

**SoMeIL: A social media infodemic listening for public health behaviours
conceptual framework**

by

Shu-Feng Tsao

A thesis

presented to the University of Waterloo

in fulfillment of the

thesis requirement for the degree of

Doctor of Philosophy

in

Public Health Sciences

Waterloo, Ontario, Canada, 2023

© Shu-Feng Tsao 2023

Examining Committee Membership

The following served on the Examining Committee for this thesis. The decision of the Examining Committee is by majority vote.

External Examiner.

Dr. Lu Xiao

Associate Professor, School of Information Studies,
Syracuse University

Supervisors

Dr. Helen Chen

Professor of Practice, School of Public Health
Sciences, University of Waterloo

Dr. Zahid A. Butt

Assistant Professor, School of Public Health
Sciences, University of Waterloo

Committee Member

Dr. Samantha B. Meyer

Associate Professor, School of Public Health
Sciences, University of Waterloo

Dr. Plinio Morita

Associate Professor, School of Public Health
Sciences, University of Waterloo

Internal-External Committee Member

Dr. Aimée Morrison

Associate Professor, Department of English
Language and Literature, University of Waterloo

Author's Declaration

This thesis consists of material all of which I authored or co-authored: see Statement of Contributions included in the thesis. This is a true copy of the thesis, including any required final revisions, as accepted by my examiners.

I understand that my thesis may be made electronically available to the public.

Statement of Contributions

As this thesis adopts a manuscript-based approach, there are three manuscripts written for publication. Shu-Feng Tsao was the sole author of Chapters 1 and 5, which were written under the supervision of Dr. Helen Chen and Dr. Zahid A. Butt and were not written for publication. References for Chapters 1 and 5 are included in the final references starting from page 152.

The manuscript in Chapter 2 has been published in a peer-reviewed journal. Manuscripts in Chapter 3 and Chapter 4 have been submitted to peer-reviewed journals and are currently under peer review. Exceptions to sole authorship of material are as follows:

Chapter 2: Tsao S-F, Chen H, Tisseverasinghe T, Yang Y, Li L, Butt ZA. What social media told us in the time of COVID-19: a scoping review. *Lancet Digit Health* [Internet]. 2021;3(3):e175–94. Available from: [http://dx.doi.org/10.1016/s2589-7500\(20\)30315-0](http://dx.doi.org/10.1016/s2589-7500(20)30315-0)

The study was conducted at the University of Waterloo by Shu-Feng Tsao under the supervision of co-authors Dr. Helen Chen and Dr. Zahid A. Butt. Other co-authors have helped review articles included in the study and proofread the manuscript.

Chapter 3: Tsao S-F, Chen H, Meyer SB, Butt ZA. Proposing a conceptual framework: social media listening to public health behavior. Currently under journal review.

Shu-Feng Tsao was responsible for framework conceptualization, data collection, analysis, and drafting, submitting, and revising the manuscript under supervision from co-authors Dr. Helen Chen, Dr. Zahid A. Butt, and Dr. Samantha B. Meyer. They have provided methodological guidance and feedback on the draft manuscript.

Chapter 4: Tsao S-F, Chen H, Meyer SB, Butt ZA. A computer-assisted qualitative validation and demonstration of the SoMeIL conceptual framework. Currently under journal review.

Shu-Feng Tsao was responsible for data collection and analysis, as well as drafting, submitting, and revising the manuscript under supervision from co-authors Dr. Helen Chen and Dr. Zahid A. Butt. They have provided methodological guidance and feedback on the draft manuscript.

Chapter 5: Tsao S-F, Chen H, Meyer SB, Butt ZA. Validating social media infodemic listening conceptual framework using structural equation modelling. Currently under journal review..

Shu-Feng Tsao designed the study, collected, and analyzed data, as well as drafted and submitted the manuscript with supervision from Dr. Helen Chen and Dr. Zahid A. Butt. They have contributed to the papers by offering overall theoretical and statistical guidance and detailed suggestions on the papers' structure and writing.

Abstract

Introduction

The coronavirus disease 2019 (COVID-19) pandemic has escalated health infodemics given substantially digitalized daily life since the pandemic began. The number of social media users has skyrocketed. However, this has brought issues given misleading health information circulating on social media platforms that can lead to undesirable behaviours compromising individual or public health in real life. One long-lasting health issue is vaccine hesitancy, which has been further compounded by health infodemics on social media. According to the World Health Organization, health infodemics occur when too much information that makes true information competes with misinformation for people's attention, understanding, and adherence to recommended health interventions. Existing theories and theoretical constructs have been applied to study public behaviours influenced by health infodemics on social media. However, these theories have limited to individual behaviours and ignored other critical factors. Furthermore, the current theories have rarely reflected the nature of social media as information can be disseminated instantly and massively without geographical restrictions regardless of information quality. Therefore, this dissertation aimed to address these limitations by proposing a solution that can listen to public discourse on social media and infer their behavioural intentions in real life.

Methods

The scoping review (Study I) was conducted by following the methods of Arksey and O'Malley as well as Levac et al. to identify and synthesize literature related to the research question. The theory construction methodology was used in the conceptual paper (Study II) to review existing theories and propose a new conceptual framework. Next, the Latent Dirichlet allocation topic modelling and qualitative thematic analysis were applied in the preliminary and partial qualitative validation study (Study III). The last study (Study IV) applied structural equation modeling (SEM) to infer people's intentions toward COVID-19 vaccination in real life from Twitter amid the pandemic as a preliminary and partial validation for the proposed conceptual framework.

Results

A total of 2,405 articles published between November 1, 2019, and November 4, 2020, were retrieved from PubMed, Scopus, and PsycINFO. After removing duplicates, non-empirical literature, and irrelevant studies, a total of 81 articles written in English published in peer-reviewed journals were included in the scoping review (Study I). Six themes were found and reported: (1) surveying public attitudes, (2) identifying infodemics, (3) assessing mental health, (4) detecting or predicting COVID-19 cases, (5) analyzing government responses to the pandemic, and (6) evaluating quality of health information in prevention education videos. The findings also suggested knowledge gaps in real-time COVID-19 surveillance using social media data and limited machine learning or artificial intelligence techniques used in overall COVID-19 research using social media data except the first theme. In the conceptual paper (Study II), a new conceptual framework—social media infodemic listening for public health behaviors (SoMeIL) —was proposed to address limitations in existing theories given lacking systematic and theoretical foundation for such research. After the SoMeIL was proposed, validations were needed. A preliminary qualitative validation and demonstration using Twitter data about the Canadian Freedom Convoy were conducted to partially validate and illustrate how the SoMeIL conceptual framework could be applied (Study III). Finally, the findings from SEM in the last study (Study IV) showed statistically significant associations between the latent variable and the observed variables derived from Twitter. This study provided preliminary evidence to validate partial components in the proposed SoMeIL conceptual framework that could be used as a proxy to infer people’s vaccination intentions in real life. It also demonstrated the feasibility of using Twitter data in SEM research besides typical surveys.

Conclusion

The scoping review (Study I) was important since it identified various roles that social media data have played in research related to the COVID-19 pandemic. It also informed us of knowledge gaps to be bridged. This led us to the conceptual paper (Study II) since we identified limitations in existing theories when the current theories or theoretical constructs were applied in health research that analyzed social media data. A new conceptual framework—SoMeIL—was proposed accordingly. A preliminary qualitative study was followed to validate and demonstrate partial components of the SoMeIL conceptual framework. The last study (Study IV) showed preliminary evidence to show that parts of the SoMeIL conceptual framework was workable given statistically significant relationships

found among certain constructs. As a result, Twitter data in this dissertation could be used as a proxy to infer people's vaccination behavior in real life as suggested by the proposed conceptual framework. Yet more research is needed to further validate and improve the proposed SoMeIL conceptual framework. If social media listening can be integrated into future pandemic preparedness as the proposed conceptual framework suggests, it can help health authorities and governmental agencies promptly shape public perception, disseminate more scientific information, and influence behaviors during a health crisis in a timely fashion.

Acknowledgements

Since I started my PhD program in fall, 2020, in the middle of the COVID-19 pandemic, there have been many people who helped me along the way on this journey. I would like to thank them because this journey would not have been possible without them.

First and foremost, I am extremely grateful to my supervisors Dr. Helen Chen and Dr. Zahid A. Butt for their invaluable advice, continuous support, and patience during my PhD study. Their immense knowledge and plentiful experience have encouraged me in all the time of my academic research and daily life. In addition, I'd like to express my gratitude to my committee members for their time, effort, and understanding in helping me succeed in my studies throughout the pandemic. Words cannot express my gratitude to my supervisors' and committee members' generosity and encouragement, my time spent studying and living in the University of Waterloo has been truly rewarding.

I gratefully acknowledge the funding received towards my PhD from the 2023-24 Ontario Graduate Scholarship for international students.

I am also thankful to my fellows, Yang Yang, Lianghua Li, and Therese Tisseverasinghe, who have helped me collect, review, and summarize articles for my scoping review. My gratitude extends to Alexander MacLean, Shih-Hsio Huang, and Gaya Bin Noon, who were involved in other research projects using social media data. Although some papers are not included in my dissertation, they have helped me gain further knowledge in my overall doctoral research.

Last but not least, I want to thank my parents for their unconditional love and support throughout my life, and my best friends for their ongoing emotional and intellectual support for the past few years when we were all isolated during the pandemic lockdowns.

Dedication

I dedicate this dissertation as a tribute to my parents and best friends. I am incredibly grateful to my loving parents, whose inspiration and drive for persistence continue to reverberate in my ears. I also dedicate this dissertation to my best friends for their help and encouragement during my PhD journal and the pandemic.

Table of Contents

Examining Committee Membership.....	ii
Author’s Declaration	iii
Statement of Contributions.....	iv
Abstract	vi
Acknowledgements	ix
Dedication	x
List of Figures	xv
List of Tables.....	xvi
Chapter 1 : Background.....	1
1.1 Problem Statement	1
1.2 Research questions and objectives	3
1.3 Dissertation Structure	4
1.4 Methodological considerations.....	6
Summary of Chapter 1	7
Chapter 2 : What social media told us in the time of COVID-19: a scoping review	9
2.1 Summary	9
2.2 Introduction	10
2.3 Methods.....	11
2.3.1 Overview	11
2.3.2 Data Sources.....	11
2.3.3 Screening procedure	12
2.4 Results	12

2.4.1 Social media as contagion and vector.....	46
2.4.2 Social media for surveillance and monitoring.....	47
2.4.3 Social media as disease control.....	54
2.5 Discussion.....	59
2.6 Conclusion.....	61
2.7 References.....	63
Chapter 3 : Proposing a conceptual framework: social media infodemic listening (SoMeIL) for public health behaviours.....	73
3.1 Abstract.....	74
3.2 Introduction.....	75
3.3 Methods.....	77
3.4 Results.....	78
3.4.1 Synthesis of Theories.....	78
3.4.2 Proposed Conceptual Framework.....	81
3.5 Discussion.....	88
3.6 References.....	90
Chapter 4 : A computer-assisted qualitative validation and demonstration of the SoMeIL conceptual framework.....	97
4.1 Abstract.....	99
4.2 Introduction.....	100
4.3 Methods.....	102
4.3.1 Data Collection and Preprocessing.....	102
4.3.2 Data Analysis.....	102

4.3.3 Ethical Approval.....	103
4.3.4 Funding.....	103
4.4 Results	103
4.5 Discussion	106
4.6 Conclusion.....	109
4.7 References	110
Chapter 5 : Validating part of the social media infodemic listening conceptual framework using structural equation modelling	113
5.1 Abstract	115
5.2 Introduction	116
5.3 Methods	120
5.3.1 Data collection.....	120
5.3.2 Measures.....	122
5.3.3 Statistical Analysis	122
5.3.4 Ethical approval.....	123
5.3.5 Funding.....	123
5.4 Results	123
5.5 Discussion	129
5.6 Conclusion.....	132
5.7 Declaration of interests.....	132
5.8 Data sharing statement	132
5.9 References	133
Chapter 6 Conclusion	138

6.1 Summary of key findings	138
6.2 Study limitations and strengths	140
6.3 Directions for future research.....	141
6.4 Implications for public health practices and policies	142
References	144
Appendix A : Chapter 3 Supplementary Materials	150

List of Figures

Figure 2-1: Preferred Reporting Items for Systematic Reviews and Meta-Analyses flow chart of article extraction from the literature search.....	13
Figure 2-2: Modified Social Media and Public Health Epidemic and Response framework.....	14
Figure 3-1: Social media infodemic listening (SoMeIL) for public health behavior conceptual framework	82
Figure 4-1: Qualitative thematic analysis results	105
Figure 5-1: Proposed conceptual framework.....	120
Figure 5-2: Model 1.....	126
Figure 5-3: Model 2.....	128

List of Tables

Table 2-1: Summary of chosen article.....	45
Table 2-2 Summaries of findings and research gaps identified in each theme.....	58
Table 3-1 Theories and models used in health infodemic research in the context of the COVID-19 pandemic.	79
Table 3-2 Attributes of each components in the SoMeIL conceptual framework.....	84
Table 4-1: Latent Dirichlet allocation topic modeling results	104
Table 5-1: Keywords and hashtags for data collection.....	122
Table 5-2: Descriptive statistics of measured variables	124
Table 5-3: Spearman correlations for measured variables	124
Table 5-4: Fit statistics for each latent variable and full measurement model	125
Table 5-5: Model fit indices for two SEM models.....	127
Table 5-6: Model 1 estimates after conversion	129
Table 4--5-7: Model 2 estimates after conversion.....	129

Chapter 1: Background

1.1 Problem Statement

Throughout the COVID-19 pandemic, social media has rapidly become a crucial tool for various purposes, including but not limited to information distribution and consumption, social connections, and service utilization.¹ The number of social media users has surged since the pandemic began as people were confined to their houses given lockdowns and other interventions.^{2,3,4} A global survey reported an increase of 61% on social media.² Facebook's overall usage increased by 37%, with the biggest increase from users in the 18-34 age group.² The same age group also contributed to an over 40% increase in usage on Instagram and WhatsApp besides Facebook.² In Canada, a 2020 survey revealed that the number of daily users on YouTube and Instagram increased by 16% and 8%, respectively.³ Furthermore, approximately 94% of Canadian adults had at least one social media account.³ A follow-up study was conducted in 2022 and its findings showed that TikTok had the largest increase in the number of users (11%).⁴ Facebook still had the highest percentage of daily users, but it decreased from 77% in 2020 to 70% in 2022.⁴

Literature has suggested that social media has implicitly or explicitly shaped people's perceptions or attitudes toward the COVID-19 pandemic and related interventions, such as the COVID-19 vaccination. Although social media has been useful for faster information disseminations, it has also caused issues, especially health infodemics.^{1,5,6} According to the World Health Organization (WHO), an infodemic is defined as "too much information—including false or misleading information—in digital and physical environments during a disease outbreak."⁷ Health infodemics can lead to many undesirable outcomes, such as detrimental behaviors undermining individual or public health, and mistrusts in health professionals, authorities, or governments.⁷ When the COVID-19 vaccines became available in late 2020, health infodemics have further exacerbated vaccine hesitancy, which has been one of top ten threats in global health since 2019 according to WHO.⁸ Literature has shown that social media can contribute to vaccine hesitancy in several ways.⁹⁻¹⁴ To begin with, social media platforms has been breeding grounds for the spread of misinformation about COVID-19 vaccines. Conspiracy theories, myths, fake news, and misleading information about vaccine ingredients, safety, and efficacy can easily go viral and create doubts among users.⁹⁻¹⁸ Although personal stories can be powerful in building empathy, social media can also share misleading anecdotes of adverse reactions or vaccine failures, which may not be representative of the general vaccine experience.⁹⁻¹² Influencers

or celebrities on social media with large followers can sway public opinion.^{13, 14} Throughout the COVID-19 pandemic, some influencers have promoted vaccine hesitancy based on personal beliefs or misinformation, leading their followers to question vaccination.^{13, 14} Furthermore, social media algorithms actively show users content that aligns with their existing beliefs and interests.^{9, 15} This can create echo chambers, where users are exposed to information that reinforces their doubts about vaccines, further entrenching vaccine hesitancy.^{16, 17} Additionally, social media can amplify antivaccine sentiment by providing a platform for vocal individuals or against vaccination^{9, 16, 17}. These groups can spread their views widely and attract like-minded followers even though they may account for a very small proportion of the general public in the real world.^{9, 16, 17} Last but not least, some antivaccine content on social media uses emotional appeals and fearmongering to dissuade people from getting vaccinated.^{11, 18} Such tactics can heighten anxieties and uncertainties about the vaccines' safety and efficacy.^{11, 18} It is crucial to recognize that social media's role in vaccine hesitancy is complex, and not all users are influenced by misinformation or become hesitant toward vaccination due to social media exposures. Yet social media platforms have been identified as substantial contributors to the spread of vaccine misinformation and the shaping of public opinions on the COVID-19 vaccination.

Given social media's contribution to the COVID-19 vaccine hesitancy, it is necessary to better understand public discourse regarding the COVID-19 pandemic and vaccines on social media at larger scales. Machine learning (ML) or artificial intelligences (AI) techniques have been adopted to collect and analyze hundreds of thousands or even millions of social media data. WHO has coined such approach as "social listening" and deployed its "Early AI-supported Response with Social Listening Platform" (EARS) to monitor discussions related to the COVID-19 pandemic and vaccine hesitancy.^{19, 20} Following the WHO's approach and calls for more social listening research, countless studies have been conducted and published with existing theories or constructs derived from these theories. A systematic review indicated that the health belief model (HBM) has been applied frequently, followed by the theory of planned behavior (TPB) and the social cognitive theory (SCT).²¹ Another theoretical review investigated how social media has shaped public risk perceptions of the pandemic through lens of fear drive model, self-determination theory, perceived locus of causality, and cultivation theory.²² The other scoping review has summarized 26 theories and 51 theoretical constructs in different disciplines that have also been used to investigate the COVID-9 vaccine

hesitancy,²³ with new theories or framework being developed, such as the syndemic conceptual framework for COVID-19 vaccine hesitancy²⁴ and WHO's behavioural and social drivers (BeSD) of vaccination.²⁵ It is common that the same theories or their constructs have been continuously used in health research, from the vaccine hesitancy to the social listening.²¹⁻²⁵ However, there are limitations in existing theories given the increasingly multidimensional and evolving information ecosystems, including social media, in modern society. Since most existing theories were developed before social media has become popular, they do not really reflect the current nature of social media, such as instant dissemination of user-generated content or message without physical restrictions with reliable Internet and smartphone coverage. Unlike conventional televisions and newspapers, social media users can create and distribute their content instantly and massively, as well as receive content generated by other users on these platforms instantly. This can be compounded by social media algorithms that actively "recommend" content to users according to their digital footprints and connections, or a viral event at a moment, regardless of content quality. Therefore, despite abundant theories and advanced ML or AI techniques, there has been a lack of systematic approaches to conduct social listening on social media. In addition, limited health research has directly used social media data to investigate how people's online behaviors on social media can reflect their behaviors in real life as theories have suggested.

1.2 Research questions and objectives

To address these limitations and gaps in literature and research, this dissertation addressed the following research questions:

- How has social media data been used in research related to the COVID-19 pandemic?
- Are there existing theories which adequately address people's health behaviours using social media data?

