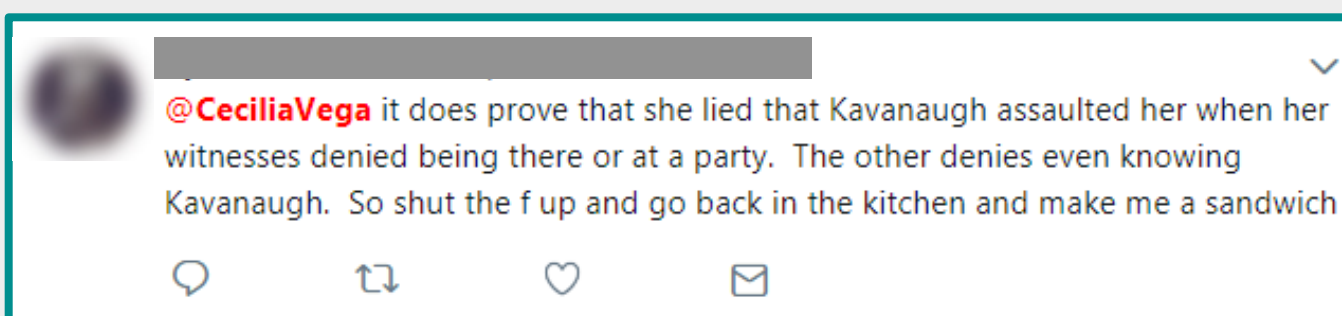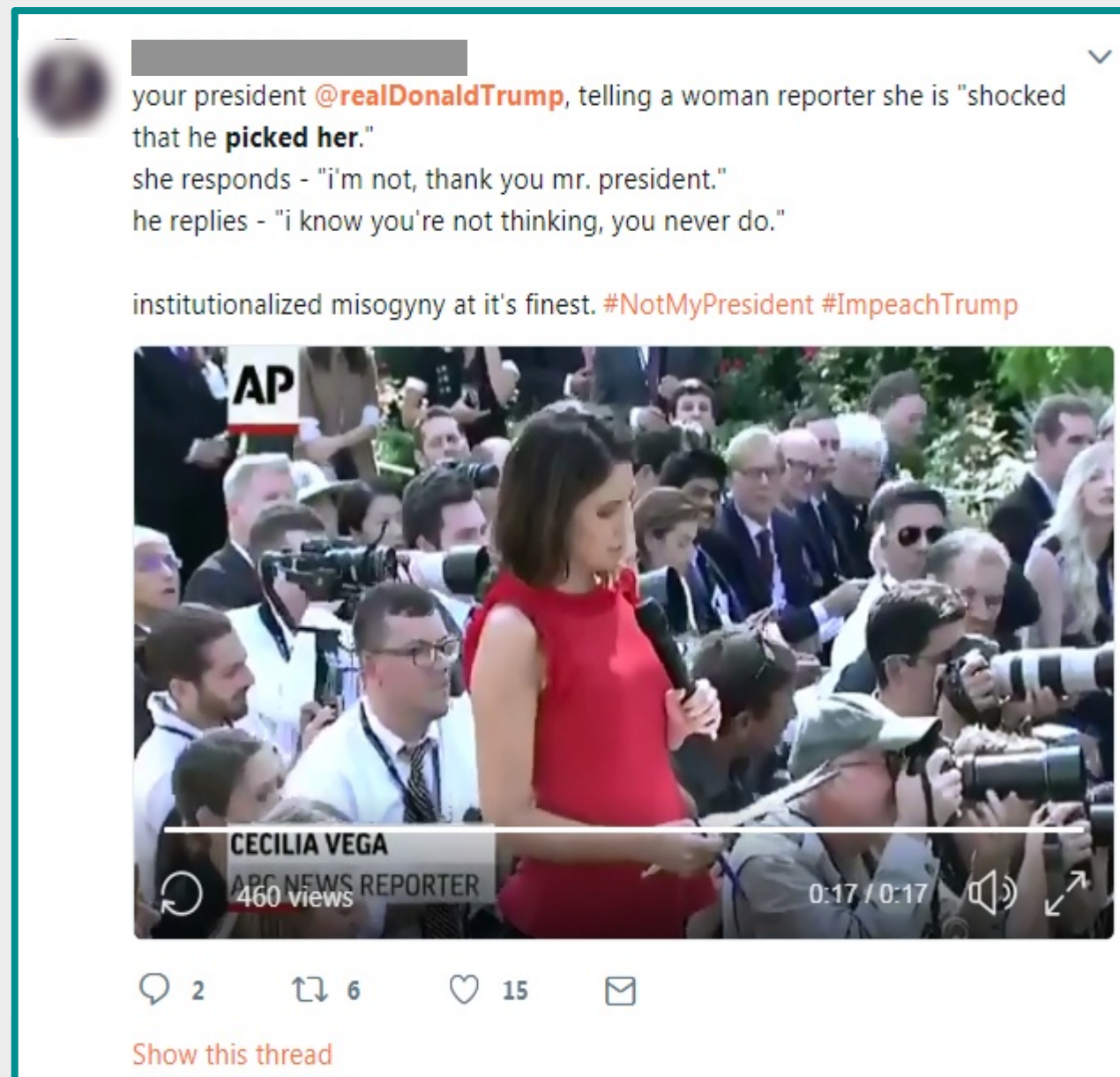
# Monitoring and Tagging Hate Speech in Social Media

