# Querying and Visualizing Air Pollution Data Using Linked Open Data and Python

Leyan CHENG & Yuxin GONG & Jiani XU

# Catalogue

# Introduction

**01**

# Introduction

- In recent years, air pollution has become a major environmental issue around the world. It doesn't just affect the air we breathe—it also harms our health, nature, and even the climate. As cities grow and industries develop, sources of pollution like cars, factories, and wildfires are becoming more common.

- We now have better ways to search and understand data about air pollution from different places and topics. In this project, we use SPARQL to pull structured information from Wikidata, helping us explore and analyse.

- By combining this kind of smart data search with Python visualizations, we want to show a clearer picture of air pollution in open data and find patterns that can help raise awareness and guide better decisions.

# Exploratory Analysis of the Database

02

# Origin and Structure of the Data

## Wikidata
a collaborative, open knowledge base

## SPARQL
using SPARQL queries through the Wikidata Query Service

## Customized queries
instead of relying on pre-built dataset

## The data structure
a triple model: subject – predicate – object

## Results
The results of the queries were returned in tabular format

## Analysis and visualization
This format enabled further analysis and visualization using Python

# Data Related to Pollution



We started by searching for "air pollution" in the wikidata search bar to find the representation code.

We selected different types of information for better understanding of air pollution. (Causes、Effects、Policies, laws and Scientific studies)

These linked data help us explore and visualise air pollution directly from multiple perspectives.

# Data Analysis: SPARQL Queries
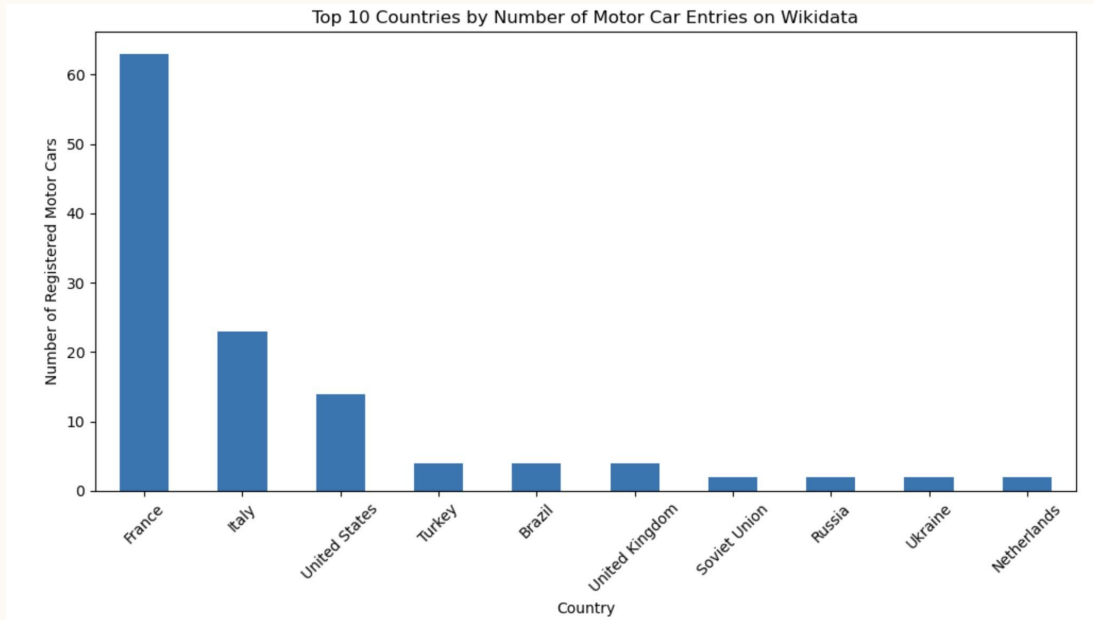
03

# Sources of air pollution

```
1  SELECT ?cause ?causeLabel WHERE {
2    wd:Q131123 wdt:P828 ?cause.
3    SERVICE wikibase:label { bd:serviceParam wikibase:language "en". }
4  }
```

we can see that the two main causes of air pollution are motor car and wildfire which correspond to the two types of human activities and natural sources of pollution respectively

| cause | causeLabel |
| --- | --- |
| 🔍 wd:Q1420 | motor car |
| 🔍 wd:Q2025 | carbon monoxide |
| 🔍 wd:Q165632 | dust |
| 🔍 wd:Q169950 | wildfire |
| 🔍 wd:Q7692360 | volcanic eruption |
| 🔍 wd:Q20962970 | NOx |
| 🔍 wd:Q50429805 | air pollutant |

