



UNIVERSITÉ
DE GENÈVE



CSIC
CONSEJO SUPERIOR DE INVESTIGACIONES CIENTÍFICAS

POLITECNICO
MILANO 1863

unitar
United Nations Institute
for Training and Research



Université
de Paris



Crowd4SDG

Extracting actionable information from unconventional data sources: social media analysis with VisualCit

Carlo Bono, Politecnico di Milano - carlo.bono@polimi.it
Barbara Pernici, Politecnico di Milano - barbara.pernici@polimi.it

May 3, 2022



This project has received funding from the European Union's Horizon 2020
research and innovation programme under grant agreement No 872944



Who am I

PhD candidate @ Politecnico di Milano (Department of Electronics, Information and Bioengineering)

My (super!) supervisor is [Barbara Pernici](#)

Background: Philosophy and Computer Science

Main research interests: data-driven applications, adaptive information systems, machine learning

Currently working on EU H2020 project Crowd4SDG

For support with the tools and/or movie recommendations: carlo.bono@polimi.it



Why social media content analysis?

Social media contents are about «everything that is happening right now». So convenient when you need **information!**

But they're **huge**, and they're **noisy**. We want to retain **relevant** data only (in our dreams).

So we need **scalable** tools and algorithms that **understand** the contents of social media posts, in order to **filter** and **augment** them



Why social media content analysis?

Crowd4SDG

Most of the functionalities that we are going to discuss are aimed at **image** data (hence the «Visual» in VisualCit), but many of the concepts are general

Often one (implicit?) goal will be to **reduce** the number of elements to be manually processed, while keeping **quality** high



- <https://crowd4sdg.eu/>
- <https://pernici.faculty.polimi.it/crowd4sdgpolimi/>
- V. Negri, D. Scuratti, S. Agresti, D. Rooein, G. Scalia, J.L. Fernandez-Marquez, A. Ravi Shankar, M. Carman and B.Pernici, *Image-based Social Sensing: Combining AI and the Crowd to Mine Policy-Adherence Indicators from Twitter*, accepted at ICSE, Track Software Engineering in Society, May 2021 [link](#)
- Barbara Pernici, *CROWD4SDG: Crowdsourcing for sustainable developments goals*, 248-253, in Book of Short Papers, SIS 2020, Pearson, 2020 ([link](#))



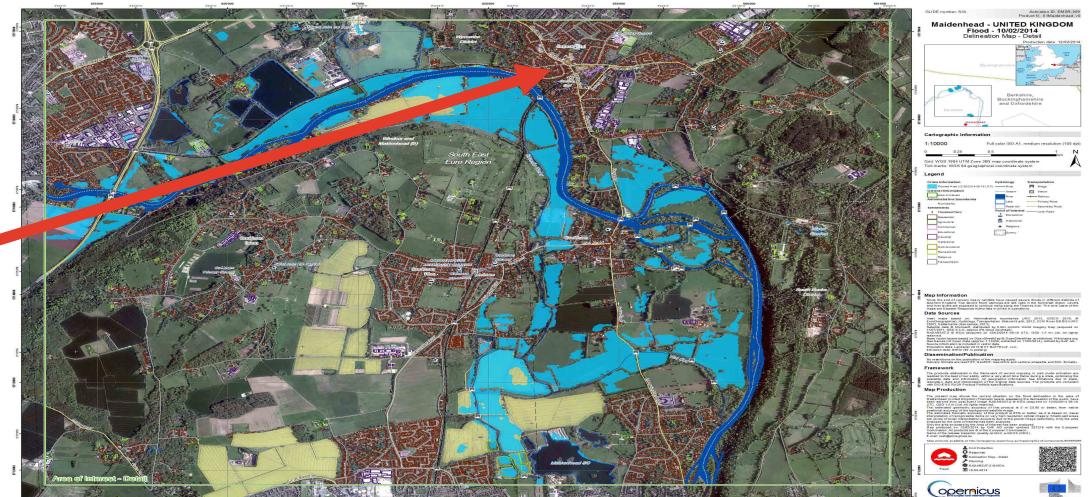
Getting evidence from tweets

In a nutshell:
ingest social media, filter/enrich data, aggregate

Flood levels up overnight in Datchet Need wellies for the pavement past Spices. High St and Queens Rd blocked off.