These questions led to the development of a framework and partial validation, meeting the following objectives:

- To identify roles of social media have played since the COVID-19 pandemic began and knowledge gaps (Study I Chapter 2).

- To review current theories and identify limitations and gaps in their applications to social media listening (Study II Chapter 3).
- To propose a conceptual framework that can address identified gaps (Study II Chapter 3).
- To validate the SoMeIL conceptual framework.
- To qualitatively validate part of the SoMeIL conceptual framework using Twitter data about the Canadian Freedom Convoy as an example (Study III Chapter 4).
- To quantitatively validate the partial SoMeIL conceptual framework using Twitter data regarding the first COVID-19 vaccine behaviors as an example (Study IV Chapter 5).

1.3 Dissertation Structure

This dissertation adopts a manuscript-based approach. It includes four manuscripts (Chapter 2 already published in peer-reviewed journals, Chapters 3, 4, and 5 currently under journals' peer reviews) for which I am the first author. For all studies, I was responsible for study design conceptualization, data collection, analysis, and drafting, submitting, and revising manuscripts. My co-authors provided methodological guidance and feedback on the draft manuscripts. To meet the research objective, the first study (Study I) in Chapter 2 conducted was a scoping review designed to understand how social media data have been used in research related to the COVID-19 pandemic. Extreme interventions, such as lockdowns, have made people isolated and more digitally connected. In other words, people have relied on digital information channels, such as social media, to receive information and stay connected with their friends or families since the pandemic began. As a result, social media data has become one of valuable data sources for researchers to study the pandemic given the restrictions. This scoping review has not only summarized important ways that social media has been used during the pandemic, but also identified gaps in current infodemic research using social media data.

Study II (Chapter 3) proposed a conceptual framework addressing how public health behaviors, such as the COVID-19 vaccination, could be inferred using social media data with advanced ML or AI techniques, especially natural language processing (NLP). It was designed to understand the concepts derived from theories related to health behavior changes and communications. Through this conceptual paper, I documented critical concepts hypothetically attributed to people's vaccination

intentions or behaviors, including but not limited to attitudes, motivations, perceptions, needs, abilities, and actions/behaviors. This investigation of theories has helped to refine my understanding of the concepts and to clarify how these three intermediate variables may shape people's intentions or acceptance of the COVID-19 vaccines. Furthermore, it has helped me design and propose a conceptual framework that uses ML or NLP techniques and parameters retrieved from social media platforms as proxies to figure the vaccination intentions or behaviors. This conceptual paper was critical in identifying limitations in existing theories and the lack of systematic ways for such social media listening research. This study added value to the existing literature in terms of proposing a solution; namely, the conceptual framework. A qualitative study was included in this paper to demonstrate how the conceptual framework could be applied.

Study III (Chapter 4) aimed to preliminarily validate partial components of the SoMeIL conceptual framework by using Latent Dirichlet allocation (LDA) topic modelling and qualitative thematic analysis. The health information in this study was the massive COVID-19 vaccination campaigns by the Canadian governments to encourage their residents to take the first dose of COVID-19 vaccines. The partial components included online reaction behaviours, emotion, and self-reported COVID-19 vaccination as the offline reaction behaviours. These components helped investigate associations between people's self-reported offline behaviours and their emotions and online reaction behaviours. Another objective of this study was to demonstrate how the SoMeIL framework could be applied.

Study IV (Chapter 5) aimed to validate the proposed conceptual framework by using structural equation modeling (SEM). Another objective of this study was to demonstrate that social media data could be used in SEM analysis since surveys have been primarily used in typical SEM.

Since the sensitive Twitter data was retrieved from Twitter's application programming interface (API), ethics approval was obtained from the Office of Research Ethics (ORE), University of Waterloo. To achieve this, all identifiers that could be used to identify Twitter users from the data were stored separately and later deleted permanently, except unique identifier for each tweet. In our studies, aggregated tweets were analyzed, and individual tweets were quoted anonymously.

Although the included studies have been undertaken independently, together they provide empirical knowledge to the broader goal of the dissertation study. Their sequential presentation herein highlights my doctoral research. Since infodemic research is a relatively young field and has been

evolving, there is a scarcity of appropriate theoretical and systematic foundations in social media listening. Hence, following the development from the scoping review to the conceptual article and its validation study, this research overall provided a more comprehensive understanding regarding how social media data can be used as proxy for people's vaccination intentions, as suggested by the proposed conceptual framework.

1.4 Methodological considerations

For Chapter 2 Study I, it was guided by the scoping review methods of Arksey and O'Malley²⁶ and Levac et al.²⁷ The authors followed the five-step scoping review protocol²⁸ and the Preferred Reporting Items for Systematic Reviews and Meta-Analyses extension for scoping reviews.²⁹ The five steps included: (1) defining research questions, (2) identifying relevant literature, (3) selecting studies, (4) charting the data, and finally (5) collating, summarizing, and reporting the results.²⁹ Instead of systematic reviews, the scoping review was conducted because we aimed to identify the types of available evidence and knowledge gaps regarding what roles social media has played since the COVID-19 pandemic began. This was a relatively general research question compared to very specific research questions needed in systematic reviews.³⁰ Furthermore, it was conducted during the beginning of the COVID-19 pandemic when insufficient literature was published related to our research questions.³⁰ Therefore, a scoping review was more appropriate than a systematic review given the research question and the study timing.

For Chapter 3 Study II, the theory construction methodology (TCM) was adopted to develop the conceptual framework³¹. There are five steps in TCM: (1) identification of relevant phenomena, (2) development of a proto theory, (3) development of a formal model, (4) adequacy evaluation of the formal model, and (5) assessment of overall worth of the formal model.³¹ TCM was proposed to address lack of scientific methodology to develop a theory with explanatory theory formation instead of hypothesis testing in psychology.³¹ Like psychology, hypothesis testing has also dominated in quantitative health studies. Although qualitative research, such as grounded theory, can be used to develop a theory with both quantitative and qualitative data,^{32, 33, 34} the development of the proposed conceptual framework in the second study does not fit into the philosophies of grounded theory. One major difference is that NLP techniques, including topic modeling and sentiment analysis, are used to code data instead of iterative coding process by human researchers in grounded theory research.

Additionally, we did not approach the data inductively without prior theories or frameworks. On the contrary, we aimed to develop a conceptual framework addressing known gaps and limitations in existing theories. Although TCM is a relatively new methodology with limited research applications, it was more appropriate to apply the TCM to develop the conceptual framework since it combines both quantitative and qualitative principles for theory development. TCM also provides an overall structural guidelines and flexibility for the development and validation of the conceptual framework.

For Chapter 4 Study III, LDA topic modelling³⁵ and qualitative thematic analysis³⁶⁻³⁷ were used to preliminarily validate partial components of the SoMeIL conceptual framework.

For the last study in Chapter 5, structural equation modeling (SEM) was used to partially validate the SoMeIL conceptual framework. SEM has been commonly used to test associations when latent variables are involved.³⁸ A latent variable refers to an unmeasured or unobserved variable but is assumed to exist based on theories or other observable or measurable variables in statistical models.³⁸ In the SoMeIL conceptual framework, the latent variable was people's self-reported intention to become vaccinated against COVID-19 on Twitter, which could be used as a proxy for people's vaccination behaviors in real life. This latent variable could be inferred by other observed variables derived from the Twitter platform, such as sentiment scores and the numbers of likes and shares. Hence, SEM was used to test associations between the latent variable and other measured variables, thus preliminarily validating some components of the SoMeIL conceptual framework.

Summary of Chapter 1

In conclusion, this chapter provided an overview of a complex problem investigated throughout this dissertation. In the problem statement, a literature review was conducted to identify the knowledge gaps known so far that this dissertation has addressed, while providing an outline of the research objectives that each chapter achieves. This chapter also highlights methodological considerations relevant to the studies that were beyond the space allocated in an article written for peer-reviewed publication. Therefore, by conducting a scoping literature review, by proposing a conceptual framework, and by preliminarily validating part of the SoMeIL conceptual framework, this dissertation has identified and addressed the knowledge gaps regarding research using social media data to infer public health behaviors, such as people's vaccination intentions during the COVID-19 pandemic and beyond. The last chapter of this dissertation will discuss the contributions and impacts

of this dissertation, providing recommendations drawn from the findings to initiate a more robust social media listening research that caters to a wide variety of researchers addressing health infodemics and vaccine hesitancy.

Chapter 2: What social media told us in the time of COVID-19: a scoping review

Status: Published

Citation: Tsao S-F, Chen H, Tisseverasinghe T, Yang Y, Li L, Butt ZA. What social media told us in the time of COVID-19: a scoping review. *Lancet Digit Health* [Internet]. 2021;3(3):e175–94.

Available from: [http://dx.doi.org/10.1016/s2589-7500\(20\)30315-0](http://dx.doi.org/10.1016/s2589-7500(20)30315-0)

In the wake of the global COVID-19 pandemic, the landscape of communication has transformed substantially, with social media emerging as a pivotal channel for the creation, distribution, and consumption of information. This scoping review dived into a comprehensive exploration of relationship between COVID-19 and social media during its initial outbreak spanning from November 2019 to November 2020. Through an examination of 81 peer-reviewed studies, six overarching themes emerged that illuminated different utilizations of social media data in research related to the COVID-19 pandemic. This chapter begins with the study’s abstract, followed by the full-text manuscript.

2.1 Summary

With the onset of the COVID-19 pandemic, social media has rapidly become a crucial communication tool for information generation, dissemination, and consumption. In this scoping review, we selected and examined peer-reviewed empirical studies relating to COVID-19 and social media during the first outbreak from November 2019 to November 2020. From an analysis of 81 studies, we identified five overarching public health themes concerning the role of online social media platforms and COVID-19. These themes focused on surveying public attitudes, identifying infodemics, assessing mental health, detecting or predicting COVID-19 cases, analysing government responses to the pandemic, and evaluating quality of health information in prevention education videos. Furthermore, our Review emphasises the paucity of studies on the application of machine learning on data from COVID-19-related social media and a scarcity of studies documenting real-time surveillance that was developed with data from social media on COVID-19. For COVID-19, social

media can have a crucial role in disseminating health information and tackling infodemics and misinformation.

2.2 Introduction

Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2), and the resulting COVID-19, is a substantial international public health issue. As of Jan 18, 2021, an estimated 95 million people worldwide had been infected with the virus, with about 2 million deaths.¹ As a consequence of the pandemic, social media is becoming the platform of choice for public opinions, perceptions, and attitudes towards various events or public health policies regarding COVID-19.² Social media has become a pivotal communication tool for governments, organisations, and universities to disseminate crucial information to the public. Numerous studies have already used social media data to help to identify and detect outbreaks of infectious diseases and to interpret public attitudes, behaviours, and perceptions.^{3, 4, 5, 6} Social media, particularly Twitter, can be used to explore multiple facets of public health research. A systematic review identified six categories of Twitter use for health research, namely content analysis, surveillance, engagement, recruitment, as part of an intervention, and network analysis of Twitter users.⁵ However, this review included only broader research terms, such as health, medicine, or disease, by use of Twitter data and did not focus on specific disease topics, such as COVID-19. Another article analysed tweets on COVID-19 and identified 12 topics that were categorised into four main themes: the origin, source, effects on individuals and countries, and methods of decreasing the spread of SARS-CoV-2.⁷ In this study, data were not available for tweets that were related to COVID-19 before February, 2020, thereby missing the initial part of the epidemic, and the data for tweets were limited to between Feb 2 and March 15, 2020.

Social media can also be effectively used to communicate health information to the general public during a pandemic. Emerging infectious diseases, such as COVID-19, almost always result in increased usage and consumption of media of all forms by the general public for information.⁸ Therefore, social media has a crucial role in people's perception of disease exposure, resultant decision making, and risk behaviours.^{9, 10} As information on social media is generated by users, such information can be subjective or inaccurate, and frequently includes misinformation and conspiracy theories.¹¹ Hence, it is imperative that accurate and timely information is disseminated to the general public about emerging threats, such as SARS-CoV-2. A systematic review explored the major

approaches that were used in published research on social media and emerging infectious diseases.¹² The review identified three major approaches: assessment of the public's interest in, and responses to, emerging infectious diseases; examination of organisations' use of social media in communicating emerging infectious diseases; and evaluation of the accuracy of medical information that is related to emerging infectious diseases on social media. However, this review did not focus on studies that used social media data to track and predict outbreaks of emerging infectious diseases.

Analysing and disseminating information from peer-reviewed, published research can guide policy makers and public health agencies to design interventions for accurate and timely knowledge translation to the general public. Therefore, keeping in view the limitations of existing research that we have previously mentioned, we did a scoping review with the aim of understanding the roles that social media has had since the beginning of the COVID-19 crisis. We investigated public attitudes and perceptions towards COVID-19 on social media, information about COVID-19 on social media, use of social media for prediction and detection of COVID-19, the effects of COVID-19 on mental health, and government responses to COVID-19 on social media. Our objective was to identify and analyse studies on social media that were related to COVID-19 and focused on five themes: infodemics, public attitudes, mental health, detection or prediction of COVID-19 cases, government responses to the pandemic, and quality of health information in videos.

2.3 Methods

2.3.1 Overview

Studies exploring the use of social media relating to COVID-19 were reviewed by use of the scoping review methods of Arksey and O'Malley¹³ and Levac and colleagues.¹⁴ We followed the five-step scoping review protocol and the Preferred Reporting Items for Systematic Reviews and Meta-Analyses extension for scoping reviews.

2.3.2 Data Sources

Exploratory searches were done on COVID-19 Open Research Dataset Challenge and Google Scholar in April 2020. These searches helped to define the Review scope, develop the research questions, and determine eligibility criteria. After such activity, MEDLINE and PubMed, Scopus, and PsycINFO

were selected for this Review because they include peer-reviewed literature in the fields of medicine, behavioural sciences, psychology, health-care systems, and clinical sciences. Variations of the key search terms can be found in the panel. Since the start of the current pandemic, COVID-19 articles were reviewed and published at an unprecedentedly rapid rate, with numerous publications that were available ahead of print referred to as preprints or articles in press. In this Review, we consider peer-reviewed preprints to be equivalent to published peer-reviewed articles, and relevant articles were screened accordingly.

2.3.3 Screening procedure

Mainly, the primary reviewer (S-FT) screened title and abstract for each article to decide whether an article met the inclusion criteria. If the criteria were confirmed, then the article was included; otherwise, it was excluded. Paragraphs in articles were assigned a code representing one of the five themes (e.g., I for infodemic), then a code was assigned to the article on the basis of the majority of paragraph codes. Next, quotes were sorted under each code, applying Ose's method.¹⁵ Braun and Clark's thematic analysis method was used and involved searching for the text that matched the identified predictors (i.e., codes) from the quantitative analysis and discovering emergent codes that were relevant to either the study objective or identified in the relevant literature review.¹⁶ Finally, we categorised the codes into main themes. These codes and themes were compared and clarified by S-FT, ZAB, and YY to draw conclusions around the main themes. S-FT is fluent in English and Mandarin. The secondary reviewer (ZAB) is fluent in English, and the tertiary reviewer and domain expert (YY and HC) are both fluent in English and Mandarin. Any discrepancies among reviewers were discussed with the research team to reach consensus.

2.4 Results

With the application of appropriate search filters, a total of 2405 articles were retrieved from the identified databases: PubMed (1,084 articles), Scopus (1021 articles), and PsycINFO (300 articles). Among these, 670 duplicates were excluded. Of the remaining 1,735 articles, 1,434 were deemed to be non-empirical, such as comments, editorial essays, letters, opinions, and reviews. These exclusions left 301 articles for a full-text review on the basis of the screening results of titles and abstracts. After the full-text review, 81 articles were included in this scoping review (Figure 2-1).

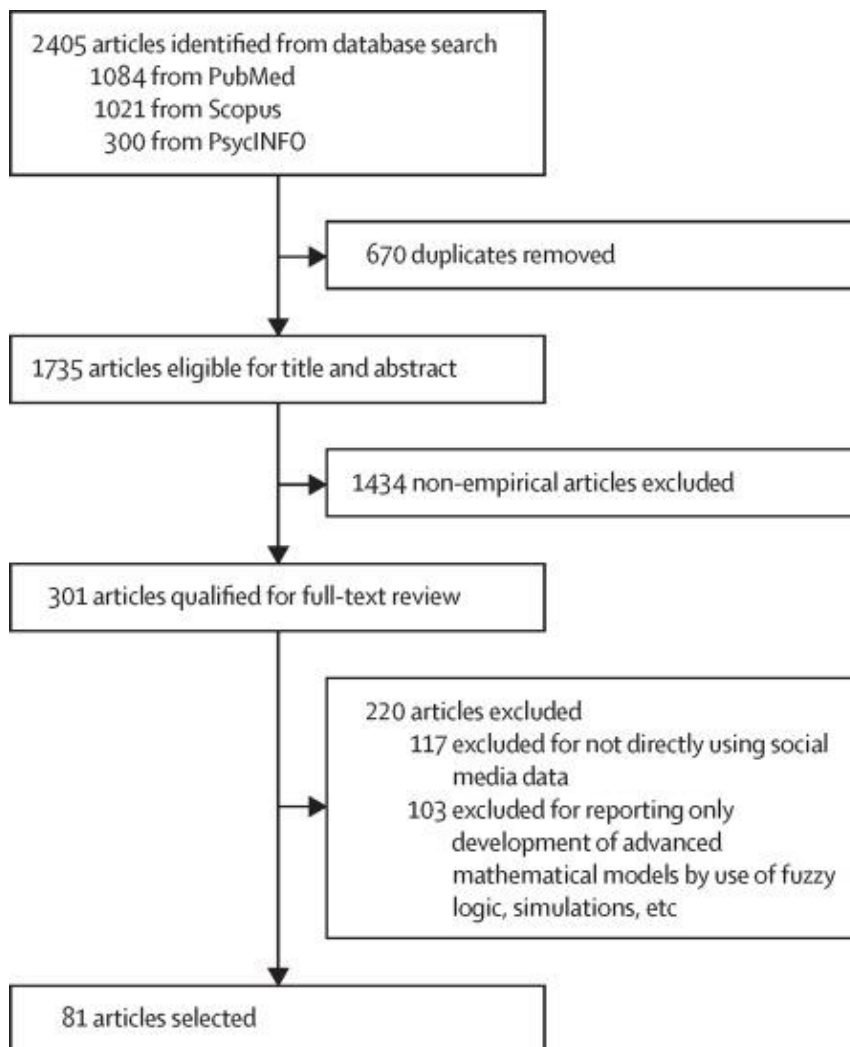


Figure 2-1: Preferred Reporting Items for Systematic Reviews and Meta-Analyses flow chart of article extraction from the literature search

Table 2-1 summarises the 81 articles that were selected on COVID-19 and social media. All articles were written in English. Data from Twitter (45 articles) and Sina Weibo (16 articles) were undoubtedly the most frequently studied. To categorise these chosen articles, we adopted a novel framework called Social Media and Public Health Epidemic and Response (SPHERE) and developed a modified version of SPHERE framework to organise the themes for our scoping review (Figure 2-2).⁹⁸ Themes were identified through reviewers' consensus based on our modified SPHERE

framework. We identified six themes: infodemics, public attitudes, mental health, detecting or predicting COVID-19 cases, government responses, and quality of health information in prevention education videos.

