# Technologie en P/CVE

—

Olivier Cauberghs

**textgain**

—

# EOOH dashboard tech roadmap

**Olivier Cauberghs**
**www.eooh.eu**
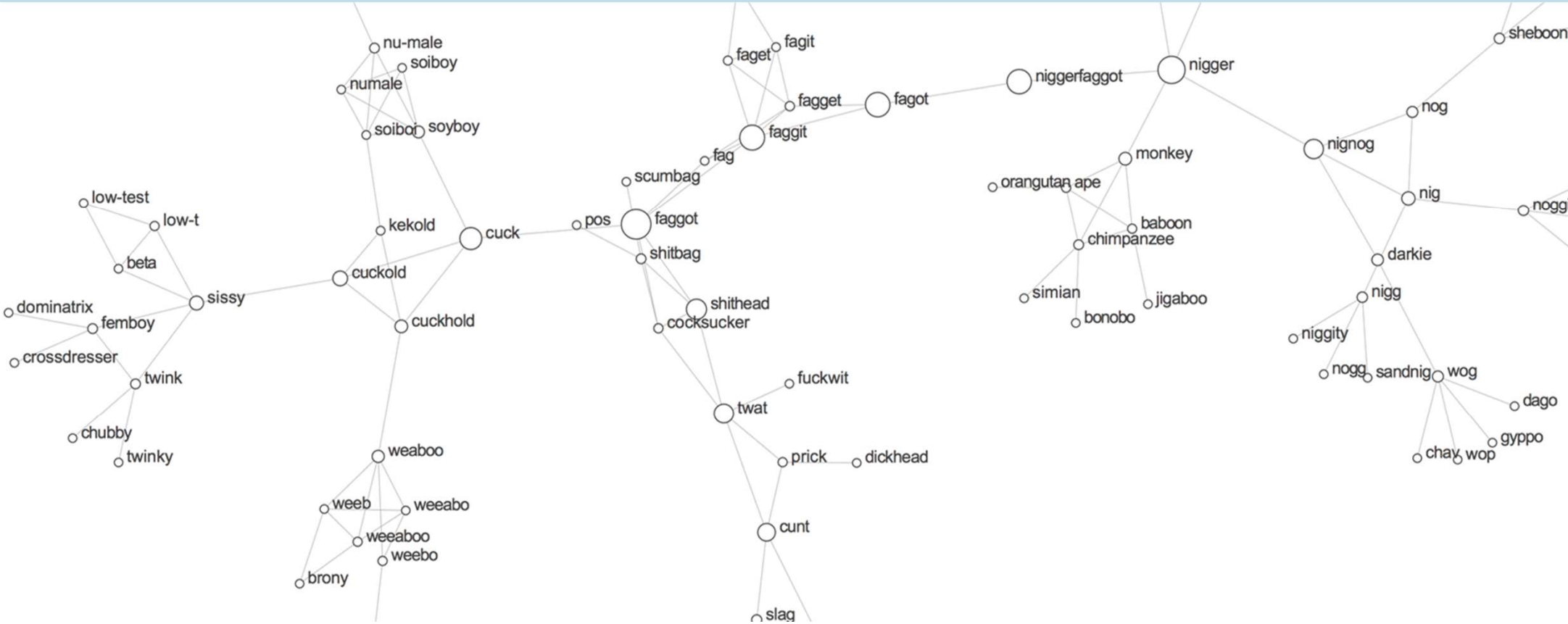**DG JUST** • REC-RRAC-AG-2020 • PANORAMA 963801

# TRANSPARENT AI

- Technology that detects problematic content can help moderators

- Technology that learns by itself (AI) can also be prejudiced however

- That's why Textgain will develop new transparent AI for EOOH
  - compliant with the EU's GDPR privacy regulations
  - compliant with the EU's AI & Ethics recommendations

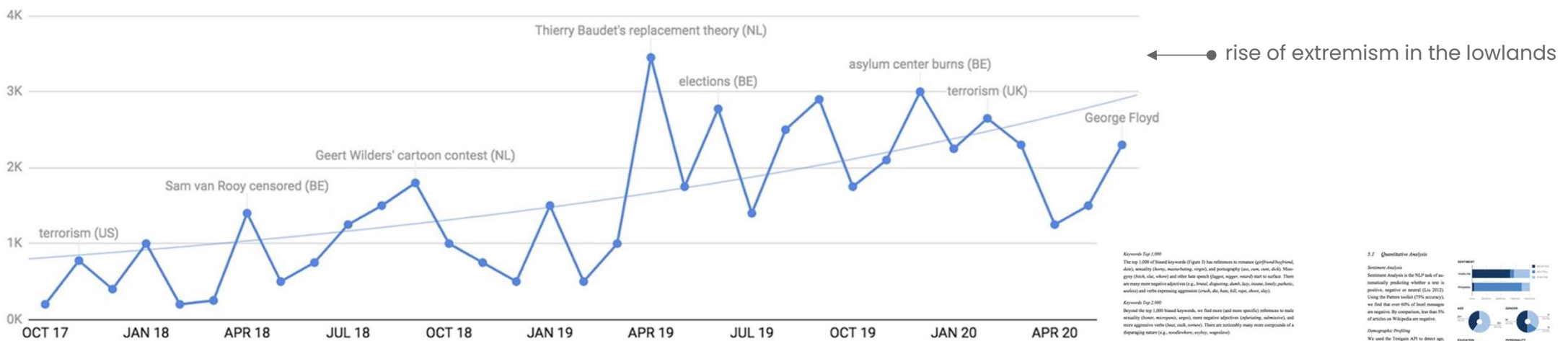- For all **24 languages** in the EU + Arabic, Turkish and Russian

