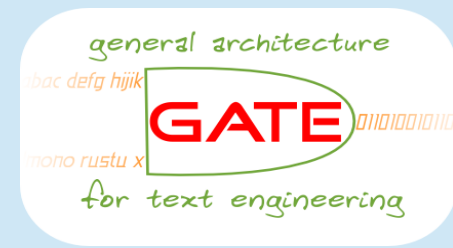# Adding value to NLP:
## More semantic technologies in use

# 15 years ago in NLP, we were…

- Using ontologies for adding reasoning/additional information to text

- Developing linguistic tools for conceptual modelling

- Semantic enrichment for ontology mapping

- Developing ways to get from text to ontologies and back

- Trying to persuade the Semantic Web community that NLP was important

# Mapping semantic relations

Ford, Henry Ford -- (United States manufacturer of automobiles who pioneered mass production (1863-1947))

car, auto, automobile, machine, motorcar -- (4-wheeled motor vehicle; usually propelled by an internal combustion engine; "he needs a car to get to work")

*A hyponym of car (from a total of 30):*

   => Model T -- (the first widely available automobile powered by a gasoline engine; mass-produced by Henry Ford from 1908 to 1927)

- Finding "hidden" semantic links between concepts using gloss information from dictionaries
- (We want to know that Henry Ford makes cars)
- These days we could just check Wikipedia ☺

# Many of the underlying technologies are still relevant

- The techniques may have changed a bit: we now have
    - lots of lovely big data to play with
    - open source knowledge bases
    - powerful GPUs (crack cocaine for NLP students)
    - deep learning (all the cool kidz ☺ )

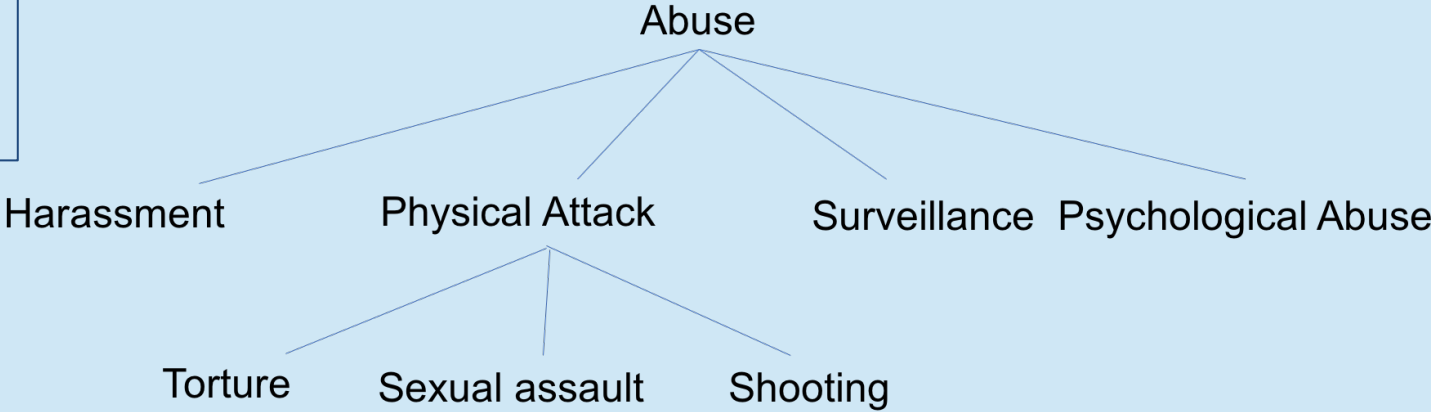# How to impress a "freedom of the media" researcher



Journalist and radio host Luka Popov from northern Serbia was found dead in his home in Srpski Krstur on Friday 17 June 2016. According to Serbian daily Blic, Popov's body was found with visible injuries and he had been presumably "tortured and murdered". The police are currently investigating the circumstances behind the death. The Serbian Journalist associations NUNS and UNS, and journalist association of Vojvodina DNV have urged the authorities to thoroughly investigate the circumstances around the journalist's death. The OSCE has condemned the murder. The OSCE's representative on freedom of the media, Dunja Mijatovi_, said that the investigation must be carried ... authorities must do their utmost ... murder to justice." Luka Popov w... of Coka and Novi Knezevac.

Incidents
- ✔ Incident Date
- ✔ Job Title
- ✔ Location
- ✔ Media
- ✔ Person
- ✔ Violence

Violence›

| | deliberate | ▼ | yes | ▼ | ✕ |
| C | kind | ▼ | death | ▼ | ✕ |

"What is this wizardry?"

Abuse

Harassment    Physical Attack    Surveillance    Psychological Abuse

Torture    Sexual assault    Shooting

# Feral databases

The problem: too much unstructured information

On 26 June 1996, while driving her red Opel Calibra, Guerin stopped at a red traffic light on the Naas Dual Carriageway  near Newlands Cross, on the outskirts of Dublin, unaware she was being followed. She was shot six times, by one of two men sitting on a motorcycle.

How many journalists died in 1996?

The solution: categorise it and stick it in a spreadsheet

"Everybody loves spreadsheets!"