## Hate speech against journalists



- ▶ Special targets of hate speech
- ▶ Public personas
- ▶ Report on somebody else
- ▶ Fake news runs faster than the truth

## Impact of hate speech on journalists

▸ Undermines the role of journalism

▸ Impacts trust

▸ Incites criticism – instigates violence

▸ Intimidation and harassment

▸ Self-censorship, mental well-being

▸ Fear of repercussion

– jailed, disappeared, abducted, killed



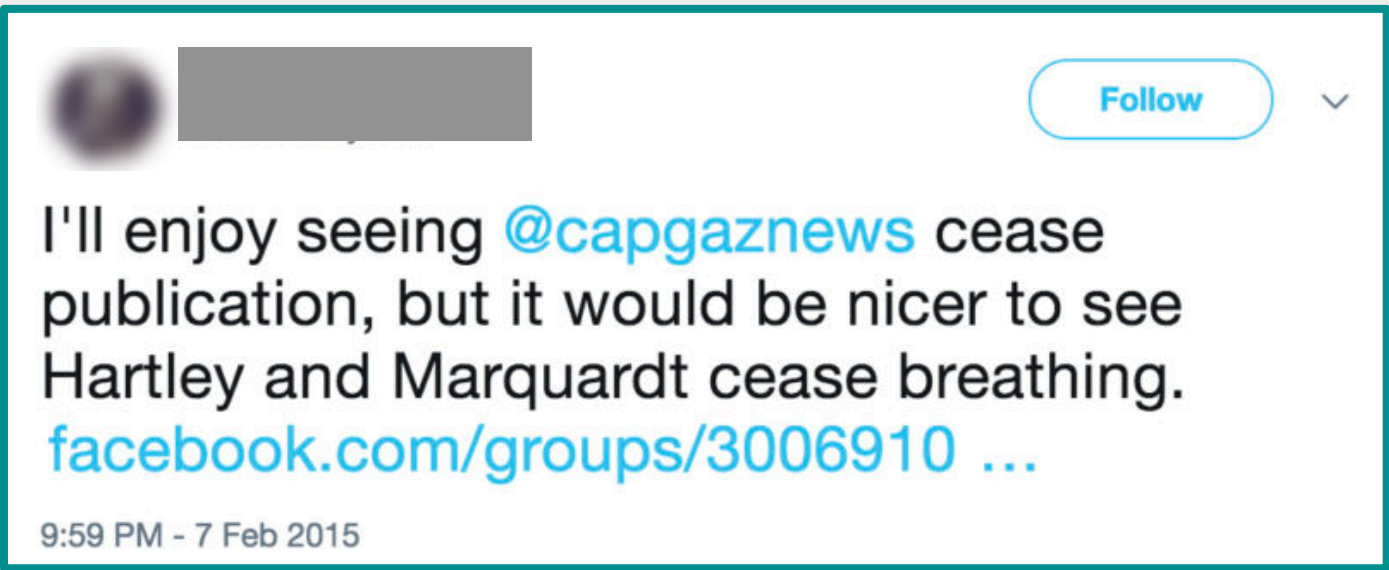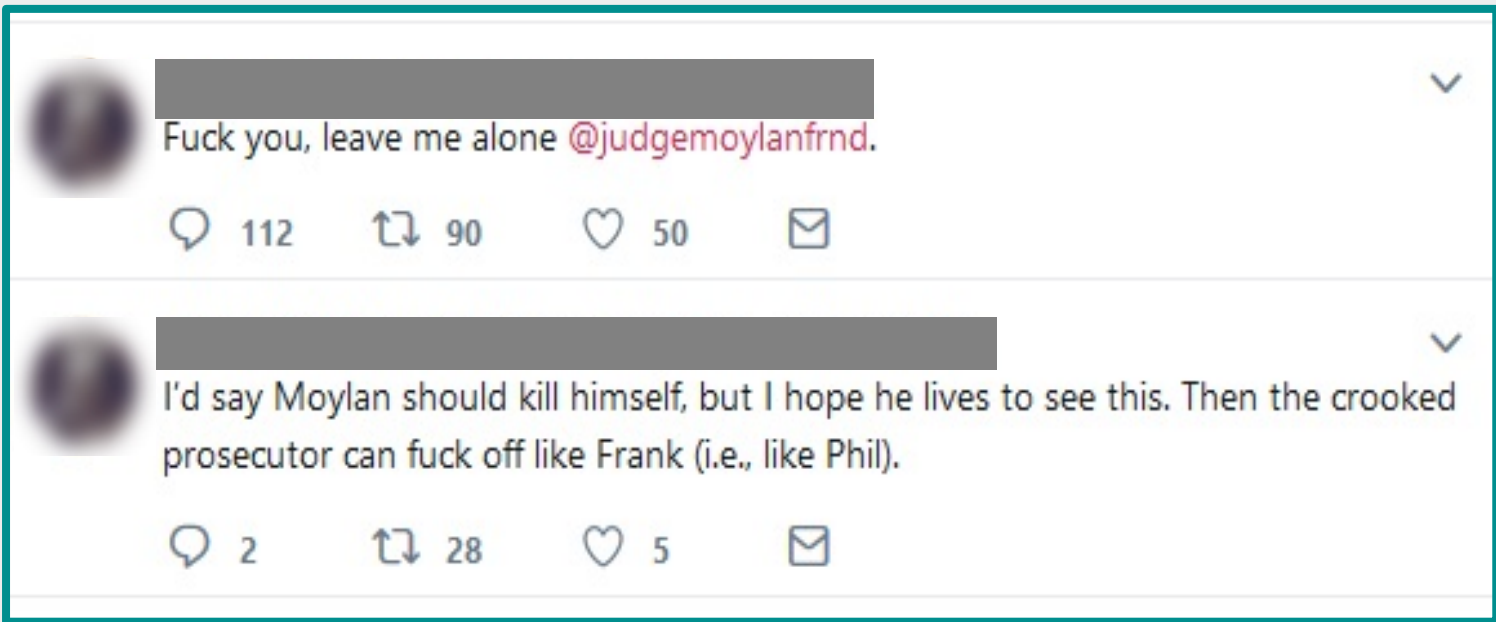**"Kill these journalists once and for all"**

In the summer of 2017, a fearful message was spread on social networks in Togo: *"Kill all these journalists once and for all."* It was accompanied by photos of four journalists, pasted onto images of pigs. The journalists' personal data was disseminated and they were accused by their critics of supporting the Lomé regime.

