# SOCIAL MEDIA TRENDS, 3<sup>RD</sup> OF OCTOBER

# Health of the Media Conversation: Using AI to Detect and Understand Harmful Speech

**Eric Karstens**

Funding Consultant

European Journalism Centre

**Sophia Karakeva**

Communications and Marketing Executive

Datascouting

**Stavros Vologiannidis**

Professor for the Technological Educational Institute of Central Macedonia and Founder

Datascouting

SOCIAL MEDIA TRENDS, 3RD OF OCTOBER

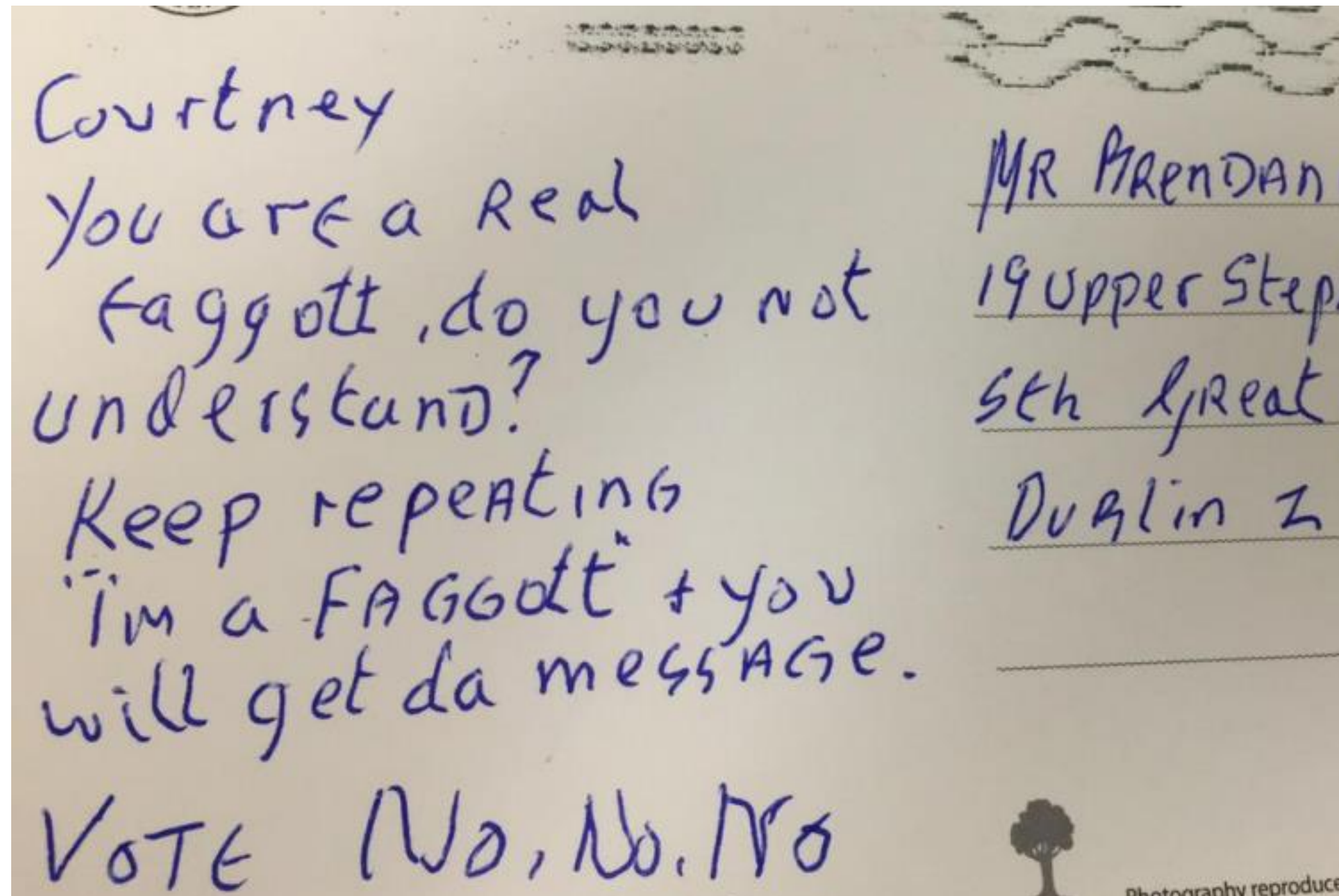# Health of the Media Conversation: Using AI to Detect and Understand Harmful Speech