# That's great, we have lots of NLP tools to do that!

Date: 26-06-1996    Location: other    Person    Location: city

On 26 June 1996, while driving her red Opel Calibra, Guerin stopped at a red traffic light on the Naas Dual Carriageway near Newlands Cross, on the outskirts of Dublin, unaware she was being followed. She was shot six times, fatally, by one of two men sitting on a motorcycle.

Event: Killing: shooting

| Name | Date | Location City | Location Country | Event type | Type of death |
|------|------|---------------|------------------|------------|---------------|
| Guerin | 26061996 | Dublin | Ireland | Killing | Shooting |

# Woe betide anyone who opens the sacred spreadsheet!

↓

# Mass panic

↓

# Let's start a new spreadsheet!

↓



Aaaaggghhhh!!!

# Harry Potter and the Spreadsheet of Doom

# This was once the latest in calendar technology…

# Case studies: how can NLP and semantics help?

- Violations against journalists
- Disaster relief
- The B-word

# Global monitoring of violations against journalists

- The situation: *cases of killing, kidnapping, enforced disappearance, arbitrary detention and torture of journalists, associated media personnel, often with impunity*

- The task:  UN SDG agenda (indicator 16.10.1): we need a comprehensive monitoring system capturing the scope & nature of violations (lethal & non-lethal)

- The problem: the data is a horrible mess!
  - accessibility of reliable information on violations
    - gaps in the data beyond killings; cannot compile data on a wide range of violations; practical challenges of collection
  - systematisation and compiling of information
    - conceptual inconsistency; much information is unstructured
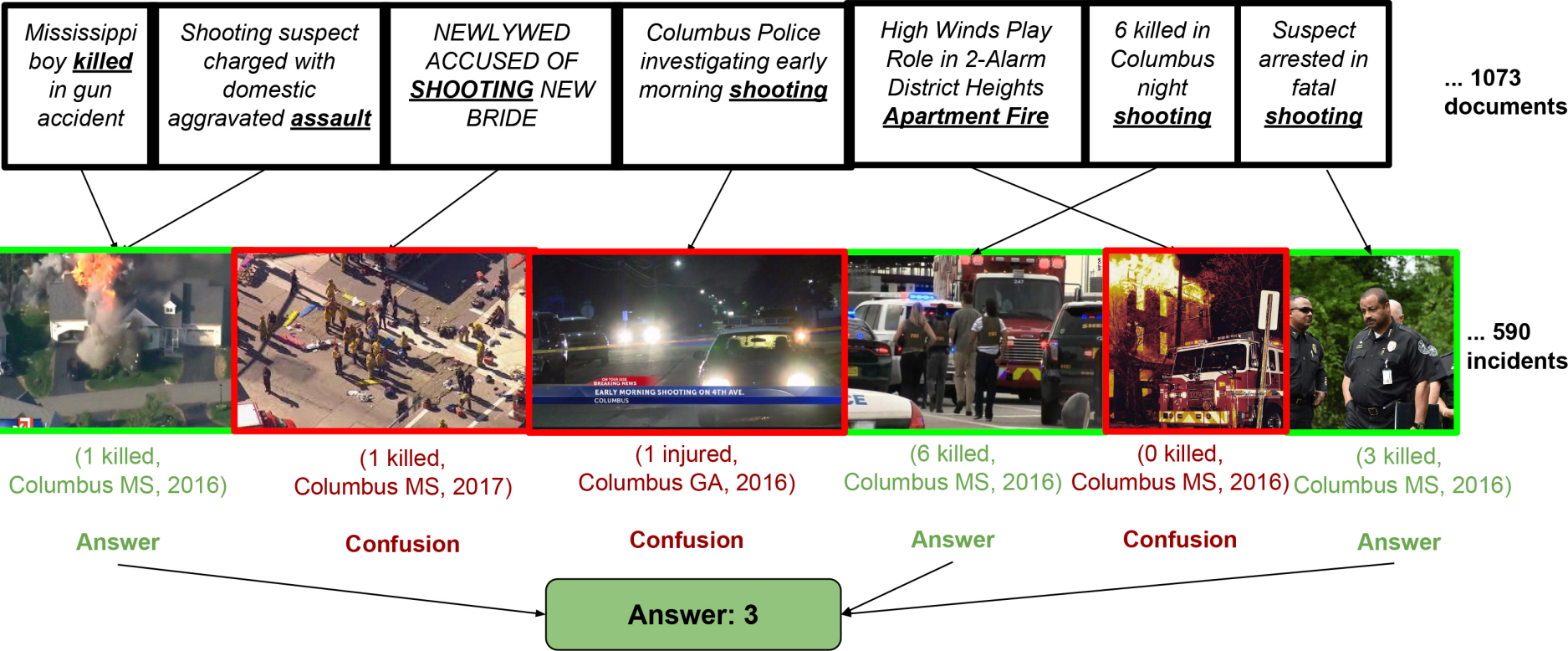
# Events are complex

*The family of murdered Maltese anti-corruption journalist Daphne Caruana Galizia is demanding an independent public inquiry because she had suffered years of intimidation.*

