# Covid-19 available corpora

**Novel Coronavirus (2019-nCoV) [Multilingual]**
A collection created by the Content Development Group of the International Internet Preservation Consortium in collaboration with Archive-It to preserve web content related to the ongoing Novel Coronavirus (Covid-19) outbreak. Identification of seed websites and initial web crawling began in February 2020, and the collection will continue to add new content as needed during the course of the outbreak and its containment. High priority subtopics include: coronavirus origins; information about the spread of infection; regional or local containment efforts; medical and scientific aspects; social aspects; economic aspects; and political aspects. Websites from anywhere in the world and in any language are in scope : https://archive-it.org/collections/13529


**COVID-19 Open Research Dataset (CORD-19) [English]**
Free resource of over 51,000 scholarly articles, including over 40,000 with full text about COVID-19 and the Coronavirus family of viruses for use by the global research community. https://www.semanticscholar.org/cord19

**Coronavirus Corpus [English]**
The Coronavirus Corpus is designed to be a comprehensive record of the social, cultural, and economic impact of the Coronavirus (COVID-19) in 2020 and beyond, and it is part of the English-Corpora.org suite of corpora, which offer insight into genre-based, historical, and dialectal variation in English. Early May 2020 the corpus was  about 270 million words in size, and it continues to grow by 3-4 million words each day. https://www.english-corpora.org/corona/


**Social Media for Public Health [English]**
Social Media for Public Health collects links to resources for academics and other investigators, including datasets pertinent to the COVID-19 pandemic. http://www.socialmediaforpublichealth.org/covid-19/resources/

**Links to other resources :**
https://github.com/bigheiniu/awesome-coronavirus19-dataset
https://www.covid19-archive.com/
list of broadcast news related to cover-19 (on archives.org) : http://data.gdeltproject.org/blog/2020-coronavirus-narrative/live_tvnews/MASTERFILELIST.TXT
https://www.clarin.eu/covid-19


**NLP / Corpus analysis tools applied to covid-19 corpora**
analysis of covid-19 lexical dynamics by OED : https://public.oed.com/blog/corpus-analysis-of-the-language-of-covid-19/
SKetchEngine access to COVID-19 Open research Dataset : https://app.sketchengine.eu/#dashboard?corpname=preloaded%2Fcovid19
IBM Knowledge graphs from corpora : https://ds-covid19.res.ibm.com/about