Journalism and hate speech

Eric Karstens
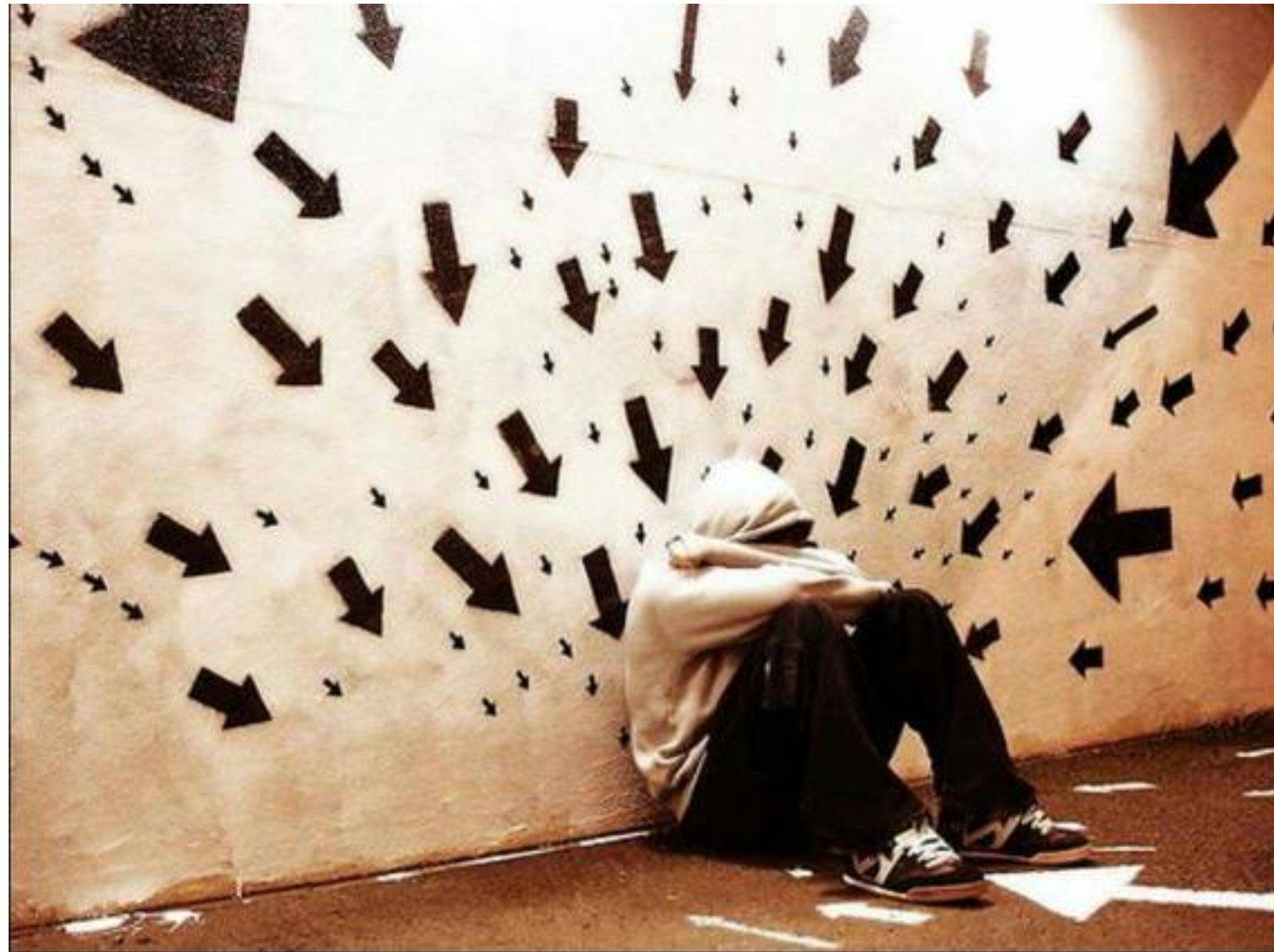
Funding Consultant

European Journalism Centre

50 FIBEP WORLD MEDIA INTELLIGENCE CONGRESS

COPENHAGEN OCTOBER 1-3 2018 MARRIOTT HOTEL

European Journalism Centre

DATASCOUTING
Actionable Information

@EricKarstens | @ejcnet | Eric Karstens | European Journalism Centre