*She was killed by a car bomb near her home in October. Her widely-read blog accused top politicians of corruption.*

*One of her sons, Paul, said three pet dogs were killed and attempts were made to burn down the journalist's home.*

- Several separate events, but
  - are they related?
  - was it the same perpetrator(s)?
- Traditional methods use a person-centric approach with no relation between events

# But we know how to understand complex events

Question: How many killing incidents happened in 2016 in Columbus, Mississippi?

| *Mississippi boy **killed** in gun accident* | *Shooting suspect charged with domestic aggravated **assault*** | *NEWLYWED ACCUSED OF **SHOOTING** NEW BRIDE* | *Columbus Police investigating early morning **shooting*** | *High Winds Play Role in 2-Alarm District Heights **Apartment Fire*** | *6 killed in Columbus night **shooting*** | *Suspect arrested in fatal **shooting*** |

**... 1073 documents**

**... 590 incidents**

(1 killed, Columbus MS, 2016)   (1 killed, Columbus MS, 2017)   (1 injured, Columbus GA, 2016)   (6 killed, Columbus MS, 2016)   (0 killed, Columbus MS, 2016)   (3 killed, Columbus MS, 2016)

**Answer**   **Confusion**   **Confusion**   **Answer**   **Confusion**   **Answer**

**Answer: 3**

SemEval 2018 Task on Counting Events and Participants in the Long Tail

# Information categories in current monitoring efforts

| Assimilated categories | Press Freedom Tracker | Council of Europe | Mapping Media Freedom | CPJ | Media Governance and Industries Lab |
|---|---|---|---|---|---|
| Death | Physical attack | Death | Death | Death | Killed |
| | | | | | |
| Missing | | | | Missing | Missing/ Disappeared |
| | | | | | |
| Arrest/ Imprisonment | Arrests and criminal charges | Detention or Imprisonment | Arrest/Detention | Imprisoned | Imprisoned |
| | | | | | Arrested |
| | | | | | |
| Interrogation | | | Interrogation | | |
| | | | | | |
| Torture | | | | Tortured | |
| | | | | | |
| Assault/Attack | Physical Attack | Attacks on physical safety and integrity | Physical assaults | | Physical assault |
| | Prior Restraint | | | | |

# ICCS Crime Classification Scheme

## Section 02 Acts causing harm or intending to cause harm to the person

| LEVEL 02 | LEVEL 03 | LEVEL 04 | CRIME |
|---|---|---|---|
| 0201 | | | **Assaults and threats** |
| | 02011 | | Assault |
| | | 020111 | Serious assault |
| | | 020112 | Minor assault |
| | 02012 | | Threat |
| | | 020121 | Serious threat |
| | | 020122 | Minor threat |
| | 02019 | | Other assaults or threats |
| 0202 | | | **Acts against liberty** |
| | 02022 | | Deprivation of liberty |
| | | 020221 | Kidnapping |
| | | 020222 | Illegal restraint |
| | | 020223 | Hijacking |
| | | 020229 | Other deprivation of liberty |

# HURIDOCS classification scheme

| | |
|---|---|
| **03** | **Violations against the right to liberty** |
| 03 01 | Direct actions which violate the right to liberty |
| 03 01 01 | Arrest |
| 03 01 02 | Detention, imprisonment |
| 03 01 02 01 | Denial of release in case of unlawful arrest |
| 03 01 02 02 | Denial of release pending trial |
| 03 01 02 03 | Denial of release on bail |
| 03 01 03 | Disappearance |
| 03 01 04 | Abduction; kidnapping |
| 03 01 05 | House arrest |
| 03 01 06 | Slavery |
| 03 01 07 | Mass roundup |
| 03 01 08 | Curfew |

# Putting it all together



**CPJ database**

Dangerous Assignment

Abu Hailu

Was refused medical treatment in prison and died as a result

Record 01

**HURIDOCS scheme**

0101 Direct actions which violate the right to life

010106 Death in certain locations

010103 Killings in the context of conflict

01010601 Death in detention or police custody

**ICCS scheme**

01 Acts leading to death or intending to cause death

0107 Unlawful killing associated with armed conflict

# Reconciling information from different sources

| Relation | Example | Issue | Solution | Same Event? |
|----------|---------|-------|----------|-------------|
| **Exact match** | Location: London Location: London | Identical facts | Two events can be merged | Very likely |
| **Equivalent match** | Type: murder Type: assassination | Semantically equivalent facts | Two events can be merged | Very likely |
| **No-conflict** | Location: London Date: 2019 | Different but semantically compatible facts | Facts from two events can be merged | Probably |
| **Specificity Conflict** | Location: London Location: UK | One fact is more specific than the other, but both are semantically compatible | Facts can be merged at either specific or generic level | Probably |
| **Direct conflict** | Location: London Location: Paris | Facts conflict semantically: both cannot be correct | More info needed to resolve conflict | Unlikely |

# Reconciling our information

| Name | Date | City | Country | Event type | Type of death |
|---|---|---|---|---|---|
| Guerin | 26061996 | Dublin | Ireland | Killing | Shooting |

➕

| Name | Date | Organisation | Location | Event type |
|---|---|---|---|---|
| Veronica Guerin | 26 June 1996 | Sunday Independent | Ireland | Murder |

