# HATE
METER

Hate speech tool for monitoring, analysing
and tackling Anti-Muslim hatred online

REC Action Grant (REC-DISC-AG-2016-04)
24 months (01.02.2018 - 31.01.2020)

## Training module B for stakeholders



eCrime
ICT, law & criminology

UNIVERSITÀ
DI TRENTO

FBK
FONDAZIONE
BRUNO KESSLER

UNIVERSITÉ
TOULOUSE 1
CAPITOLE

Teesside
University

AMNESTY
INTERNATIONAL
ITALIA

STOP HATE UK

CCIF
COLLECTIF CONTRE
L'ISLAMOPHOBIE
EN FRANCE

# D19 – Training Module B for stakeholders

| WP5 | Training, dissemination and sustainability strategy |
|---|---|
| Due Date: | 30/11/2019 |
| Submission Date: | 30/11/2019 |
| Responsible Partner: | CCIF |
| Version: | 1.0 |
| Status: | Final |
| Author(s): | Isis Koral (CCIF), Andrea Di Nicola, Daniela Andreatta, Gabriele Baratto and Elisa Martini (UNITRENTO), Francesca Cesarotti and Giulia Pirozzi (Amnesty Italy), Bill Howe (StopHate UK), Marco Guerini and Sara Tonelli (FBK), Georgios A. Antonopoulos and Parisa Diba (TEES) |
| Reviewer(s): | Isis Koral (CCIF), Andrea Di Nicola, Daniela Andreatta, Gabriele Baratto and Elisa Martini (UNITRENTO), Francesca Cesarotti and Giulia Pirozzi (Amnesty Italy), Bill Howe (StopHate UK), Serena Bressan, Marco Guerini and Sara Tonelli (FBK), Jérôme Ferret and Mario Laurent (UT1-Capitole), Georgios A. Antonopoulos and Parisa Diba (TEES) |
| Deliverable Type: | R |
| Dissemination Level: | CO |

# Glossary

| | |
|---|---|
| **CAP** | Computer Assisted Persuasion |
| **CCIF** | Collectif contre l'islamophobie en France |
| **LEA** | Law Enforcement Agency (plural: LEAs) |
| **MS** | Member State (plural: MSs) |
| **NGO** | Non-governmental organisation (plural: NGOs) |
| **NLP** | Natural Language Processing |
| **SD** | Standard Deviation |
| **WP** | Work Package (plural: WPs) |

# Table of contents

# Executive summary

This document is the **Deliverable "D19 – Training Module B"** of the European project **"Hatemeter - Hate speech tool for monitoring, analysing and tackling anti-Muslim hatred online"** (hereafter referred to as the "Hatemeter", project reference: 764583), which aims at systematising, augmenting and sharing knowledge of anti-Muslim hatred online, and at increasing the efficiency and effectiveness of non-governmental organisations (NGOs) in preventing and tackling Islamophobia at the EU level.

This deliverable is meant to be used by stakeholders outside the Hatemeter project (i.e., NGO/CSO representatives; civil servants; Muslim community leaders; journalists/media; LEAs) as a **manual** to understand the main goals of **Hatemeter Platform** (also referred to as the Hatemeter Tool), how it works and what it can achieve. Together with Training Module A for academics/research organizations, the present document is one of the main deliverables of the "**Activity 5.1 - Networking, capacity building and training activities for target stakeholder groups**" within the Hatemeter **WP5 "Training, dissemination and sustainability strategy"**.

Hatemeter addresses a **strategic challenge towards NGOs' needs to tackle anti-Muslim hate speech online**, by developing and testing an **ICT tool** (i.e., Hatemeter Platform) that **automatically monitors and analyses Internet and social media data on the phenomenon and produces computer-assisted responses and hints to support counter-narratives and awareness raising campaigns**. In this regard, the Hatemeter Platform uses a combination of natural language processing (NLP), machine learning, and big data analytics/visualisation to: a) identify and systematise in real-time actual "red flags" of anti-Muslim hate speech and/or possible related threats online (**Real-time identification**); b) understand and assess the sets of features and patterns associated with trends of Islamophobia online (**In-depth understanding**); c) facilitate effective tactical/strategic planning against anti-Muslim hatred online through the adoption of an innovative Computer Assisted Persuasion (CAP) approach (**Tactical/strategic response**); d) produce an accurate counter-narrative framework for preventing and tackling Islamophobia online, and building knowledge-based and tailored awareness raising campaigns (**Counter-Narratives Production**).

As such, this **Training Module B** outlines **practical uses** and also **provides examples of the location and analysis of hate speech against Muslims on Twitter** from the three NGOs involved in the Hatemeter project, namely Amnesty International (AMN) in Italy, Collectif Contre l'Islamophobie en France (CCIF) in France, and Stop Hate UK (STOPHATE) in the United Kingdom, using the Hatemeter Platform developed by FBK and tested by the three NGOs. In these countries, the magnitude of the phenomenon is significant, but no systematic responses have been implemented. This document is tailored to the scientific and technical requirements and needs of specific key actors and is designed of the location and analysis of hate speech against. The Training Module B will be employed during the training seminar for professionals in Toulouse on the 18th December 2019 and made available through the Hatemeter **website** (www.hatemeter.eu).

This document is organised as follows. The "Introduction" starts by giving an idea of the problem addressed by the project Hatemeter and, consequently, it briefly describes **what Islamophobia is** (subsection 1.1). The second subsection (1.2) presents a **general overview** of the project and its aims. The second section of the document is a technical description of the Hatemeter Platform and demonstrates the variety of data analysis that can be performed (subsection 2.1, and more specifically: 2.1.1 Recent trends Functionality, 2.1.2 Hashtag trends Functionality, and 2.1.3 Hate speakers Functionality) and the **Computer-Assisted Persuasion tool** (subsection 2.2, and more specifically: 2.2.1 Alerts Functionality, and 2.2.2 Counter-narratives Functionality). The third section summarises **evidence of online Islamophobia** in the three countries involved in the project: namely, Italy (subsection 3.1), France (subsection 3.2) and the United Kingdom (subsection 3.3).

Each subsection starts with a brief description of the **context and background** of online Islamophobia and anti-Muslim hatred in that country and then specifically presents **evidence** of hate speech collected through the Hatemeter Platform. The fourth and last section proposes some **suggestions and insights** on the use of the Platform and of the Hatemeter methodology for stakeholders in future work.

# 1. Introduction

## 1.1 What is Islamophobia?

Islamophobia is defined as "all acts of discrimination or violence against institutions or individuals because of their real or supposed affiliation to Islam" (CCIF, 2019), evinced as **feelings of anxiety** or **perceptions of fear** and **hatred**. Additionally, Islamophobia does not merely entail anxious awareness or perceptions rooted in apprehension and contempt, but also the discriminatory attitudes and hostile practices through which it is manifest and expressed; such as harassment, verbal and physical abuse as well as hate crimes, perpetrated in both **offline** and **online contexts**.

The European Islamophobia Report (Bayrakli and Hafez, 2019a) recorded that Muslims are among the **first victims** of the rise of far-right extremism in Europe. Below are some reported examples. In Austria, the Office for Documenting Islamophobia and Anti-Muslim Racism recorded a 74% increase of anti-Muslim racist acts in its 2018 report. In France, the Collectif contre l'islamophobie en France (CCIF) recorded a 52% increase.[1] In the UK, the number of cases recorded in official statistics rose by 17% in 2017-18 and religion-specific incidents multiplied by 40%. In Italy, between 2017-2018 there was a (significant) increase in hateful posts on social media. Finally, in the Netherlands, the Anti-discrimination Agencies announced that Muslims were the target of 91% of reported cases of religious discrimination.

In the last decade, Islamophobia has gained momentum through the Internet, which, along with new media technologies including social media platforms and global digital networks, has enabled the spread of polemic and anti-Islamic and anti-Muslim discourse to a worldwide audience (Larsson, 2007), (Horsti, 2017). With the advent of the Internet, online or cyber Islamophobia has seen a **large increase**, with spaces on the Internet now becoming a Platform for the spread of this rhetoric, with xenophobic viewpoints and racist attitudes towards Muslims being easily disseminated into public debate (Ekman, 2015). Online Islamophobia takes place primarily through **blogs** and **social media**, as well as through traditional media outlets seen online (Aguilera-Carnerero and Azeez, 2016).

In 2018, The Collective Against Islamophobia in Belgium underlined that 29% of reported Islamophobic incidents in 2018 pertain to Islamophobia both in the media and online (Bayrakli and Hafez, 2019b). However, as Faytre (2019: 18) points out, "Islamophobic controversies often originate from social media before being debated in mainstream media and triggering reactions among politicians."

According to Oboler (2016), anti-Muslim hate, much like many other forms of hate, is unlikely to remain purely virtual, with online Islamophobia likely to incite religious hatred and xenophobia, leading to **real world crimes** and a **rise in political extremism** both on the far-right and from the radicalisation of Muslim youth in response to such messages of exclusion. The outcome is a vicious circle that is difficult to break. Thus, as Larsson (2007) points out, it is important to examine to what extent the Internet is being used to spread and foster anti-Muslim and anti-Islamic opinions in contemporary society.

---

[1] Additionally, it is necessary to underline that according to the French government (Gouvernement.fr, 2019), the official number of anti-Muslim acts is the lowest since 2010. This statement recalled the CNCDH report (2018) on racism, antisemitism and xenophobia, which affirms that the number of anti-Muslim hate incidents is lower than previous years because many acts are not reported to the police. On the contrary, the CCIF's number reported in the text indicates an increase. This difference has two possible explanations: i) CNCDH and CCIF utilise different definition of Islamophobic acts and consequently they register different information; ii) the phenomenon may suffer under-reporting via the police, because people feel more comfortable reporting these types of acts to CCIF.

## 1.2 Project Hatemeter

Project "Hatemeter - Hate speech tool for monitoring, analysing and tackling anti-Muslim hatred online" aims at **systematising, augmenting and sharing knowledge of anti-Muslim hatred online,** and at increasing the efficiency and effectiveness of NGOs in **preventing and tackling Islamophobia** at the EU level, by developing and testing an **ICT tool (i.e., Hatemeter Platform)** that **automatically monitors and analyses Internet and social media data** on the phenomenon, and **produces computer-assisted responses and hints** to support **counter-narratives** and **awareness raising campaigns.**

More specifically, backed by a strong interdisciplinary effort (criminology, social sciences, computer sciences, statistics, law), the Hatemeter Platform uses a combination of **natural language processing (NLP), machine learning, and big data analytics/visualization** to:

A. identify and systematise in real-time actual "red flags" of anti-Muslim hate speech and/or possible related threats online (**Real-time Identification**);

B. understand and assess the sets of features and patterns associated with trends of Islamophobia online (**In-depth Understanding**);

C. develop an effective tactical/strategic planning against anti-Muslim hatred online through the adoption of the innovative Computer Assisted Persuasion (CAP) approach (**Tactical/Strategic Response**);

D. produce an accurate counter-narrative framework for preventing and addressing Islamophobia online and building knowledge-based and tailored awareness raising campaigns (**Counter-Narratives Production**).

The Hatemeter Platform has been **piloted and tested by three NGOs of EU MSs** where the magnitude of the problem is considerable but no systematic responses have been implemented (**France, Italy and the United Kingdom**), thus enabling Project Hatemeter to address several objectives of the **Annual Colloquium on Fundamental Rights "Tolerance and respect: preventing and combating anti-Semitic and anti-Muslim hatred in Europe"** and the **European Agenda on Security** (2015), as well as the **priorities** of the REC call for proposals.