# HATE MAIL THEN AND NOW

# HATE SPEECH AGAINST JOURNALISTS

# IMPACTS OF HATE SPEECH

# RECLAIMING AGENCY

# Health of the media conversation: using AI to detect and understand and harmful speech

Mapping and monitoring hate speech
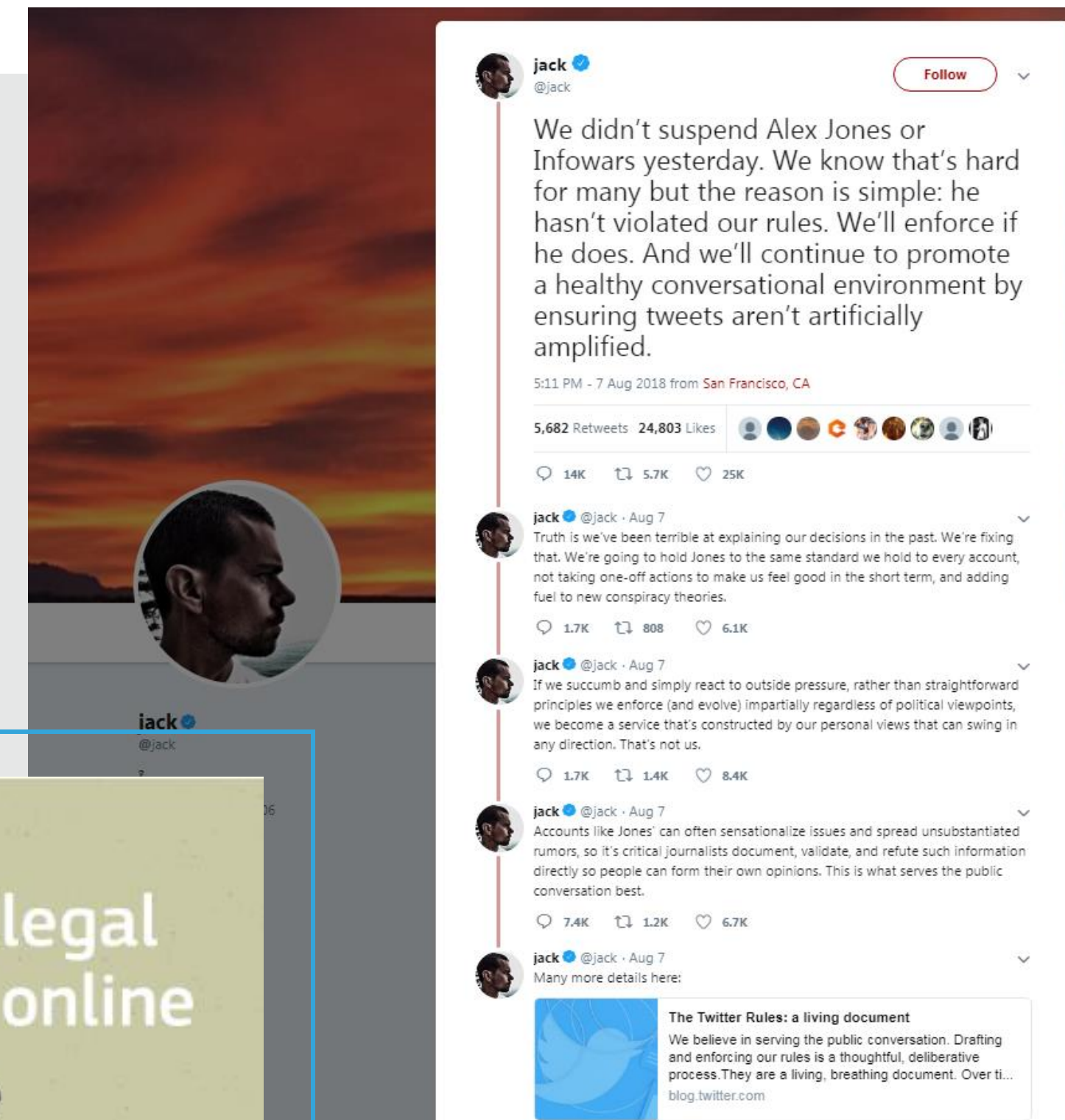
Sophia Karakeva