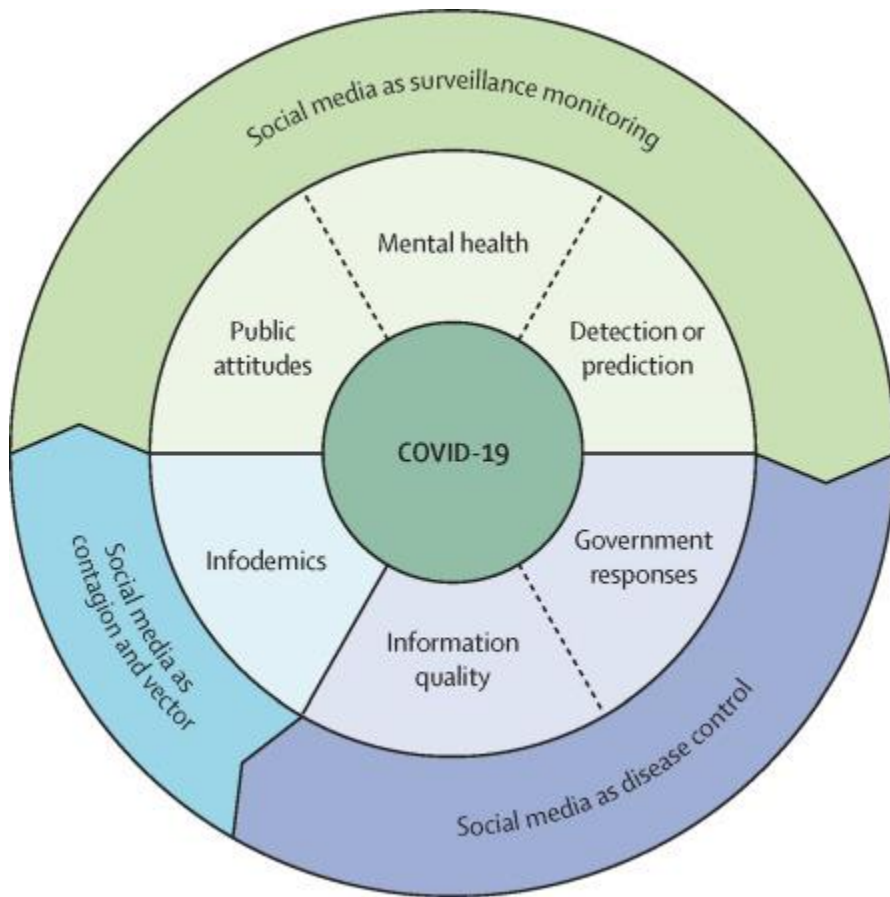


Figure 2-2: Modified Social Media and Public Health Epidemic and Response framework

	Publication month	Origin	Social media	Study population and sample size	Methods	Key findings
Detection or prediction of COVID-19 cases						
Li et al ¹⁷	March	China	Google Trends, Baidu Search Index, and Sina Weibo Index	Keywords of coronavirus and pneumonia were searched, and trend data was collected from Google Trends, Baidu Search Index, and Sina Weibo Index from Jan 2 to Feb 20, 2020	Lag correlation	Lag correlations showed a maximum correlation between trend data and the number of diagnoses at 8–12 days before for laboratory-confirmed cases and 6–8 days before for suspected cases
Liu et al ¹⁸	August	China	Sina Weibo	Sina Weibo messages between Jan 20 and Feb 15, 2020; 599 participants	Gathered data via Sina Weibo, then followed up with telephone call; statistical analysis taken with Fisher exact test; rates of death calculated with Kaplan-Meier method; multivariate Cox regression used to establish risk factors for mortality	Older age (i.e., >69 years), diffuse pneumonia, and hypoxaemia are factors that can help clinicians to identify patients with COVID-19 who have poor prognosis; aggregated data from social media can also be comprehensive, immediate, and informative in disease prognosis
O'Leary and Storey ¹⁹	September	USA	Google Trends, Wikipedia, and Twitter	Google Trends searches for coronavirus and COVID-19 between Jan 21 and April 5, 2020;	Regression analysis	To model the number of cases, the current Wikipedia page views, tweets from 1 week before, and Google Trends

				Wikipedia page views for coronavirus and COVID-19 between Jan 12 and April 5, 2020; number of Twitter original tweets between Jan 27 and April 5, 2020; numbers of COVID-19 cases and deaths in the USA ²⁰		searches from 2 weeks before were used; to model of the number of deaths, each variable was taken from 1 week earlier than for cases
Peng et al ²¹	June	China	Sina Weibo	1,200 records	Spatiotemporal distribution of COVID-19 cases in the main urban area of Wuhan, China; kernel density analysis; ordinary least square regression	Older people (i.e., >60 years) are at high risk of severe symptoms and have high prevalence in the COVID-19 outbreak, and they account for >50% of the total number of Sina Weibo help seekers; early transmission of COVID-19 in Wuhan, China, could be divided into three phrases: scattered infection, community spread, and full-scale outbreak
Qin et al ²²	March	China	Baidu Search Index	Social media search index for dry cough, fever, chest distress, coronavirus, and pneumonia from Dec 31, 2019, to Feb 9, 2020; data for new suspected	Subset selection; forward selection; lasso regression; ridge regression; elastic net	Case numbers of new suspected COVID-19 correlated significantly with the lagged series of social media search index; social media search index could detect new suspected COVID-19 cases 6–9

				cases of COVID-19 from Jan 20 to Feb 9, 2020		days earlier than could laboratories
Zhu et al ²³	April	China	Sina Weibo	1,101 Sina Weibo posts related to COVID-19 from Dec 31, 2019, to Feb 12, 2020	Descriptive statistics: numbers and percentage; time series analysis	Attention to COVID-19 was low until China openly admitted human-to-human transmission on Jan 20, 2020; attention quickly increased and remained high over time
Government responses						
Basch et al ²⁴	April	USA	YouTube	100 most widely viewed videos uploaded in January, 2020	Descriptive analysis: frequency, percentage, mean, and standard deviation	Percentage of each of the seven key prevention behaviours that are listed on the US Centers for Disease Control and Prevention website that were covered in the 100 videos varied from 0% (eg, use a face mask for protection if you are caring for the ill) to 31% (avoid close contact with people who are sick); overall, videos that covered at least one prevention behaviour accounted for less than one-third of the 100 videos
Basch et al ²⁵	April	USA	YouTube	100 most widely viewed YouTube videos as of Jan 31, 2020, and March 20, 2020, with keyword of coronavirus in English, with English subtitles, or in Spanish	Descriptive analysis: frequency, percentage, mean, and standard deviation	<50% of videos in either sample covered any of the prevention behaviours that are recommended by the US Centers for Disease Control and Prevention

Khatri et al ²⁶	March	Singapore	YouTube	150 videos collected on Feb 1–2, 2020, with keywords of 2019 novel coronavirus (50 videos), and Wuhan virus in English (50 videos) and Mandarin (50 videos)	Descriptive analysis: percentage and mean; DISCERN score; Medical Information and Content Index score	Mean DISCERN score for reliability was 3.12 of 5.00 for English and 3.25 of 5.00 for Mandarin videos; mean cumulative Medical Information and Content Index score of useful videos was 6.71 of 25.00 for English and 6.28 of 25.00 for Mandarin
Li et al ²⁷	March	China	Sina Weibo	36,746 Sina Weibo data from Dec 30, 2019, to Feb 1, 2020; a random sample of 3000 Sina Weibo posts as training dataset	Linear regression; support vector machine; Naive Bayes; natural language processing	Classified the information related to COVID-19 into seven types of situational information and their predictors
Merkley et al ²⁸	April	Canada	Twitter and Google Trends	33,142 tweets from 292 social media accounts of federal members of parliament from Jan 1 to March 28, 2020; 87 Google search trends for the search term coronavirus in the first half (i.e., days 1–14) and second half (i.e., days 15–31) of March 2020; a survey of 2499 Canadian citizens ≥ 18 years from April 2 to April 6, 2020	Linear regression	No members of parliament from any party downplaying the pandemic; no association between Conservative Party vote share and Google search interest in the coronavirus

Rufai and Bunce ²⁹	April	USA	Twitter	203 viral tweets from G7 world leaders from Nov 17, 2019, to March 17, 2020 with keywords COVID-19 or coronavirus and a minimum of 500 likes	Qualitative design; content analysis	166 of 203 of tweets were informative; 9.4% (19) were morale-boosting; 6.9% (14) were political
Sutton et al ³⁰	September	USA	Twitter	690 accounts representing public health, emergency management, and elected officials and 149,335 tweets	χ^2 analyses; negative binomial regression modelling	Systematic changes were made in message strategies over time and identified key features that affect message passing, both positively and negatively; results have the potential to aid in message design strategies as the pandemic continues, or in similar future events
Wang et al ³¹	September	USA	Twitter	13,598 tweets related to COVID-19 from Jan 1 to April 27, 2020	Temporal analysis and networking analysis	16 categories of message types were manually annotated; inconsistencies and incongruencies were identified in four critical topics (i.e., wearing masks, assessment of risks, stay at home order, and disinfectant and sanitizer); network analysis showed increased communication coordination over time
Infodemics						

Ahmed et al ³²	October	UK	Twitter	22,785 tweets and 11,333 Twitter users with #FilmYourHospital from April 13 to April 20, 2020	Social network analysis; user analysis	The most important drivers of the #FilmYourHospital conspiracy theory are ordinary citizens; YouTube was the information source most linked to by users; the most retweeted post belonged to a verified Twitter user
Ahmed et al ³³	May	UK	Twitter	A subsample of 233 tweets from 10 140 tweets collected from 19:44 h UTC on Friday, March 27, 2020, to 10:38 h UTC on Saturday, April 4, 2020, were used for content analysis	Descriptive statistics: numbers, percentage; social network analysis; content analysis	34.8% (81 of 233) of tweets linked 5G and COVID-19; 32.2% (75) of tweets denounced the conspiracy theory
Brennen et al ³⁴	October	UK	Digital visual media	96 samples of visuals from January to March, 2020	Qualitative coding	Organised all findings into six trends: authoritative agency, virulence, medical efficacy, intolerance, prophecy, satire; a small number of manipulated visuals, all were produced by use of simple tools; no examples of so-called deepfakes (i.e., techniques that are used to make synthetic videos that closely resemble real videos) or other techniques

						that were based on artificial intelligence
Bruns et al ³⁵	August	Australia	Facebook	89,664 distinct Facebook posts from Jan 1 to April 12, 2020	Time series; network analysis	Substantially increased number of posts about 5G rumours on Facebook after March 19, 2020; network analysis showed that coalitions of various groups were brought together by conspiracy theories about COVID-19 and 5G technology
Galhardi et al ³⁶	October	Brazil	WhatsApp, Instagram, and Facebook	Fake news collected from March 17 to April 10, 2020, on the basis of data from the Eu Fiscalizo app (version 5.0.5)	Quantitative content analysis	WhatsApp is the main channel for sharing fake news, followed by Instagram and Facebook
Gallotti et al ³⁷	October	Italy	Twitter	>100 million Tweets	Developed an Infodemic Risk Index	Before the rise of COVID-19 cases, entire countries had measurable waves of potentially unreliable information, posing a serious threat to public health
Islam et al ³⁸	October	Bangladesh	Fact-checking agency websites, Facebook, Twitter, and websites for television networks and newspapers	2,311 infodemic reports related to COVID-19 between Dec 31, 2019, and April 5, 2020	Descriptive analysis; spatial distribution analysis	Misinformation that is fuelled by rumours, stigma, and conspiracy theories can have potentially severe implications on public health if prioritised over scientific guidelines; governments and other agencies should understand the patterns of rumours, stigma, and

						conspiracy theories that are related to COVID-19 and circulating globally so that they can develop appropriate messages for risk communication
Kouzy et al ³⁹	March	Lebanon	Twitter	673 English tweets collected on Feb 27, 2020; 617 tweets after exclusion of tweets that were humorous or not serious	Descriptive statistics; bar chart; χ^2 statistic to calculate p value (2-sided; $p=0.05$ significance threshold) for the association between account or tweet characteristics and the presence of misinformation or unverifiable information about COVID-19	153 (24.8%) of 617 tweets had misinformation; 107 (17.3%) had unverifiable information; misinformation rate higher in informal individual or group accounts than in formal individual or group accounts (33.8% [123 of 364] vs 15.0% [30 of 200], $p<0.001$)
Moscadelli et al ⁴⁰	August	Italy	Fake news and corresponding verified news that was circulated in Italy	2,102 articles between Dec 31, 2019, and April 30, 2020	Social media trend analysis by use of BuzzSumo	Links containing fake news were shared 2,352,585 times, accounting for 23.1% (2,352,585 of 10,184,351) of total shares of all reviewed articles

Pulido et al ⁴¹	April	Spain	Twitter	942 valid tweets between Feb 6 and Feb 7, 2020	Communicative content analysis	Misinformation was tweeted more but retweeted less than tweets based on scientific evidence; tweets based on scientific evidence had more engagement than misinformation
Rovetta and Bhagavathula ⁴²	August	Italy	Google Trends and Instagram	2 million Google Trends queries and Instagram hashtags from Feb 20 to May 6, 2020	Classification of infodemic monikers (i.e., a term, query, hashtag, or phrase that generates or feeds fake news, misinterpretations, or discrimination); computed the mean peak volume with a 95% CI	Globally, growing interest exists in COVID-19, and numerous infodemic monikers continue to circulate on the internet
Uyheng and Carley ⁴³	October	USA and Philippines	Twitter	12 million tweets from 1.6 million users from the USA and 15 million tweets from 1 million users from the Philippines between March 5 and March 19, 2020	Hate speech score assigned to each tweet by use of machine learning algorithm; bot scores were assigned to each user via BotHunter algorithm; social	Analysis showed idiosyncratic relationships between bots and hate speech across datasets, emphasising different network dynamics of racially charged toxicity in the USA and political conflicts in the Philippines; bot activity is linked to hate in both countries, especially in communities that

					media analysis via ORA software; network analysis via centrality analysis; cluster analysis via Leiden algorithm	are dense and isolated from others
Mental health						
Gao et al ⁴⁴	April	China	Sina Weibo	Online survey on Wenjuanxing platform from Jan 31 to Feb 2, 2020; with 4872 Chinese citizens aged ≥ 18 years from 31 provinces and autonomous regions in China	Multivariable logistic regression	Social media exposure was frequently positively associated with high odds of anxiety (odds ratio 1.72, 95% CI 1.31–2.26) and combination of depression and anxiety (odds ratio 1.91, 95% CI 1.52–2.41)
Li et al ⁴⁵	March	China	Sina Weibo	Sina Weibo posts from 17,865 active Sina Weibo users between Jan 13 and Jan 26, 2020	Sentiment analysis; paired sample t-test	Negative emotions and sensitivity to social risks increased; scores of positive emotions and life satisfaction decreased after outbreak declaration
Prevention education in videos						
Hakimi and Armstrong ⁴⁶	September	USA	YouTube	49 of the first 100 videos on YouTube with the most views that were identified by the search term DIY hand sanitizer; 51 videos were excluded because they were not in	Codified video content; assessed by use of Cohen's κ ; descriptive statistics calculated; assessed by χ^2 test	Most videos did not describe labelling storage containers, 69% (34 of 49) of videos encouraged the use of oils or perfumes to enhance hand sanitizer scent, and 2% (1) of videos promoted the use of

				English or not related to the search term	with 2-sided p value <0.05 as the threshold for significance	colouring agents to be more attractive for use among children specifically; significantly increased mean number of daily calls to poison control centres regarding unsafe paediatric exposure to hand sanitiser since the first confirmed patient with COVID-19 in the USA (p<0.001); significantly increased mean number of daily calls in March, 2020, compared with the previous 2 years (p<0.001)
Hernández-García and Giménez-Júlvez ⁴⁷	June	Spain	YouTube	129 videos in Spanish with the terms “prevention,” “coronavirus,” and “prevention COVID19”	Univariate analysis; multiple logistic regression model	Information from YouTube in Spanish on basic measures to prevent COVID-19 is usually not complete and differs according to the type of authorship (i.e., mass media, health professionals, individual users, or others)
Moon and Lee ⁴⁸	August	South Korea	YouTube	105 most viewed YouTube videos from Jan 1 to April 30, 2020	Modified DISCERN index; Journal of the American Medical Association Score benchmark criteria; Global Quality Score;	37.14% (39 of 105) of videos contained misleading information; independent user-generated videos showed the highest proportion of misleading information at 68.09% (32 of 47); misleading videos had more likes, fewer

					Title–Content Consistency Index; Medical Information and Content Index	comments, and longer running times than did useful videos; transmission and precautionary measures were the most frequently covered content
Ozdede and Peker ⁴⁹	July–August	Turkey	YouTube	The top 116 English language videos with at least 300 views	Precision indices and total video information and quality index scores were calculated	High number of views on dentistry YouTube videos related to COVID-19; quality and usefulness of these videos are moderate
Yüce et al ⁵⁰	July	Turkey	YouTube	55 English videos about COVID-19 control procedures for dental practices collected on March 31, 2020, between 9:00 h and 18:00 h	Modified DISCERN instrument; descriptive statistics	Only two (3.6%) of 55 videos were good quality, whereas 24 (43.6%) videos were poor quality
Public attitudes						
Abd-Alrazaq et al ⁷	April	Qatar	Twitter	2.8 million English tweets (167,073 unique tweets from 160,829 unique users) from Feb 2 to March 15, 2020	Word frequencies of single (i.e., unigrams) and double words (i.e., bigrams); sentiment analysis; mean number of retweets, likes, and followers for each topic;	Identified 12 topics and grouped into four themes; average sentiment positive for ten topics and negative for two topics

					interaction rate per topic; LDA for topic modelling	
Al-Rawi et al ⁵¹	November	Canada	Twitter	Over 50 million tweets referencing #Covid-19 and #Covid19 for more than 2 months in early 2020	Mixed method: analysed emoji use by each gender category; the top 600 emojis were manually classified on the basis of their sentiment	Identified five major themes in the analysis: morbidity fears, health concerns, employment and financial issues, praise for front-line workers, and unique gendered emoji use; most emojis are extremely positive across genders, but discussions by women and gender minorities are more negative than by men; when discussing particular topics (e.g., financial and employment matters, gratitude, and health care), there are many differences; use of several unique gender emojis to express specific issues (e.g., coffin, skull, and siren emojis were used more often by men than by other genders when discussing fears and morbidity, whereas the use of the folded hands emoji as a thankful gesture for front-line workers was found more often in discussions by women than by

						other genders and the bank emoji was noted only in women's discussions)
Arpaci et al ⁵²	July	Turkey	Twitter	43 million tweets between March 22 and March 30, 2020	Evolutionary clustering analysis	Unigram terms appear more frequently than bigram and trigram (i.e., triple words) terms; during the epidemic, many tweets about COVID-19 were distributed and attracted widespread public attention; high-frequency words (e.g., death, test, spread, and lockdown) indicated that people were afraid of being infected and people who were infected were afraid of death; people agreed to stay at home due to fear of spread and called for physical distancing since they became aware of COVID-19
Barrett et al ⁵³	August	USA	Twitter	188 tweets about Governor Dan Patrick's statement on March 23, 2020, about generational self-sacrifice.	Thematic analysis	90% (169 of 188) of tweets opposed calculated ageism, whereas only 5% (9) supported it and 5% (10) conveyed no position; opposition centred on moral critiques, political-economic critiques, assertions of the worth of older adults (e.g., >60 years), and public health arguments; support

						centred on individual responsibility and patriotism
Boon-Itt and Skunkan ⁵⁴	November	Thailand	Twitter	107,990 English tweets related to COVID-19 between Dec 13, 2019, and March 9, 2020	Sentiment analysis; topic modelling by use of LDA	Sentiment analysis showed a predominantly negative feeling towards the COVID-19 pandemic; topic modelling revealed three themes relating to COVID-19 and the outbreak: the COVID-19 pandemic emergency, how to control COVID-19, and reports on COVID-19
Budhwani and Sun ⁵⁵	May	USA	Twitter	16,535 tweets about Chinese virus or China virus between March 9 and March 15, 2020, 177,327 tweets between March 19 and March 25, 2020	Descriptive analysis; spatial analysis	Nearly 10 times increase at the national level; all 50 states had an increase in the number of tweets exclusively mentioning Chinese virus or China virus instead of coronavirus disease, COVID-19, or coronavirus; mean 0.38 tweets referencing Chinese virus or China virus were posted per 10,000 people at the state level in the pre-period (i.e., March 9–15, 2020), and 4.08 of these stigmatising tweets were posted in the post-period (i.e., March 19–25, 2020), also indicating a 10 times increase

