5th International Conference on AI in Computational Linguistics

# AraCOVID19-MFH: Arabic COVID-19 Multi-label Fake News & Hate Speech Detection Dataset

Mohamed Seghir Hadj Ameur[a], Hassina Aliane[a]

[a]*Research and Development in Digital Humanities Division, Research Centre on Scientific and Technical Information (CERIST)/Algiers, Algeria*

## Abstract

Along with the COVID-19 pandemic, an "infodemic" of false and misleading information has emerged and has complicated the COVID-19 response efforts. Social networking sites such as Facebook and Twitter have contributed largely to the spread of rumors, conspiracy theories, hate, xenophobia, racism, and prejudice. To combat the spread of fake news, researchers around the world have and are still making considerable efforts to build and share COVID-19 related research articles, models, and datasets. This paper releases "AraCOVID19-MFH"[1] a manually annotated multi-label Arabic COVID-19 fake news and hate speech detection dataset. Our dataset contains 10,828 Arabic tweets annotated with 10 different labels. The labels have been designed to consider some aspects relevant to the fact-checking task, such as the tweet's check worthiness, positivity/negativity, and factuality. To confirm our annotated dataset's practical utility, we used it to train and evaluate several classification models and reported the obtained results. Though the dataset is mainly designed for fake news detection, it can also be used for hate speech detection, opinion/news classification, dialect identification, and many other tasks.

*Keywords:* Arabic COVID-19 Multi-label Dataset; Annotated Dataset; Fake News Detection; Hate Speech Detection; Misinformation; Social Media; Arabic Language

## 1. Introduction

Coronavirus disease (COVID-19) is an infectious respiratory disease caused by the "Sars-CoV-2" virus [17]. It was discovered in Wuhan, China, in December 2019, and declared a global pandemic by the World Health Organization (WHO) in March 2020 [10]. To reduce the spread of the virus, governments have adopted several measures such as closing borders, travel restrictions, quarantine, and containment. As of late January 2021, COVID-19 has caused more than 100 million confirmed cases and 2 million deaths worldwide[2]. The World Health Organization has reported that along with the COVID-19 pandemic, an "infodemic" of false and misleading information has emerged and has complicated the COVID-19 response efforts[3]. Indeed, with over 4.2 million active users[4], social networking sites such as Facebook and Twitter have contributed largely to the spread of rumors, conspiracy theories, hate,

---

[1] https://github.com/MohamedHadjAmeur/AraCOVID19-MFH

*E-mail addresses:* mhadjameur@cerist.dz (Mohamed Seghir Hadj Ameur)., ahassina@cerist.dz (Hassina Aliane).

[2] https://www.worldometers.info/coronavirus/
[3] https://www.who.int/news-room/feature-stories/detail/immunizing-the-public-against-misinformation
[4] https://datareportal.com/social-media-users