=

| Name | Date | Organisation | Location | Event type | Sources |
|---|---|---|---|---|---|
| Veronica Guerin | 26061996 | Sunday Independent | Dublin, Ireland | Shooting: deliberate, fatal | 2 |

# Automatic categorization of free text



Entities and events are extracted, along with features, and linked to other knowledge sources, e.g. Wikipedia, other categorization schemes

# Social Media and Disaster Relief

- Hurricane Sandy in the US:
  - 1.1 million tweets in the first day; over 20 million in total
  - > 800K photos with #Sandy hashtag on Instagram
- Haze in Singapore: > 23 million
- Nepal earthquake in 2015: more than half a million posts

# Aid workers in Nepal discussing strategy

# Tools to help disaster victims get aid quickly: pinpointing geographic locations

- Many NGOs are not local to the disaster area and may not have a good grasp of the geography
- Place names are ambiguous
- Find mentions of locations in the text, match them to a knowledge base, and plot them on a map

## Location



## Annotation from YODIE

A huge fire has engulfed a tower block in West London. The fire broke out shortly before 1am on Wednesday at Grenfell Tower in Latimer Road near Notting Hill.The tower is at least 24 storeys high and contains 120 apartments. 200 firefighters have been tackling the blaze with 40 engines. A number of people have been treated for a range of injuries according to the fire brigade There have been multiple reports of people trapped in the blaze. These have not been confirmed by police or the fire brigade. Streets around the tower have been sealed off and residents in ther houses evacuated. London

# CREES Google Sheets Add-on

Adding semantic information (from BabelNet) improves cross-lingual crisis event classification

| Notes | CREES | | |
|---|---|---|---|
| Text of request | RELATED | EVENT | INFO |
|  | non-related | none | sympathy_and_support |
| Young nurse needing rescue! | related | floods | donations_and_volunteering |
|  | non-related | none | sympathy_and_support |
| Together we will rebuilt page - Celia Torres | related | floods | infrastructure_and_utilities |
| Only have bank address from  Twitter | non-related | floods | other_useful_information |
| #disabled lady is #stranded in #LeagueCity #dickinsontexas #needrescue #help #hurricaneharveyanimalrescue | related | floods | affected_individuals |
|  | non-related | none | sympathy_and_support |
|  | non-related | none | sympathy_and_support |
| Been trying numbers for hrs - busy signal. Elderly couple desperately #NeedWaterRescue:10202 Willowgrove, 77035. Near S Post Oak & W Belfort | related | floods | other_useful_information |
|  | non-related | none | sympathy_and_support |
| Can't get through on the coast guard number reported 2:19 pm CT | non-related | floods | affected_individuals |

# GATE MIMIR and Prospector

- Built on top of annotated text in GATE

- can be used to index and search over text, annotations, semantic metadata

- allow queries that arbitrarily mix full-text, structural, linguistic and semantic annotations

- open source

- you can issue ridiculously complex queries

```
{
{Person gender=male sparql = "SELECT ?inst WHERE {
   ?inst :party
<http://dbpedia.org/resource/Labour_Party_%28UK%29> .
 ?inst :almaMater
<http://dbpedia.org/resource/University_of_Edinburgh>  }"
 }
 [0..3] root:say
 [0..5] {Money currency="£" value > 2000000}
}

IN

{Document date > 20050000 topic=uk_health}
```