Communications and Marketing Executive

Datascouting

50 FIBEP
WORLD MEDIA
INTELLIGENCE
CONGRESS

COPENHAGEN
OCTOBER 1-3 2018
MARRIOTT HOTEL

DATASCOUTING
Actionable Information

European
Journalism
Centre

# 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

# The legal and discursive characteristics of hate speech

▸ Legal liability

▸ Danger of over-regulation
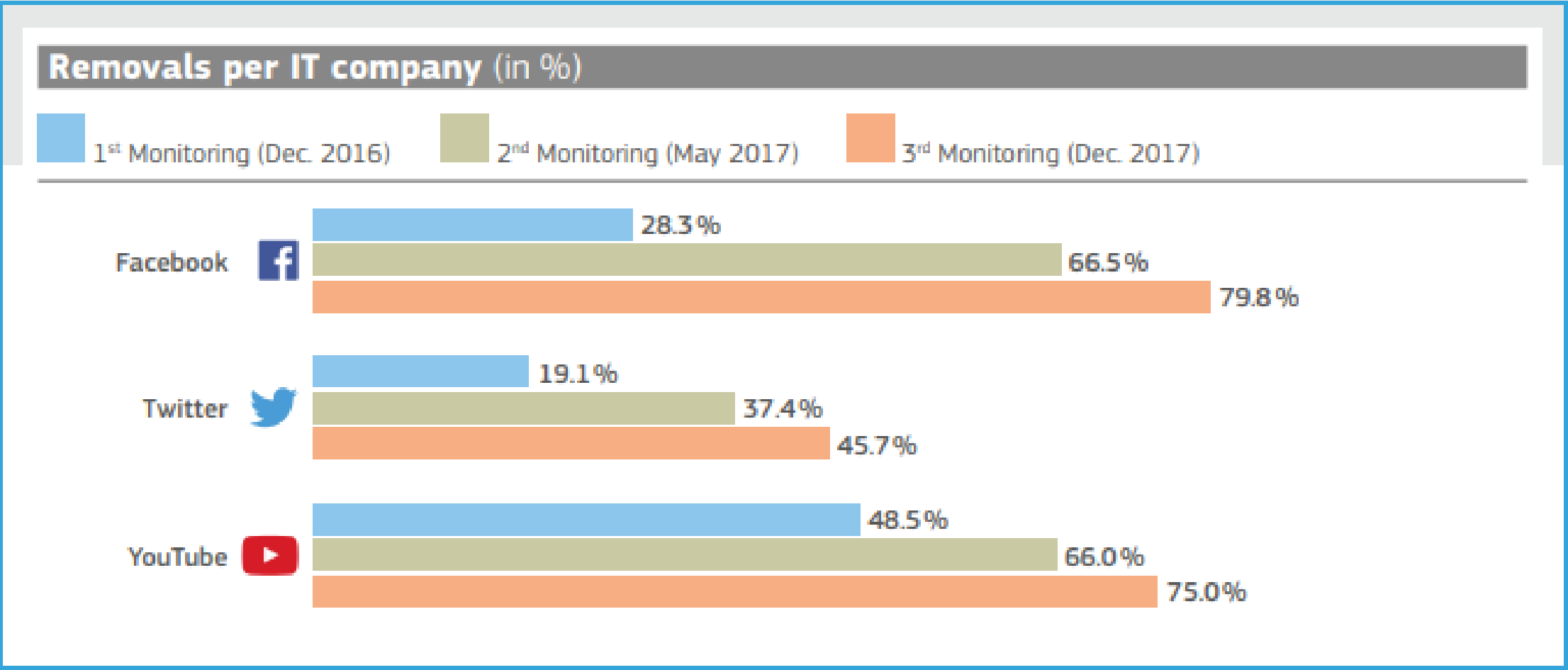
Hate Speech

Free Speech

## ARTICLE 19

"Everyone has the right to freedom