In order to strengthen **cooperation** between key actors and to ensure the **widest circulation** and **long term impact** of project results upon future research streams and operational strategies, the project favours **capacity building and training** and the **sustainability and transferability** of the Hatemeter Platform amongst **other target stakeholder groups** (e.g., LEAs, journalists/media, etc.) across the EU and for **other forms of hate speech**, through the building of an **"EU laboratory for countering online anti-Muslim hate speech on the Internet and social media" (i.e., Hatemeter Lab).**

The present document is one of the main deliverables of the "Activity 5.1 - Networking, capacity building and training activities for target stakeholder groups" within the Hatemeter WP5 "Training, dissemination and sustainability events", and precisely, the "D19 - Training Module B". This **Training Module** is intended as a **manual for academics and research organizations** outside the Hatemeter project to understand the main goals of the Hatemeter Platform (or Hatemeter Tool), its main functions and possible uses.
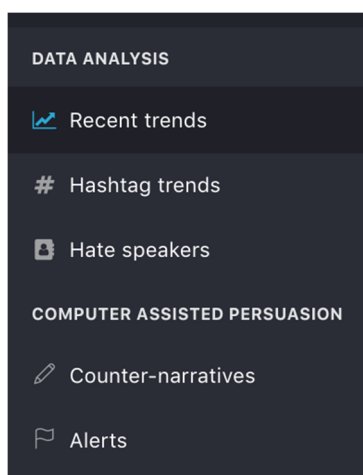
# 2. Hatemeter Platform

The HATEMETER Platform is a **web-based tool designed to support researchers and stakeholders** (e.g. NGO operators, researchers, law enforcement agencies) in analysing and countering anti-Muslim speech online. **Twitter** and **YouTube** are continuously monitored and, when keywords or hashtags related to anti-Muslim discourse are detected, the corresponding messages are retrieved and stored in the project database. Here we **focus on the analytics applied** to tweets to better understand the Islamophobic messages, which are displayed within the Platform utilising different perspectives. However, because the analytics relating to YouTube were not tested during the pilot phases, they have not been included in the training kit, although they are visible in the final release of the Platform (v.3). From a technical perspective the Platform relies on a relational database and a tomcat application server. The interface is based on existing javascript libraries such as C3.js (https://c3js.org), D3.js (https://d3js.org) and Sigma.js (http://sigmajs.org).

**The Platform is password protected**, and each account gives access to data in a **specific language** (Italian, English or French) and a version of the underlying database. The analytic tools are consistent across the three versions, since they have been designed so to be language-independent.

The **Platform functionalities** are divided into three main groups, displayed on the left of the homepage view: 1) DATA ANALYSIS, 2) COMPUTER ASSISTED PERSUASION, and 3) PROJECT HATEMETER (see Figure 1). The first item includes all analysis concerned with Islamophobic hashtags and keywords and the spread of networks of users. The second item displays the Platform which, taking a hate message as input, automatically generates potential counter-narratives. The third item includes general information on the project and a link to the website. The first two groups are described in more detail below (section 2.1, 2.2).

**Figure 1 - Platform menu with main functionalities**



*SOURCE: Screenshot from the Hatemeter Platform*

## 2.1 Data Analysis

The "**DATA ANALYSIS**" item presents three views: "*Recent trends*", "*Hashtag trends*" and "*Hate speakers*". There is a significant difference between the first item and the others: "*Recent trends*" performs **real-time monitoring** and therefore calls Twitter APIs on the fly. In contrast, the others, rely upon the **HATEMETER database**, i.e. the outcome of monitoring Islamophobic discourse, which has been in operation since October 2018. This means that, while the analysis of past data is stable, those in real-

time may be subject to change whenever Twitter's Application Programming Interface (API) policy undergoes revision.
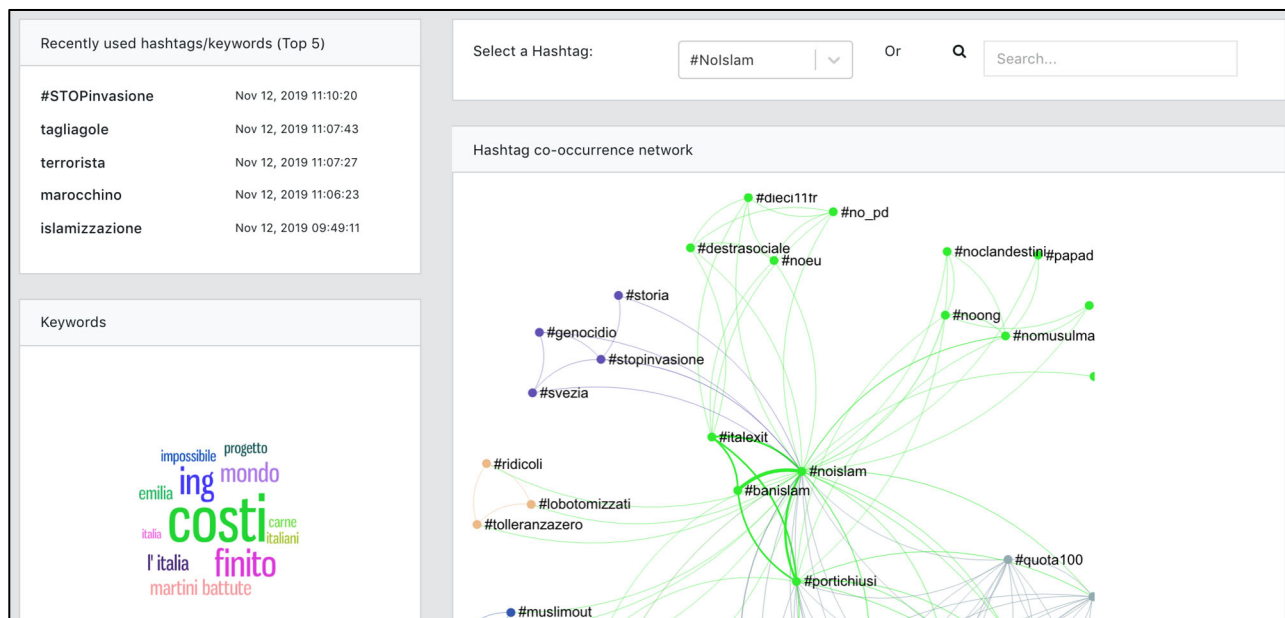
In the early stages of the project a pool of academics and activists defined, for each language, a set of **hashtags and keywords** that are directly associated with Islamophobic messages (e.g. #STOPIslam, #Muslimshit[2], #BanIslam). These hashtags and keywords were used as query terms to access Twitter APIs on a regular basis (twice weekly) and collect all messages containing the query term. The collected tweets were analysed using **text processing tools** to extract the most relevant information related to anti-Muslim hatred online. These could be the **metadata** connected to the messages (i.e., user, date, frequency), the **content popularity** (number of replies, i.e. answers to a message or tweet, and retweets, i.e. broadcasting a tweet or message posted by another person), and the **network** in which the discourse is spread (i.e. nodes that had most interactions involving the hashtags or keywords of interest).

The information distilled and structured in the previous steps is made available to final users through an advanced visualisation Platform. This provides functionalities for the **visual exploration and analysis of the data**, enabling content monitoring, synchronic and diachronic comparisons, close and distant reading, data clustering, network analysis, etc. A pictorial and graphic format is used as much as possible to ensure the tool is language and country-independent.

## 2.1.1 The "Recent trends" Functionality

Under DATA ANALYSIS, the "*Recent trends*" view allows users to **monitor current Twitter activities around pre-defined Islamophobic keywords and hashtags** or utilise search terms defined on-the-fly by the user. After selecting a hashtag or writing a term in the field, each box displays related information: in the "*Recently used hashtags/keywords*", the system ranks the most recently used Islamophobic terms by the date and time they were last posted on Twitter.

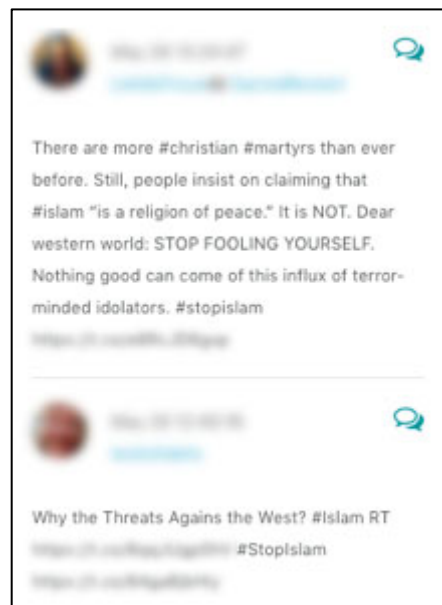**Figure 2 - View of the "Recent trends" tab**



*SOURCE: Screenshot from the Hatemeter Platform*

---

[2] Original language: Musulmerda

For example, Figure 2 shows what would have been displayed for the Italian pilot on 12 November 2019: the five most recently tweeted Islamophobic hashtags and terms, the network of hashtags currently being used in conjunction with #NoIslam (selected by the user) and the tag cloud extracted from the recent tweets containing #NoIslam and extracted with Keyphrase Digger, a keyword extractor tool (Moretti et al., 2015). In the hashtag network, **each node is a hashtag co-occurring with the central one**, and the arcs connecting nodes represent co-occurrence. **The thicker the arc, the more frequently two hashtags have been used together**. Finally, on the right it is possible to see the list of tweets, ranked by date, containing the hashtag or keyword selected by the user (Figure 3). This information is displayed to provide a **more fine-grained view of the Islamophobic content** currently circulating online. By clicking on one of these tweets, it is possible to open the message of interest within Twitter, for example, in order to check the replies or the number of retweets. Each of the tweets displayed inside the feed also contains the link to the user originating the post, as well as that of the re-tweeting user. These messages are also **linked to the counter-narratives tool**, which is activated by clicking on the speech bubble icon on the right. Further details of this tool are outlined in Section 2.2.

**Figure 3 – Feed of tweets containing a user-defined hashtag or keyword, ranked by date**



*SOURCE: Screenshot from the Hatemeter Platform – Tweet feed*

Overall, the information shown in this tab can be **utilised by researchers and NGO operators** to *i)* discover new hashtags to be monitored, and *ii)* get an overview of the main topics discussed in Islamophobic tweets.

## 2.1.2 The "Hashtag trends" Functionality

By clicking on the "Hashtag Trends" item under the "DATA ANALYSIS" menu, users can see a **more fine-grained analysis of hashtags**, focused upon their **trend *over time***. In this case, the information is not retrieved on the fly using Twitter APIs, but is gathered from the underlying database. In particular, users can select a hashtag, and a pertinent temporal snapshot. The system with then provide **general statistics** reporting the number of tweets, retweets and replies containing the hashtag of interest within the selected snapshot. In addition, the system displays the **hashtag co-occurrence network** and the **day-by-day statistics** showing the hashtag presence on Twitter over time and enabling the visualisation of the message popularity day-by-day as well as comparing different snapshots.

Figure 4 demonstrates an example: users select a hashtag (in this case #NoIslam) and a temporal snapshot of interest (10 days ending on April 10th). In the upper right of the figure, it is possible to read the general statistics for these 10 days, reporting the number of tweets, retweets and replies containing the selected hashtag. On the left of the figure, the system shows the hashtag co-occurrence network (similar to the one displayed in Figure 2), whereas in the lower right there are the day-by-day statistics showing the hashtag presence on Twitter over time.

**Figure 4 – Feed of tweets containing a user-defined hashtag or keyword, ranked by date**
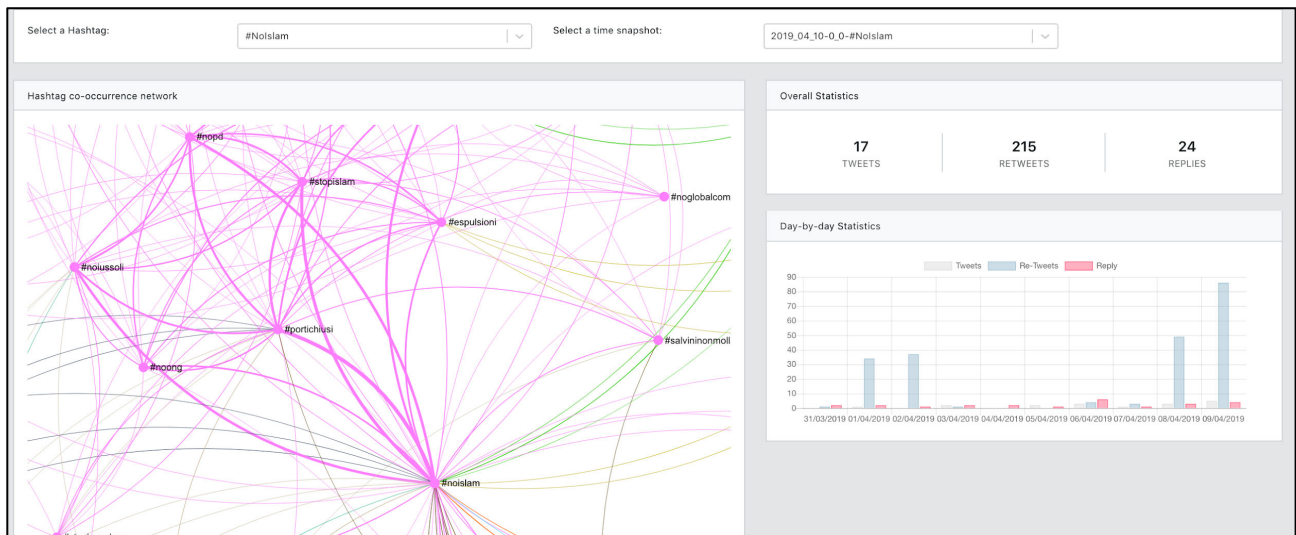


*SOURCE: Screenshot from the Hatemeter Platform*

Furthermore, a box at the bottom of the tab displays a **list of the most retweeted messages amongst those collected through the query of interest** (see Figure 5). In particular, each tweet is displayed along with the date posted and the number of retweets (ranked in decreasing order). Note that this box is not always present, since for some hashtags and some temporal snapshots no retweets with a frequency > 1 are present.