# Prospecting Biomedical Literature

# Choosing A Specific Instance

# What diseases are in these documents?

# What pathogens?

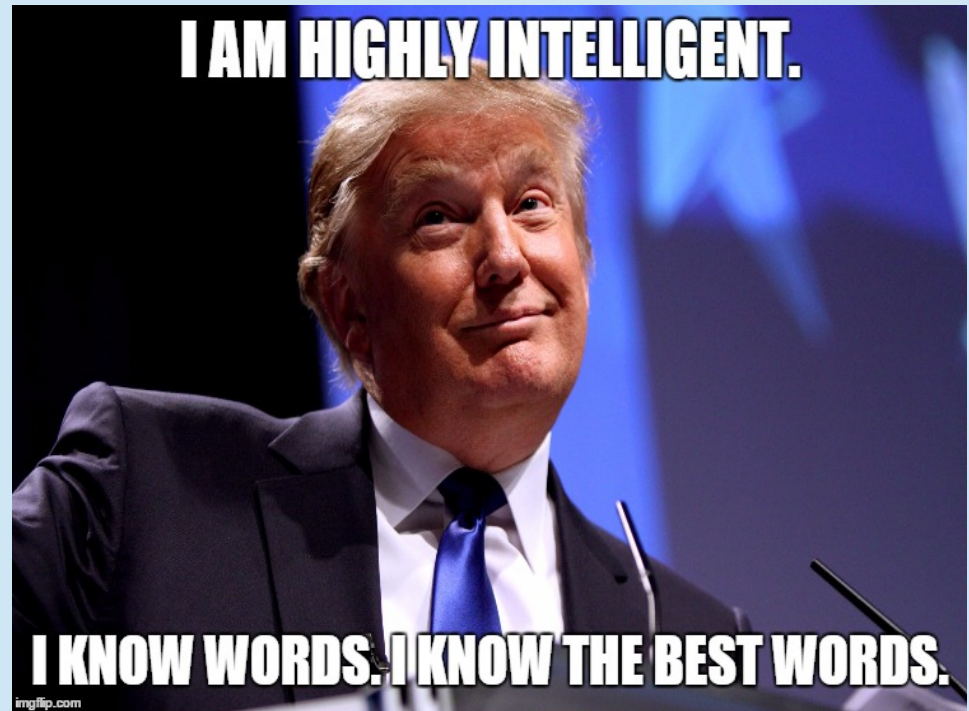# Disease vs Disease Co-ocurrences

# Diseases vs Pathogens
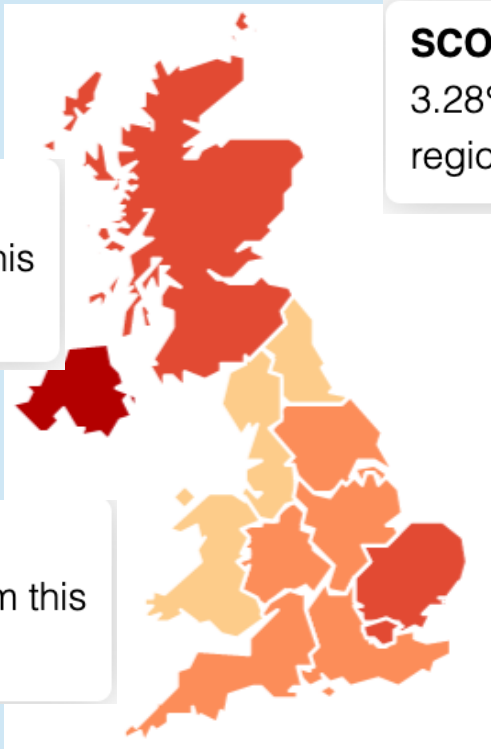
# Social media and politics

- How do people in different parts of the country talk about elections and political events?

- How do the MPs talk about different topics?

- How does the public respond to them?

- What kind of people send hate speech, who to, and what about?

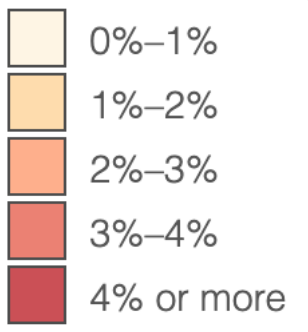# Hate speech against politicians Jan-Feb 2019