Chang et al ⁵⁶	November	Taiwan	10 news websites, 11 discussion forums, 1 social network, 2 principal media sharing networks	1.07 million Chinese texts from Dec 30, 2019, to March 31, 2020	Deductive analysis	Online news promoted negativity and drove emotional social posts; stigmatising language that was linked to the COVID-19 pandemic showed an absence of civic responsibility that encouraged bias, hostility, and discrimination
Chehal et al ⁵⁷	July	India	Twitter	29,554 tweets during the second lockdown (i.e., April 15–May 3, 2020); 47,672 tweets during the third lockdown (May 4–17, 2020)	Sentiment analysis by use of the National Research Council of Canada Emotion Lexicon	A positive approach in the second lockdown but a negative approach in the third lockdown
Chen et al ⁵⁸	September	China	Sina Weibo	1,411 posts pertinent to COVID-19 taken from Healthy China, an official Sina Weibo account of the National Health Commission of China, from Jan 14 to March 5, 2020	Descriptive analysis; hypothesis testing	Media richness (i.e., potential information load, where low richness is only text and high richness is not only text) negatively predicted citizen participation via government social media, but dialogic loop (i.e., stimulation of public dialogue, provision of the dialogue channel, and response to public feedback in a timely manner) facilitated engagement
Damiano and Allen Catellier ⁵⁹	August	USA	Twitter	600 English tweets from the USA were selected: 300 from February 2020,	Frequencies; χ^2 statistics	Neutral sentiment; tweets about COVID-19 risks and emotional outrage accounted for <50%

				and 300 from March, 2020		(135 of 600); few tweets were related to blame
Darling-Hammond et al ⁶⁰	September	USA	Twitter	339,063 tweets from non Asian respondents of the Project Implicit Asian Implicit Association Test from 2007–20 and were broken into two datasets: the first dataset was from Jan 1, 2007, to Feb 10, 2020; the second data set was from Feb 11 to March 31, 2020	Local polynomial regression; interrupted time-series analyses	Implicit Americanness Bias steadily decreased from 2007 to 2020; when media entities began using stigmatising terms, such as Chinese virus, starting from March 8, 2020, Implicit Americanness Bias began to increase; such bias was more pronounced among conservative individuals than among non-conservative individuals
Das and Dutta ⁶¹	July	India	Twitter	410,643 tweets with #IndiaLockdown and #IndiafightsCorona from March 22 to April 21, 2020	National Research Council of Canada lexicon for corpus-level emotion mining; sentimentr from open-source R software for sentiment analysis to create additional sentiment scores; LDA for topic models; Natural Language Toolkit to develop	For the broad corpus-level analysis, the context of positiveness was substantially higher than were negative sentiments; however, positive sentiment trends were similar to negative sentiment trends in terms of topics covered when the analysis was done at individual tweet level; the results showed that the discussion of COVID-19 in India on Twitter contains slightly more positive sentiments than negative sentiments

					sentiment-based topic models	
De Santis et al ⁶²	July	Italy	Twitter	1,044,645 tweets	A general purpose methodological framework, grounded on a biological metaphor and on a chain of NLP and graph analysis techniques	Energy evolution through time was monitored; daily hot topics were identified (e.g., COVID-19, Walter Ricciardi's retweet of an anti-Trump tweet from Michael Moore, Gabriele Gravina's argument against suspension of Italian football, increased COVID-19 cases in Italy, high case numbers in Lombardy, Italy, and an interview of Matteo Salvini about COVID-19 topics by Massimo Giletti)
Dheeraj ⁶³	May–June	India	Reddit	868 posts related to COVID-19	Fetching the articles: Python Reddit Application Programming Interface Wrapper; data preprocessing: Reddit Application Programming Interface and Natural Language Toolkit library	Of 868 posts on Reddit that were related to COVID-19 articles, 50% (434) were neutral, 22% (191) were positive, and 28% (243) were negative

Essam and Abdo ⁶⁴	August	Egypt	Twitter	1,920,593 tweets with corona, coronavirus, or COVID-19 keywords from Feb 1 to April 30, 2020	Thematic analysis	The dominant themes that were closely related to coronavirus tweets included the outbreak of the pandemic, metaphysics responses, signs and symptoms in confirmed cases, and conspiracies; the psycholinguistic analysis showed that tweeters maintained high amounts of affective talk (i.e., expression of feelings), which was loaded with negative emotions and sadness; Linguistic Inquiry and Word Count's psychological categories of religion and health dominated the Arabic tweets discussing the pandemic situation
Yin FL et al ⁶⁵	March	China	Sina Weibo	Sina Weibo posts from Dec 31, 2019, to Feb 7, 2020	Multiple-information susceptible-discussing-immune model	Model reproduction ratio declined from 1.78 to 0.97, showing that the peak of posts had passed but the topic was still on social media afterwards with a decreased number of posts
Gozzi et al ⁶⁶	October	Italy, UK, USA, and Canada	News, YouTube, Reddit, and Wikipedia	227,768 web-based news articles from Feb 7 to May 15, 2020; 13,448 YouTube videos from	Linear regression; topic modelling by use of LDA	Collective attention was mainly driven by media coverage rather than epidemic progression, rapidly became saturated, and

				Feb 7 to May 15, 2020; 107,898 English user posts and 3,829,309 comments on Reddit from Feb 15 to May 15, 2020; 278,456,892 views of Wikipedia pages that were related to COVID-19 from Feb 7 to May 15, 2020		decreased despite media coverage and COVID-19 incidence remaining high; Reddit users were generally more interested in health, data regarding the new disease, and interventions needed to halt the spreading with respect to media exposure than were users of other platforms
Green et al ⁶⁷	July	USA	Twitter	19,803 tweets from Democrats and 11,084 tweets from Republicans between Jan 17 and March 31, 2020	Random forest	Democrats discussed the crisis more frequently—emphasising public health and direct aid to US workers—whereas Republicans placed greater emphasis on national unit, China, and businesses
Han et al ⁶⁸	April	China	Sina Weibo	1,413,297 Sina Weibo messages, including 105,330 texts with geographical location information, from 00:00 h on Jan 9, 2020, to 00:00 h on Feb 11, 2020	Time series analysis; kernel density estimation; Spearman correlation; LDA model; random forest algorithm	Public response was sensitive to the epidemic and notable social events, especially in urban agglomerations
Jelodar et al ⁶⁹	June	China	Reddit	563,079 English comments related to COVID-19 from Reddit between Jan 20 and March 19, 2020	Topic modelling by use of LDA and probabilistic latent semantic analysis;	The results showed a novel application for NLP based on a long short term memory model to detect meaningful latent topics and sentiment–comment

					sentiment classification by use of recurrent neural network	classification on issues related to COVID-19 on social media
Jimenez-Sotomayor et al ⁷⁰	April	Mexico	Twitter	A random sample of 351 of 18,128 tweets were analysed from March 12 to March 21, 2020	Qualitative content classification	The most common types of tweets were personal opinions (31.9% [112 of 351]), followed by informative tweets (29.6% [104]), jokes or ridicule (14.2% [50]), and personal accounts (13.4% [47]); 72 of 351 tweets were most likely intended to ridicule or offend someone and 21.1% (74) had content implying that the life of older adults (i.e., referred to in tweets as “elderly”, “older”, and “boomer”) was less valuable than that of younger people or downplayed the relevance of COVID-19
Kim ⁷¹	August	South Korea	Twitter	27,849 individual tweets about COVID-19 between Feb 10 and Feb 14, 2020	27 849 individual tweets about COVID-19 between Feb 10 and Feb 14, 2020	Social network size was a negative predictor of incivility
Kurten and Beullens ⁷²	August	Belgium	Twitter	373,908 tweets and retweets from Feb 25 to March 30, 2020	Time series; network bigrams; emotion lexicon; LDA	Notable COVID-19 events immediately increased the number of tweets; most topics focused on the need for EU

						collaboration to tackle the pandemic
Kwon et al ⁷³	October	USA	Twitter	259,529 unique tweets containing the word coronavirus between Jan 23 and March 24, 2020	Trending analysis; spatiotemporal analysis	Early facets of physical distancing appeared in Los Angeles (CA, USA), San Francisco (CA, USA), and Seattle (WA, USA); social disruptiveness tweets were most retweeted, and intervention implementation tweets were most favoured
Lai et al ⁷⁴	October	USA	Reddit	522 comments from an Ask Me Anything session on COVID-19 on March 11, 2020, from 14:00 h to 16:00 h EST	Content analysis	The highest number of posts were about symptoms (27% [141 of 522]), followed by prevention (25% [131]); symptoms was the most common intended topic for further discussions (28% [94 of 337])
Li et al ⁷⁵	April	China	Sina Weibo	115,299 Sina Weibo posts from Dec 23, 2019, to Jan 30, 2020; 11,893 of them were collected from Dec 31, 2019, to Jan 20, 2020, for qualitative analysis; total daily cases of COVID-19 in Wuhan, China, were obtained from the	Linear regression model; qualitative content analysis	Positive correlation between the number of Sina Weibo posts and the number of reported cases, with ten COVID-19 cases per 40 posts; posts grouped into four themes

				Chinese National Health Commission		
Li et al ⁷⁶	September	USA	Twitter	155,353 unique English tweets related to COVID-19 that were posted from Dec 31, 2019, to March 13, 2020	Content analysis	Peril of COVID-19 was mentioned the most often, followed by content about marks (i.e., cues to identify members of a stigmatised group: flu-like symptoms, personal protective equipment, Asian origin, and health-care providers and essential workers), responsibility, and group labelling; information on conspiracy theories was more likely to be included in tweets about group labelling and responsibility than in tweets about COVID-19 peril
Lwin et al ⁷⁷	May	Singapore	Twitter	20,325,929 tweets from 7,033,158 unique users from Jan 28 to April 9, 2020	Sentiment analysis	Public emotions shifted strongly from fear to anger over the course of the pandemic, while sadness and joy also surfaced; anger shifted from xenophobia at the beginning of the pandemic to discourse around the stay-at-home notices; sadness was emphasised by the topics of losing friends and family members, whereas topics that

						were related to joy included words of gratitude and good health; emotion-driven collective issues around shared public distress experiences of the COVID-19 pandemic are developing and include large-scale social isolation and the loss of human lives
Ma et al ⁷⁸	July	China	WeChat	Top 200 accounts from Jan 21 to Jan 27, 2020	Simple linear regression; multiple linear regression; content analysis	For non-medical institution accounts in the model, report and story types of articles had positive effects on whether users followed behaviours; for medical institution accounts, report and science types of articles had a positive effect
Medford et al ⁷⁹	June	USA	Twitter	126,049 English tweets from 53 196 unique users with matching hashtags that were related to COVID-19 from Jan 14 to Jan 28, 2020	Temporal analysis; sentiment analysis; topic modelling by use of LDA	The hourly number of tweets that were related to COVID-19 starkly increased from Jan 21, 2020, onwards; fear was the most common emotion and was expressed in 49.5% (62,424 of 126,049) of all tweets; the most common predominant topic was the economic and political effect
Mohamad ⁸⁰	June	Brunei	Twitter, Instagram, and TikTok	30 individual profiles from Instagram, Twitter, and TikTok	Qualitative content analysis	Five narratives of local responses to physical distancing practices were apparent: fear,

						responsibility, annoyance, fun, and resistance
Nguyen et al ⁸¹	September	USA	Twitter	3,377,295 US tweets that were related to race from November 2019, to June, 2020	Support vector machine was used for sentiment analysis	Proportion of negative tweets referencing Asians increased by 68.4%; proportion of negative tweets referencing other racial or ethnic minorities was stable; common themes that emerged during the content analysis of a random subsample of 3,300 tweets included: racism and blame, anti-racism, and effect on daily life
Odlum et al ⁸²	June	USA	Twitter	2,558,474 Tweets from Jan 21 to May 3, 2020	Clustering algorithm; NLP; network diagrams	15 topics (in four themes) were identified; positive sentiments, cohesively encouraging online discussions, and behaviours for COVID-19 prevention were uniquely observed in African American Twitter communities
Park et al ⁸³	May	South Korea	Twitter	43,832 unique users and 78,233 relationships on Feb 29, 2020	Network analysis; content analysis	Spread of information was faster in the COVID-19 network than in the other networks; tweets containing medically framed news articles were more popular than were tweets that included news articles adopting non-medical frames

Pastor ⁸⁴	April	Philippines	Twitter	Tweets were collected on three Tuesdays in March 2020, since lockdown in Philippines	NLP for sentiment analysis	Negative sentiments increased over time in lockdown
Samuel et al ⁸⁵	June	USA	Twitter	900,000 tweets from February to March, 2020	Sentiment analysis packages; textual analytics; machine learning classification methods: Naive Bayes and logistic regression	For short tweets, classification accuracy was 91% with Naive Bayes whereas accuracy was 74% with logistic regression; both methods showed weaker performance for longer tweets
Samuel et al ⁸⁶	August	USA	Twitter	293,597 tweets, 90 variables	Textual analytics to analyse public sentiment support; sentiment analysis by use of R package Syuzhet (version 1.0.6)	For the reopening of the US economy, there was more positive sentiment support than there was negative support; developed a novel sentiment polarity based public sentiment scenarios framework
Su et al ⁸⁷	June	China and Italy	Sina Weibo and Twitter	850 Sina Weibo users with posts published from Jan 9 to Feb 5, 2020; 14,269 tweets from 188 unique Twitter users from Feb 23 to March 21, 2020	Wilcoxon tests	Individuals focused more on home and expressed a high level of cognitive process after a lockdown in both Wuhan, China, and Lombardy, Italy; level of stress decreased, and the attention to leisure increased in Lombardy, Italy, after the lockdown; attention to group, religion, and emotions became

						more prevalent in Wuhan, China, after the lockdown
Thelwall and Thelwall ⁸⁸	May	UK	Twitter	3,038,026 English tweets from March 10 to March 23, 2020	Word frequency comparison; χ^2 analysis	Women were more likely to tweet about the virus in the context of family, physical distancing, and health care, whereas men were more likely to tweet about sports cancellations, the global spread of the virus, and political reactions
Wang et al ⁸⁹	July	China	Sina Weibo	999,978 randomly selected Sina Weibo posts that were related to COVID-19 from Jan 1 to Feb 18, 2020	Unsupervised Bidirectional Encoder Representations from Transformers model: classify sentiment categories; Term Frequency-Inverse Document Frequency model: summarise the topics of posts; trend analysis; thematic analysis	People were concerned about four aspects regarding COVID-19: the virus origin, symptoms, production activity, and public health control
Wicke and Bolognesi ⁹⁰	September	Ireland	Twitter	203,756 tweets	Topic modelling	Although the family frame covers a wider portion of topics, among the figurative frames,

						war (a highly conventional one) was the frame used most frequently; yet, this frame does not seem to be appropriate to elaborate the discourse around some aspects that are involved in the situation
Xi et al ⁹¹	September	China	Sina Weibo	188 unique topics, their views, and comments from Jan 20 to April 28, 2020	Thematic analysis; temporal analysis	Six themes were identified: the most prominent theme was older people contributing to the community (46 [24%] of 188) followed by older patients (defined by keywords—e.g., “older people”, “old-aged people”, “grandmother”, “grandfather”, “old grandmother”, “old grandfather”, “old woman”, and “old man”) in hospitals (43 [23%]); the theme of contributing to the community was the most dominant in the first phase (Jan 20–Feb 20, 2020; period of COVID-19 outbreak in China); the theme of older patients in hospitals was most dominant in the second (Feb 21–March 17, 2020; turnover period) and third phase (March 18–April 28,

						2020; post-peak period in China)
Xie et al ⁹²	August	China	Baidu Search Index and Google Trends	Number of cases by Feb 29, 2020: 79,968 cumulative confirmed cases, 41,675 cured cases, 2,873 dead cases	Kendall's T _b rank test	Both the Baidu Search Index and Google Trends indices showed a similar trend in a slightly different way; daily Google Trends were correlated to seven indicators, whereas daily Baidu Search Index was correlated to only three indicators; these indexes and rumours are statistically related to disease-related indicators; information symmetry was also noted
Xue et al ⁹³	November	Canada	Twitter	1,015,874 tweets from April 12 to July 16, 2020	LDA	Nine themes about family violence were identified
Yigitcanlar et al ⁹⁴	October	Australia	Twitter	96,666 tweets from Australia in Jan 1 to May 4, 2020	Descriptive analysis; content analysis; sentiment analysis; spatial analysis	Social media analytics is an efficient approach to capture attitudes and perceptions of the public during a pandemic; crowdsourced social media data can guide interventions and decisions of the authorities during a pandemic; effective use of government social media channels can help the public to follow the introduced measures and restrictions

Yu et al ⁹⁵	July	Spain	Twitter	22,223 tweets	Topic modelling; network analysis	Identified eight news frames for each newspaper's Twitter account; the entire pandemic development process is divided into three periods: precrisis, lockdown, and recovery period; understanding of how Spanish news media cover public health crises on social media platforms
Zhao et al ⁹⁶	May	China	Sina Weibo and microblog hot search list	4,056 topics from Dec 31, 2019, to Feb 20, 2020	Word segmentation; word frequency; sentiment analysis	The trend of public attention could be divided into three stages; the hot topic keywords of public attention at each stage were slightly different; the emotional tendency of the public towards the COVID-19 pandemic-related hot topics changed from negative to neutral between January and February, 2020, with negative emotions weakening and positive emotions increasing overall; COVID-19 topics with the most public concern were divided into five categories: the situation of the new cases of COVID-19 and its effects, front-line reporting of the pandemic and the measures of prevention and control, expert

						interpretation and discussion on the source of infection, medical services on the front line of the pandemic, and focus on the pandemic and the search for suspected cases
Zhu et al ⁹⁷	July	China	Sina Weibo	1,858,288 microblog data	LDA	A so-called double peaks feature appeared in the search curve for epidemic topics; the topic changed over time, the fluctuation of topic discussion rate gradually decreased; political and economic centres attracted high attention on social media; the existence of the subject of rumours enabled people to have more communication and discussion

All studies were published in 2020. LDA=latent Dirichlet allocation. NLP=natural language processing.