of opinion and expression;
this right includes
freedom to hold opinions without interference
and to seek, receive and impart information
and ideas through any media
and regardless of frontiers."
- Article 19; Universal Declaration of Human Rights

# Where does hate speech spread the most?

▸ 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
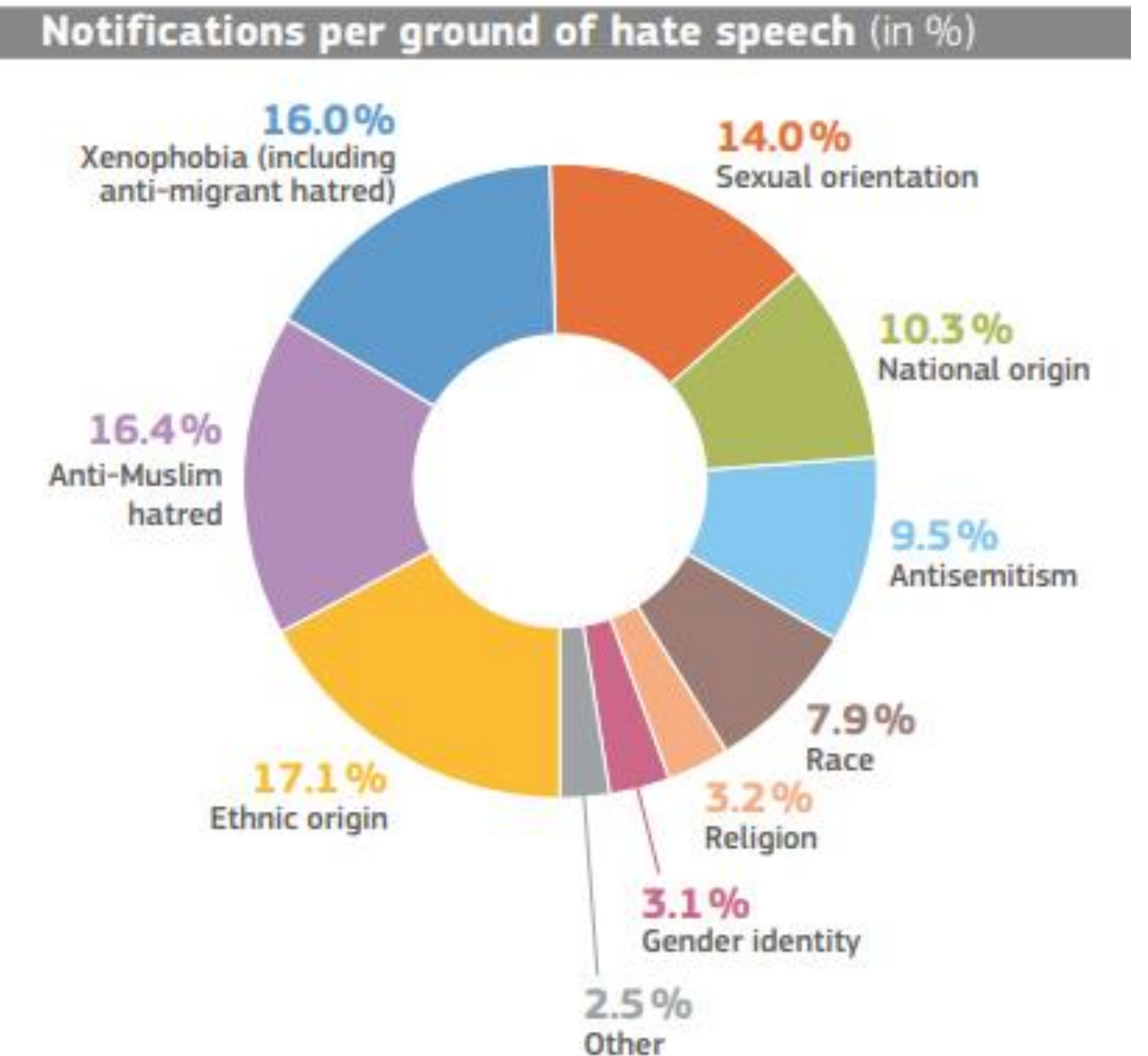
▸ 70% of hate speech content was removed in 2017