**SCOTLAND**
3.28% of the replies sent to MPs from this region were abusive

**NORTHERN IRELAND**
4.79% of the replies sent to MPs from this region were abusive

**WALES**
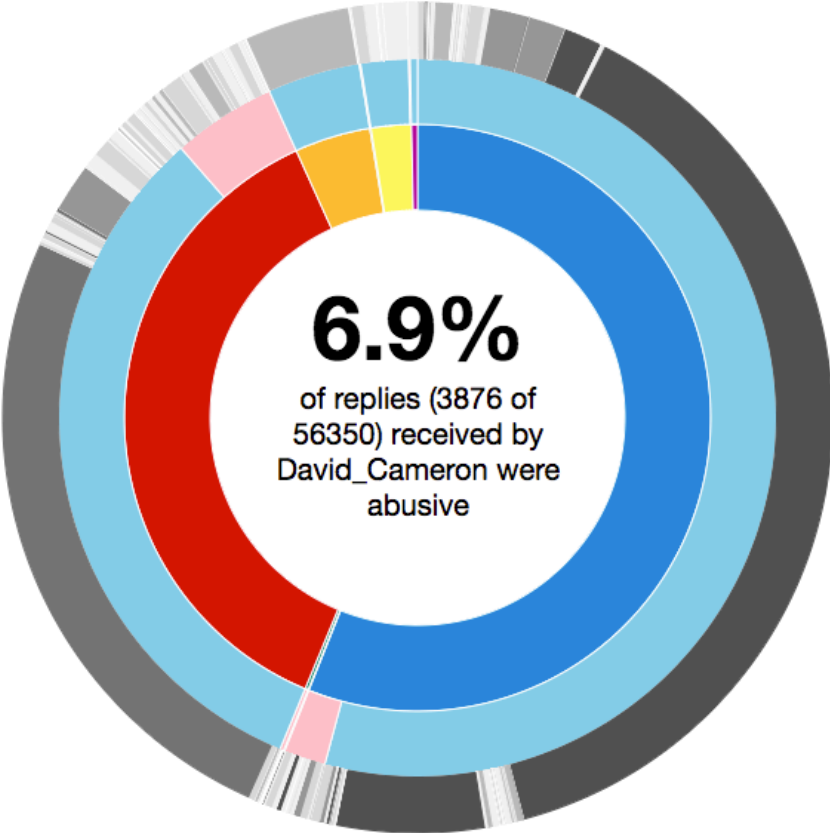1.51% of the replies sent to MPs from this region were abusive