**Figure 5 – View of the most retweeted messages for a given hashtag/keyword in a given time period. The number next to the date shows how many times the message was retweeted.**



*SOURCE: Screenshot from the Hatemeter Platform*

### 2.1.3 The "Hate speakers" Functionality

By clicking on the "Hate Speakers" item under the "DATA ANALYSIS" menu, users of the Platform can display an analysis related to the **most active users** within a community spreading Islamophobic messages online. As in the previous tab, data is extracted from the Hatemeter database and not collected on the fly.

After selecting a hashtag and a time snapshot, the Platform displays the network of users who posted messages containing the selected hashtag. This is called the "**User co-occurrence network**", in which colours are automatically assigned by the network analysis algorithm to identify communities of users, i.e. those who interact more often with each other through replies or retweets. In the "**Most connected users**" frame, the Platform displays a ranked list of users with the most connections inside the network, i.e. those that are more likely to give visibility to Islamophobic messages. Finally, a tag cloud representing the messages exchanged inside the identified community of users is also created, again using the KD keyword extraction tool.
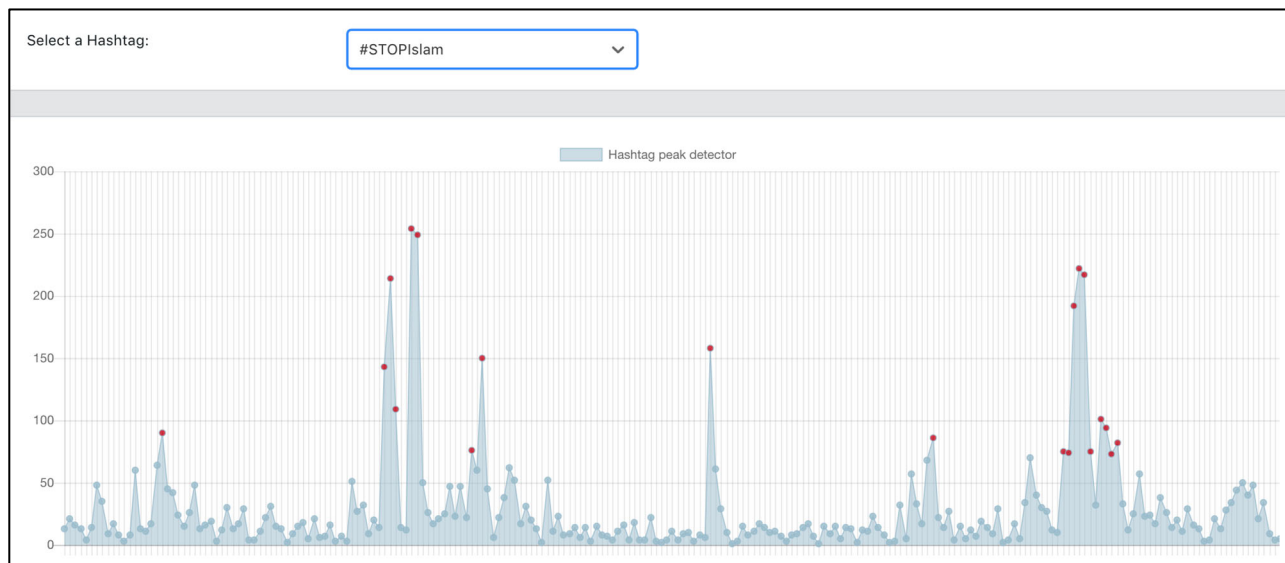
## 2.2 Computer-Assisted Persuasion (CAP)

The goal of the "Computer-Assisted Persuasion" CAP item of the Hatemeter Platform is to provide a **suite of tools for stakeholders to support actual intervention online**, i.e. responding to hatred online to prevent and combat it. To this end, the CAP approach is twofold: i) providing timely, context-sensitive and incident/crime-specific alerts based on complex data analytics; ii) producing an accurate counter-narrative framework on the basis of a "best experiences" repository containing real, positive examples of counter-narratives to Islamophobia.

### 2.2.1 The "Alerts" Functionality

The "**ALERTS**" view was created to increase awareness of the scale of Islamophobic messages, monitoring the trend of hashtags and keywords over time without a focused, time-bound framework such as the one described in the "DATA ANALYSIS" section of the Platform. In this view, users are asked to select one of the pre-defined **hashtags or keywords** monitored since the beginning of the project. The system displays the trend of the selected term over time, based on the collation of the number of tweets, retweets and replies. The visualisation displays a **peak** when the Islamophobic hashtag or keywords have been particularly present on Twitter, enabling a possible alert for operators. This is computed on the fly by the system, taking into account the average frequency of messages plus one standard deviation. This value is dynamically computed for each hashtag or keyword, since some of them may be generally frequent and more present online than others. As an example, we report the analysis for the #STOPIslam hashtag in Figure 6.

**Figure 6 – Hashtag trend for the #STOPIslam hashtag. Red dots signal that the hashtag appeared frequently on that day (more than average + 1 SD)**



*SOURCE: Screenshot from the Hatemeter Platform*

## 2.2.2 The "Counter-narratives" Functionality

The main functionality for CAP intervention relies on a **chatbot-like application** that, given a short Islamophobic text as input, provides **five automatically generated suggestions** that could be used **to counter the hate speech and deescalate the argument**. The counter-narrative suggestions are provided in the specific language of the account in use (Italian, English or French). Figure 7 displays a screenshot of the English version of the interface, where the message above is the input and those below are suggestions provided by the chatbot that could be used to build counter-narratives.

Operators using the Platform can set the input according to different strategies: if operators want to reply to one of the Islamophobic messages displayed in the **list of tweets** retrieved by the Twitter API, they can click on the speech bubble icon on the upper right of the message (see Figure 3 above). In this way, the content of the tweet is displayed as input in the counter-narratives tab, while **five possible answers** are automatically generated by the application. As an alternative, operators can also write an **input text**, and request the application provide possible replies, or copy and paste it from other sources.

**Figure 7 – Screenshot of chatbot-like tool for the creation of counter-narratives**



*SOURCE: Screenshot from the Hatemeter Platform*

After checking the answers, operators can choose to **use one of them to reply to an Islamophobic message**. Each answer can also be **edited and modified** (using the 'pencil' icon) and saved (using the 'floppy disk' icon), before copying and pasting it into Twitter or any other social media to post the reply. If none of the responses proposed by the system are effective, operators can also draft their own utilising the interface. Although having the possibility to post the reply directly into Twitter would be handier for operators, this would require them to link the Platform to their social media account (or that of the NGO), making it necessary to share personal data to enable this connection. Since we want to avoid this, and protect operators who often interact online using fake accounts, a different solution was implemented, requiring copying and pasting the replies from the Platform.

The **suggestion tool** has been implemented utilising a **data-driven approach**. In particular, it relies on a pool of pre-existing "Islamophobic message – counter-argument" pairs that are used by the tool as examples to select and rank possible replies given an input message. A language specific pool of pre-existing "Islamophobic message – counter-argument" pairs is used for each account. We implemented a tf-idf[3] response retrieval mode, which is built by calculating the tf-idf word-document matrix of the message pairs mentioned above. The suggested responses for a new input message are obtained by finding the hate message in the **pool of examples** that is most similar to the input one and presenting in the interface the five most relevant responses.

---

[3] Base concept in NLP. For a technical explanation, see Spärck Jones (1972).

The example pairs needed to build the tool were created **with the help of NGOs operators** following the **same data collection procedure for each language** (i.e. English, French and Italian). The data collection was conducted using the following steps:

1. **Hate speech preparation.** For each language, we asked two native speaker experts (NGOs' trainers) to write around 50 prototypical short Islamophobic hate texts.

2. **Counter-narratives collection forms.** We prepared three online forms (one per language) in which users can respond to hate text prepared by NGOs' trainers, operators were asked to compose up to 5 counter-narratives for each hate text.

3. **Counter-narratives collection sessions.** For each language, we performed three data collection sessions on different days. Each session lasted roughly three hours and had a variable number of operators -- usually around 20. Operators were gathered in NGOs' premises with a computer and received a brief introduction to the counter-narrative collection task.

4. **Data annotation and augmentation.** After the data collection phase, we employed the services of three non-expert annotators to perform additional work. In particular, they were asked to (i) paraphrase original hate content to augment the number of pairs per language, (ii) annotate hate speech sub-topics and counter-narrative types (iii) translate content from French and Italian to English to have parallel data across languages.

In total, we had **more than 500 hours of data collection with NGOs** involved in the Hatemeter Partnership, during which we collected **4,081 hate speech/counter-narrative pairs**; specifically, 1,288 pairs for English, 1,719 pairs for French, and 1,074 pairs for Italian. At least **111 operators** participated in the **9 data collection sessions** carried out face-to-face and remotely. Each counter-narrative needed about **8 minutes** on average to be composed. The paraphrasing of hate messages and the translation of French and Italian pairs to English brought the total number of pairs to more than **14 thousand.**

**If operators modify a counter-narrative message** or write their own messages and save them within the CAP interface, they are automatically added to the pool of examples to increase its variability and size, so that **they can be used to suggest additional counter-arguments in the future**.

A **final evaluation experiment** showed that, as expected, using the **Hatemeter counter-narrative (CN) suggestion tool** drastically **reduces the time needed to obtain a CN**: with the CN suggestion tool the time needed for composing a new CN or for modifying a suggestion is almost halved as compared to the one without suggestion tool (**8 minutes vs 4.5 on average**).

# 3. Hate speech against Muslims on social media: evidence and analysis

This section presents **evidence of online Islamophobia** from the three countries involved in the project, namely Italy, France and the UK, subdivided into three parts. Each subsection starts with a brief description of the context and background of online Islamophobia and anti-Muslim hatred in the country, and then specifically presents **evidence of hate speech collected through the Hatemeter Platform**, and a **description of the information gathered** by the visualisation tools available within the Platform. The three subsections are the result of collaboration between the three NGOs and two of the three universities involved in the project. In more detail, Amnesty International Italy and the University of Trento wrote section "3.1 Italy"; Collectif Contre l'Islamophobie en France (CCIF) wrote section "3.2 France"; and Stop Hate UK (STOPHATE) wrote section "3.3 UK" with a contribution from Teesside University for the subsection 3.3.2.

## 3.1 Italy

### 3.1.1 Islamophobia in Italy

The situation of online hate speech against Muslim communities in Italy presents a distressing picture (see deliverable **D7 "Guidelines on the socio-technical requirements of the Hatemeter Platform"[4]**). The target groups for online hatred in Italy tend to be **migrants from all nationalities**: hatred is associated with the idea that immigration policies are economically and socially unsustainable. Islamophobia becomes particularly noticeable in that migrants are usually linked with Muslims and Muslims are then linked with terrorists.