# Enquire about the number of motor cars' model in each country

```
1  SELECT ?car ?carLabel ?countryLabel WHERE {
2    ?car wdt:P31 wd:Q1420.
3    ?car wdt:P17 ?country.
4    SERVICE wikibase:label { bd:serviceParam wikibase:language "en". }
5  }
```


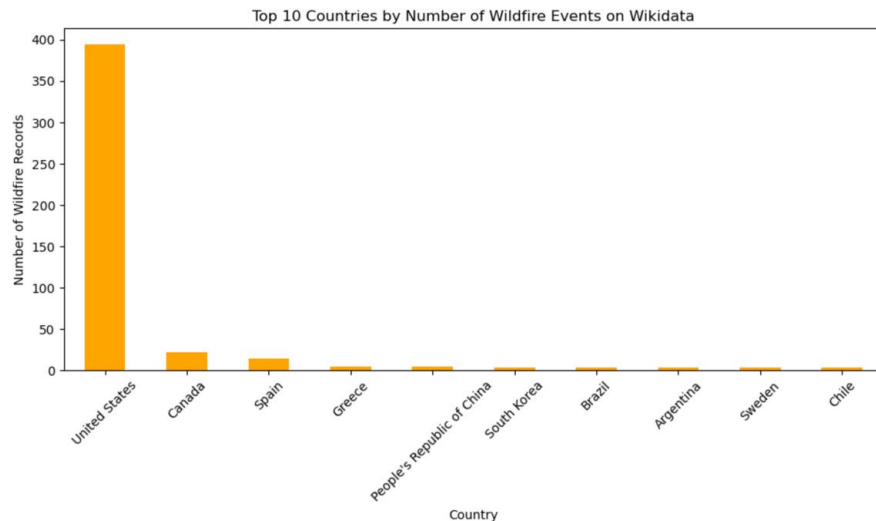Top 10 Countries by Number of Motor Car Entries on Wikidata

Motor vehicles from France are well recorded in semantic data sources. This suggests that transport-related emissions might be high. Countries like Italy and the United States also have many entries, which matches their known industrial growth and urban development. These countries often face problems like air pollution, greenhouse gas emissions, and city smog.

On the other hand, countries like Turkey, Brazil, and Russia have fewer records. This could mean that Wikidata has less data from these places, or that these countries have fewer motor vehicles.

# Enquire about the number of wildfire in each country

```
1  SELECT ?fire ?fireLabel ?countryLabel WHERE {
2    ?fire wdt:P31 wd:Q169950.
3    ?fire wdt:P17 ?country.
4    SERVICE wikibase:label { bd:serviceParam wikibase:language "en". }
5  }
6  LIMIT 500
```



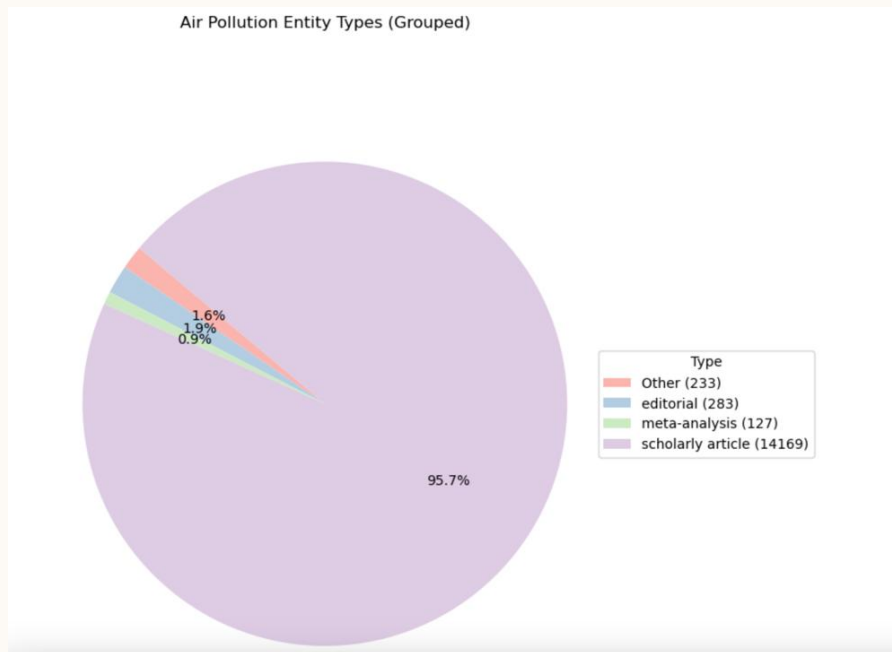Top 10 Countries by Number of Wildfire Events on Wikidata

The large number of recorded wildfires in the United States indicates that wildfire seasons have the potential to have an even greater impact on regional and even global air pollution.