Table 2-1: Summary of chosen article

2.4.1 Social media as contagion and vector

According to WHO, the term infodemic, a combination of information and epidemic, refers to a fast and widespread dissemination of both accurate and inaccurate information about an epidemic, such as COVID-19.⁹⁹ 12 articles studied infodemics that were related to COVID-19 that were circulating on social media platforms. Rovetta and Bhagavathula⁴² analysed over 2 million queries from Google Trends and Instagram between Feb 20 and May 6, 2020. Their findings showed that as global interest for COVID-19 information increased, so did its infodemic.⁴² Gallotti and colleagues analysed over 100 million tweets and identified that, even before the onset of the COVID-19 pandemic, infodemics threatened public health, although not to the same extent.³⁷ Pulido and colleagues sampled and analysed 942 tweets, which revealed that although false information had a higher number of tweets, it also had less retweets and lower engagement than did tweets comprising scientific evidence or factual statements.⁴¹ Kouzy and colleagues³⁹ investigated the extent to which misinformation or unverifiable information about the COVID-19 pandemic was spread on Twitter by analysing 673 English tweets. Their results showed that misinformation accounted for 24.8% (153 of 617) of all serious tweets (i.e., not humour-related posts). Healthcare or public health accounts had the lowest amount of misinformation; yet still 12.3% (7 of 57) of their tweets included unverifiable information. Moscadelli and colleagues⁴⁰ collected and reviewed 2,102 news articles that were circulated on the internet. Their analysis showed that fake news was shared over 2 million times, which accounted for 23.1% (2,352,585 of 10,184,351) of total shares between Dec 31, 2019, and April 30, 2020.⁴⁰ Similarly, another quantitative study by Galhardi and colleagues comparing the proportion of fake news shared on WhatsApp, Instagram, and Facebook in Brazil showed that fake news was mainly shared on WhatsApp.³⁶ A UK study by Ahmed and colleagues analysed 22,785 tweets posted by 11,333 Twitter users with #FilmYourHospital to identify and evaluate the source of the conspiracy theory on Twitter.³² Their work uncovered that ordinary people were the major driver behind the spread of conspiracy theories.³² Another study investigated the 5G and COVID-19 conspiracy theory that was circulating on Twitter with a random subsample of 233 tweets. The content analysis showed that 34.8% (81) of tweets linked 5G and COVID-19 and 32.2% (75) condemned such theory.³³ Similar research by Bruns and colleagues investigated 89 664 distinct Facebook posts in Australia that were related to this conspiracy from Jan 1 to April 12, 2020, by use of time series and network analysis.³⁵ The results showed that this conspiracy went viral after March 19, 2020, with unusual coalition

among various groups on Facebook. Islam and colleagues analysed 2,311 infodemic reports that were related to COVID-19 from Dec 31, 2019, to April 5, 2020, and showed that misinformation was mainly driven by rumours, stigma, and conspiracy theories that were circulating on various social media and other online platforms.³⁸ Associations between infodemic and bot activities on social media are another important research direction. One study analysed 12 million tweets from the USA and 15 million tweets from the Philippines from March 5 to March 19, 2020, and both countries showed a positive relation between bot activities and rate of hate speech in communities that are denser and more isolated than others.⁴³ Brennen and colleagues qualitatively analysed 96 samples of visuals (i.e., image or video) from January to March, 2020, and categorised misinformation into six trends, noting that, fortunately, there has been no involvement of artificial intelligence deepfake techniques (i.e., techniques used to make synthetic videos that closely resemble real videos) so far.³⁴

2.4.2 Social media for surveillance and monitoring

Three themes emerged under this category: public attitudes, mental health, and detection or prediction of COVID-19 cases. Public attitudes and mental health are reflections regarding the public perceptions and mental health effects of the pandemic; detection or prediction of COVID-19 cases includes typical surveillance studies aiming to propose ways to detect or predict COVID-19 cases.

48 selected articles gauged the attitudes and emotions that were expressed by social media users regarding the COVID-19 pandemic, mainly by use of content and sentiment analysis. Twitter accounted for 33 articles and Sina Weibo accounted for 8 articles. Public attitude can be further divided into the following sub-themes: public sentiment towards the COVID-19 pandemic and interventions, stigma and racism, and ageism.

To learn about the public sentiment towards the overall COVID-19 pandemic and its interventions, Abd-Alrazaq and colleagues⁷ analysed 167,073 unique English tweets that were divided into four categories: origin, source, regional and global effects on people and society, and methods to reduce transmission of SARS-CoV-2. Tweets regarding economic loss had the highest mean number of likes, whereas travel bans and warnings had the lowest number of likes.⁷ Kwon and colleagues investigated 259,529 English tweets in the USA, using trending and spatiotemporal analyses, and noted that tweets about social disruptiveness had the highest number of retweets, whereas tweets about COVID-19 interventions had the highest number of likes.⁷³ A content analysis of 522 Reddit comments showed

that the topic of symptoms accounted for 27% (141) of all comments, followed by the topic of prevention (25% [131]).⁷⁴ Likewise, another content analysis of 155,353 unique English tweets showed that the most mentioned topic was “peril of COVID-19”.⁷⁶ Additionally, a study that examined 126,049 English tweets by use of sentiment analysis and latent Dirichlet analysis for topic modelling showed that the most common emotion that was mentioned was fear, and the most common topic that was mentioned was the economic and political effects.⁷⁹ Al-Rawi and colleagues studied emojis in over 50 million tweets and identified five primary subjects: morbidity fears, health concerns, employment and financial issues, praise for front-line workers, and unique gendered emoji use.⁵¹ Samuel and colleagues investigated 293,597 tweets with sentiment analysis and noted more positive emotions than negative emotions towards the US economy reopening.⁸⁶ Analysing 2,558,474 English tweets by use of clustering and network analyses, Odlum and colleagues identified that African Americans shared positive sentiments and encouraged virtual discussions and prevention behaviours.⁸² A study investigated gender differences in terms of topics by analysing 3,038,026 English tweets.⁸⁸ The results showed that tweets from women were more likely to be about family, physical distancing, and health care, whereas tweets from men were more likely to be about sports cancellations, pandemic severity, and politics. In Canada, Xue and colleagues analysed 1,015,874 tweets via latent Dirichlet analysis to identify nine themes about family violence.⁹³ In Australia, Yigitcanlar and colleagues analysed 96,666 tweets and identified that the public's attitude could be captured efficiently through social media analytics.⁹⁴ One qualitative content analysis of 30 profiles from Instagram, Twitter, and TikTok in Brunei identified five types of attitudes towards physical distancing: fear, responsibility, annoyance, fun, and resistance.⁸⁰ In Turkey, to show the effects of social media on human psychology and behaviour, Arpaci and colleagues⁵² used evolutionary clustering analysis on 43 million tweets between March 22 and March 30, 2020. The study suggested that high-frequency word clusters, such as death, test, spread, and lockdown denoted the public's underlying fear of infection and death from the virus, whereas terms such as stay home and social distancing corresponded to behavioural shifts.⁵² A study in Luzon, Philippines,⁸⁴ in which sentiment analysis was done by use of natural language processing, showed that most Filipino Twitter users expressed negative emotions towards COVID-19, and the negative mood grew stronger over time in lockdown.⁸⁴ Sentiment analysis of 107,990 English tweets uncovered that a negative feeling towards the COVID-19 pandemic dominated, and topic modelling showed three major themes in people's

concerns: the COVID-19 pandemic emergency, how to control COVID-19, and reports on COVID-19.⁵⁴ Another study analysed 373,908 Belgian tweets and retweets, which showed that the public relied on the EU coalition to tackle the pandemic.⁷² De Santis and colleagues analysed 1,044,645 tweets to identify daily hot topics in Italy that were related to the COVID-19 pandemic and developed a framework for prospective research.⁶² One thematic analysis study of 1,920,593 Arabic tweets in Egypt showed that negative emotions and sadness were high in tweets showing affective discussions, and the dominant themes included the outbreak of the pandemic, metaphysics responses, signs and symptoms in confirmed cases, and conspiracism.⁶⁴ In Singapore, Lwin and colleagues examined 20,325,929 tweets using sentiment analysis and showed that public emotions shifted over time: from fear to anger and from sadness to gratefulness.⁷⁷ Chang and colleagues examined over 1.07 million Chinese texts from various online sources in Taiwan using deductive analysis and identified that negative sentiments mainly came from online news with stigmatising language linked with the COVID-19 pandemic.⁵⁶ In India, one study investigated 410 643 tweets via sentiment analysis and latent Dirichlet analysis and showed that positive emotions were overall substantially higher than negative sentiments, but this observation diminished at individual levels.⁶¹ Another study analysed 29,554 tweets from the second lockdown (i.e., April 15–May 3, 2020) and 47,672 tweets from the third lockdown (i.e., May 4–May 17, 2020) via sentiment analysis uncovered positive attitudes towards the second lockdown but negative attitudes towards the third lockdown in India.⁵⁷ One study analysed 868 posts from Reddit and noted sentiments to be 50% (434) neutral, 22% (191) positive, and 28% (243) negative in India.⁶³ A study in South Korea examined 43,832 unique users and their relations on Twitter by use of content and network analyses and showed that tweets including medical news were more popular than tweets containing non-medical news.⁸³ A study from Ireland analysed 203,756 tweets through topic modelling and identified that war was the most frequently used frame for the pandemic.⁹⁰ In the USA, Damiano and colleagues qualitatively analysed 600 English tweets and showed neutral sentiment across most tweets.⁵⁹ Politics also had an essential role in shaping people's opinion.⁵⁹ A study of 19,803 tweets from Democrats and 11,084 tweets from Republicans by use of random forest in the USA showed that Democrats put more emphasis on public health and direct aid to US workers, whereas Republicans put more emphasis on national unity, China, and businesses.⁶⁷ Results of a study involving various online data sources from Italy, the UK, the USA, and Canada showed that media was the major driver of the public's attention, but attention decreased

with saturation of the media with news about COVID-19.⁶⁶ Compared with other users, Reddit users focused more on health, data related to new disease, and preventative interventions. Researchers in Spain studied 22,223 tweets by use of topic modelling and network analysis.⁹⁵ They identified eight frames and noted that the entire pandemic could be divided into three periods: precrisis, lockdown, and recovery periods. Using 563,079 English Reddit posts that were related to COVID-19, Jelodar and colleagues proposed a novel method to detect meaningful latent topics and sentiment–comment classification.⁶⁹ Samuel and colleagues examined over 900,000 tweets to study the accuracy of tweet classifications among logistic regression and Naive Bayes methods.⁸⁵ They identified that Naive Bayes had 91% of accuracy compared with 74% from the logistic regression model.⁸⁵

Han and colleagues analysed 1,413,297 Sina Weibo posts and observed that the public paid attention to information regarding the epidemic, especially in metro areas.⁶⁸ Zhao and colleagues studied 4,056 topics from the Sina Microblog hot search list and noted that the public emotions shifted from negative to neutral to positive over time and that five major public concerns existed: the situation of the new cases of COVID-19 and its effects, front-line reporting of the pandemic and the measures of prevention and control, expert interpretation and discussion on the source of infection, medical services on the front line of the pandemic, and focus on the pandemic and the search for suspected cases.⁹⁶ Li and colleagues⁷⁵ did an observational infoveillance study with a linear regression model by analysing 115,299 Sina Weibo posts. The results showed that the number of Sina Weibo posts positively correlated with the number of reported cases of COVID-19 in Wuhan. Additionally, the qualitative analysis classified the topics into the following four overarching themes: cause of the virus, epidemiological characteristics of COVID-19, public responses, and others.⁷⁵ Chen and colleagues examined relationships between citizen engagement through government social media and media richness, dialogic loop, content type, and emotion valence.⁵⁸ Citizen engagement through government social media refers to sum of shares, likes, and comments in this study, so the higher the sum, the greater the citizen engagement through government social media. Media richness quantifies how much information that a sender transfers to a receiver via a medium and is based on the media richness theory (i.e., “the potential information load of communication media, emphasising the abilities of promoting shared meaning”).¹⁰¹ Dialogic loop, or dialogic communication theory, is defined as an approach that promotes a dialogue between a speaker and audience. According to the American Psychological Association, emotion valence refers to “the value associated with a stimulus, expressed

on a continuum from pleasant to unpleasant or from attractive to aversive".¹⁰⁰ For instance, happiness is typically considered to be pleasant valence. Chen and colleagues analysed 1,411 posts that were related to COVID-19 from Healthy China, an official account of the National Health Commission of China on Sina Weibo. Findings showed an inverse association between media richness and citizen engagement through government social media, indicating that posts with plain texts had higher citizen engagement through government social media than did posts with pictures or videos. A positive association between dialogic loop and citizen engagement through government social media was noted, as evidenced by 96% (1,355 of 1,411) of responses to these posts having hashtags and 25% (353 of 1,411) containing questions. In terms of media richness, when posts had both a high media richness and positive emotion, citizen engagement through government social media increased, whereas when posts had a high media richness and negative emotion, citizen engagement decreased. Regarding content type, when posts were related to the latest news about the pandemic, stronger negative emotions led to increased citizen engagement through government social media.⁵⁸ Yin and colleagues⁶⁵ proposed a new multiple-information susceptible-discussing-immune model to analyse the public opinion propagation of COVID-19 from Sina Weibo posts that were collected from Dec 31, 2019, to Feb 27, 2020. The researchers reported that the reproduction rate of this proposed model reached 1.78 in the early stage of COVID-19 but decreased to around 0.97 and was maintained at this level. Such a result showed that the information on COVID-19 would continue to increase slowly in the future until it stabilises. However, this stability would depend on how much information is received on COVID-19. Wang and colleagues⁸⁹ analysed 999,978 randomly selected Sina Weibo posts that were related to COVID-19 through an unsupervised Bidirectional Encoder Representations from Transformers model for sentiments and a term frequency-inverse document frequency model for topic modelling. The authors identified four public concerns: the virus origin, symptom, production activity, and public health control in China.⁸⁹ Xi and colleagues examined 241 topics with their views and comments via thematic and temporal analysis and noted that older adults contributing to the community was the most frequent theme in the first phase of COVID-19 in China (i.e., Jan 20–Feb 20, 2020).⁹¹ The theme of older patients in hospitals was most frequent in the second (i.e., Feb 21–March 17, 2020) and third phase (i.e., March 18–April 28, 2020). Using Wilcoxon tests, Su and colleagues examined posts from 850 Sina Weibo users and 14,269 tweets from Italy.⁸⁷ The findings showed that Italian people paid more attention to leisure, whereas Chinese people paid more attention

to the community, religion, and emotions after lockdowns. Analysing the top 200 accounts from WeChat via regressions and content analysis, Ma and colleagues showed that both non-medical and medical reports had positive effects on people's behaviours.⁷⁸ Using Kendall's Tau-B rank test, Xie and colleagues investigated relations among the Baidu Attention Index, daily Google Trends, and numbers of COVID-19 cases and deaths.⁹² Daily Google Trends were correlated to seven indicators, whereas daily Baidu Search Index was correlated only to three indicators.⁹² Zhu and colleagues analysed 1,858,288 Sina Weibo posts and noted that topics changed over time but political and economic posts attracted greater attention than did other topics.⁹⁷

Regarding stigma and racism, Kim⁷¹ analysed 27,849 individual tweets in South Korea by use of a binary logistic regression to gauge network size and semantic network analysis to capture contextual and subjective factors. The results indicated that size of personal social network was inversely correlated with impolite language use. Namely, users with larger social networks were less likely to post uncivil messages on Twitter than were users with smaller social networks. This study suggested that the size of the social network influenced the language choice of social media users in their postings.⁷¹ Research compared public stigma before and after the introduction of the terms Chinese virus or China virus in 16,535 English tweets from before introduction and 177,327 tweets from after introduction.⁵⁵ The results showed an almost 10 times increase, nationwide and statewide and in the USA, from 0.38 tweets posted per 10,000 people referencing the two terms before introduction to 4.08 tweets posted per 10,000 after introduction. A similar study examined 339,063 tweets from non-Asian respondents via local polynomial regression and interrupted time-series analysis.⁶⁰ The findings showed that, when stigmatising terms, such as Chinese virus, were used by media (starting from March 8, 2020), the bias index (i.e., Implicit Americanness Bias) began to increase, and such bias was more profound in conservatives than in members of any other political subgroup. Nguyen and colleagues analysed 3,377,295 tweets that were related to race in the USA using sentiment analysis and uncovered a 68.4% increase in negative tweets referring to Asian people, whereas tweets referring to other races remained stable.⁸¹

Regarding ageism, a study⁷⁰ investigating Twitter content that was related to both COVID-19 and older adults analysed a random sample of 351 English tweets. 21.1% (74) of the tweets implied diminished regard for older adults by downplaying or dismissing concerns over the high fatality of

COVID-19 in this population.⁷⁰ Similar research examined 188 tweets via thematic analysis and showed that 90% (169) of tweets opposed ageism, whereas 5% (9) of tweets favoured ageism, and 5% (10) of tweets were neutral.⁵³

Two of 81 reviewed studies, both based in China, focused on assessing the mental health of social media users.^{44, 45} A cross-sectional study⁴⁴ investigated the relationship between anxiety and social media exposure, which is theoretically defined as “the extent to which audience members have encountered specific messages”.¹⁰² The researchers distributed an online survey based on the Chinese version of WHO-Five Well-Being Index for depression and the Chinese version of Generalized Anxiety Disorder Scale for anxiety. Respondents included 4,872 Chinese citizens aged 18 years and older from 31 provinces and autonomous regions in China. After controlling for all covariates through a multivariable logistic regression, the study showed that frequent social media exposure increased the odds ratio of anxiety, showing that frequent social media exposure is potentially contributing to mental health problems during the COVID-19 outbreak.⁴⁴ To explore how people's mental health was influenced by COVID-19, Li and colleagues⁴⁵ analysed posts from 17,865 active Sina Weibo users to compare sentiments before and after declaration of COVID-19 outbreak by the National Health Commission in China on Jan 20, 2020. The researchers identified increased negative sentiments, including anxiety, depression, and indignation, after the declaration and decreased positive sentiments expressed in the Oxford happiness score. Additionally, cognitive indicators showed increased sensitivity to social risks but decreased life satisfaction after the declaration.⁴⁵

Six of 81 studies investigated the detection or prediction of COVID-19 outbreaks with social media data. Qin and colleagues²² attempted to predict the number of newly suspected or confirmed COVID-19 cases by collecting social media search indexes for symptoms (e.g., dry cough, fever, and chest distress), coronavirus, and pneumonia. The data were analysed by use of subset selection, forward selection, lasso regression, ridge regression, and elastic net. Results showed that the optimal model was constructed via the subset selection. The lagged social media search indexes were a predictor of new suspected COVID-19 cases and could be detected 6–9 days before confirmation of new cases.²² To evaluate the possibility of early prediction of COVID-19 cases via internet searches and social media data, Li and colleagues¹⁷ used the keywords coronavirus and pneumonia to retrieve corresponding trend data from Google Trends, Baidu Search Index, and Sina Weibo Index. By use of

the lag correlation, the results showed that the correlation between trend data with the keyword coronavirus and number of laboratory-confirmed cases was highest 8–12 days before increase in confirmed COVID-19 cases in the three platforms. Similarly, the correlation between trend data for the keyword coronavirus and new suspected COVID-19 cases was highest 6–8 days before increase in new suspected cases. The correlation between trend data for the keyword pneumonia and new suspected cases was highest 8–10 days before increase in new suspected COVID-19 cases across the three platforms.¹⁷ Peng and colleagues studied 1,200 Sina Weibo records using spatiotemporal analysis, kernel density analysis, and ordinary least square regression and noted that scattered infection, community spread, and full-scale outbreak were three phases of early COVID-19 transmission in Wuhan, China.²¹ Older people are at high risk of severe COVID-19 and accounted for over 50% of help seeking on Sina Weibo. To identify COVID-19 patients with poor prognosis, Liu and colleagues analysed Sina Weibo messages from 599 patients along with telephone follow-ups.¹⁸ The findings suggested risk factors involving older age, diffuse distribution of pneumonia, and hypoxaemia. A regression study analysed Google Trends searches, Wikipedia page views, and tweets and showed that current Wikipedia page views, tweets from a week before, and Google Trends searches from two weeks before can be used to model the number of COVID-19 cases. To model the number of deaths, all three variables should be one week earlier than for cases.¹⁹