# ONLINE TOOL CONNECTS TO PUBLIC APIs (25)

Telegram too

# CONTINUOUSLY LEARNS NEW UNKNOWN EXPRESSIONS

# ONLINE TOOL GENERATES TREND REPORTS



rise of extremism in the lowlands

QAnon: researchgate.net/publication/The-QAnon-superconspiracy
Islamic State: arxiv.org/pdf/1803.04596.pdf
German far-right: organisms.be/downloads/jaki2018.pdf
Incels: organisms.be/downloads/incels.pdf

# DASHBOARD
# European Observatory of Online Hate

# Dashboard



SK-UA O     PFIZER | MESSAGE

- SK-UA on Facebook
- Messages to Black actress Moses Ingram in Obi-Wan Kenobi Star Wars
- pfizer documents

15/9   16/9   17/9   18/9   19/9   20/9

300
201
102
3

## Navigation (sidebar)

- Dashboard
- Community
- Manage
- Monitor

Collapse |<

- Settings
- Support
- guydepauw

## Announcements

| | New | 03 August 15:32 |
The visibility of toxic words has been improved.

| | New | 26 July 16:02 |
A bug affecting the 'Not interested' button has been fixed.

| | New | 13 June 14:30 |
New platforms added: TikTok, Minds, Steam and Google News

| | New | 12 May 14:30 |
Addition of filters for platforms, keywords and languages.

| | New | 28 January 15:30 |
"path", "to" & "by" now accept multiple entries.

| | New | 28 January 15:30 |
"path" display bug is now fixed.

| | New | 28 January 12:47 |
YouTube was added as platform.

v1

# Dashboard
# Community
# Manage
# Monitor

Collapse |<

Settings
Support
guydepauw

Paula

0 followers   DUPLICATE   FOLLOW

GO TO CHANNEL

**Platforms**

Twitter

**Status**

Urgent

**Languages and keywords**

> EN

| | category | # | score | keywords |
|---|---|---|---|---|
| 🤡 | RIDICULE | ●●○ | 0.25 | baby, stupid, degeneracy |
| 👲 | CONTEMPT | ●●○ | 0.35 | disgusting, illegal, nazi |
| 👳 | RACISM | ●●○ | 0.40 | racist, jews, replacement theory |
| 💁 | SEXISM | ●●○ | 0.20 | gay, hooters, ho |
| ✊ | POLITICS | ●●○ | 0.30 | nazis, woke, illegal |
| ☝ | RELIGION | ●●○ | 0.35 | jews, replacement theory, hell |
| 💣 | THREAT | ●●○ | 0.20 | murder, killing, rape |

# Coming up...

- User manual
- Speed boosts through better hardware
- Secure email pipelines for automated email alerts and weekly digests
- Trending words and phrases

Judaism

- *hassidic zionist*
- *join telegram*
- *conquer nwo*
- *exposed reptilian*
- *gassed millions*
- *politics demonic*
- *starved camps*
- *systematic effort*
- *camps typhus*
- *camps de-loused*

Islam

- *shirk*
- *insult lgbtq+*
- *global agenda*
- *global fascism*
- *jewish massacres*
- *abused girls*
- *men jailed*

Refugees

- *grateful ukr*
- *good kills*
- *independence day*
- *battalion delivering*
- *ukie escape*
- *national army*
- *chechens trophies*
- *malkachan love*
- *globohomo art*

Ukraine

- *#crimea*
- *taiwan*
- *western globalist*
- *sovereign development*
- *zaporizhzhia nuclear*
- *#visabanforrussians*
- *nuclear plant*
- *visa ban*

**Coming up...**

- Theme taxonomy
  - 118 experts
  - 727 channels (532 active)
  - Develop a living theme taxonomy based on the interests of the experts
    - Hands-on pan-European trend analysis

- Local upload of data for analysis

## Dataset description

- The dataset has over 550K messages (*n*=558,918).
- These were posted between July 01, 2022 and August 31, 2022.
- Most data was collected from Twitter (80%), Facebook and 4plebs.
- Most messages are written in English (EN, 65%), Dutch (NL, 10%) and French (FR, 5%).
- The **average toxicity score** is **0.35**, which is quite high.
- About 50K messages (10%) are very toxic, with a score > 0.8.
- Some of these are written in English (EN, 55%).
- For example: *">>391799441 >Chinks vs kikes Chinks are nothing but kike slaves. Globohomo wouldn't be able to exist without chink slave labor. If you want to defeat the Jews, then china must be destroyed."* (source: 4plebs).
- Messages that are **very toxic** (> 0.8) are more frequent on Twitter (35%).
- Toxicity often involves RACISM (10%), POLITICS, RELIGION and/or CONTEMPT.
- About 5% of messages involve THREATENING.
- About 2% of messages involve DISINFORMATION.
- The most frequent combination is RACISM + RELIGION.
- Some common toxic keywords include: *holohoax*, *holokauszt*, *nagyon*, *tant*, *scum*, *jews control*, *jews hate*, *latvijā*, *ježiš*, and *shitskin*.

# JOIN US AND BECOME A EOOH PARTNER!

Q&A

Textgain
Olivier Cauberghs
Olivier@Textgain.com