

# Data Sets

---

## General Data Sets

---

- [UCSF Industry Documents](#)  
Massive digital archive of documents created by industries which influence public health, such as the tobacco, chemical, drug and fossil fuel industries.
- [LuminDatabase](#)  
A searchable database that collects and analyzes legal complaints and requests for removal of online materials (DMCA take-downs), helping Internet users to know their rights and understand the law.
- [National Center for Biotechnology Information](#)  
Allows you to search 39 different scientific databases such as Pubmed, SRA, OMIM, MedGen and more from a single page.
- [Afrobarometer](#)  
Large archive of sociological surveys conducted in African countries over the last ~20 years.
- [Arabbarometer](#)  
Large archive of sociological surveys conducted in the Arab countries of Africa and the Middle East from 2007 to 2021.
- [CensoredPlanet](#)  
A censorship measurement platform that collects data using multiple remote measurement techniques in more than 200 countries. Provides reports and offers their raw data sets which are available for download.
- [DomainsProject](#)  
Massive dataset of over 600 million domains. Total size is ~16 GB.
- [Face Recognition Datasets](#)  
A large collection of face datasets for training facial recognition systems and other things of that nature.
- [Kaggle](#)  
Offers over 50,000 public datasets for all kinds of various things.
- [Common Crawl](#)  
An open repository of web crawl data that can be accessed and analyzed by anyone.
- [OCCRP Catalogue of Research Databases](#)  
A massive collection of public data sources compiled by OCCRP researchers that are the most useful for investigative reporting.

## Government Data Sets

---

- [FBI Vault](#)  
The US FOIA library. Contains over 6,700 scanned FOIA documents.
- [CIA Reading Room](#)  
Search and view documents released through the FOIA and other CIA release programs.
- [US Department of State Records](#)  
Search through 221,373 documents reviewed and released to the public.
- [US National Archives](#)  
An independent agency of the United States government charged with the preservation and documentation of government and historical records. It is also tasked with increasing public access to those documents which make up the National Archive.

- [UK National Archives](#)  
One of the worlds largest archives, containing over 11 million historical government and public records. From Domesday Book to modern government files. Includes paper records, digital records, websites, photographs, posters, maps, drawings and paintings.
- [Canada Declassified](#)  
A digital repository of government records declassified under the Canadian Access to Information Act. Spans from 1945 through 1991.
- [Archives Canada](#)  
A gateway to over 800 archival repositories across Canada.
- [Australian National Archives Search](#)  
A tool to search the national archives of Australia. Includes an advanced search function.

## Leaked Data Sets

---

- [WikiLeaks](#)  
An international non-profit organization that publishes leaks and classified media from governments, companies and organizations alike. All data is provided by anonymous sources. #FreeAssange.
- [CryptoMe](#)  
Publishes documents that are prohibited by governments worldwide. Particularly material on freedom of expression, privacy, cryptology, dual-use technologies, national security, intelligence, secret governance, open, secret and/or classified documents.
- [ICIJ Offshore Leaks](#)  
Data from more than 785,000 offshore companies, foundations and trusts from the [Panama Papers](#), [Offshore Leaks](#), [Bahamas Leaks](#), and the [Paradise Papers](#).
- [ICIJ Luxembourg Leaks](#)  
Also known as the "LuxLeaks", is a collection of over 350 documents about Luxembourg's tax rulings set up by PricewaterhouseCoopers from 2002 to 2010 to the benefits of its clients.