2.4.3 Social media as disease control

To inoculate the public against misinformation, public health organisations and governments should create and spread accurate information on social media because social media has had an increasingly important role in policy announcement and health education. Six of 81 articles were categorised as government responses because they examined how government messages and health education material were generated and consumed on social media platforms. Two studies analysed data from Sina Weibo,^{23, 27} and the other four studies analysed data from Twitter.^{28, 29, 30, 31}

Zhu and colleagues²³ measured the attention of Chinese netizens—i.e., citizen of the net—to COVID-19 by analysing 1,101 Sina Weibo posts. They noted that Chinese netizens paid little attention to the disease until the Chinese Government acknowledged and declared the COVID-19 outbreak on Jan 20, 2020. Since then, high levels of social media traffic occurred when Wuhan, China, began its quarantine (Jan 23–Jan 24, 2020), during a Red Cross Society of China scandal (Feb

1, 2020), and following the death of Li Wenliang (Feb 6–Feb 7, 2020).²³ Li and colleagues²⁷ collected 36,746 Sina Weibo posts to identify and categorise the situational information using support vector machines, Naive Bayes, and random forest as well as features in predicting the number of reports using linear regression. Except for posts that were categorised as counter rumours (i.e., used to oppose rumours), they identified that the higher the word count, the more reposts there were. Likewise, posts from unverified users had more reposts for all categories than did posts from verified users, excluding the counter rumours. For counter rumours, reposts increased with the number of followers and if the followers were from urban areas.²⁷ A qualitative content analysis was done to investigate how G7 leaders used Twitter for matters concerning the COVID-19 pandemic by collecting 203 tweets.²⁹ The findings showed that 166 of 203 tweets were informative, 48 tweets were linked to official government resources, 19 (9.4%) tweets were morale-boosting, and 14 (6.9%) tweets were political.²⁹ To assess the political partisan polarisation in Canada regarding COVID-19, Merkley and colleagues²⁸ randomly sampled 1,260 tweets from the social media of 292 federal members of parliament and collected 87 Google Trends for the search term coronavirus. 2,499 Canadian respondents aged 18 years and above were also surveyed. The results showed that, regardless of party affiliation, members of parliament emphasised the importance of measures for physical distancing and proper hand-hygiene practices to cope with the COVID-19 pandemic, without tweets exaggerating concerns or misinformation about COVID-19. Search interest in COVID-19 among municipalities was strongly determined by socioeconomic and urban factors rather than Conservative Party vote share.²⁸ Sutton and colleagues studied 149,335 tweets from public health, emergency management, and elected officials and observed that the underlying emotion of messages changed positively and negatively over time.³⁰ Wang and colleagues investigated 13,598 tweets that were related to COVID-19 via temporal and network analyses.³¹ They categorised 16 types of messages and identified inconsistent and incongruent messages expressed in four crucial prevention topics: mask wearing, risk assessments, stay at home order, and disinfectants or sanitizers.

Eight chosen studies investigated the quality (i.e., the number of recommended prevention behaviours that were covered in the videos—e.g., wearing a facemask, washing hands, physical distancing, etc.) of YouTube videos with COVID-19 prevention information. Basch and colleagues²⁴ did a cross-sectional study and retrieved the top 100 YouTube videos with the most views that were uploaded in January 2020, with the keyword of coronavirus in English, with English subtitles, or in

Spanish. These 100 videos generated over 125 million views in total. However, fewer than 33 videos included any of the seven key prevention behaviours that are recommended by the US Centers for Disease Control and Prevention.²⁴ A follow-up study with the same criteria and a successive sampling design gathered the top 100 YouTube videos that were most viewed in January and March, 2020.²⁵ Findings showed that, in total, the January sample generated over 125 million views, and the March sample had over 355 million views. Yet, fewer than 50 videos in either sample contained any of the prevention behaviours that are recommended by the US Centers for Disease Control and Prevention.²⁵ Additionally, a study investigated the top 100 YouTube videos about do-it-yourself hand sanitizer with the most views and showed that the average number of daily calls about paediatric poisoning increased substantially in March 2020, compared with the previous 2 years.⁴⁶

To analyse the information quality on YouTube about the COVID-19 pandemic and to compare the contents in English and Chinese Mandarin videos, Khatri and colleagues²⁶ collected 150 videos with the keywords 2019 novel coronavirus and Wuhan virus in English and Mandarin. The DISCERN score and the medical information and content index were calculated as a reliable way to measure the quality of health information. The mean DISCERN score for reliability was low: 3.12 of 5.00 for English videos and 3.25 for Mandarin videos. The mean cumulative medical information and content index score of useful videos was also undesirable: 6.71 of 25.00 for English videos and 6.28 for Mandarin videos.²⁶ In Spain, a similar study of 129 videos in Spanish identified that information in videos about preventing COVID-19 was usually incomplete and differed according to the type of authorship (i.e., mass media, health professionals, individual users, and others).⁴⁷ Likewise, one study in South Korea noted that misleading videos accounted for 37.14% (39 of 105) of most-viewed videos and had more likes, fewer comments, and longer viewing times than did useful videos.⁴⁸ Two studies in Turkey investigated the quality of YouTube videos regarding COVID-19 information in dentistry.^{49, 50} One of these studies analysed the top 116 English videos with at least 300 views and showed moderate quality and useful information from these videos.⁴⁹ The other study, however, showed poor quality for 24 of 55 (43.6%) English videos, whereas good quality accounted for only 2 (3.6%) videos.⁵⁰

Table 2-2 has summarised general findings and identified research gaps of each theme. It is acknowledged that except the “public attitudes” and “infodemics” themes, the other four themes

together have accounted for 21, or 26%, of the 81 reviewed articles. Therefore, the other themes naturally need to be investigated further in the future.

Themes	Overall Findings	Overall Gaps
Public attitudes	<ul style="list-style-type: none"> • Sentiment analysis and topic modelling have been commonly applied to investigate people’s attitudes towards COVID-19 related events on social media. • Qualitative analyses, such as content analysis and thematic analysis, have also been widely used for similar purposes. • Twitter and Weibo are the mostly investigated social media platform. 	<ul style="list-style-type: none"> • Other social media platforms also need to be investigated. • Limited studies have applied theories in understanding public attitudes. • Public attitudes or sentiments have not been incorporated into many intervention studies to decide if an intervention is effective.
Mental health	<ul style="list-style-type: none"> • Social media data from Weibo can be useful to detect mental health issues at the population level. 	<ul style="list-style-type: none"> • More studies are needed to use different social media data when investigating mental health issues.
Detection or prediction of COVID-19 cases	<ul style="list-style-type: none"> • Various methods, from statistical correlations to more advanced machine learning techniques, have been used to forecast or predict the number of COVID-19 cases by incorporating social media data. 	<ul style="list-style-type: none"> • Real-time surveillance that incorporates various social media data and other data are needed. • Machine learning techniques need large amount of data, which can be a disadvantage in

		early pandemic or infectious disease outbreaks.
Government responses	<ul style="list-style-type: none"> • Studies have shown people’s reactions to government responses to the COVID-19 pandemic. • Researchers think timely government responses is critical, but more studies are needed. 	<ul style="list-style-type: none"> • More research is needed to investigate how efficient and effective these official responses can lead to public belief or behavioural changes. • There is a need to compare impacts of infodemics with that of government responses.
Information quality	<ul style="list-style-type: none"> • Low information quality found in COVID-19 education or prevention YouTube videos. 	<ul style="list-style-type: none"> • Videos on other social media platforms need to be investigated.
Infodemics	<ul style="list-style-type: none"> • COVID-19 related information has substantially increased, but information quality has not been consistent. • Health misinformation have appeared less on health organization official accounts, but some of them still have shared unverified information. • Fake news, conspiracies, and other misinformation have been shared widely on social media regardless of physical borders. 	<ul style="list-style-type: none"> • More research is needed to understand how misinformation can undermine public health preventions. • More studies are needed to investigate how bots on social media have played in sharing misinformation.

Table 2-2 Summaries of findings and research gaps identified in each theme

2.5 Discussion

Studies on social media data showed our attitudes and mental state to some extent during the COVID-19 crisis. These studies also showed how we generated, consumed, and propagated information on social media platforms when facing the rapid spread of the SARS-CoV-2 and extraordinary measures for the containment. In our Review, public attitudes accounted for nearly 59% (48 of 81) of the reviewed articles. In terms of social media platforms, 56% (45 of 81) of the chosen articles used data from Twitter, followed by Sina Weibo (20% [16 of 81]). Machine learning analyses, such as latent Dirichlet analysis and random forest, were applied in research that studied public attitudes.

We identified six themes on the basis of our modified SPHERE framework, including infodemics, public attitudes, mental health, detection or prediction of COVID-19 cases, government responses to the pandemic, and quality of prevention education videos. However, a common limitation in all chosen studies on social media data is the comparison of data due to differences in quality, such as formats, metrics, or even the definition of common variables (e.g., the amount of time required for a post to be on an individual's screen to be counted as a view). For instance, the definition of a view on one social media platform is likely to be different from another. Besides, not every social media platform offers accessible data, like Twitter and Sina Weibo. To address these challenges, the selected studies have controlled for many factors, including social media platforms, languages, locations, time, misspellings, keywords, or hashtags. However, such search strategies resulted in many study limitations, such as non-representative sample sizes, selection bias, cross-sectional study design, or retrospective study design. We also observed that, given the large amount of available data, most studies across all domains sampled small data size for analyses, except for four studies under the theme of public attitudes that analysed over one million posts via machine learning methods. Additionally, data from Twitter and Sina Weibo accounted for over 70% (59 of 81) of our selected studies. Research examining other social media platforms, including Facebook, Instagram, TikTok, Snapchat, and WhatsApp, is scarce due to barriers of data availability and accessibility. We also identified future research topics that are needed for each category during the COVID-19 pandemic as shown in Table 2-2. From an infodemics perspective, additional research is needed to investigate how misinformation, rumours, and fake news (e.g., anti-mask wearing reports) undermine preventions and compromise public health, although social media companies, such as Twitter and Facebook, have

started to remove accounts that are based on misinformation. Bot posts are another topic to be addressed and studies evaluating effective counter-infodemic interventions are also needed.

Articles regarding public attitudes towards the COVID-19 pandemic have shown sentiments that shifted over time. Yet, this theme can be a useful indicator when evaluating interventions, such as physical distancing and wearing masks, that aim to reduce the risk of COVID-19 infection. However, public sentiments had not been incorporated into many intervention studies by the time that we did this Review. When a disease, such as COVID-19, starts spreading and causing negative sentiments, timely, proper, and effective risk communication is needed to help ease people's anxiety or negative attitudes regarding the COVID-19 pandemic, especially through social media.

Mental health is another issue that requires further investigation. Our chosen studies did not address mental health issues on the basis of age, as symptoms and interventions tend to vary with age. Public health measures, such as physical distancing, that were implemented in the COVID-19 pandemic exacerbated risk factors and adverse health behaviours at the individual and population levels. Studies showed that social media data were useful to detect mental health issues at the population level. Due to the early outbreak of COVID-19 and the prevalence of social media use (e.g., Sina Weibo and WeChat) in China, two studies reported increased issues of mental health among the Chinese population.^{44, 45} A similar trend of deteriorating mental health could happen in other regions. At the time of writing, British Columbia has recorded the highest number of overdose deaths in Canada (May 2020).¹⁰³

In terms of the surveillance of the COVID-19 pandemic, six chosen studies showed methods to detect or predict the number of COVID-19 cases by use of social media data. According to our Review, unlike other infectious diseases, such as influenza and malaria, COVID-19 has not had real-time monitoring surveillance developed with social media data. It is possible that the pandemic has evolved so rapidly that finding COVID-19 vaccinations or therapies has been prioritised over real-time monitoring surveillance with social media. Besides, scarcity of accurate and reliable data sources might discourage the development of the COVID-19 real-time surveillance. Moreover, whether COVID-19 is a one-time event or will become seasonal, like influenza, is unknown. If COVID-19 becomes seasonal, then it might be meaningful and useful to establish a real-time model to monitor the disease by use of social media data.

Government responses that were distributed via social media have been increasingly crucial in combating infodemics and promoting accurate and reliable information for the public. However, little has been studied about how efficient and effective these official responses are at leading to public belief or behavioural changes. It also remained unknown whether government posts would reach greater numbers of social media users or have greater effects on them than would infodemics.

YouTube has served as one of the major platforms to spread information concerning the control of COVID-19. Nonetheless, our chosen studies showed that most YouTube videos were of undesirable quality because they contained few recommended preventions from governments or public health organisations. The undesirable quality is a worrisome observation if accurate and reliable videos and other types of information are not created and disseminated in a timely manner. Therefore, videos, especially from public health authorities, should include accurate and reliable medical and scientific information and use relevant hashtags to reach a large audience, generate a high number of views, and increase responses. Moreover, our selected studies were limited to YouTube videos only. Additionally, a substantial proportion of the studies were done using Sina Weibo, which, although used by many people, is exclusive to China and might lead to an over-representation of a single country in this Review.

In summary, although our Review has limitations that are embedded from the chosen studies, we recognised six themes that have been studied so far and identified future research directions. Our adopted framework can serve as a fundamental and flexible guideline when studying social media and epidemiology.

2.6 Conclusion

Our Review identified various topics, themes, and methodological approaches in studies on social media and COVID-19. Among the six identified themes, public attitudes comprised most of the articles. Among the selected studies, Twitter was the leading social media platform, followed by Sina Weibo. Few studies included machine learning methods, whereas most studies used traditional statistical methods. Unlike influenza, we were not able to find studies documenting real-time surveillance that was developed with social media data on COVID-19. Our Review also identified studies that were related to COVID-19 on infodemics, mental health, and prediction. For COVID-19, accurate and reliable information through social media platforms can have a crucial role in tackling

infodemics, misinformation, and rumours. Additionally, real-time surveillance from social media about COVID-19 can be an important tool in the armamentarium of interventions by public health agencies and organisations.

2.7 References

1. Johns Hopkins University and Medicine Coronavirus resource center: world map. <https://coronavirus.jhu.edu/map.html> Date accessed: January 18, 2021.
2. Pérez-Escoda A, Jiménez-Narros C, Perlado-Lamo-de-Espinosa M, Pedrero-Esteban LM. Social networks' engagement during the COVID-19 pandemic in Spain: health media vs. healthcare professionals. *Int J Environ Res Public Health*. 2020; **17**: 5261
3. Jordan SE , Hovet SE, Fung IC, Liang H, Fu K, Tse ZTH. Using Twitter for public health surveillance from monitoring and prediction to public response. *Data (Basel)*. 2019; **4**: 6.
4. Shah Z, Surian D, Dyda A, Coiera E, Mandl KD, Dunn AG. Automatically appraising the credibility of vaccine-related web pages shared on social media: a Twitter surveillance study. *J Med Internet Res* 2019; **21**: e14007.
5. Sinnenberg L, Buttenheim AM, Padrez K, Mancheno C, Ungar L, Merchant RM. Twitter as a tool for health research: a systematic review. *Am J Public Health* 2017; **107**: e1–8.
6. Steffens MS, Dunn AG, Wiley KE, Leask J. How organisations promoting vaccination respond to misinformation on social media: a qualitative investigation. *BMC Public Health* 2019; **19**: 1348.
7. Abd-Alrazaq A, Alhuwail D, Househ M, Hamdi M, Shah Z. Top concerns of tweeters during the COVID-19 pandemic: infoveillance study. *J Med Internet Res* 2020; **22**: e19016.
8. Freberg K, Palenchar MJ, Veil SR. Managing and sharing H1N1 crisis information using social media bookmarking services. *Public Relat Rev* 2013; **39**: 178–84.
9. Giustini D, Ali SM, Fraser M, Kamel Boulos MN. Effective uses of social media in public health and medicine: a systematic review of systematic reviews. *Online J Public Health Inform* 2018; **10**: e215.
10. Al-Dmour H, Masa'deh R, Salman A, Abuhashesh M, Al-Dmour R. Influence of social media platforms on public health protection against the COVID-19 pandemic via the mediating effects of public health awareness and behavioral changes: integrated model. *J Med Internet Res* 2020; **22**: e19996.

11. Bridgman A, Merkley E, Loewen PJ, et al. The causes and consequences of COVID-19 misperceptions: understanding the role of news and social media. June 18, 2020. <https://misinforeview.hks.harvard.edu/article/the-causes-and-consequences-of-covid-19-misperceptions-understanding-the-role-of-news-and-social-media> (accessed Sept 15, 2020).
12. Tang L, Bie B, Park SE, Zhi D. Social media and outbreaks of emerging infectious diseases: a systematic review of literature. *Am J Infect Control* 2018; **46**: 962–72.
13. Arksey H, O'Malley L. Scoping studies: towards a methodological framework. *Int J Soc Res Methodol* 2005; **8**: 19–32.
14. Levac D, Colquhoun H, O'Brien KK. Scoping studies: advancing the methodology. *Implement Sci* 2010; **5**: 69.
15. Ose SO. Using Excel and Word to structure qualitative data. *J Appl Soc Sci (Boulder)* 2016; **10**: 147–62.
16. Braun V, Clarke V. Using thematic analysis in psychology. *Qual Res Psychol* 2006; **3**: 77–101.
17. Li C, Chen LJ, Chen X, Zhang M, Pang CP, Chen H. Retrospective analysis of the possibility of predicting the COVID-19 outbreak from Internet searches and social media data, China, 2020. *Euro Surveill* 2020; **25**: 2000199.
18. Liu D, Wang Y, Wang J, et al. Characteristics and outcomes of a sample of patients with COVID-19 identified through social media in Wuhan, China: observational study. *J Med Internet Res* 2020; **22**: e20108.
19. O'Leary DE, Storey VC. A Google–Wikipedia–Twitter model as a leading indicator of the numbers of coronavirus deaths. *Intell Syst Account Finance Manag* 2020; **27**: 151–58.
20. Worldometer. COVID-19 coronavirus pandemic. 2020. <https://www.worldometers.info/coronavirus> (accessed Sept 15, 2020).
21. Peng Z, Wang R, Liu L, Wu H. Exploring urban spatial features of COVID-19 transmission in Wuhan based on social media data. *ISPRS Int J Geoinf* 2020; **9**: 402.
22. Qin L, Sun Q, Wang Y, et al. Prediction of number of cases of 2019 novel coronavirus (COVID-19) using social media search index. *Int J Environ Res Public Health* 2020; **17**: 2365.