As in other countries, public discussion has changed with the birth of the **Internet and social media.** There is an increasing connection between **alternative information websites** (e.g. blogs, informal webpages not connected to a specific institution/journal)**, social networks** and **traditional mass media** (especially newspapers). Social networks facilitate the quick and easy movement of information on 'hate news' between alternative information websites and traditional mass media and vice-versa. 'Hate preachers' tend to be individuals, rather than groups and the use of social networks has replaced their use of websites and blogs. Social networks have a greater capacity to convey messages, while open platforms make individuals accountable for the content of their messages.

**Fake news** and **inflammatory statements against Muslims proliferate** on the Internet and via social media platforms (Giancalone 2017). The Internet offers immediacy, pervasiveness, amplification, replicability, social validation and the persistence of specific messages. At the same time, social media platforms offer a polycentric proliferation of hate speech and promote the diffusion of demagogic and propagandistic messages. Importantly, the **online and the offline worlds** are increasingly **connected** and the impact of one upon the other is often underestimated (this is the so-called 'prejudice of the digital dualism') (Giovannetti and Minicucci, 2015). A recent study on Islamophobia in Italy reports an increase in discriminatory articles in newspapers and in hate speech against Islam by **Internet-based neo-fascist** and **Catholic fundamentalist groups** (Alietti and Padovan, 2018). According to another study, Muslims are the **fourth most targeted** group on Twitter (Vox, 2019). Amongst the criticisms and attacks there are many that conceptualise Islam as a violent, absolutist, anti-democratic religion that is against and incompatible with Western values (Malchiodi, 2016). Moreover, the United Nations refers to the existence of a

---

[4] The contributors of the deliverable are: Andrea Di Nicola, Stefano Bonino and Elisa Martini (UNITRENTO), Jérôme Ferret, Mario Laurent and Jen Schradie (UT1-Capitole), Georgios Antonopoulos and Parisa Diba (TEES).

dangerous prejudice against immigrants in Italy, especially emanating from politics and the media (Osservatorio sulle Discriminazioni, 2010).

Hate speech is soaring alongside xenophobia, Islamophobia, anti-Semitism and racism, as a result of both **terrorist attacks** and **migration flows** (Bortone and Cerquozzi, 2017).

### 3.1.2 Evidence of Islamophobia in Italy through the Hatemeter Platform

In order to address the phenomenon of online hate speech, **Amnesty International Italy** has developed a program based on multiple levels of intervention and targeting different demographics. Among these activities, there is the "**Hate Barometer**" (Amnesty International Italy, 2019), i.e. a social media-monitoring tool. The Barometer measures the level of hate on **Facebook and Twitter** through the implementation of an **innovative methodology**, which combines the **use of algorithms** and the **involvement of activists** in the role of "content evaluators". In parallel, Amnesty International Italy also promoted the foundation of a permanent network of experts from universities, institutes of research, civil society organizations and institutions who share skills, good practice and tools on a regular basis to tackle hate speech. The data collected through these activities does not only have an impact on public opinion, but is also used lobby for change to prevent online hate. In order to achieve these changes, Amnesty International Italy conducts an ongoing dialogue with the relevant institutions and stakeholders. Within this framework, the Hatemeter Partnership and its Platform helped in advancing the work conducted on social media, with new functionalities such as an innovative monitoring and counter-narrative tool, with a specific focus on anti-Muslim hate speech.

Amnesty International Italy tested the Hatemeter Platform during the **two Piloting Sessions** (see Deliverables no. 14 and 15 on the platform evaluation), thanks to a simultaneous project called "**Task Force Hate Speech**". The "Task Force Hate Speech" project was realized after an experimental phase that took place between April and November 2016, in which the pilot project "Web Task Force" was tested. Based on a group of 15 young activists', the pilot project had the strategic goal of improving Amnesty International Italy's knowledge of, and responsiveness to hate speech.

Due to the successful results of the pilot project, the Italian Section of Amnesty International decided to continue the work on countering hate speech with a specific group of activists: the **Task Force Hate Speech**, on monitoring, preventing and combatting the spread of online hate speech. The group now has over 200 trained activists who intervene in the comments sections of online platforms (Facebook and Twitter) where are likely to find explicit hate comments targeting minorities or vulnerable groups. The **main objectives** being to re-focus upon the factual basis of the **information** presented in order to promote, an **objective discussion** on the specific topic, and to promote the use of a **specific, less adversarial and respectful dialogue and tone/spirit**. Utilising the initial Hate Barometer, Amnesty International Italy choose to monitor social media content generated during the **election campaign** of February-March 2018, aiming at "tackling the violent, aggressive, discriminatory discourse and spreading an appropriate use of words". Amnesty International Italy monitored declarations and comments posted on the social profiles (i.e. Facebook and Twitter) of political candidates to verify the "level of hate" contained within the political arena and addressed at vulnerable groups, such as migrants, Roma, LGBTI, and members of the Jewish and Muslim communities. These were candidates listed as nominated for the Chamber of Deputies and Senate of the Republic of the main four Italian parties and coalitions (i.e. centre-right, centre-left, 5 Stars Movement, and Free and Equal) and of the candidates to the presidency of the regions of Lazio and Lombardy.

In November 2018, the Amnesty International Task Force Hate Speech started **working jointly with the Project Hatemeter** focusing on the contrast of the online Islamophobia. The collaboration with the training sessions and continued during the piloting of the Hatemeter Platform, as well as in the development of

an awareness raising campaign. The Task Force gave support to the creation of the Hatemeter Platform by providing examples of hate speech and Islamophobic discourse followed by counter narrative examples. The initial phase of the Hate Barometer (February-March 2018) has also been refined as a result of the work done during the piloting of the Hatemeter Platform, and the outcomes resulted not only in a second version of the Hatemeter Platform (see Deliverable no. 11 – released in May 2019), but also a new, advanced Amnesty Hate Barometer, also in May 2019.

During the **Hatemeter piloting sessions** from January to March 2019 and from June to September 2019, as well as in May 2019 Amnesty International Italy collected, through specific algorithms, more than 4 million pieces of content. Over **180 trained activists** assessed 100,000 pieces of content, with the aim of detecting possible correlations between tone/spirit and the political rhetoric of politicians, and sentiments of social media users toward specific topics and groups of people. On this occasion, the Task Force utilised the **Hatemeter Platform**, both for the collection of hateful content and the rapid generation of alternative replies to messages posted on Twitter.

Thanks to the Platform, we were able to shed light upon rise of hateful messages **directed at religious minorities** where these are connected to "terrorism" (Islam). Anti-Muslim sentiments find their roots not only in the connections between Islam and "invasion", "terrorism", and "barbarities", but also in the idea that Islam presents an obstacle to the advancement of the feminist and LGBTI movements. Moreover, these hateful messages are more likely to have recourse to acrimony and personal attacks.

By way of example, below are reported some tweets collected by Amnesty International Italy as comments to **posts from Italian politicians**. These tweets were collected and presented in the *"Barometro dell'odio. Elezioni europee 2019"* report (2019) and for this reason the screenshots are not available.

More specifically, the texts of examples 1 and 2 were collected from the official page of the former Italian Minister of the Interior (right wing) under a tweet displaying a video starting with an argument between a train conductor and a foreign passenger without a ticket, which escalates into a violent quarrel. The post reached more than 1 million views, more than 80.000 likes and reactions, and more than 22.000 replies. In the examples below (1 and 2) the concept of "invasion" is frequently cited, and is also referenced within the political debate. As shown by the figures, the **counter-narrative functionality of the Hatemeter Platform** can assist in responding to these tweets, by providing suggestions and insights to NGOs operators and thereby accelerate the response.

**Figure 8 – Example 1: hate tweet and answers provided by the "counter-narrative" Functionality of the Hatemeter Platform.**



*SOURCE: UNITN elaboration – Screenshot from the Hatemeter Platform*

| English translation of the tweet: The African and Islamic invasion that the European Union is imposing upon us will be the end of the Western world. |
|---|
| English translation of the first counter-narrative: Where do you get these conclusions about this alleged 'invasion' from? |
| English translation of the second counter-narrative: In Italy, Muslims are less than 5%, personally I would not define this an invasion. |
| English translation of the third counter-narrative: I do not think so. Muslim population in Italy represents only the 4% of the total population. I would not define this an invasion. |

**Figure 9 - Example 2: hate tweet and answers provided by the "counter-narratives" Functionality of Hatemeter Platform.**



*SOURCE: UNITN elaboration – Screenshot from the Hatemeter Platform*

| English translation of the tweet: We must end this African and Islamic invasion … they must stay in Africa … we do not need them. They are dangerous for us and all Europe. |
|---|
| English translation of the first counter-narrative: Where do you get these conclusions about this alleged 'invasion' from? |
| English translation of the second counter-narrative: In Italy, Muslims are less than 5%, personally I would not define this an invasion. |
| English translation of the third counter-narrative: I do not think so. Muslim population in Italy represents only the 4% of the total population. I would not define this an invasion. |
| English translation of the fourth counter-narrative: Hello, could you please explain more clearly what you mean by Islamic invasion? And can you provide me with data to support your statement? |

In both cases (Figures 8 and 9), **the counter-narratives functionality suggests that users reflect upon the concept of "invasion"**, primarily by providing objective data on the presence of Muslims in Italy (i.e. the Platform affirms that Muslims represent between 4 and 5% of the total population). The aim is to present a basis for reflection, which may then curtail the spread of the hate speech.

On the page of **another Italian politician** of the right wing, one image reporting "No to the Islamisation of Europe" had been shared more than seven thousand times and it generated more than 53,000 likes and reactions, and more than 2,000 comments. Among these, Amnesty International Italy retrieved examples of hate speech (reported in Figure 10 and 11), from which can be observed the emphasis placed upon the conservation of national identity.

**Figure 10 - Example 3: hate tweet and answers provided by the "counter-narratives" Functionality of the Hatemeter Platform.**



*SOURCE: UNITN elaboration – Screenshot from the Hatemeter Platform*

| English translation of the tweet: They come here to dominate us, with the help of the left wing that wants globalization |
|---|
| English translation of the first counter-narrative: Being devoted to one's religion does not mean proselytizing. |
| English translation of the second counter-narrative: Such an indiscriminate answer does not seem a solution to the migratory phenomenon; people who come to Italy do not come to "conquer us", they come to try to have a better future. |
| English translation of the third counter-narrative: A person who lives in conditions of profound social distress and migrates to another country, exports his or her own tradition and religion, and wants to be able to maintain it, just as the country that hosts his/her wants to keep their own. This does not mean submission, it means cohabitation, exchange, mutual enrichment. |

**Figure 11 - Example 4: hate tweet and answers provided by the "counter-narratives" Functionality of the Hatemeter Platform.**



*SOURCE: UNITN elaboration – Screenshot from the Hatemeter Platform*

| English translation of the tweet: CLOSED PORTS AND NO ISLAM |
| --- |
| English translation of the first counter-narrative: As regards diseases, it is a false news that those who arrive in Italy bring new ones: as soon as they arrive, they are subject to a medical examination. |
| English translation of the second counter-narrative: Sure enough, it's true that Islam is an Abrahamic religion like Christianity and Judaism. However, religious fanaticism does exist, which distorts religion and makes it a tool for manipulation |
| English translation of the third counter-narrative: I would like to point out that Islam and Isis are different. Islam is a peaceful religion, while ISIS is a military and terrorist organization that supports Islamic fundamentalism. |

Figures 10 and 11 display the **suggestions of the counter-narratives functionality** of the Hatemeter Platform for examples 3 and 4. In these cases, the proposals are more general, trying to assure the users that migrants escape from difficult situations in their countries and even if they want to maintain their religious beliefs and habits in Italy, they do not want to impose them upon the local population. Nevertheless, not all sentences can be suitable answers to the tweet (e.g. Figure 11, the first counter-narrative refers to "diseases" that are not cited in the original tweet) and **NGO operators should read them carefully to choose the most efficient counter-narrative.**

In conclusion, starting from the Piloting Sessions of the Hatemeter Platform and the simultaneous activities of the Amnesty Hate Barometer, it has emerged that the adjective **"Islamic"** is frequently

associated with other terms, such as: **"Islamic extremism"**, **"Islamic terrorism"**, **"Islamic fundamentalism"**, **"problem"**, **"immigration"**, and **"danger"** (Faloppa 2019). Such monitoring activities confirm some trends already registered with the online monitoring in 2018 realised by Amnesty International Italy (i.e. during the political campaign). Even if politicians seem to maintain more moderate tones/spirit in comparison to their followers, candidates to the European elections wrote almost one in five tweets that were considered negative because of discriminatory content directed at minor religious groups in Europe (Vitullo 2019). As in the Hate Barometer of 2018, **the Muslims community is confirmed as the religious group most targeted by political discourse in Italy** (Vitullo 2019).