0%–1%
1%–2%
2%–3%
3%–4%
4% or more

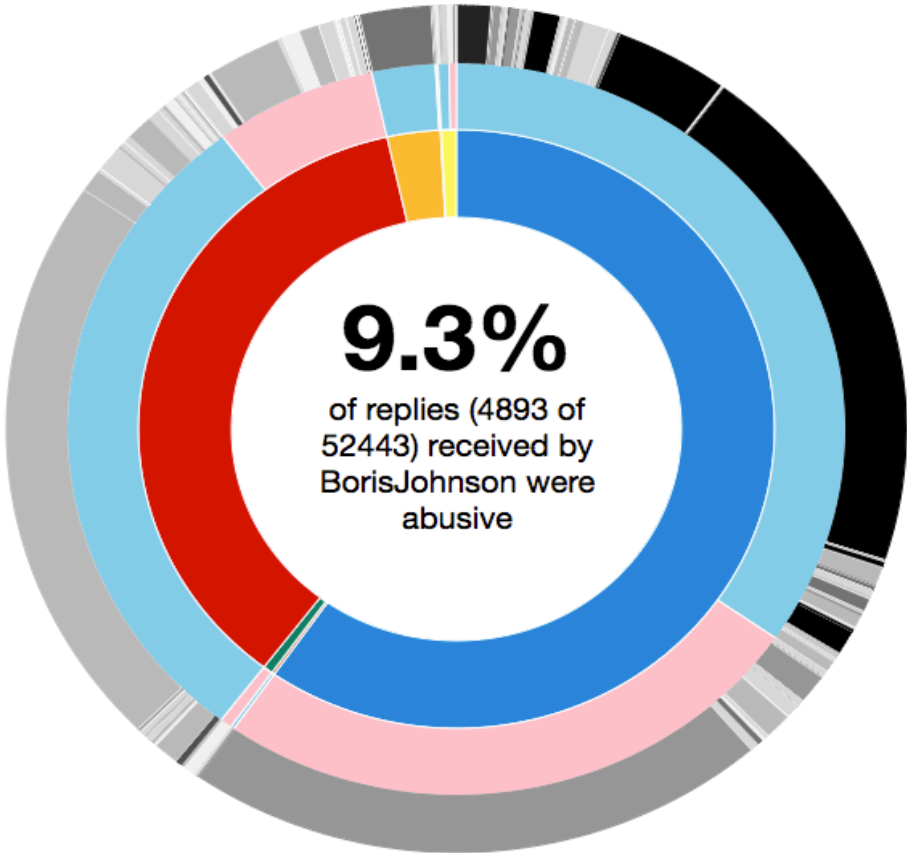Annotation of tweets with simple DBpedia and NUTS linking lets us index and query our data in interesting ways
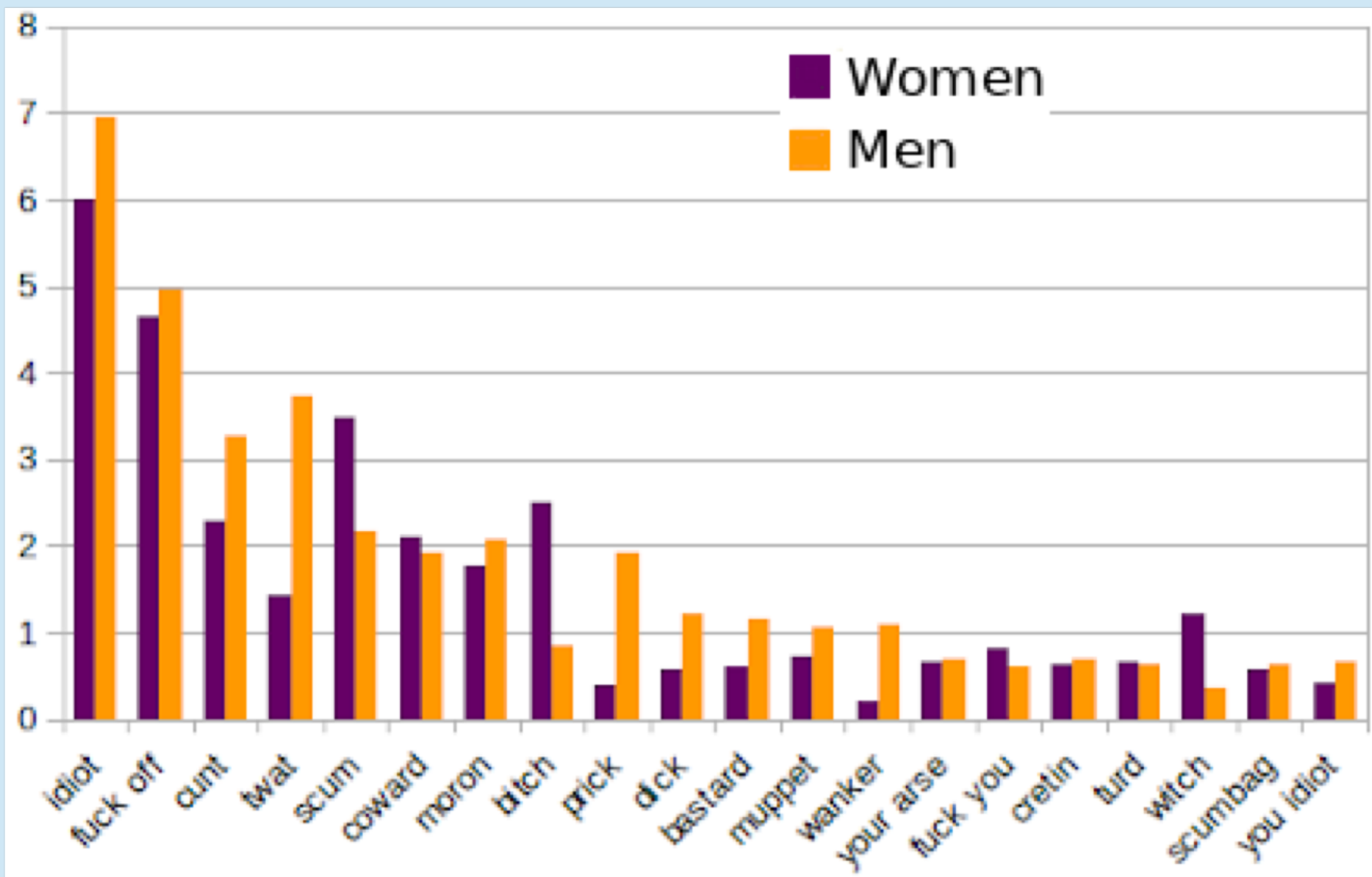
# Who gets the abuse?