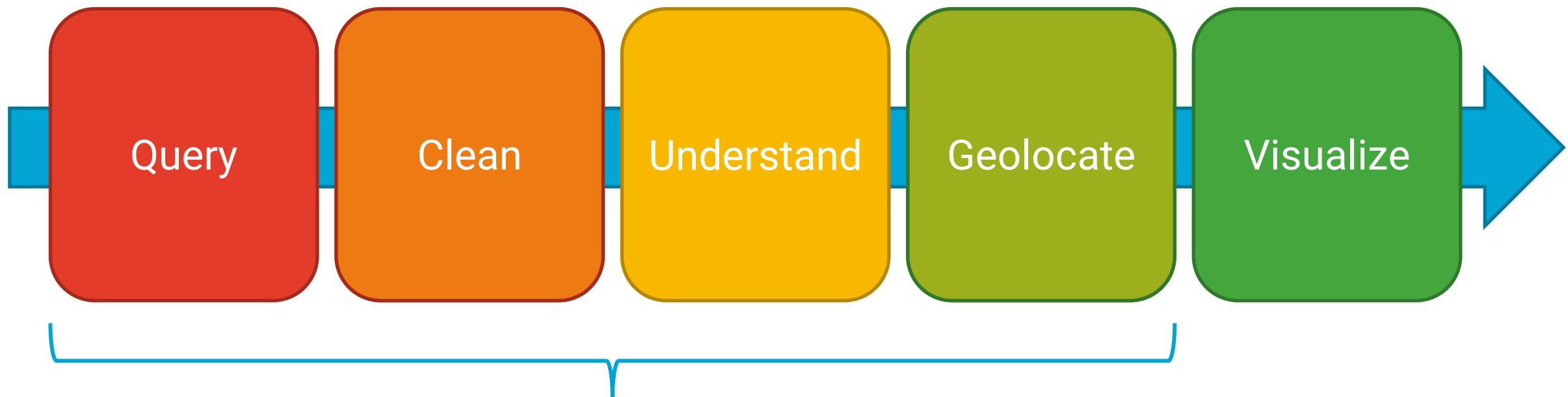


00:28 - 11 feb 2014





A data analysis pipeline in practice



In some sense they all are ways of filtering



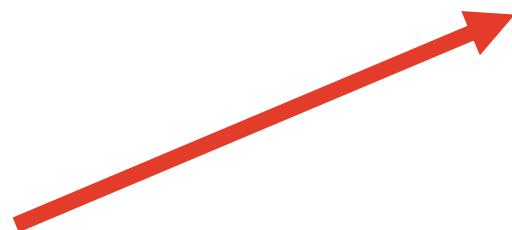
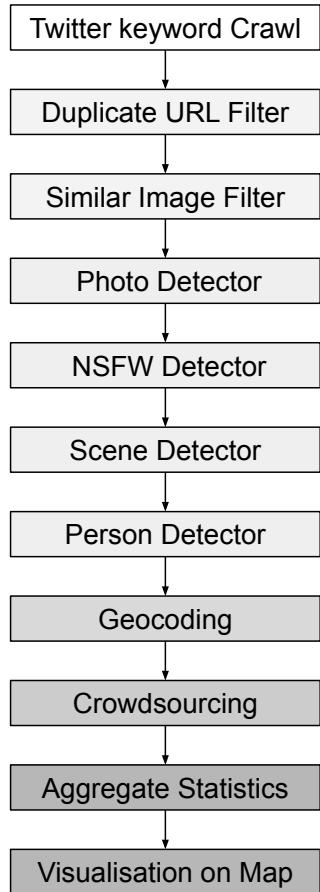
Crowd4SDG

«Helping machines do their work»



Crowd4SDG

PIPELINE



CROWDSOURCING

BETA CITIZEN SCIENCE CENTER ZURICH

HOME DISCOVER ABOUT FORUM EN LOGIN

[Go back to the project](#)

Is this a photo (rather than a cartoon, graph, meme, etc.)?

Does it look like it has been taken recently (in the last three months)?

Are there people in this image?

Are the people wearing masks?

If so, which type?