## 3.2 France

### 3.1.1 Islamophobia in France

Article 1 of the French constitution of 1958 sets out a **"principle of equality"**, essentially forbidding all forms of discrimination: *"La France est une République indivisible, laïque, démocratique et sociale. Elle assure l'égalité devant la loi de tous les citoyens sans distinction d'origine, de race ou de religion. Elle respecte toutes les croyances. Son organisation est décentralisée"*, roughly translated as: "France shall be an indivisible, secular, democratic and social Republic. It shall ensure the equality of all citizens before the law, without distinction regardless of origin, race or religion. It shall respect all beliefs". As such, equality of treatment and non-discrimination principles are supposed to be at the heart of the French national identity. However, **equality is interpreted as uniformity** and France essentially implements the assimilation of minorities and differences. This is the French concept of the **"melting pot"**: French people are all equal because everyone is the same and treated as such. However, this concept erases each citizen's individuality in the eyes of the state, even though it is initially supposed to protect them from the state making distinctions and preferences amongst them. In truth, this has been used increasingly in past years as a means to force citizens to present themselves as religiously "neutral".

Clear-cut cases of **racism** are in theory therefore **punishable by law**, even though **cultural forms of Islamophobic** racism are still gaining traction. In recent years, the influence of far-right political groups has spread across Europe on a broader spectrum, and tainted other more "mainstream" political groups. In France, Islamophobic hate speech is justified by the "defence of secularism", "national identity" and "preserving France's Judeo-Christian roots", to cite only a few. According to the CCIF's annual report this has contributed to a 52% increase in Islamophobic acts committed in 2018 in comparison with 2017,

Instead of understanding and respecting the Faith of Muslims and the freedom to choose their own lifestyles, **political and mass media** discourse in France is mainly focused upon the objective of **educating Muslims** that presents **Western ideals** as synonymous with freedom. Paradoxically, this need for Muslims to "integrate into society" is used to **ban the practice of their faith in society**. Although there is no legal or logical incompatibility between French values and Islam, French Muslims face increasing demands that they disavow various aspects of their religious practice to demonstrate their allegiance to their country.

Recently, **"societal concerns"** have grown over the **religious headwear** of Muslim women, known as the hijab. The sheer amount of political debate over women's practice of their faith has concrete and severe consequences: **70% of all Islamophobic acts are committed against women.** As identifiable people of Muslim faith, they are targeted and harmed because of it. In 2018, **676 cases of Islamophobia** were reported to the CCIF. This number is an under-representation, since many acts go unreported. This is where one starts to see that online controversies have real life consequences. **Cross-partisan Islamophobia** is gradually being normalized, and as such is obviously increasingly concerning, putting various human rights in clear danger. Islamophobia is not a form of "freedom of expression". It is a crime that harms and murders, and destroys entire communities. The **link between online debates,**
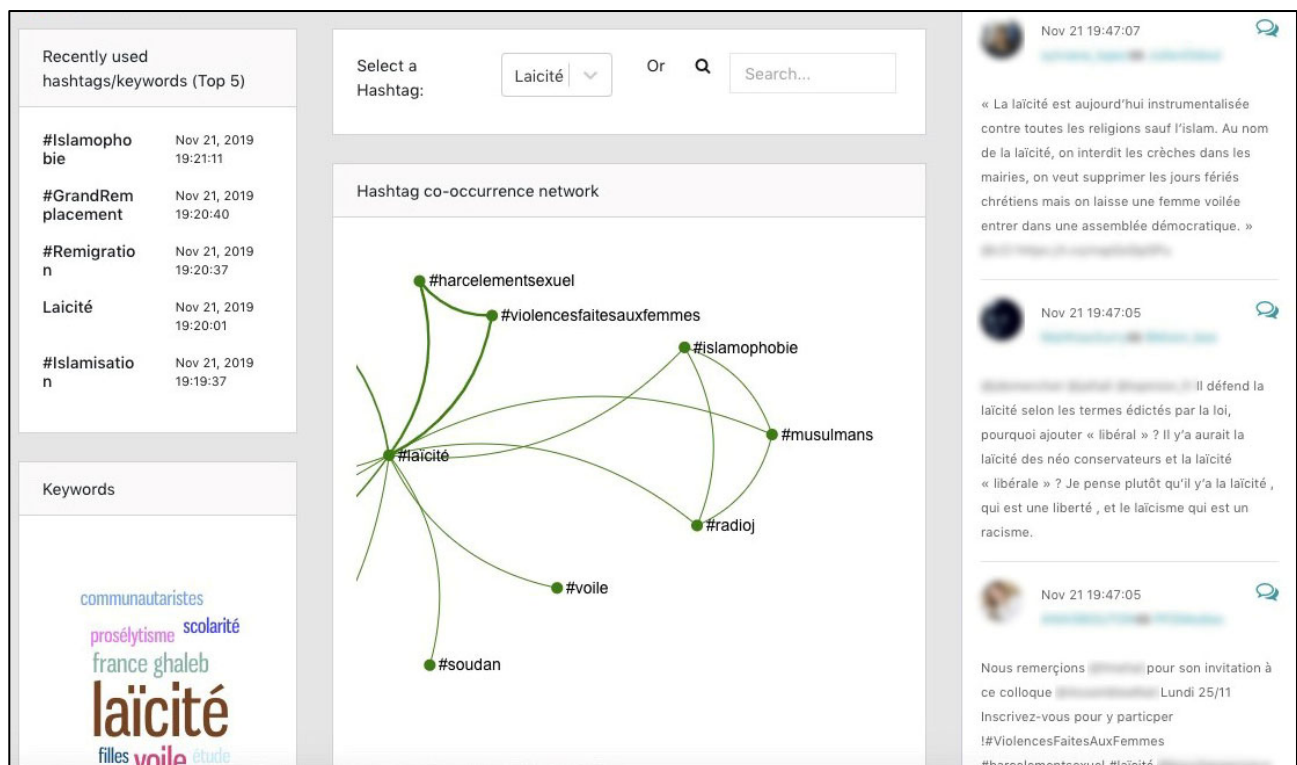
Islamophobic arguments being spread on national television, and **tangible hate crimes** is not hard to identify.

**This is why combatting hate-speech is so crucial and targeting online hatred so essential;** the normalisation of Islamophobia online has proven to have real life effects and violent outcomes. It is therefore essential to understand the phenomenon, collect data on when, why and how it is happening, and analyse it to generate strategies to stop its proliferation.

### 3.1.2 Evidence of Islamophobia in France through the Hatemeter Platform

The **Hatemeter Platform** also allows French NGOs to gather evidence of online Islamophobia. For example, a Twitter user misrepresented the concept of French "laïcité" on his twitter feed to promote an interview he gave. In the following tweet displayed by Figure 12, gathered thanks to the Hatemeter tool through the *"Recent trends"* functionality, he completely negated the reality of Islamophobia.

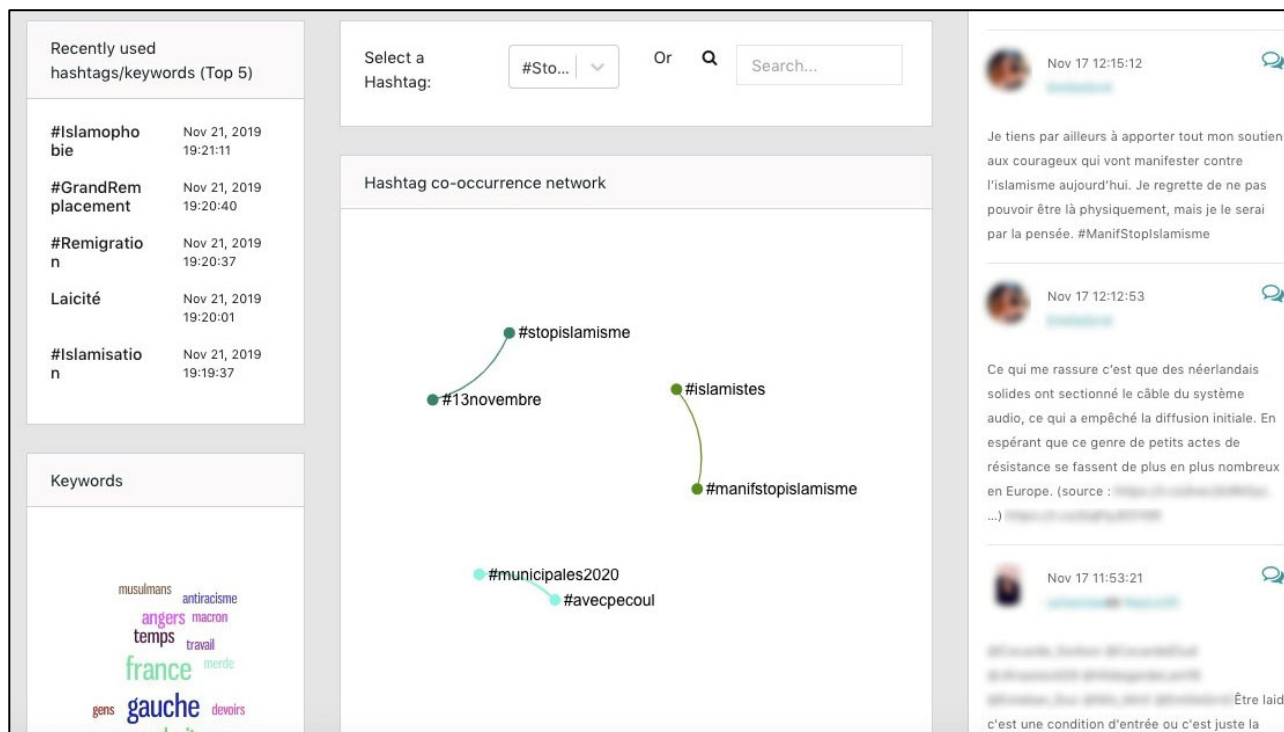**Figure 12 - The "Recent Trends" Functionality of Hatemeter Platform V.2 utilising #Laicité**



*SOURCE: CCIF elaboration – Screenshot from the Hatemeter Platform*

Utilising the Hatemeter Platform, any NGO that wishes to prove the usage of **#laïcité** in an Islamophobic context, can search the hashtag #laïcité and instantly flip through all the most recent tweets that used that key word, and see the links made with other topics, usually indicating Islamophobic discourse. For example, here we can see that laïcité was linked to **"musulmans"** (muslims) and **"voile"** (headscarf).