Countries such as Canada, Spain, Greece and Brazil are also known for the occurrence of seasonal or climate-induced wildfires, although they show fewer entries.

Wildfires are becoming more frequent as a result of climate change and deforestation, and their role in worsening air quality and contributing to global warming cannot be underestimated.

# Air Pollution Topic Count number of typelabel

```
1  SELECT ?entity ?entityLabel ?typeLabel WHERE {
2    ?entity wdt:P921 wd:Q131123.
3    ?entity wdt:P31 ?type.
4    SERVICE wikibase:label { bd:serviceParam wikibase:language "en". }
5  }
```

Air Pollution Entity Types (Grouped)



1.6%
1.9%
0.9%

95.7%

**Type**
- Other (233)
- editorial (283)
- meta-analysis (127)
- scholarly article (14169)

This pie chart shows the types of items related to air pollution. Most of them are academic articles, which make up 95.7% of the total. This means that air pollution is mostly studied in universities and research institutions, and many of the articles are peer-reviewed.
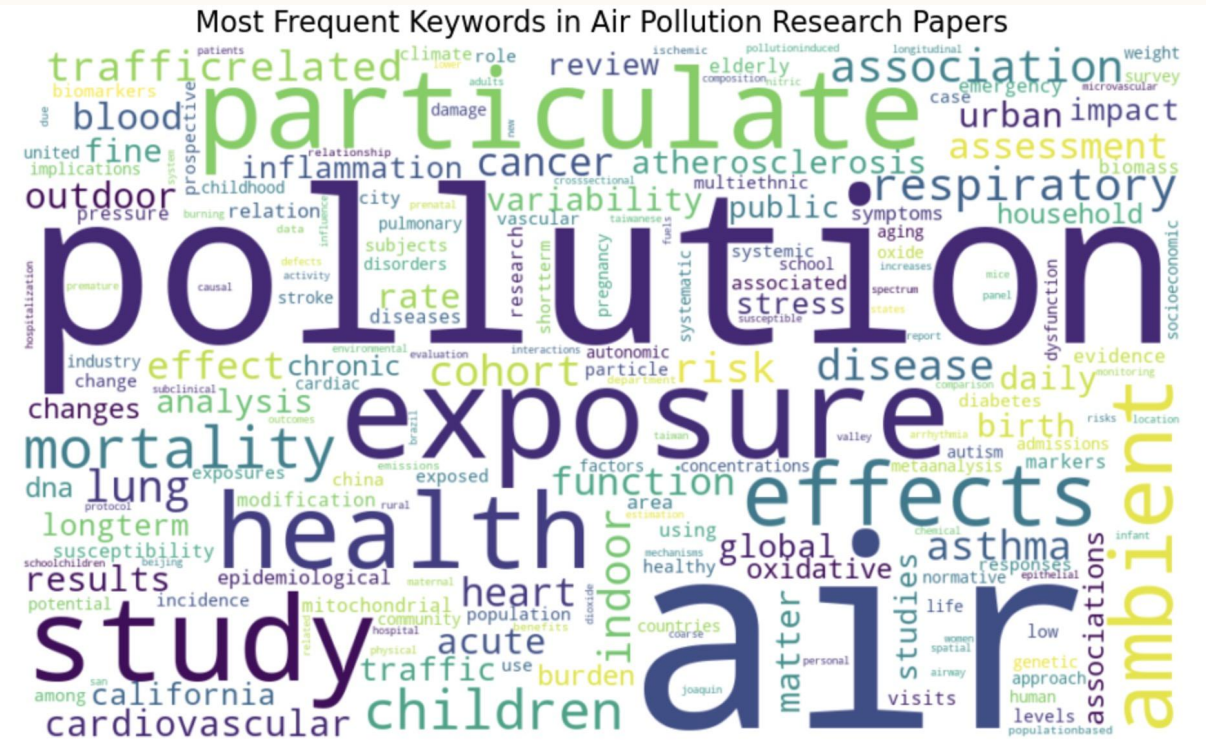
Other types, like editorials (1.9%), meta-analyses (0.9%), and others (1.6%), appear much less often.

In general, the chart shows that academic articles are the main way people share research about air pollution, while other ways of sharing scientific information are not used as much.