Tweet

Coronavirus en Argentina: el Gobierno amplió la definición de caso sospechoso. <https://t.co/0ef8KjyAsH>
<https://t.co/1fqPhEQaqJ>

Country/Territory: Argentina

<https://t.co/1fqPhEQaqJ>





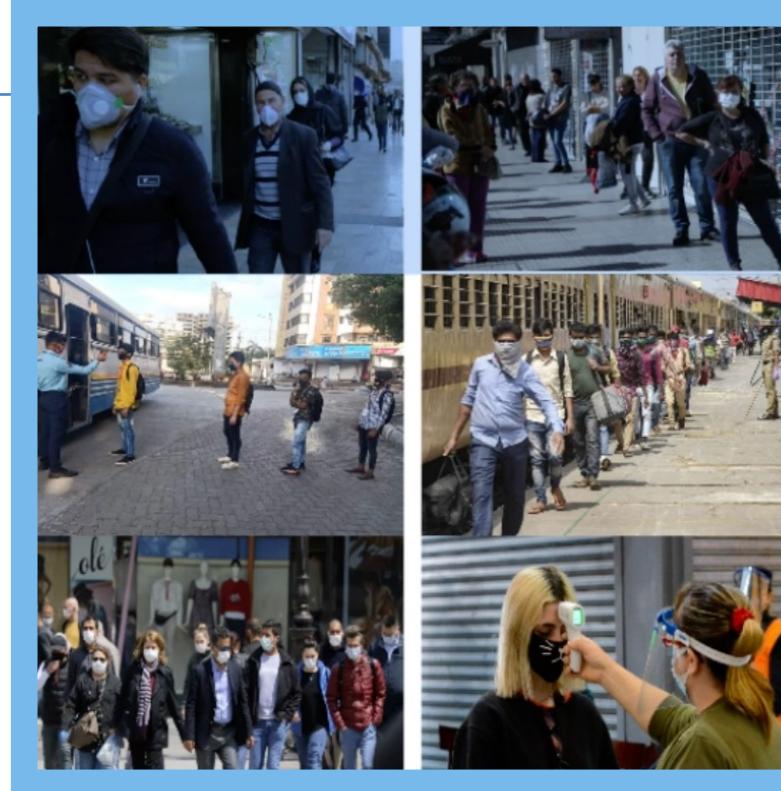
Crowd4SDG

VisualCit pipeline step by step





- Which keywords?
- Do keywords change over time?
- How much data?
- Usually done via an API



Crawling Social Media
(e.g. by keywords)



- Don't do the same work again and again!
- Different meaning for image and text



Duplicated URL
Identical / Similar Image Filter



Crowd4SDG

- Drawings
- Documents
- Memes
- Graphs
- Screenshots

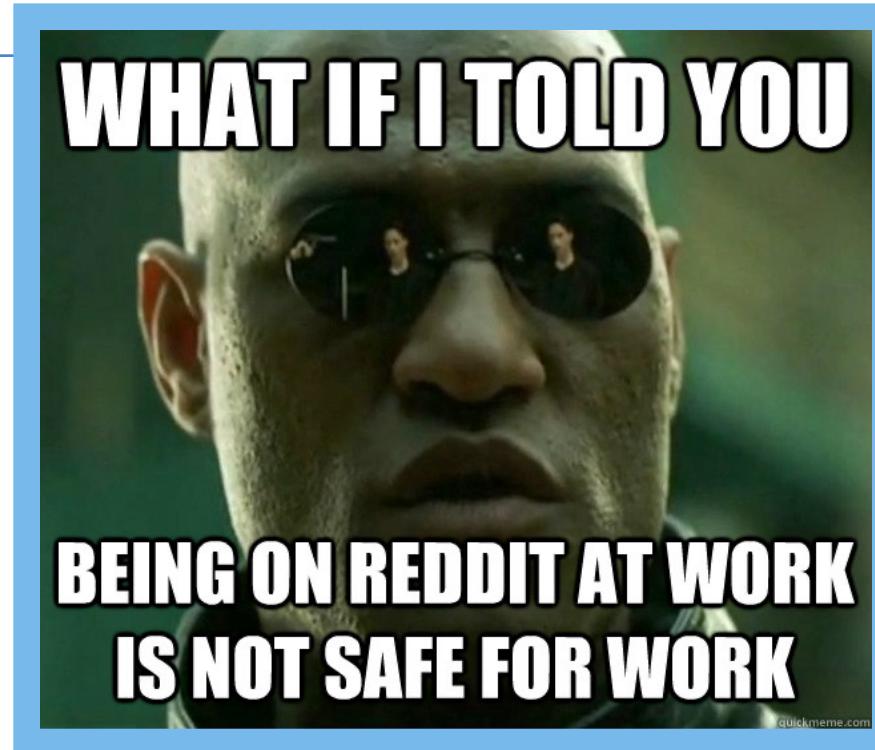


Remove “noisy” images
(e.g. non-photos)