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



**Removals per IT company** (in %)

| | 1st Monitoring (Dec. 2016) | 2nd Monitoring (May 2017) | 3rd Monitoring (Dec. 2017) |
|---|---|---|---|
| Facebook | 28.3% | 66.5% | 79.8% |
| Twitter | 19.1% | 37.4% | 45.7% |
| YouTube | 48.5% | 66.0% | 75.0% |

# Who is most targeted?

▸ Ethnic origin, anti-Muslim hatred, & xenophobia

▸ Gender matters to social media

▸ The developing world

> " *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

$#*%!

#ToxicTwitter

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

16.0% Xenophobia (including anti-migrant hatred)

14.0% Sexual orientation

10.3% National origin

16.4% Anti-Muslim hatred

9.5% Antisemitism

17.1% Ethnic origin

7.9% Race

3.2% Religion

3.1% Gender identity

2.5% Other

Source: 3rd monitoring exercise of the European Union's Code of Conduct, released in January 2018

50 FIBEP
WORLD MEDIA
INTELLIGENCE
CONGRESS

COPENHAGEN
OCTOBER 1-3 2018
MARRIOTT HOTEL

DATASCOUTING
Actionable Information

European
Journalism
Centre

# Hate me, hate me not

*"It's like public lynching... It has made be frightened for my physical safety when I am out in the streets."*
Turkish journalist Amberin Zaman

▸ Hate me not: not allowed on our platform
▸ Hate me: often the most engaged content

▸ Hate me not: strong measures to take it down
▸ Hate me: not everything has the same value

▸ Hate me: highly reactive, high error rate
▸ Hate me not: automation and AI tools

▸ Hate me: "new normal"
▸ Hate me not: education, campaigns, regulation



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.

**facebook**

**This post goes against our Community Standards**

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

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

# 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")



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.

# What's being done to fight it?

▸ Expanded hate speech rules

▸ More human annotation

▸ Academic & scientific research

▸ Sophisticated technology

SOCIAL MEDIA TRENDS, 3RD OF OCTOBER

# HEALTH OF THE MEDIA CONVERSATION: USING AI TO DETECT AND UNDERSTAND AND HARMFUL SPEECH

## AUTOMATICALLY TAGGING HATE SPEECH

Stavros Vologiannidis

ASSISTANT PROFESSOR

FOUNDER OF DATASCOUTING

@SVOLOGIA | @DATASCOUTING | #FIBEP | #WMIC18

# HATE SPEECH AND TECHNOLOGY

▸ Hate speech is nothing new

▸ But the problem is much more evident in the age of **social media**

▸ Several attempts have been made to automatically identify hate speech…

▸ **Challenge: separating hate speech from offensive speech**

  ▸ Hate speech changes forms, uses different expressions

  ▸ Keyword based identification techniques do not work well

# MACHINE LEARNING APPROACH

- Use **Natural Language Processing** and **Machine Learning** techniques in order to identify hate speech