The Platform allows analysis of the most recent trends and use of semantics by Islamophobic accounts and tweets. However, in our findings, we can often deduce other conclusions than the use of specific words of justification. For example, a student activist closely linked with extreme right militia in French universities posted a tweet of support towards a hate crime committed against a Mosque in Amsterdam (see Figure 13). The user employed the hashtag **"StopIslamisme"**, showing the endless possibilities to voice support of Islamophobic acts across borders thanks to social media, and the instantaneous
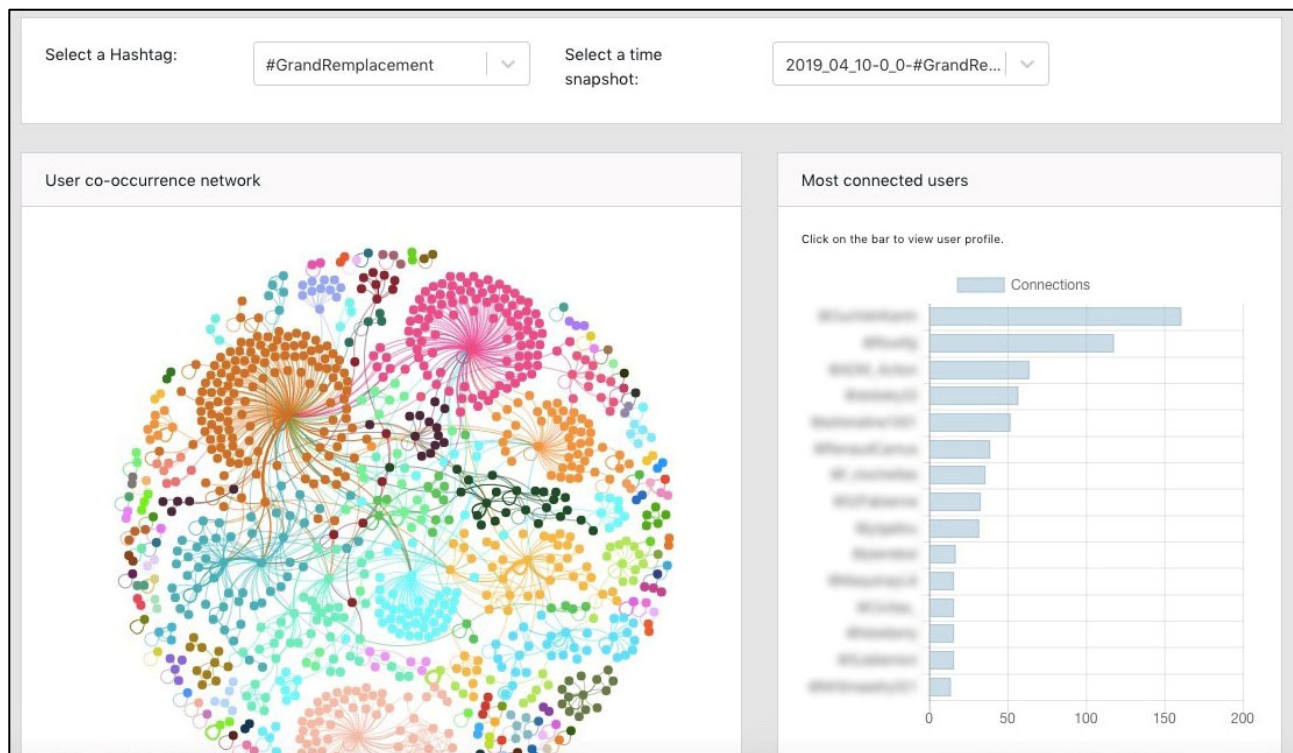
influence this may have on various countries. In the example of Figure 13, one tweet is clearly calling for such "acts of resistance" (hate crimes) to be multiplied across Europe.

**Figure 13 - The "Recent Trends" Functionality of Hatemeter Platform V.2 utilising #Stopislamisme**



*SOURCE: CCIF elaboration – Screenshot from the Hatemeter Platform*

Moreover, the **"Hate Speakers" functionality** facilitates delving deeper into a specific account to analyse its influence or any other specific reason for the data analysis. Here, it is possible to research a specific topic and see all the accounts that gained the most traction due to the use of a hashtag. For instance, by searching **"#GrandRemplacement"** in the most recent timeframe the platform offers, one of the twitter accounts that mostly use this hashtag is highlighted (see Figure 14).

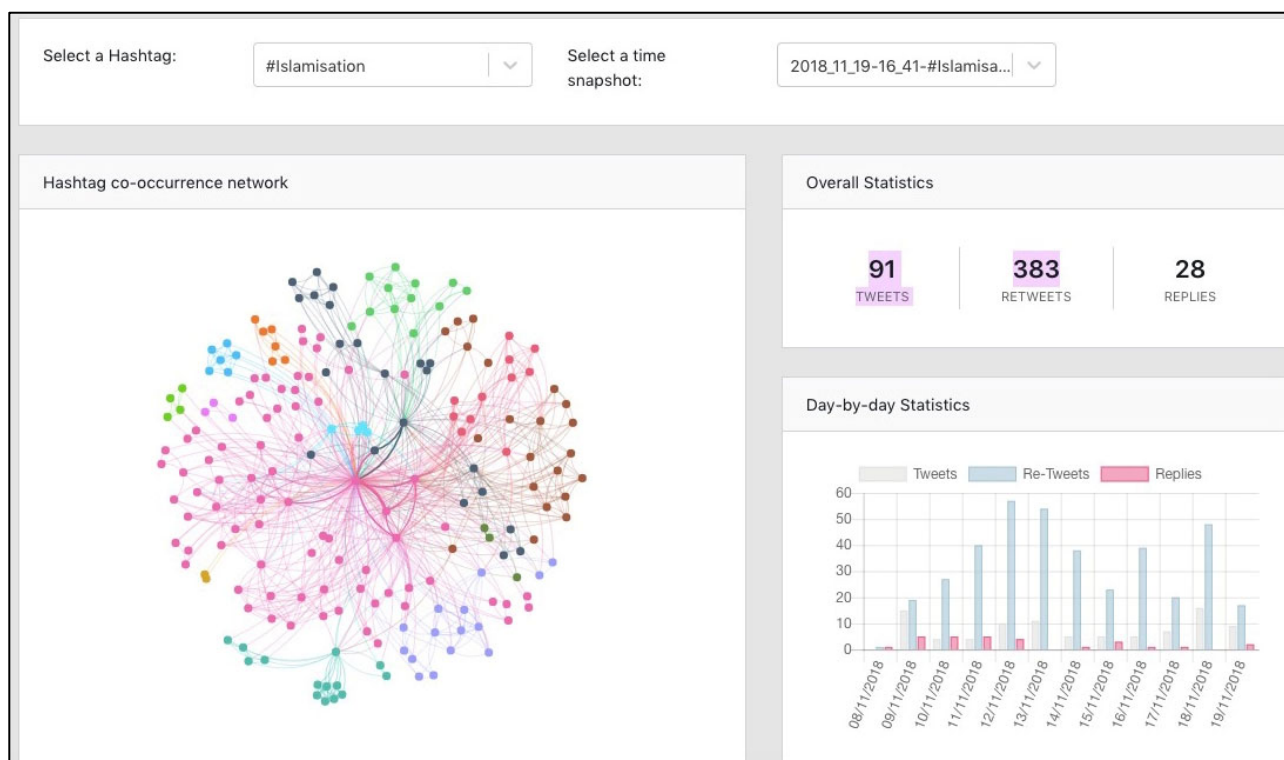**Figure 14 - The "Hate Speakers" Functionality of the Hatemeter Platform V.2, using #GrandRemplacement**



*SOURCE: CCIF elaboration – Screenshot from the Hatemeter Platform*

By clicking on the user's account, the NGO can realise the impact this Twitter user might have. Here, the user has 161 connections to the hashtag #GrandRemplacement. Their tweets range from declaring a "cultural war" to protesting in the streets against Islam with extremist groups like Génération Identitaire.

The Platform offers the possibility of **searching for keywords with the use of timeframes**; this allows linking the rise of a specific point of view to dated events. In France, activists are aware that elections are often accompanied by a rise of Islamophobic sentiment. However, thanks to the data that can be collected by the Hatemeter project, NGOs can now evidence the online effects of political statements. This will be **crucial for the NGOs' work in fighting Islamophobia**, because they can acquire greater legitimacy by **providing facts through data analysis**.

Furthermore, the cumulative use of these timeframes and keywords after televised events allow us to see the online response to political debates. For instance, by using **the "Hashtag Trends" functionality** with **"#Islamisation"** from 8th to 19th November 2019, we could analyse, compare, and draw conclusions from overall and day-by-day statistics (see Figure 15).

**Figure 15 - The "Hashtag Trends" Functionality of the Hatemeter Platform V.2, using #Islamisation**



*SOURCE: CCIF elaboration – Screenshot from the Hatemeter Platform*

By scrolling down to see **"most retweeted messages"**, it is possible to notice the most indicative reactions to the controversial political debates of that timeframe. For instance, the second tweet in the list shown below shows the momentum gained by an Islamophobic response to the **Gilets Jaunes protests** t (see Figure 16).

**Figure 16 - The "Hashtag Trends" Functionality of Hatemeter Platform V.2, using #islamisation, most retweeted messages**



*SOURCE: CCIF elaboration – Screenshot from the Hatemeter Platform*

This analytical use of trends and time frames can also be used as a tool to see **the tone of different discourses on social media after terrorist attacks**. How does public opinion of Muslims change after a terrorist attack committed by someone in the "name of Islam"? How do people respond to online Islamophobia? What are the numbers on each side? Hopefully**, these questions will be able to be answered thanks to the Platform**. Furthermore, in addition to answering these questions, the Hatemeter Platform will provide NGOs with enough quantifiable data to hopefully prevent such hate speech: because the same key words and arguments appear repeatedly, thanks to the **"Counter-narratives" functionality,** activists will be able to provide ready-made responses to online hate **speech.** These counter-narratives will cover as many topics as possible, and are created to save time for activists, preventing them from having to repeat themselves. As the screenshot below shows (see Figure 17), any operator of the Platform can answer multiple points raised by a form of hate speech with a click of a button, thanks to the counter-narrative tool. Platform users (e.g. NGO operators) can copy and paste a tweet into the Hatemeter tool for counter-narratives, and have available **multiple potential answers from which to select a reply**. Figure 17 provides an example of the types of answer the Platform can offer in response to the hate-tweet reported in the first rectangle, affirming *"Le voile est un symbole de soumission, l'islam politique envahit la France"* (English translation: the hijab is a sign of submission, political Islam is invading France).

**Figure 17 - The "counter-narratives" Functionality of the Hatemeter Platform V.2**



*SOURCE: CCIF elaboration – Screenshot from the Hatemeter Platform*

| English translation of the tweet: The veil is a symbol of submission, political Islam invades France. |
|---|
| English translation of the first counter-narrative: Islam is a religion not a political party. The veil is a fabric not a symbol. |
| English translation of the second counter-narrative: Islam is a religion. The veil is a sign of devotion and not a political flag. |

## 3.3 UK

### 3.3.1 Islamophobia in the UK

In the UK, the Equality Act 2010, amalgamated over 116 separate pieces of legislation into one single Act. The Act provides a legal framework to protect the rights of individuals and advance equality of opportunity for all. The Act covers the same groups that were protected by pre-existing equality legislation (i.e. age, disability, gender reassignment, race, religion or belief, sex, sexual orientation, marriage and civil partnership and pregnancy and maternity), that are now designated as "protected characteristics". In addition, the Public Order Act 1986 makes it an offence to intentionally stir up hatred on the grounds of race, religion and sexual orientation, and the Crime and Disorder Act 1988 creates racially and religiously aggravated offences. By referring to both pieces of legislation, the term **"Hate Crime"** is used to describe a range of criminal behaviours where the perpetrator is motivated by or demonstrates hostility towards the victim's "protected characteristics" (see above). Where an offence is designated as a "hate crime" an "enhanced" or more severe penalty can be applied. While the legislation seeks to offer protections for followers of all faiths, this has to be considered in the context of the **increasing normalisation of "Islamophobia" in public and political discourse**, the print and broadcast media and the online space.

As such, this is reflected in the Government's recently published statistics on hate crimes recorded by the police and information on hate crime from the **Crime Survey for England and Wales** (Home Office 2019). It indicates that in 2018/19, where the perceived religion of the victim was recorded, just **under half (47%) of religious hate crime offences were targeted against Muslims (3,530 offences)**. The next most commonly targeted group were Jewish people, who were targeted in 18% of religious hate crimes (1,326 offences). In 2017/18, where the perceived religion of the victim was recorded, just over half (52%) of religious hate crime offences were targeted against Muslims (2,965 offences) (Home Office 2018). This is a much greater proportion than the proportion of the population of England and Wales that identify as Muslims. In the 2011 Census, **4.8% of the population of England and Wales** identified as **Muslims**. While the figures do indicate a **rise in the number of recorded offences**, this also has to be considered in the context of concerns that a **large number of incidents are simply not reported to the authorities**.

Recent attempts to agree and adopt a definition of Islamophobia, assert that it "is rooted in racism and is a type of racism that targets expressions of **Muslimness or perceived Muslimness**" (All Party Parliamentary Group on British Muslims 2018). While adopted by parties including Labour, the Liberal Democrats and the Scottish Conservatives, the definition was rejected by the Government because the definition is currently too broad and might be used to hinder free speech and criticism of the historical and theological actions of Islamic States (Elgot 2019).