# Monitoring and Tagging Hate Speech in Social Media

## From malicious threats to actual intent



From left up: John McNamara, Rob Hiaasen, Gerald Fischman
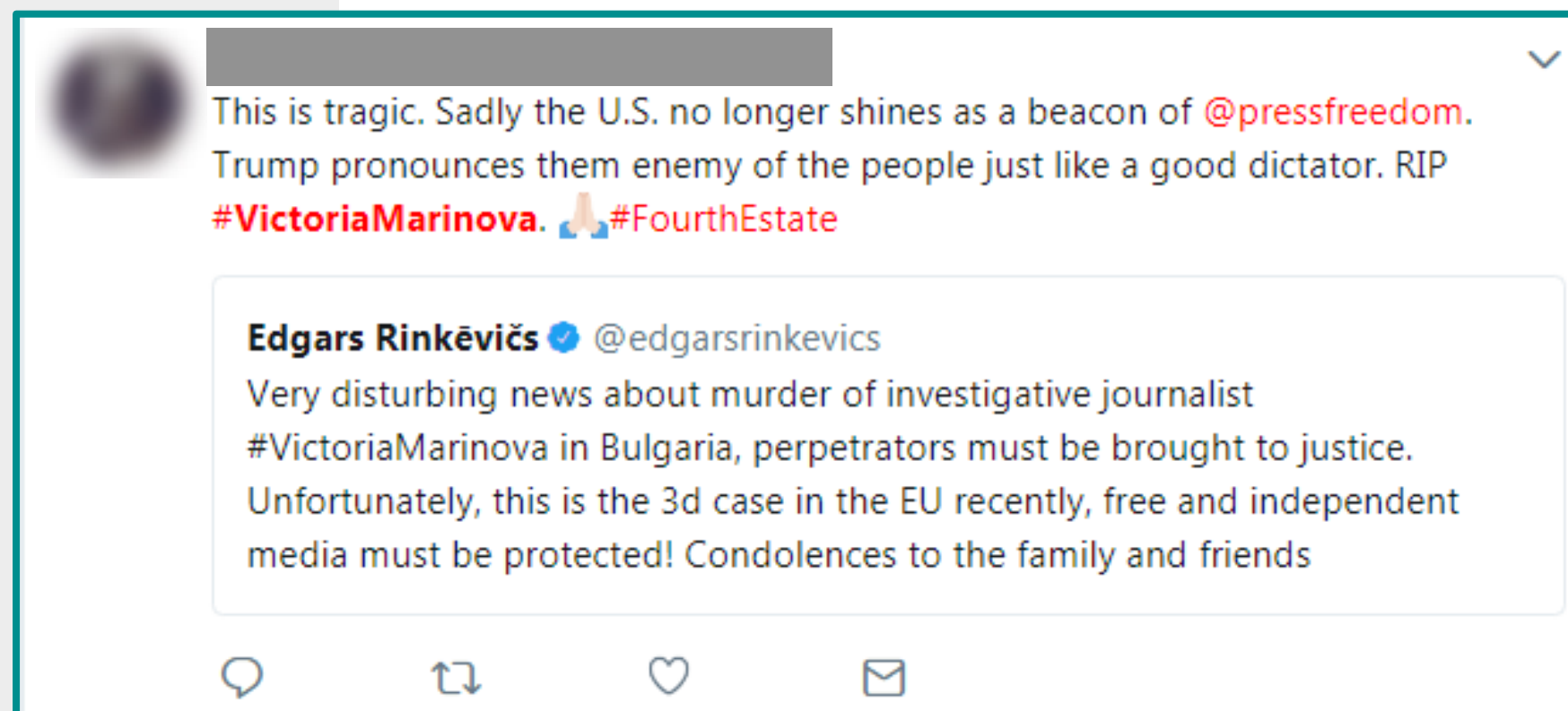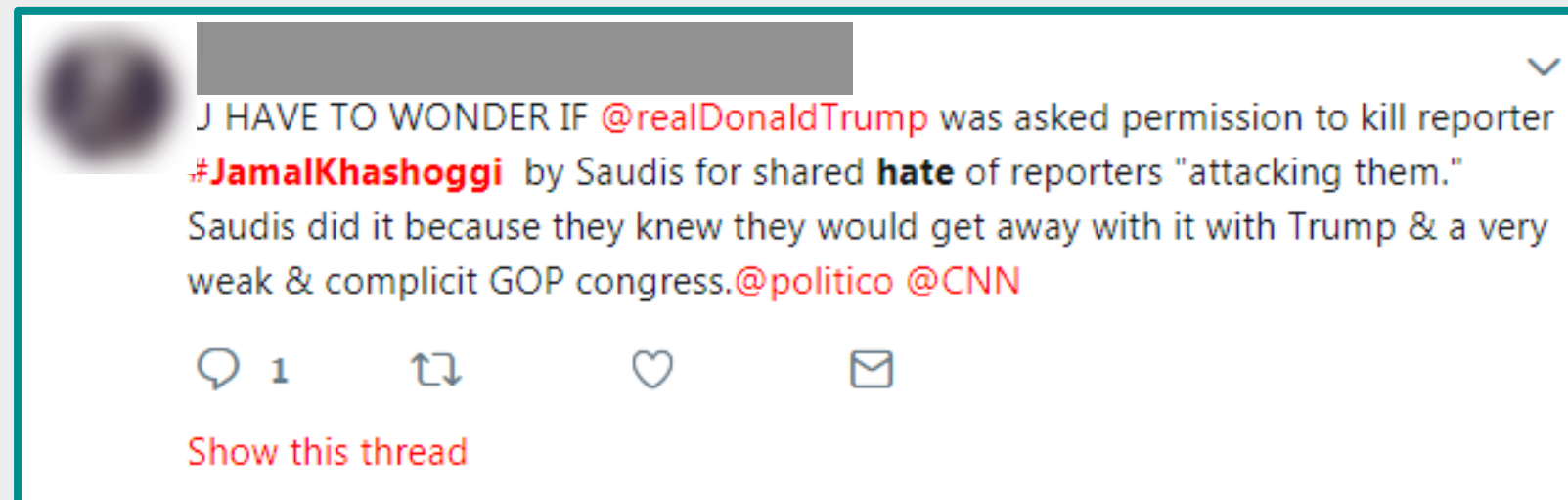From left down: Rebecca Smith, Wendi Winters