## 2015



6.9%
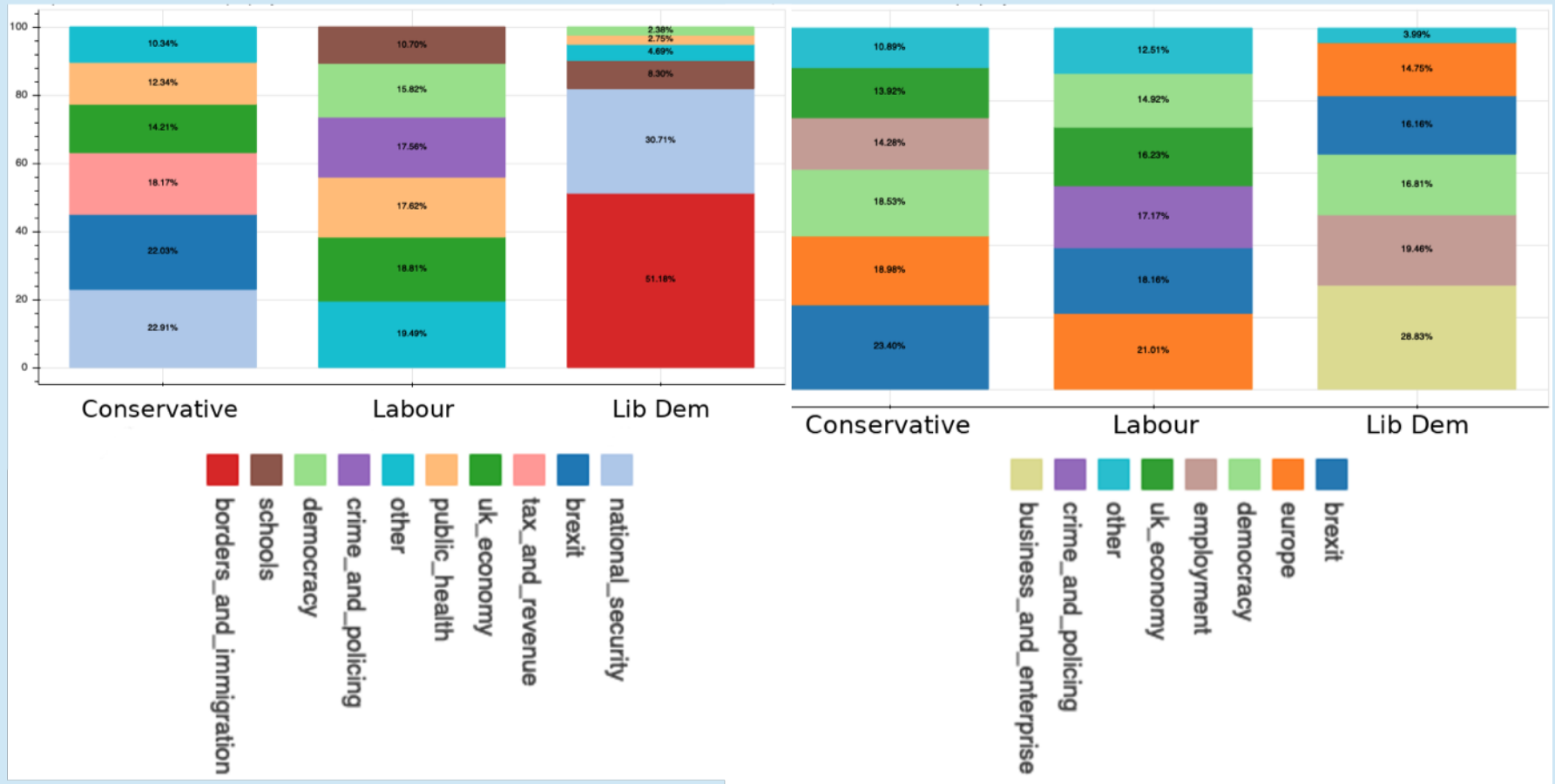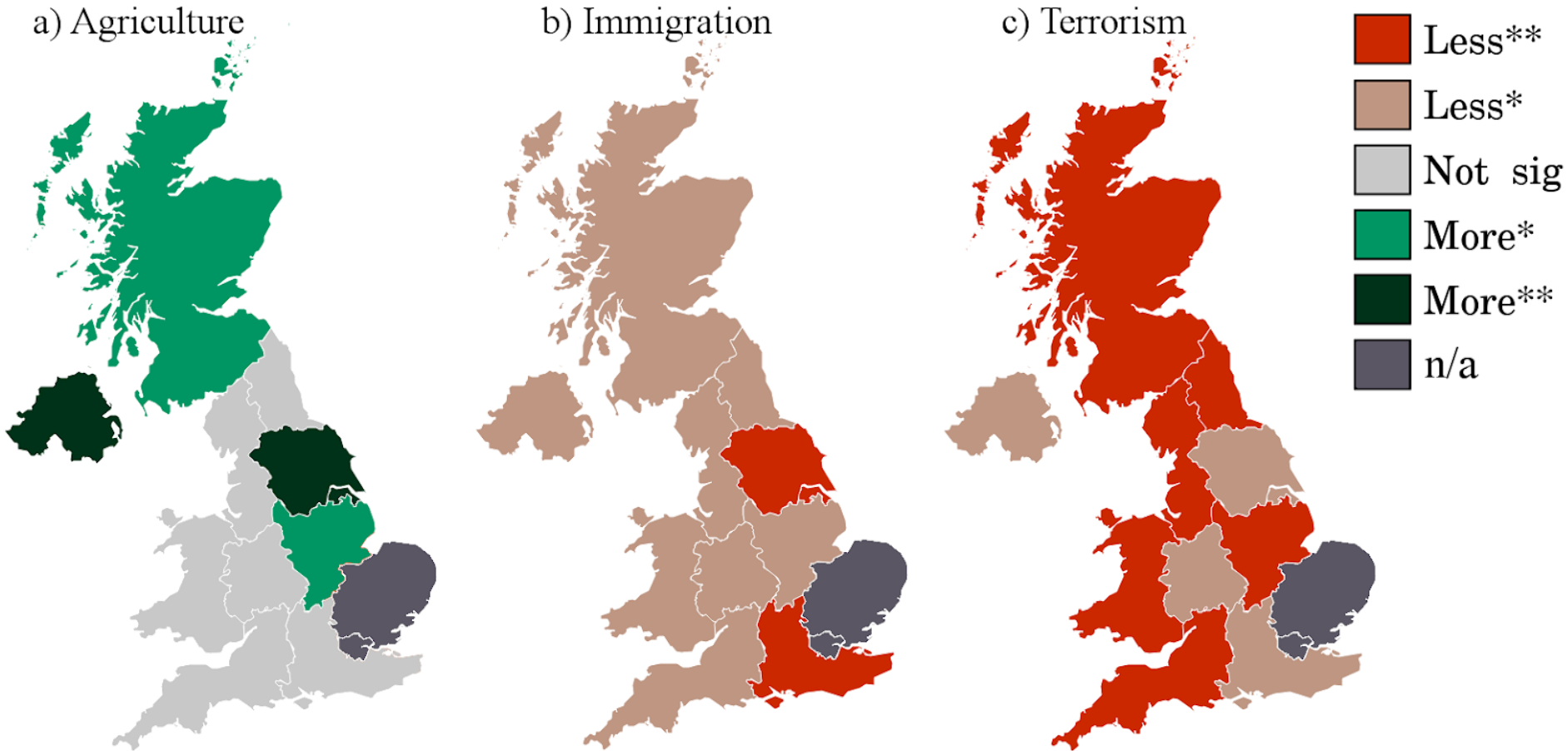of replies (3876 of 56350) received by David_Cameron were abusive

## 2017



9.3%
of replies (4893 of 52443) received by BorisJohnson were abusive

# Top abusive words

# Words vs people

# Topics of abuse per party

## 2017

## 2019

# Local and national news coverage of Brexit



a) Agriculture    b) Immigration    c) Terrorism

Legend:
- Less** (red)
- Less* (light brown)
- Not sig (light grey)
- More* (green)
- More** (dark green)
- n/a (dark grey/blue)

- Green/blue indicates areas where local coverage outweighed national coverage
- Across the regions, Twitter remainers showed closer congruence with local press than Twitter leavers

# So where are we at?

- We have lots of cool technology and data
- There is massive demand for its use in the real world
- It doesn't always need to be complex, it just needs to be applied
- Sometimes even simple techniques can have really useful outputs
- Don't neglect interdisciplinary "borrowing"
- The usual disclaimers about accuracy – real-world evaluation

# Acknowledgements

And all the GATE team in Sheffield!