By contrast, we now also have at least **two political parties actively campaigning on "Anti-Islam" issues** as part of their respective manifestos (i.e. For Britain and United Kingdom Independence Party). They indicate the growing acceptance of a range of anti-Islamic narratives and stereotypes, which have gained increasing traction across the UK but are also clearly proliferating in the online space, both on a national and trans-national basis.

In closing this section and as a means of demonstrating **how these narratives are shaping public opinion** and beliefs, and, in some cases, provoking individual actions, below is a report taken by **Stop Hate UK** from a concerned member of the public who contacted our **Helpline** on 25/11/2019:

> *I am getting in touch as I have just witnessed an incident of racial abuse in our local post office.*
>
> *There were several people in the queue and a woman at the back, suddenly went towards the counter and asked the post-master 'Did you post my signed for letter?' She then accused him of not posting it. He responded, calmly trying to explain that she could check it by entering the code and asked when she sent it. She replied it was months ago, she then accused him of being a 'Muslim, thief'. This was followed by 'you don't belong round here, we do' (she is likely to be a British person but I can't say what her ethnicity is). She told him that he 'should go home'.*
>
> *Several of us asked her to stop, it wasn't right etc., but she ignored it.*
>
> *She left the post office, but returned almost straight away, this time to accuse the post-master of being a Muslim paedophile and he was to stay away from her kid. She was swearing and called him a f\*cking dirty Muslim several times.*
>
> *The whole incident was between 2-3 minutes. The post-master remained calm, he didn't respond in an aggressive way.*

### 3.3.2 Evidence of Islamophobia in the UK through the Hatemeter Platform

In considering the evidence submitted in relation to **Stop Hate UK's trialling of both versions of the Platform**, it is perhaps important to note several points:

- As an established Hate Crime organisation, providing **'third party' (independent) reporting facilities across the UK,** we routinely receive reports of online Islamophobia from members of the public, both as 'victims' and 'bystanders', which are also used to identify potential 'targets' for both reporting and counter messaging by staff and volunteers involved in our 'No Hate Speech' team.

- In **testing the Hatemeter Platform**, we are operating within a project established in 2016, with pre-existing practices and methodologies. We would therefore acknowledge that, at this stage we have not relied solely upon usage of the Platform as our primary mechanism for identifying and responding to instances of online Islamophobia.

However, **our usage of the Hatemeter platform has produced positive results** in terms of assisting new volunteers / staff members **to quickly and accurately identify sources of Islamophobic hate speech on Twitter**. Specifically, operators have been able to observe various **user networks**, including central users and seed accounts, with the Tool enabling NGO operators to discern who these users are within a network, not just 'infamous' accounts acting as central nodes, but also unknown accounts and consistent accounts across time and space.

The Platform has provided confirmation of our own tentative conclusions and, in particular, in relation to the graphic representations of account interactions, aided an understanding of the **international / transnational aspects and interrelationships of online Islamophobia.** Given our observations of the range of interactions between UK and US based accounts, and between European/Eastern European based accounts, we anticipate that this function will become a major asset, enabling much more accurate identification and evidencing of these interrelationships. This in turn will assist in the development of **more effective and targeted strategies for combatting online Islamophobia**

Our utilisation of the **Counter-message generation aspect of the Platform** has been minimal, and largely confined to the Deployment day itself, for a number of obvious reasons which themselves are hopefully reflected in the feedback we have provided as part of the evaluation process. Nevertheless, the responses suggested by the version utilised during Deployment Day 2, clearly demonstrated their value as a **'catalyst' or starting point for the generation of a counter-message**, (which clearly aids speed and efficiency in the generation of a response) and their major potential as a training tool,

The evidence below is a set of separate instances, in which the Hatemeter Platform, and its various functions, have been utilised to identify and combat hate speech pertaining to online Islamophobia. Where possible, the outcome of engagement has also been included.

**Instance 1**

| Date: | 30/01/2019 |
|---|---|
| Issue: | Utilising the Hashtag search function, an account has been located and subsequently reported for Islamophobia, specifically a post that advocates 'eradicating' Muslims from the USA. |
| Outcome | "Account found in violation of Twitter rules", leading to a ban. |

**Figure 18 – First instance about online Islamophobia retrieved through the Hatemeter Platform**



*SOURCE: Stop Hate UK's research elaboration – Screenshot from Twitter*

## Instance 2

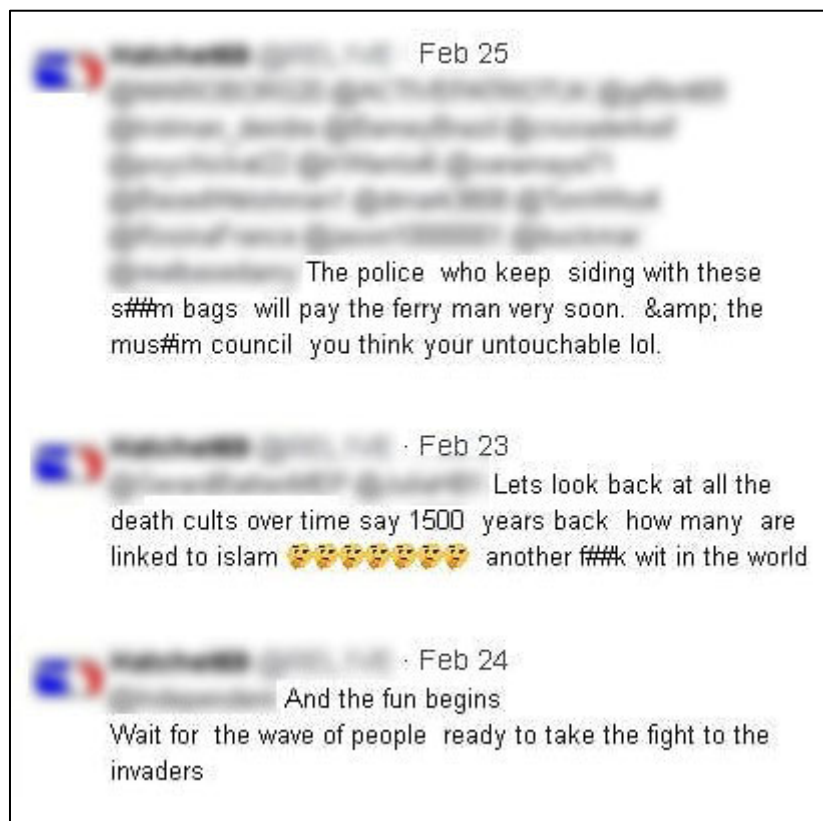| Date: | 11/02/2019 |
|---|---|
| Issue: | Utilising the Hashtag search function, an account has been located and subsequently reported for Islamophobic post, directed at US Member of Congress Ilhan Omar |
| Outcome | "Account found in violation of Twitter rules." |

**Figure 19 – Second instance about online Islamophobia retrieved through the Hatemeter Platform**



*SOURCE: Stop Hate UK's research elaboration – Screenshot from Twitter*

## Instance 3

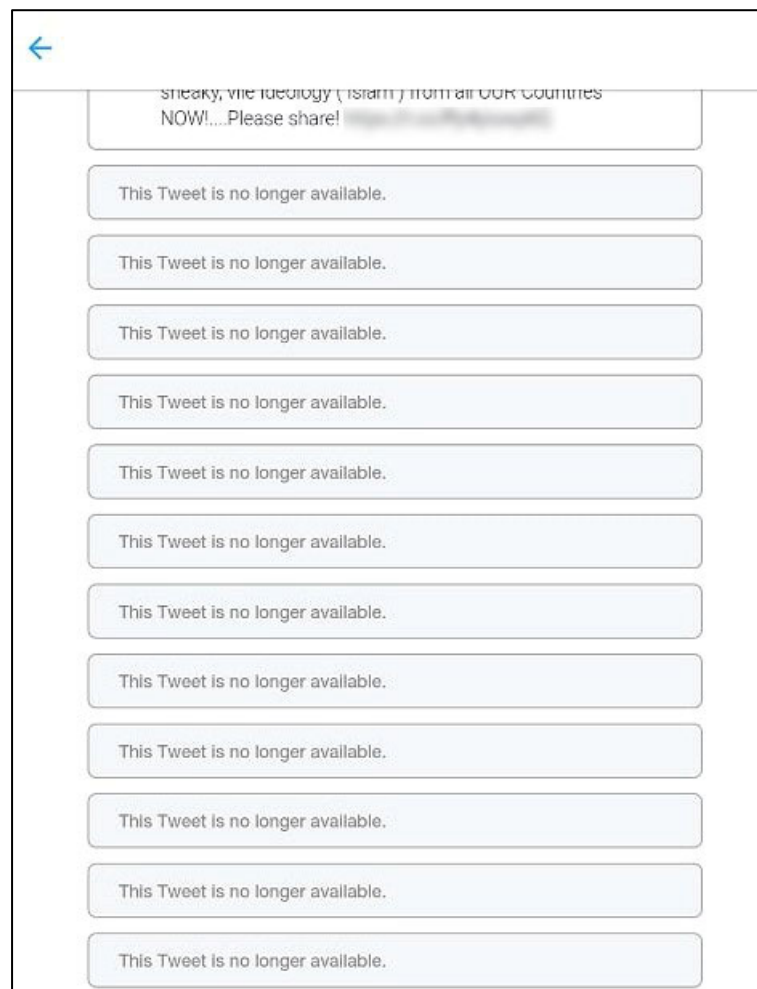| Date: | 28/02/2019 |
|---|---|
| Issue: | Using the Hatemeter Platform, an account has been located, via its interaction with several more prolific accounts. Reported a range of (5) Islamophobic posts. |
| Outcome | On the, 13/03/2109, informed: "Account found to have breached rules on abusive behaviour". |

**Figure 20 – Third instance about online Islamophobia retrieved through the Hatemeter Platform**



*SOURCE: Stop Hate UK's research elaboration – Screenshot from Twitter*

## Instance 4

| Date:    | 28/02/2019 |
|----------|------------|
| Issue:   | (Note: Included only to demonstrate the large number of individual tweets reported.) An account has been reported, for a range of breaches of the 'Hateful Conduct Policy, including Islamophobia etc. |
| Note:    | This was an alternative account established after this high profile figure had had their personal account suspended. Given the account's profile and large number of followers, evidenced by use of the Platform, we were able to prioritise reporting activity and thereby contribute to the suspension and removal of this account. |
| Outcome  | Received note on the 08/03/2019: "Account found in violation of Twitter rules." This report appears to have resulted in the removal of 29 Posts from the account. Account suspended / removed on the 10/03/2019. |

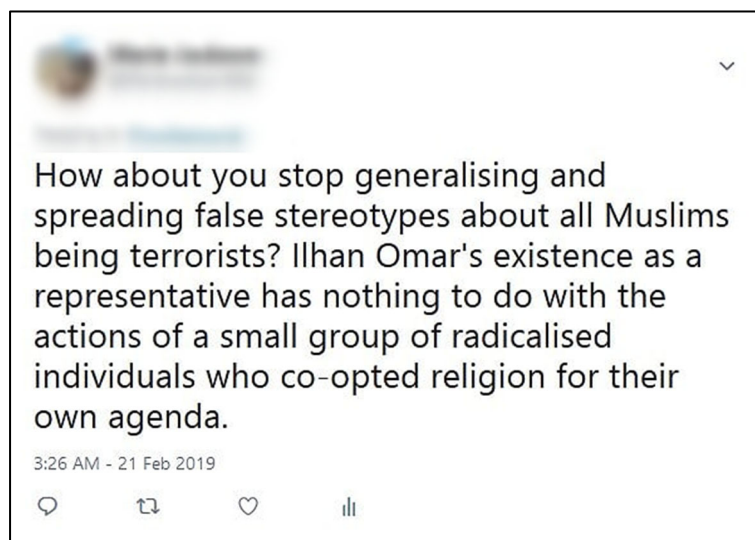**Figure 21 – Fourth instance about online Islamophobia retrieved through the Hatemeter Platform**