# Monitoring and Tagging Hate Speech in Social Media

## Reclaiming Agency



> J HAVE TO WONDER IF @realDonaldTrump was asked permission to kill reporter #JamalKhashoggi by Saudis for shared **hate** of reporters "attacking them." Saudis did it because they knew they would get away with it with Trump & a very weak & complicit GOP congress.@politico @CNN
>
> 💬 1   ⟲        ♡        ✉
>
> Show this thread



> This is tragic. Sadly the U.S. no longer shines as a beacon of @pressfreedom. Trump pronounces them enemy of the people just like a good dictator. RIP #VictoriaMarinova. 🙏#FourthEstate
>
> **Edgars Rinkēvičs** ✔ @edgarsrinkevics
> Very disturbing news about murder of investigative journalist #VictoriaMarinova in Bulgaria, perpetrators must be brought to justice. Unfortunately, this is the 3d case in the EU recently, free and independent media must be protected! Condolences to the family and friends
>
> 💬        ⟲        ♡        ✉

- ▸ Hate speech detection
- ▸ Early warning mechanism
- ▸ Counter measures



Bulgarian investigative journalist Viktoria Marinova and Saudi journalist Jamal Khashoggi. | Photo: Reuters

# Ευχαριστώ

Sophia Karakeva

Communications and Marketing Executive | Datascouting

Vice President | FIBEP