23. Zhu Y, Fu KW, Grépin KA, Liang H, Fung IC. Limited early warnings and public attention to coronavirus disease 2019 in China, January–February, 2020: a longitudinal cohort of randomly sampled Weibo users. *Disaster Med Public Health Prep* 2020; published online April 3. <https://doi.org/10.1017/dmp.2020.68>.
24. Basch CH, Hillyer GC, Meleo-Erwin ZC, Jaime C, Mohlman J, Basch CE. Preventive behaviors conveyed on YouTube to mitigate transmission of COVID-19: cross-sectional study. *JMIR Public Health Surveill* 2020; **6**: e18807.
25. Basch CE, Basch CH, Hillyer GC, Jaime C. The role of YouTube and the entertainment industry in saving lives by educating and mobilizing the public to adopt behaviors for community mitigation of COVID-19: successive sampling design study. *JMIR Public Health Surveill* 2020; **6**: e19145.
26. Khatri P, Singh SR, Belani NK, et al. YouTube as a source of information on 2019 novel coronavirus outbreak: a cross-sectional study of English and Mandarin content. *Travel Med Infect Dis* 2020; **35**: 101636.
27. Li L, Zhang Q, Wang X, et al. Characterizing the propagation of situational information in social media during covid-19 epidemic: a case study on Weibo. *IEEE Trans Comput Soc Syst* 2020; **7**: 556–62.
28. Merkley E, Bridgman A, Loewen PJ, Owen T, Ruths D, Zhilin O. A rare moment of cross-partisan consensus: elite and public response to the COVID-19 pandemic in Canada. *Can J Polit Sci* 2020; published April 16. <https://doi.org/10.1017/S0008423920000311>.
29. Rufai SR, Bunce C. World leaders' usage of Twitter in response to the COVID-19 pandemic: a content analysis. *J Public Health (Oxf)* 2020; **42**: 510–16.
30. Sutton J, Renshaw SL, Butts CT. COVID-19: Retransmission of official communications in an emerging pandemic. *PLoS One* 2020; **15**: e0238491.
31. Wang Y, Hao H, Platt LS. Examining risk and crisis communications of government agencies and stakeholders during early-stages of COVID-19 on Twitter. *Comput Human Behav* 2021; **114**: 106568.

32. Ahmed W, López Seguí F, Vidal-Alaball J, Katz MS. COVID-19 and the “Film Your Hospital” conspiracy theory: social network analysis of Twitter data. *J Med Internet Res* 2020; **22**: e22374.
33. Ahmed W, Vidal-Alaball J, Downing J, López Seguí F. COVID-19 and the 5G conspiracy theory: social network analysis of Twitter data. *J Med Internet Res* 2020; **22**: e19458.
34. Brennen JS, Simon FM, Nielsen RK. Beyond (mis)representation: visuals in COVID-19 misinformation. *Int J Press/Polit* 2020; published online Oct 10. <https://doi.org/10.1177/1940161220964780>.
35. Bruns A, Harrington S, Hurcombe E. ‘Corona? 5G? or both?’: the dynamics of COVID-19/5G conspiracy theories on Facebook. *Media Int Aust* 2020; **177**: 12–29.
36. Galhardi CP, Freire NP, Minayo MCS, Fagundes MCM. Fact or fake? An analysis of disinformation regarding the COVID-19 pandemic in Brazil. *Cien Saude Colet* 2020; **25** (suppl 2): 4201–10.
37. Gallotti R, Valle F, Castaldo N, Sacco P, De Domenico M. Assessing the risks of ‘infodemics’ in response to COVID-19 epidemics. *Nat Hum Behav* 2020; **4**: 1285–93.
38. Islam MS, Sarkar T, Khan SH, et al. COVID-19-related infodemic and its impact on public health: a global social media analysis. *Am J Trop Med Hyg* 2020; **103**: 1621–29.
39. Kouzy R, Abi Jaoude J, Kraitem A, et al. Coronavirus goes viral: quantifying the COVID-19 misinformation epidemic on Twitter. *Cureus* 2020; **12**: e7255.
40. Moscadelli A, Albora G, Biamonte MA, et al. Fake news and COVID-19 in Italy: results of a quantitative observational study. *Int J Environ Res Public Health* 2020; **17**: 5850.
41. Pulido CM, Villarejo-Carballido B, Redondo-Sama G, Gómez A. COVID-19 infodemic: more retweets for science-based information on coronavirus than for false information. *Int Sociol* 2020; **35**: 377–92.
42. Rovetta A, Bhagavathula AS. Global infodemiology of COVID-19: analysis of Google web searches and Instagram hashtags. *J Med Internet Res* 2020; **22**: e20673.

43. Uyheng J, Carley KM. Bots and online hate during the COVID-19 pandemic: case studies in the United States and the Philippines. *J Comput Soc Sc* 2020; published online Oct 20. <https://doi.org/10.1007/s42001-020-00087-4>.
44. Gao J, Zheng P, Jia Y, et al. Mental health problems and social media exposure during COVID-19 outbreak. *PLoS One* 2020; **15**: e0231924.
45. Li S, Wang Y, Xue J, Zhao N, Zhu T. The impact of COVID-19 epidemic declaration on psychological consequences: a study on active Weibo users. *Int J Environ Res Public Health* 2020; **17**: 2032.
46. Hakimi AA, Armstrong WB. Hand sanitizer in a pandemic: wrong formulations in the wrong hands. *J Emerg Med* 2020; **59**: 668–72.
47. Hernández-García I, Giménez-Júlvez T. Characteristics of YouTube videos in Spanish on how to prevent COVID-19. *Int J Environ Res Public Health* 2020; **17**: 4671.
48. Moon H, Lee GH. Evaluation of Korean-language COVID-19-related medical information on YouTube: cross-sectional infodemiology study. *J Med Internet Res* 2020; **22**: e20775.
49. Ozdede M, Peker I. Analysis of dentistry YouTube videos related to COVID-19. *Braz Dent J* 2020; **31**: 392–98.
50. Yüce MÖ, Adalı E, Kanmaz B. An analysis of YouTube videos as educational resources for dental practitioners to prevent the spread of COVID-19. *Ir J Med Sci* 2020; published online July 23. <https://doi.org/10.1007/s11845-020-02312-5>.
51. Al-Rawi A, Siddiqi M, Morgan R, Vandan N, Smith J, Wenham C. COVID-19 and the gendered use of emojis on Twitter: infodemiology study. *J Med Internet Res* 2020; **22**: e21646.
52. Arpacı I, Alshehabi S, Al-Emran M, et al. Analysis of twitter data using evolutionary clustering during the COVID-19 pandemic. *Comput Mater Contin* 2020; **65**: 193–204.
53. Barrett AE, Michael C, Padavic I. Calculated ageism: generational sacrifice as a response to the COVID-19 pandemic. *J Gerontol Psychol Sci Soc Sci* 2020; published online Aug 25. <https://doi.org/10.1093/geronb/gbaa132>.

54. Boon-Itt S, Skunkan Y. Public perception of the COVID-19 pandemic on Twitter: sentiment analysis and topic modelling study. *JMIR Public Health Surveill* 2020; **6**: e21978.
55. Budhwani H, Sun R. Creating COVID-19 stigma by referencing the novel coronavirus as the “Chinese virus” on Twitter: quantitative analysis of social media data. *J Med Internet Res* 2020; **22**: e19301.
56. Chang A, Schulz PJ, Tu S, Liu MT. Blaming devices in online communication of the COVID-19 pandemic: stigmatizing cues and negative sentiment gauged with automated analytic techniques. *J Med Internet Res* 2020; **22**: e21504.
57. Chehal D, Gupta P, Gulati P. COVID-19 pandemic lockdown: an emotional health perspective of Indians on Twitter. *Int J Soc Psychiatry* 2020; published online July 7. <https://doi.org/10.1177/0020764020940741>.
58. Chen Q, Min C, Zhang W, Wang G, Ma X, Evans R. Unpacking the black box: how to promote citizen engagement through government social media during the COVID-19 crisis. *Comput Hum Behav* 2020; **110**: 106380.
59. Damiano AD, Allen Catellier JR. A content analysis of coronavirus tweets in the United States just prior to the pandemic declaration. *Cyberpsychol Behav Soc Netw* 2020; published online Dec 14. <https://doi.org/10.1089/cyber.2020.0425>.
60. Darling-Hammond S, Michaels EK, Allen AM, et al. After “The China Virus” went viral: racially charged coronavirus coverage and trends in bias against Asian Americans. *Health Educ Behav* 2020; **47**: 870–79.
61. Das S, Dutta A. Characterizing public emotions and sentiments in COVID-19 environment: a case study of India. *J Hum Behav Soc Environ* 2020; published online July 14. <https://doi.org/10.1080/10911359.2020.1781015>.
62. De Santis E, Martino A, Rizzi A. An infoveillance system for detecting and tracking relevant topics from Italian tweets during the COVID-19 event. *IEEE Access* 2020; **8**: 132527–38.
63. Dheeraj K. Analysing COVID-19 news impact on social media aggregation. *Int J Adv Trends Comput Sci Eng* 2020; **9**: 2848–55.

64. Essam BA, Abdo MS. How do Arab tweeters perceive the COVID-19 pandemic? *J Psycholinguist Res* 2020; published online Aug 14. <https://doi.org/10.1007/s10936-020-09715-6>.
65. Yin FL, Lv JH, Zhang XJ, Xia XY, Wu JH. COVID-19 information propagation dynamics in the Chinese Sina-microblog. *Math Biosci Eng* 2020; **17**: 2676–92.
66. Gozzi N, Tizzani M, Starnini M, et al. Collective response to media coverage of the COVID-19 pandemic on Reddit and Wikipedia: mixed-methods analysis. *J Med Internet Res* 2020; **22**: e21597.
67. Green J, Edgerton J, Naftel D, Shoub K, Cranmer SJ. Elusive consensus: polarization in elite communication on the COVID-19 pandemic. *Sci Adv* 2020; **6**: eabc2717.
68. Han X, Wang J, Zhang M, Wang X. Using social media to mine and analyze public opinion related to COVID-19 in China. *Int J Environ Res Public Health* 2020; **17**: 2788.
69. Jelodar H, Wang Y, Orji R, Huang S. Deep Sentiment classification and topic discovery on novel coronavirus or COVID-19 online discussions: NLP using LSTM recurrent neural network approach. *IEEE J Biomed Health Inform* 2020; **24**: 2733–42.
70. Jimenez-Sotomayor MR, Gomez-Moreno C, Soto-Perez-de-Celis E. Coronavirus, ageism, and Twitter: an evaluation of tweets about older adults and COVID-19. *J Am Geriatr Soc* 2020; **68**: 1661–65.
71. Kim B. Effects of social grooming on incivility in COVID-19. *Cyberpsychol Behav Soc Netw* 2020; **23**: 519–25.
72. Kurten S, Beullens K. #Coronavirus: monitoring the Belgian Twitter discourse on the severe acute respiratory syndrome coronavirus 2 pandemic. *Cyberpsychol Behav Soc Netw* 2020; published online Aug 27. <https://doi.org/10.1089/cyber.2020.0341>.
73. Kwon J, Grady C, Feliciano JT, Fodeh SJ. Defining facets of social distancing during the COVID-19 pandemic: Twitter analysis. *J Biomed Inform* 2020; **111**: 103601.
74. Lai D, Wang D, Calvano J, Raja AS, He S. Addressing immediate public coronavirus (COVID-19) concerns through social media: utilizing Reddit’s AMA as a framework for public engagement with science. *PLoS One* 2020; **15**: e0240326.

75. Li J, Xu Q, Cuomo R, Purushothaman V, Mackey T. Data mining and content analysis of the Chinese social media platform Weibo during the early COVID-19 outbreak: retrospective observational infoveillance study. *JMIR Public Health Surveill* 2020; **6**: e18700.
76. Li Y, Twersky S, Ignace K, et al. Constructing and communicating COVID-19 stigma on Twitter: a content analysis of tweets during the early stage of the COVID-19 outbreak. *Int J Environ Res Public Health* 2020; **17**: 6847.
77. Lwin MO, Lu J, Sheldenkar A, et al. Global sentiments surrounding the COVID-19 pandemic on Twitter: analysis of Twitter trends. *JMIR Public Health Surveill* 2020; **6**: e19447.
78. Ma R, Deng Z, Wu M. Effects of health information dissemination on user follows and likes during COVID-19 outbreak in China: data and content analysis. *Int J Environ Res Public Health* 2020; **17**: 5081.
79. Medford RJ, Saleh SN, Sumarsono A, Perl TM, Lehmann CU. An “infodemic”: leveraging high-volume Twitter data to understand early public sentiment for the coronavirus disease 2019 outbreak. *Open Forum Infect Dis* 2020; **7**: ofaa258.
80. Mohamad SM. Creative production of ‘COVID-19 social distancing’ narratives on social media. *Tijdschr Econ Soc Geog* 2020; **111**: 347–59.
81. Nguyen TT, Criss S, Dwivedi P, et al. Exploring U.S. shifts in anti-Asian sentiment with the emergence of COVID-19. *Int J Environ Res Public Health* 2020; **17**: 7032.
82. Odlum M, Cho H, Broadwell P, et al. Application of topic modeling to tweets as the foundation for health disparity research for COVID-19. *Stud Health Technol Inform* 2020; **272**: 24–27.
83. Park HW, Park S, Chong M. Conversations and medical news frames on twitter: infodemiological study on COVID-19 in South Korea. *J Med Internet Res* 2020; **22**: e18897.
84. Pastor CK. Sentiment analysis of filipinos and effects of extreme community quarantine due to coronavirus (COVID-19) pandemic. *SSRN* 2020; published online April 13. https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3574385 (preprint).
85. Samuel J, Ali NGGM, Rahman MM, Esawi E, Samuel Y. COVID-19 public sentiment insights and machine learning for tweets classification. *Information (Basel)* 2020; **11**: 314.

86. Samuel J, Rahman MM, Ali GGMN, et al. Feeling positive about reopening? New normal scenarios from COVID-19 US reopen sentiment analytics. *IEEE Access* 2020; **8**: 142173–90.
87. Su Y, Xue J, Liu X, et al. Examining the impact of COVID-19 lockdown in Wuhan and Lombardy: a psycholinguistic analysis on Weibo and Twitter. *Int J Environ Res Public Health* 2020; **17**: 4552.
88. Thelwall M, Thelwall S. COVID-19 tweeting in English: gender differences. *Relaciones Públicas* 2020; **29**: e290301.
89. Wang T, Lu K, Chow KP, Zhu Q. COVID-19 sensing: negative sentiment analysis on social media in China via BERT model. *IEEE Access* 2020; **8**: 138162–69.
90. Wicke P, Bolognesi MM. Framing COVID-19: how we conceptualize and discuss the pandemic on Twitter. *PLoS One* 2020; **15**: e0240010.
91. Xi W, Xu W, Zhang X, Ayalon L. A thematic analysis of Weibo topics (Chinese twitter hashtags) regarding older adults during the COVID-19 outbreak. *J Gerontol Psychol Sci Soc Sci* 2020; published online Sept 3. <https://doi.org/10.1093/geronb/gbaa148>.
92. Xie T, Tan T, Li J. An extensive search trends-based analysis of public attention on social media in the early outbreak of COVID-19 in China. *Risk Manag Healthc Policy* 2020; **13**: 1353–64.
93. Xue J, Chen J, Chen C, Hu R, Zhu T. The hidden pandemic of family violence during COVID-19: unsupervised learning of tweets. *J Med Internet Res* 2020; **22**: e24361.
94. Yigitcanlar T, Kankanamge N, Preston A, et al. How can social media analytics assist authorities in pandemic-related policy decisions? Insights from Australian states and territories. *Health Inf Sci Syst* 2020; **8**: 37.
95. Yu J, Lu Y, Muñoz-Justicia J. Analyzing Spanish news frames on Twitter during COVID-19—a network study of El País and El Mundo. *Int J Environ Res Public Health* 2020; **17**: 5414.
96. Zhao Y, Cheng S, Yu X, Xu H. Chinese public’s attention to the COVID-19 epidemic on social media: observational descriptive study. *J Med Internet Res* 2020; **22**: e18825.

97. Zhu B, Zheng X, Liu H, Li J, Wang P. Analysis of spatiotemporal characteristics of big data on social media sentiment with COVID-19 epidemic topics. *Chaos Solitons Fractals* 2020; **140**: 110123.
98. Schillinger D, Chittamuru D, Ramírez AS. From “infodemics” to health promotion: a novel framework for the role of social media in public health. *Am J Public Health* 2020; **110**: 1393–96.
99. WHO. Novel coronavirus (2019-nCoV): situation report—13. 2020. <https://www.who.int/docs/default-source/coronaviruse/situation-reports/20200202-sitrep-13-ncov-v3.pdf> (accessed Feb 15, 2020).
100. PsychCentral. Your emotional brain on resentment, part 2. 2020. <https://pro.psychcentral.com/your-emotional-brain-on-resentment-part-2> (accessed Sept 15, 2020).
101. Daft RL, Lengel RH. Organizational information requirements, media richness and structural design. *Manag Sci* 1986; **32**: 554–71.
102. de Vreese CH, Neijens P. Measuring media exposure in a changing communications environment. *J Commun Methods Meas* 2016; **10**: 69–80.
103. Schmunk R. B.C. records highest number of fatal overdoses in a single month, with 170 deaths. June 11, 2020. <https://www.cbc.ca/news/canada/british-columbia/overdose-deaths-bc-1.5607792> (accessed Sept 15, 2020).

Chapter 3: Proposing a conceptual framework: social media infodemic listening (SoMELL) for public health behaviours

Status: Currently under journal review.

Authors: Shu-Feng Tsao, Helen Chen, Samantha B. Meyer, Zahid A. Butt

The COVID-19 pandemic brought with it an unprecedented surge in the use of social media platforms as sources of information and communication. The rapid dissemination of information and misinformation on these social media platforms has created a unique challenge for researchers seeking to understand the impact of social media on public attitudes and the emergence of infodemics during the early stages of the pandemic. As the scoping review (Study I, Chapter 2) has shown, two predominant themes—public attitudes and infodemics—accounted for 60 out of 81 articles or 74% of the reviewed articles. This has me interested in investigating further in these two themes.

For the 48 articles related to the public attitudes, during the early stages of the COVID-19 pandemic, social media became a central platform for individuals to express their attitudes, fears, and concerns. Public attitudes were shaped by the constant flow of information and misinformation on platforms like Twitter, Facebook, and Instagram. Researchers recognized the significance of analyzing social media data to gain insights into these evolving public attitudes. Researchers have employed various methods, including sentiment analysis, content analysis, and topic modelling, to explore how social media data could be leveraged to gauge public sentiment, track shifts in attitudes, and identify influential voices within online communities. These studies have provided valuable insights into the evolving perceptions of COVID-19 among the public, shedding light on topics such as mask-wearing, social distancing, and vaccine hesitancy. However, a critical observation emerged during the review process: the limited integration of established communication and behavioral theories in these studies. Although some research drew on theories from these domains, such as the Health Belief Model or the Theory of Planned Behavior, their application often fell short of capturing the complexity of the modern information ecosystem intertwined with both online and offline channels. As a result, it has been challenging to provide comprehensive explanations for the observed changes in public attitudes.