# Most Frequent keywords in Air Pollution Research Papers Title

```
1  SELECT ?paperLabel WHERE {
2    ?paper wdt:P921 wd:Q131123.
3    ?paper wdt:P31 ?type.
4    FILTER(?type IN (
5      wd:Q13442814,
6      wd:Q13406463,
7      wd:Q179461
8    ))
9    SERVICE wikibase:label { bd:serviceParam wikibase:language "en". }
10 }
11 LIMIT 300
```

A keyword frequency analysis was conducted based on the titles of over 300 academic publications related to air pollution. The resulting word cloud highlights the most commonly recurring terms in this field of research.



Most Frequent Keywords in Air Pollution Research Papers

# Impact of air pollution

```
1  SELECT ?effect ?effectLabel WHERE {
2    wd:Q131123 wdt:P1542 ?effect.
3    SERVICE wikibase:label { bd:serviceParam wikibase:language "en". }
4  }
```
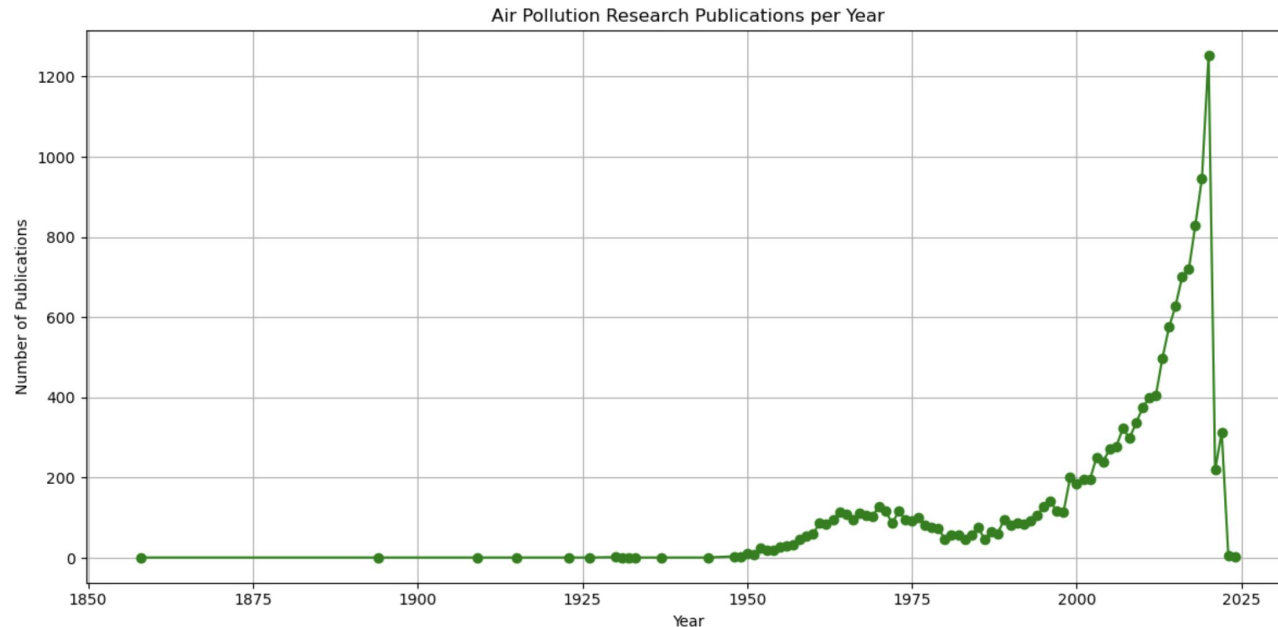
| effect | effectLabel |
|---|---|
| 🔍 wd:Q47912 | lung cancer |
| 🔍 wd:Q125928 | climate change |
| 🔍 wd:Q169994 | smog |
| 🔍 wd:Q3286546 | respiratory disease |

Real-life medical and environmental impacts are the most common consequences of air pollution.

Health impacts: lung cancer, respiratory disease → directly related to human health
Environmental impacts: climate change, smog → reflect the chain reaction of air pollution on natural systems

# Air pollution research publications per year



We analyzed the publication years of scholarly articles whose main subject is air pollution, as recorded in Wikidata. The data shows a clear upward trend in research output, especially after the year 2000.

Air pollution research has expanded rapidly, this growth reflects the rising global concern around environmental health, the development of international climate policies, and the availability of more environmental data for academic use.

**04**

# Conclusion

# Conclusion

Overall, the visualizations highlight how different data points, from human activities to natural events and academic knowledge, intertwine to shape our understanding of pollution. By combining these perspectives, this study shows that linked open data provides a rich, multidimensional view of environmental pollution – linking emission sources, event records, and knowledge representation. It also showcases various aspects of air pollution, allowing people to clearly understand the causes and effects.

# THANK YOU