Crowd4SDG

- Usually not of interest



Remove NSFW images



- Some off-the-shelf classifiers, like scene classifiers
- Or can be custom



Select images of interest by kind: scene classification



- Again, there are off-the-shelf object detectors
- Or it be custom, but it can take some effort



Select images of interest by content: object detection



- Yep, it's water!

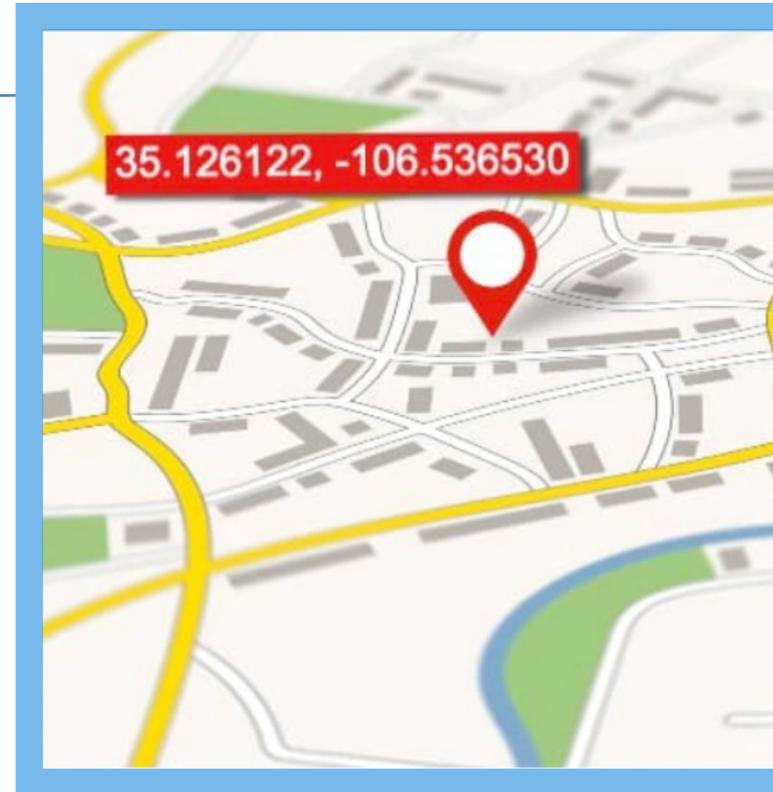
Source: floodfinder.org



Analyse image regions by
content: image segmentation



- Use some features of the posts and try to infer the location (not trivial!)



- Different applications have different constraints (e.g. granularity)

Geocoding



[Go back to the project](#)

Is this a photo (rather than a cartoon, graph, meme, etc.)?

Does it look like it has been taken recently (in the last three months)?

Are there people in this image?

Are the people wearing masks?

If so, which type?