- **Main stages** in Text Analytics with Machine Learning using **Neural Networks**

  - Transform a **document** (social media article) to **numbers** (features)

  - Create a **Machine Learning model** based on your architecture of your choice and **train** it as follows:

  - Model **looks** at documents and calculates and answer

  - Model **compares** answer against the correct one (annotated by human experts)

  - Model **tweaks** itself so that if it finds the same document again it would be more likely to calculate the correct answer

  - **Repeat**

# DOCUMENT -> NUMBERS

▸ **How to transform a document to numbers?**

  ▸ A) Convert a **word to numbers** (vector of a few hundreds of numbers)

  ▸ B) Convert a **document** = lots of words **to math**

▸ **A) Word -> numbers: Word embedding algorithms**

▸ We need to retain some **semantic properties** of the word such as similarity

  ▸ Words that share common contexts are close as mathematical objects

  ▸ King - Man + Woman = Queen

▸ We have several options of such algorithms and **pretrained** models (Wikipedia,...)

  ▸ Word2vec (Google), FastText (Facebook), .....

▸ **B) Document -> numbers: math operations on words, can include NLP**

# MODEL ACHITECTURE

▶ **Model architecture: Convolutional Neural Networks on top of character level word embeddings**

▶ **Input** : Document

▶ **Output**: Hate or not (or additional classes such as personal attack, etc)

▶ **Challenge:** we need **lots of annotated (labeled) documents to train** (and test) algorithms

# ACTIVE LEARNING APPROACH

▸ Active Learning is a methodology that can greatly **reduce the amount of annotated data** required to train a machine learning model.

▸ It does this by **prioritizing** the labeling work for the experts .

  ▸ The model **looks** at unlabeled data (cheap)

  ▸ **Identifies which data** it is **most confused about** and **requests labels** from experts for just those (human is involved - expensive).

  ▸ It **trains** on that small amount of labeled data

# ACTIVE LEARNING - BETTER DATA > MORE DATA

▶ **Minimize costs/time** by reducing expert resources

    ▶ Get the same or better classifier accuracy usually with **3 to 10 times less effort**

▶ **Better data** is **more useful** than **more data...**

▶ The best active learning based model is the one with **great accuracy** and the **minimum number of requests towards experts**

50 FIBEP
WORLD MEDIA
INTELLIGENCE
CONGRESS

COPENHAGEN
OCTOBER 1-3 2018
MARRIOTT HOTEL

DATASCOUTING

@SVOLOGIA @DATASCOUTING

Stavros Vologiannidis | DataScouting

# CURRENT STATE OF HATE SPEECH DETECTION

▸ **Social media companies** work extensively on hate speech detection

▸ Only **small datasets** are available in the research literature, even for the english language

▸ Different **hate speech definitions**, intensify the above problem

▸ Several researchers and research groups work on the subject, with a current **success rate of approximately 80%**

# OUR VISION

▸ **Hate speech**

   ▸ Create a scalable active learning based system for identifying hate speech with direct feedback from journalists

   ▸ Support multiple languages (English, French, German, Greek, Spanish)

▸ **DataScouting Media Intelligence Products**

   ▸ Include technology from the project to DataScouting's media intelligence solutions

# Thank you

**Eric Karstens**

Funding Consultant

European journalism centre

karstens@ejc.net

**Sophia Karakeva**

Marketing and Communications Executive

Datascouting

soka@datascouting.com

**Stavros Vologiannidis**

Professor for the Technological Educational Institute of Central Macedonia and Founder

Datascouting

svol@datascouting.com

European Journalism Centre

DATASCOUTING
Actionable Information

50 FIBEP WORLD MEDIA INTELLIGENCE CONGRESS

COPENHAGEN OCTOBER 1-3 2018 MARRIOTT HOTEL

European Journalism Centre

DATASCOUTING Actionable Information