*SOURCE: Stop Hate UK's research elaboration – Screenshot from Twitter*

## Instance 5

| Date: | 21/02/19 |
|---|---|
| Issue: | A Twitter user posted a picture which showed the twin towers during 9/11 with the caption '"never forget" they said', and then underneath, a picture of democrat representative Ilhan Omar with the caption 'this is proof that we have forgotten'. Note: This was an early action by a staff member who neglected to screenshot the post they were reacting to. I include it as an example of use of the Platform to identify a 'hateful account' and identify a potential target of counter-messaging activity. |
| Outcome | NGO operator replied to the image in the form of counter-messaging by saying: "How about you stop generalising and spreading false stereotypes about all Muslims being terrorists? Ilhan Omar's existence as a representative has nothing to do with the actions of a small group of radicalised individuals who co-opted religion for their own agenda." However, within 20 minutes the user had blocked the NGO operator's account. |

**Figure 22 – Fifth instance about online Islamophobia retrieved through the Hatemeter Platform**



*SOURCE: Stop Hate UK's research elaboration – Screenshot from Twitter*

## Instance 6

| Date: | 21/02/19 |
|-------|----------|
| Issue: | A Twitter user posted a picture of a gollywog with the caption 'If the ISIS bride gets let back in, surely I can come back...' in reference to Shamima Begum. The tweet received a large number of retweets and supportive replies. This example also evidenced our need to support the understanding of an NGO operator from outside of the UK, who was unaware of the cultural and historical significance of the imagery in use, in the UK context. As such, it is important for NGO operators, in any geographic context, to be made aware and possess knowledge of historical, political and socio-cultural particularities relating to race and ethnicity. |
| Outcome | NGO operator reported the tweet to Twitter, citing its breach of "Hateful Conduct policy". No recorded response from service provider. |

**Figure 23 – Sixth instance about online Islamophobia retrieved through the Hatemeter Platform**



*SOURCE: Stop Hate UK's research elaboration – Screenshot from Twitter*

## Instances 7 and 8

| Date: | 21/02/19 |
|---|---|
| Issue: | A Twitter user replied to an image I had previously countered which was posted again separately by another user. The image is of the twin towers during 9/11 and claims Americans said they would 'never forget', but insinuates that the election of democrat representative Ilhan Omar is an insult to this as she is Muslim. The user tweeted "'her face and bending the rules for her hijab makes me angry'" While this example is not in itself sourced directly via usage of the Hatemeter Platform, it is directly related to an earlier example (Figure 22) and demonstrates the practitioner refining/adapting an earlier counter-message. As such, it illustrates the stark difference between 'live interaction' with an account holder, and the safer practice of responding during a training or data-gathering exercise, and therefore the need to place emphasis upon support and supervision of staff and volunteers to address the cumulative impact of such interactions. |
| Outcome | NGO operator counter-messaged by saying: |
| | "If her face alone makes you angry you need to dial back your overt racism and Islamophobia. Ever considered that the majority of Muslims don't have anything to do with terrorism and have been suffering hatred ever since 9/11?" |
| | This provoked the user to send me a string of angry responses. |
| | From this interaction, it appears that overtly pointing out racism/Islamophobia can be quite inflammatory, but it may have been useful for other Twitter users to see their narrative exposed. |

**Figure 24 – Seventh instance about online Islamophobia retrieved through the Hatemeter Platform**



*SOURCE: Stop Hate UK's research elaboration – Screenshot from Twitter*

**Instance 9**

| Date: | Undated |
|---|---|
| Issue: | Stereotypical example of low level Islamophobic Twitter post, which was sourced using the Hatemeter Platform's 'hashtag search' function. |
| Outcome | NGO operator, no recorded comments. |

**Figure 25 – Eighth instance about online Islamophobia retrieved through the Hatemeter Platform**



*SOURCE: Stop Hate UK's research elaboration – Screenshot from Twitter*

It is salient to note that, **in 4 times out of 9** of the above aforementioned cases, there were positive outcomes, in regards to **Twitter bans and/or account suspensions**, and the **removal of inflammatory posts** produced and disseminated by hate accounts. The **ability of the Hatemeter Platform to showcase re-tweets of seed accounts** – main anti-Muslim hate accounts – was highly appreciated by NGO operators, who felt that they are able make inferences, giving them prompts to analyse the meanings behind tweets, which allows them to feel out where and what content to counter in the online fora. Such action strongly suggests that with repeated, prolonged usage of the Hatemeter Platform and its various functions, the Hatemeter Tool can **assist significantly in actively identifying potent sources of hate speech** in terms of main accounts and highly incendiary posts, and aid in attempting to secure Twitter bans and/or suspensions.

# 4. Suggestions and insights on the use of the Platform and of the Hatemeter methodology for professionals

The Hatemeter Platform was developed to assist NGOs in their daily work by identifying hate speech and dangerous networks, monitoring them and providing adequate responses in a short amount of time. The tool can be particularly useful in **training new employees and volunteers:** on the one hand, it can explain the phenomenon of Islamophobia thanks to **concrete data and facts** and helps to communicate this form of racism thanks to accurate graphs and statistics. On the other hand, it provides **counter-narrative suggestions**, which can help formulate responses to hate speech. Nevertheless, the tool will always need rephrasing by a human being to be appropriate, but it can help overcoming the most challenging part and it provides support and inspiration to practitioners in finding suitable answers. A rough estimate of the time needed to counter a hate message demonstrates that **operators using the Hatemeter Platform significantly reduce their response time by half (from 8 minutes to around 4 for each counter-message).** This is a key indicator of the effectiveness of the Platform and its impact upon the operators' daily activities.

Evaluators of the Platform (see deliverable D16 – Pilots Execution, Validation and Evaluation Report v.2) have suggested that not only NGO operators, but also **other practitioners** involved in anti-discrimination related activities could utilise it to raise awareness on and to tackle Islamophobia. Moreover, the tool allows exploration of the phenomenon from a new perspective and, for this reason, it could also be used on **other channels**. A further expansion of the Hatemeter Platform could incorporate **YouTube or any other social media** where users write public posts, by monitoring Islamophobic hashtags and keywords to retrieve videos and other forms of content that trigger hateful comments. This is focused mainly on identifying specific channels that spread online Islamophobic messages, and that operators may be interested to follow.

This tool can also be extremely appropriate in responding to **other forms of hate-speech**. The Hatemeter methodology and Platform have been tested on Islamophobia, but they can be enlarged to the study and prevention of antisemitism, various forms of racism, homophobia, and others. Other phenomena can be investigated employing the same functionalities of the Platform and simply modifying the new keywords and counter-narratives options. Alternatively, the Platform can be improved with **more functionalities.** The addition of specific analytics related to the users network structure (e.g. centrality, in-degree and out-degree of the nodes) can help in monitoring specific hostile accounts, obviously in compliance with the EU-GDPR provisions.

# References

Aguilera-Carnerero, C. and Azeez, A.H. (2016). 'Islamonausea, not Islamophobia': The many faces of cyber hate speech. Journal of Arab & Muslim Media Research, 9(1), pp. 21-40.

All Party Parliamentary Group on British Muslims (2018) 'Report on the inquiry into a working definition of Islamophobia / anti-Muslim hatred', https://static1.squarespace.com/static/599c3d2febbd1a90cffdd8a9/t/5bfd1ea3352f531a6170ceee/1543315109493/Islamophobia+Defined.pdf

Alietti, A. and Padovan, D. (2018), 'Islamophobia in Italy: National Report 2017', in: E Bayrakli, F Hafez (eds.), European Islamophobia Report, Istanbul: SETA.

Amnesty International Italia (2019), "Barometro dell'odio. Elezioni europee 2019", https://d21zrvtkxtd6ae.cloudfront.net/public/uploads/2019/05/29202706/Amnesty-barometro-odio-2019.pdf

Bayrakli, E. and Hafez, F. (2019a) European Islamophobia Report 2018, Istanbul, SETA.

Bayrakli, E. and Hafez, F. (2019b) The State of Islamophobia in Europe, in Bayraklı E. and Hafez E., European Islamophobia Report 2018, Istanbul, SETA, 2019, pp. 7-58.

Bortone, R. and Cerquozzi, F. (2017), 'L'hate speech al tempo di Internet', *Aggiornamenti Sociali*, dicembre, pp. 818-27.

CCIF (2019), 'Our Manifesto', http://www.islamophobie.net/en/manifesto/

CNCDH (2018) Rapport 2018 sur la lutte contre le racisme, l'antisémitisme et la xénophobie. https://www.cncdh.fr/fr/publications/rapport-2018-sur-la-lutte-contre-le-racisme-lantisemitisme-et-la-xenophobie

Giacalone, C. (2017), 'Islamophobia in Italy: National Report 2016', in: Bayrakli, Enes and Farid Hafez (eds), European Islamophobia Report 2016, Istanbul, SETA: 297–319.

Giovannetti, M. and Minicucci, M. (2015), 'L'hate speech nei new social media: percezioni, esperienze, approcci, reazioni e risposte dei giovani utilizzatori e dei professionisti', Relazione al convegno Hate speech e libertà di espressione.

Gouvernement.fr (2019), 'Bilan 2018 des actes racistes, antisémites, antimusulmans et antichrétiens', https://www.gouvernement.fr/bilan-2018-des-actes-racistes-antisemites-antimusulmans-et-antichretiens

Ekman, M. (2015). 'Online Islamophobia and the politics of fear: manufacturing the green scare'. Ethnic and Racial Studies, 38(11), pp. 1986-2002.

Elgot, J. (2019). 'Government criticised for rejecting definition of Islamophobia', The Guardian, https://www.theguardian.com/news/2019/may/15/uk-ministers-criticised-rejecting-new-definition-Islamophobia

Faytre, L. (2019). 'Islamophobia in France: National Report 2018', in Bayraklı E. and Hafez E., European Islamophobia Report 2018, Istanbul, SETA, 2019, pp. 319-368.

Faloppa, F. (2019), "La mappa delle parole", in Amnesty International Italia (2019), "Barometro dell'odio. Elezioni europee 2019", https://d21zrvtkxtd6ae.cloudfront.net/public/uploads/2019/05/29202706/Amnesty-barometro-odio-2019.pdf

Home Office (2019), 'Hate crime, England and Wales, 2018 to 2019', Home Office Statistical Bulletin 24/19, https://www.gov.uk/government/statistics/hate-crime-england-and-wales-2018-to-2019

Home Office (2018), 'Hate Crime, England and Wales, 2017/18', Statistical Bulletin 20/18, https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/748598/hate-crime-1718-hosb2018.pdf

Horsti, K. (2017). 'Digital Islamophobia: The Swedish woman as a figure of pure and dangerous whiteness', New Media and Society, 19(9). Pp. 1440-1457.

Islamophobia (2019), 'Islamophobia', https://Islamophobia-definition.com/

Larsson, G. (2007). 'Cyber-Islamophobia? The case of WikiIslam'. Contemporary Islam, 1(1), pp. 53-67.

Malchiodi, M. (2016), L'islam nei social media, Pavia: Osservatorio di Pavia.

Moretti, G., Sprugnoli R. and Tonelli S. (2014). Digging in the Dirt: Extracting Keyphrases from Texts with KD. In Proceedings of the Second Italian Conference on Computational Linguistics, Trento, Italy.

Oboler, A. (2016). 'The normalisation of Islamophobia through social media: Facebook'. In: Awan, I. (Ed.), Islamophobia in Cyberspace: Hate Crimes Go Viral. Routledge, New York, pp. 41-62.

Osservatorio sulle Discriminazioni (2010), Rapporto 2010, Mantova: Osservatorio sulle Discriminazioni.

Spärck Jones, K. (1972). "A Statistical Interpretation of Term Specificity and Its Application in Retrieval". Journal of Documentation. 28: 11–21.

Vitullo, A. (2019), "Odio e religione: musulmani nel mirino", in Amnesty International Italia (2019), "Barometro dell'odio. Elezioni europee 2019", https://d21zrvtkxtd6ae.cloudfront.net/public/uploads/2019/05/29202706/Amnesty-barometro-odio-2019.pdf

Vox, Osservatorio Italiano sui Diritti (2019), La mappa dell'intolleranza, http://www.voxdiritti.it/la-nuova-mappa-dellintolleranza-4/.

November 2019

HATE