Tweet

Coronavirus definici\xcd\x81
<https://t.co/1>

Country/ Ter

<https://t.co/>

Crowdsourcing



- Based on native or augmented information

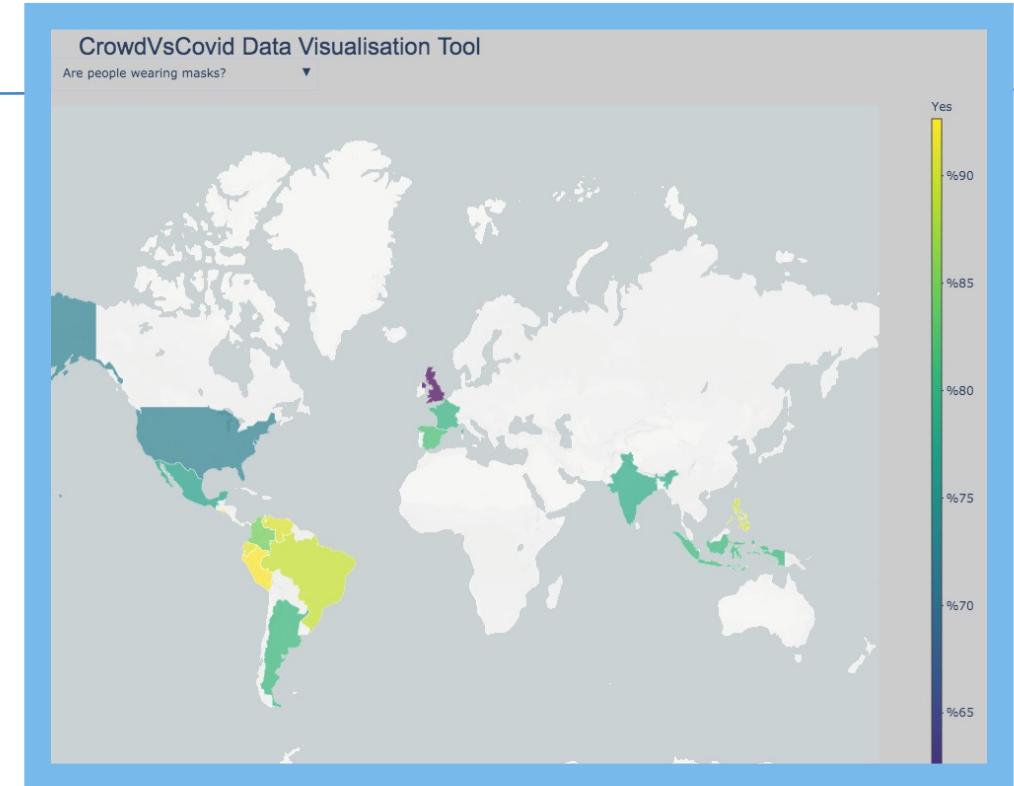
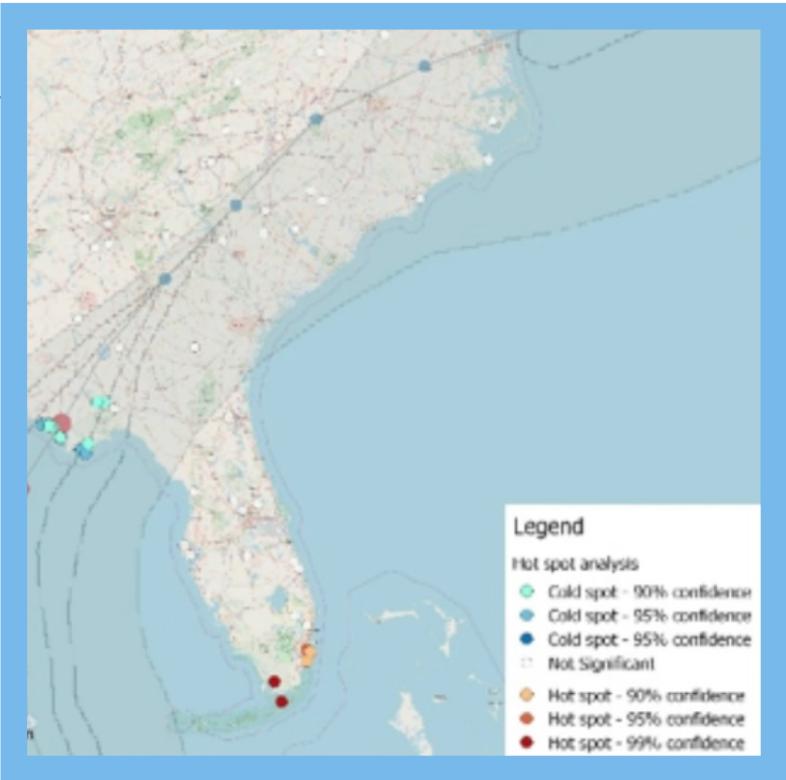


- Computed on filtered data at some stage(s)

Aggregation and
descriptive statistics



Crowd4SDG



Visualisation on Maps



UNIVERSITÉ
DE GENÈVE



CSIC
CONSEJO SUPERIOR DE INVESTIGACIONES CIENTÍFICAS

POLITECNICO
MILANO 1863

unitar
United Nations Institute
for Training and Research

Université
de Paris

Crowd4SDG

visualCit



Crowd4SDG



This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 872944

<https://pernici.faculty.polimi.it/crowd4sdgpolimi/>