Similarly, the emergence of infodemics, defined as the rapid spread of false or misleading health information, became a prominent concern during the early days of the COVID-19 pandemic. Social media platforms served as breeding grounds for the dissemination of inaccurate information, conspiracy theories, and unsubstantiated claims related to the COVID-19 virus. This phenomenon presented a unique research opportunity and challenge: understanding how infodemics on social media impact public health behavior. The scoping review has identified 12 relevant manuscripts focused on health infodemics during the COVID-19 pandemic. Scholars has used a variety of methodologies to track the spread of misinformation and identify influential sources. These studies offered valuable insights into the mechanisms driving the spread of health-related misinformation on social media platforms. However, a recurring issue emerged during the review: the limited application of existing communication and behavioral theories in the context of health infodemics. Although some studies acknowledged the relevance of theories such as the Extended Parallel Process Model (also known as Fear Appeals), these theories often fell short in capturing the dynamic and multifaceted nature of infodemics in the digital age. Therefore, researchers faced challenges in constructing comprehensive theoretical frameworks to guide their investigations.

Building upon the foundation of the scoping review, a pressing research question emerged: Are there existing theories that adequately address people's health behaviors using social media data? To answer this question, this study has two objectives: (1) to survey existing theories and identify limitations and gaps, and (2) to propose a conceptual framework that can address identified gaps. This chapter begins with the abstract of the study, followed by the full-text manuscript.

3.1 Abstract

Existing communications and behavioural theories have been adopted to address health infodemics. Although various theories and models have been used to investigate the COVID-19 pandemic, there is no framework specially designed for social listening studies using social media data and natural language processing techniques. This study aimed to propose a novel yet theory-based conceptual framework for infodemiological research. We collected theories and models used in COVID-19 related studies published in peer-reviewed journals. The theories and models ranged from health behaviours, communications, to infodemics. They are analysed and critiqued for their components, followed by proposing a conceptual framework with a demonstration. We reviewed Health Belief

Model, Theory of Planned Behaviour/Reasoned Action, Communication for Behavioural Impact, Transtheoretical Model, Uses and Gratifications Theory, Social Judgment Theory, Risk Information Seeking and Processing Model, Behavioural and Social Drivers, and Hype Loop. Accordingly, we proposed our ‘Social Media Listening for Public Health Behaviour’ Conceptual Framework by not only integrating important attributes of existing theories, but also adding new attributes. The proposed conceptual framework was demonstrated in the Freedom Convoy social media listening. The proposed conceptual framework can be used to better understand public discourse on social media, and it can be integrated with other data analyses to gather a more comprehensive picture. The framework will continue to be revised and adopted as health infodemics evolve.

Keywords: infodemic; social media; conceptual framework; social listening; machine learning; natural language processing

3.2 Introduction

The World Health Organization (WHO) has consistently reiterated the widespread and multifaced nature of health infodemics and their harmful consequences throughout the pandemic.¹ The WHO initiated and hosted infodemic conferences and trainings since early 2020 to address increasingly complex health infodemics.^{1, 2, 3} The WHO’s technical consultation has led to a framework to manage infodemics.³ Another framework that categorizes research agenda for infodemic management was developed from the first WHO’s infodemic conference.² Before infodemics can be managed, it is necessary to measure and understand them. Over the course of the COVID-19 pandemic, recent systematic reviews have shown that health infodemics, especially health misinformation, have been prevalent and far-reaching on social media before and during the pandemic.^{4, 5, 6} Depending on social media platforms, health misinformation can account for less than 1% to almost 30% of user-generated contents.⁴ Vaccine hesitancy fuelled by health misinformation has accounted for over 30% of the studies included in the systematic reviews.^{5, 6} However, given researchers from diverse backgrounds with different expertise, it is unsurprising that various theories have been used to guide studies of health infodemics.⁷ Different theories have suggested inconclusive predictors, mediators, and moderators, but scholars have constantly regarded behavioural intentions or behaviours as ultimate outcomes, yet their measurements have varied.⁷ Additionally, further research is needed to understand how online infodemics have influenced offline behavioural intentions or behaviours.² The WHO has

repeatedly called for multidisciplinary collaborations since professionals in communications, neuroscience, and digital marketing have long studied how social media have manipulated people's behaviours.^{1,8}

With the advancement in natural language processing (NLP), infodemiological research applying different NLP techniques to analyze social media data to understand public discourse—called social listening—has exponentiated. For example, the WHO has developed and deployed a “Early AI-supported Response with Social Listening” (EARS) platform to identify emerging information voids following WHO's terminologies.^{9, 10} Nonetheless, existing social listening tools, given their marketing-driven designs, need great customizations to meet the needs for infodemic social listening like the EARS platform.^{9, 10} In a public health crisis, health professionals need a tool that can efficiently harness and analyse tremendous amounts of online data to understand the public discussions in timely manners since qualitative analysis is time-consuming. Latest NLP techniques, including but not limited to topic modelling, sentiment analysis, and stance detection, have been used in infodemic social listening.^{11, 12, 13} Although improvements are still needed to decrease misclassifications in these supervised and unsupervised NLP techniques, their accuracies have been acceptable so far. These NLP techniques are commonly used as a screening layer to quickly understand public discourse at a superficial level, followed by qualitative analysis to make sense, enhance understanding, or identify information voids from the conversations. Such integrated social listening, on average, can be done on a weekly basis, along with other data sources.¹⁴

It is understandable that, in the beginning of the COVID-19 pandemic and infodemic, researchers agreed to adapting existing health theories, such as the health belief model (HBM) and social cognitive theory (SCT), and social-ecological model (SEM), and tools to overcome challenges in generating new tools given limited resources.^{2, 3} Although these health theories have been long established, most of them are developed before the existence of social media.⁸ Ubiquitous social media has changed how people consume and behave upon online health information for better or worse.⁸ Dr. Schillinger et al.'s Social media and Public Health Epidemic and Response (SPHERE) model¹⁵ and Dr. Aral's Hype Loop⁸ have demonstrated that social media have both perils and merits. That is, social media can help people make informed decisions while spreading harmful misleading

information.^{8, 15} The WHO has recommended that social listening for infodemic management should be incorporated into future pandemic preparedness.^{1, 3}

During the pandemic, social listening has mostly been reactive than proactive. Health professionals and public health organizations were rushed to debunk misinformation while competing for people's attention to urge people to follow evidence-based preventive behaviours during uncertainties.^{16, 17} Although many lessons have been learned regarding health infodemics using existing theories and tools, there is a need to carry out social listening in a systematic way based on a novel theoretical framework for health researchers. Except Dr. Aral's Hype Loop,⁸ there are limitations in current theories or frameworks developed before the existence of social media. Therefore, the objective of this paper was to propose a conceptual framework that helps monitor public discourse on social media and behaviours for future infodemiological research and possible utility of the proposed conceptual framework.

3.3 Methods

Borsboom, et al.'s theory construction methodology (TCM)¹⁸ was adapted to help develop a conceptual framework. According to TCM,¹⁸ there are five steps: (1) identification of relevant phenomena, (2) development of a proto theory, (3) formation of a formal model, (4) adequacy evaluation of the formal model, and (5) assessment of overall worth of the formal model.¹⁸ Firstly, we identified health infodemics on social media as a phenomenon of interest since we were especially interested in how online information on social media has influenced people's behavioral intentions or behaviors in a public health emergency. Next, we conducted a theory synthesis¹⁹ to develop a conceptual framework as the TCM's second and third steps were combined. We searched PubMed, Scopus, PsycINFO, and Google Scholar for theories used in reviews and original research papers written in English published in peer-reviewed journals from December 2020 to December 2022. Keywords included "social media" "online discussion," "public discourse," "behaviour," "intention," "attitude," "perception," "theory," "model," "framework" and their synonyms, but explicitly excluded "conspiracy theory" in the search. Reviews were prioritized for extractions and reading because certain theories have been commonly used in the COVID-19 related studies in health behavioral science, communications, and infodemic management. We included theories with outcomes as health behavioral intentions or behaviors at individual level and beyond. The search for relevant theories in

this process was non-exhaustive, but the results were representative of the health infodemic research conducted thus far. A total of 13 theories are included for Walker and Avant’s theory synthesis.¹⁹ After the conceptual framework was formulated, a demonstration was conducted to check and evaluate the overall conceptual framework to meet the last two steps in the TCM.¹⁸

3.4 Results

3.4.1 Synthesis of Theories

Table 3-1 shows the thirteen theories included in this study. As expected, the health belief model (HBM) has been widely employed since one systematic review reported that HBM was used in 126 quantitative studies about the COVID-19 Vaccine Hesitancy over two years.²⁰ It is also expected that some existing theories were combined or adopted by researchers to investigate complex and multifaceted health infodemics in various studies. For example, the theory of planned behaviour (TPB) itself is an extension of the theory of reasoned action (TRA).²¹ Additionally, TPB was combined with the heuristic systematic model (HSM) to create the risk information seeking and processing model (RISP),^{22, 23} or integrated with the uses and gratifications theory to investigate information-sharing behaviours.²⁴ Furthermore, Scannell et al.²⁵ weaved the social judgement theory, elaboration likelihood model of persuasion (ELM), and extended parallel process model (EPPM) to understand how persuasive COVID-19 vaccine (mis)information was to convince people, implicitly affecting their behaviours.²⁵ Overall, it has demonstrated that a theoretical approach may no longer be sufficient to address the complexity of health infodemics.

Theory/Model	Focus	Constructs
Behavioral and Social Drivers	Behavior	Confidence, Motivation, and Behavior
Capability, Opportunity and Motivation lead to Behavior	Behavior	Capability, Opportunity, Motivation, and Behavior
Elaboration Likelihood Model	Attitude or Behavior	Motivation, Ability, and Opportunity to decide Central route or Peripheral route

Extended Parallel Process Model	Behavior	Threat and Efficacy
Health Belief Model	Behavior	Perceived susceptibility, Perceived severity, Perceived benefits, Perceived barriers, Modifying variables, Cues to action, and Self-efficacy
Risk Information Seeking and Processing Model	Attitude or Behavior	Combine both theory of planned behaviour and heuristic systematic model
Social Cognitive Theory	Behavior	Behavioral capability, Observational Learning, Reinforcements, Expectations, Self-efficacy, and Reciprocal Determinism
Social Judgment Theory	Attitude	Latitude of Acceptance, Latitude of Non-commitment, and Latitude of Rejection
The Hype Loop	Behavior	Consume, Act, Sense, and Suggest
Theory of Planned Behavior	Behavior	Attitudes, Subjective norm, Perceived behavioral control, Behavioral intention, and Behavior
Theory of Reasoned Action	Behavior	Attitudes, Subjective norm, Behavioral intention, and Behavior
Transtheoretical Model	Behavior	Precontemplation, Contemplation, Preparation, Action, Maintenance, and Termination
Uses and Gratifications Theory	Behavior	Cognitive need, Affective need, Personal integrative need, Social integrative, and Tension release need

Table 3-1 Theories and models used in health infodemic research in the context of the COVID-19 pandemic.

Of these theories, several factors across theories have repeatedly been shown to affect the outcome (i.e., behaviour). Although they are described in different terms, they can be used interchangeably in most contexts. For instance, the “self-efficacy” in HBM and social cognitive theory (SCT) has shared a similar meaning with “confidence” in the behavioral and social drivers (BeSD) of vaccination, “perceived behavioral control” in TPB, and “efficacy” in EPPM. If the meaning is extended further, it can also represent “capability” in the model of capability, opportunity, and motivation lead to behavior (COM-B), “ability” in ELM, “behavioral capability” in SCT, “Act” in the Hyper Loop, and “behavioral intention” in TPB/TRA, and the Transtheoretical Model. Another group of terms—attitude, perceptions, and motivation—can also share comparable meanings, although they have different definitions in a dictionary. Five of the thirteen theories include “attitude,” another three theories consist of “motivation,” and the other two theories involve perceived variables that are associated with the outcome. In general, these words have suggested people’s views in consistent or in contrast to given health information. These terms have also suggested that there are gaps between “self-efficacy” and “(cap)ability,” “perception” and “reality,” or “subjectivity” and “objectivity.” However, it can be challenging to distinguish them because they shape each other. That is, “I believe I can do it this time (i.e., subjective self-efficacy or perception) because I did it before (i.e., an objective real action). Now I get it done (i.e., objective real action), so I know I will be able to do it next time (i.e., subjective self-efficacy or perception), with or without extra preparation or practice.” It becomes greatly interrelated and thus these two may no longer be discernible, or it is too difficult to measure them separately. Similarly, attitudes and perceptions may be indistinguishable as they both imply motivations or intentions for behavioral uptake or changes.

Although almost all theories focus on individual behaviors, factors beyond individuals are also important to be considered and yet these social determinant factors can be difficult to measure or imprecise based on self-reported measurements.²⁶⁻³² However, existing models, such as HBM, SCT, and BeSD, can incorporate variables beyond personal levels to infer the outcome. Nonetheless, unlike EPPM, these behavioral models don’t explicitly measure emotional variables, although they might be inferred in variables related to self-efficacy, perceptions, or subjective norms. One of implicit assumptions in these theories is that people can logically determine and behave to mitigate risks if they perceive greater threats or susceptibility to themselves. According to latest infodemic and social media research,^{8,9,33} unfortunately, behaviors may not be completely driven by rational reasoning;

otherwise, panic buying during the COVID-19 pandemic is not supposed to happen.³⁴ Prior studies have evidently shown how social media, given their artificial algorithm designs, can manipulate or help spread emotional posts, making it contagious at large.^{8, 35, 36} Therefore, emotion should also be taken into account when inferring behaviors, similar to perception, attitude, motivations, and others.

Given limitations and gaps identified in existing theories and frameworks, a new framework is needed to reflect the current complex infodemic issues in today's information ecosystem.^{1, 2, 3} The new conceptual framework should incorporate theories from the communication field because it will improve health professionals' understanding of public discourses. In addition, attributes measuring attitudes and emotions are included in the proposed conceptual framework: Social Media Infodemic Listening (SoMeIL) for Public Health Behavior.

3.4.2 Proposed Conceptual Framework

We propose a novel conceptual framework—SoMeIL for public health behavior (Figure 3-1)—to address these issues. Our framework aims to investigate how people's emotions and attitudes are associated with their online behaviors on social media, and eventually their offline behaviors in the real world. In other words, our proposed framework can help researchers to better understand the public discourse and to better infer collective behavioral intentions or behaviors. Double arrows illustrate potential associations these five constructs have with each other. Blurry boundaries and faded colors demonstrate that the components can happen both online and offline simultaneously. Unlike existing theories, our framework no longer assumes rational judgments and behaviors. In the following sections, we will introduce and explain each construct illustrated in our proposed conceptual framework, along with some limitations in social media data or NLP techniques when researchers use them.

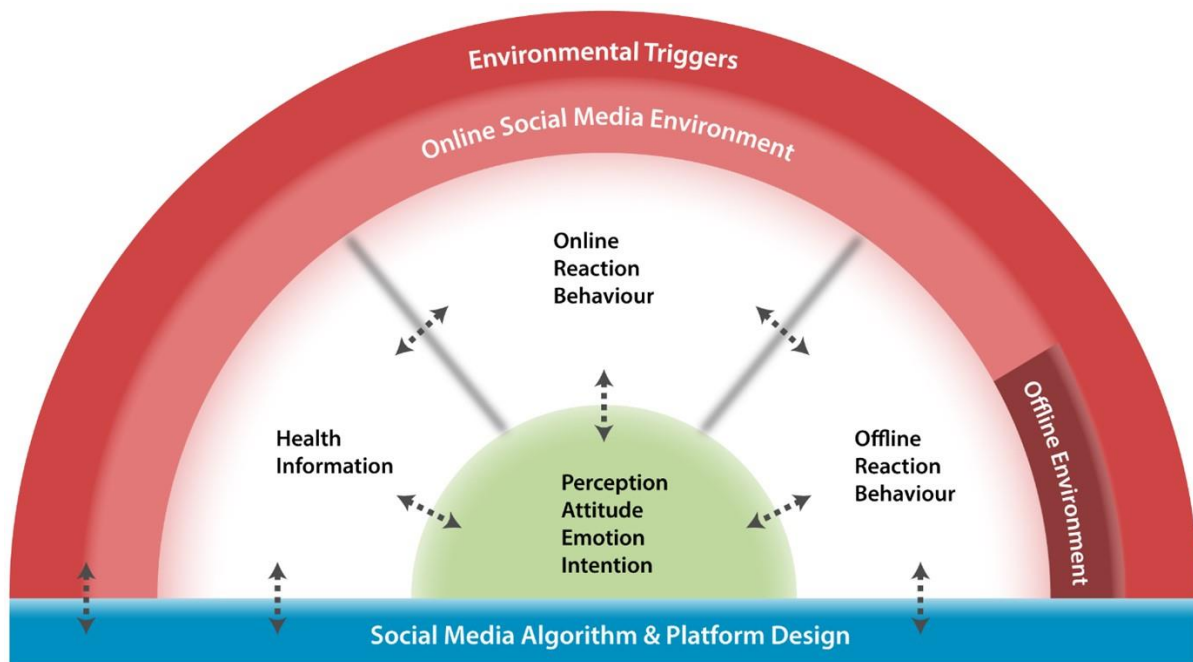


Figure 3-1: Social media infodemic listening (SoMeIL) for public health behavior conceptual framework

As Dr. Aral has demonstrated in his research,⁸ the social media algorithm is input with user attributes (Table 3-2), such as the demographics and historical behavioural data to connect friends or “recommended” posts to users based on similarity instead of heterogeneity.⁸ Studies have evidently shown that social media algorithms are intentionally designed to be addictive and affective.^{8,33} The issue is further compounded by highly personalized user experiences on social media given people’s digital footprints, encouraging echo chambers or polarization.⁸ Coupled with its engagement design, such as “like” and “follow” buttons, social media have kept their users spending more time on the platform as “engagements.”⁸ Such characteristics is defined as “user attributes on social media algorithms” in the proposed conceptual framework (Figure 3-1). However, since user attributes are voluntarily input by users when they first create their accounts, most attributes (Table 3-2) are optional, and values can be fictitious. In other words, they all can be missing data, or even untrue when values are not missing, although correct values exist. Some social media platforms require users to enter only their email and password to create an account with a username without any other details. Therefore, the users can remain primarily anonymous or unverified on the platform. Geolocations is a

special issue for researchers when modelling disease outbreaks or heat maps using Twitter data.^{37, 38, 39} For example, tweets tagged with explicit geolocation can vary from less than 1% to approximately 4% of data collected from Twitter,^{37, 38} depending on data collection methods and the amount of data collected. Although there are many ML techniques to infer geolocations for Twitter data,^{37, 39} they are not as precise or comparable as Internet Protocol (IP) addresses. Furthermore, public discussions related to vaccinations on social media have become more polarized over time, for instance.^{22, 40-43} Studies have demonstrated that user attributes, such as users' political party affiliations, religious affiliations, and who they follow (i.e., following), can potentially indicate their ideologies or attitudes toward vaccinations.⁴⁰⁻⁴⁴ Similar to the geolocation issues, researchers may not have direct access to collect these attributes. If users enter some information within these attributes, their accuracies remain uncertain. Additionally, even if researchers apply advanced ML techniques to infer these attributes, these techniques may be unable to generalize to other studies or social media platforms with different users' characteristics.⁴⁵

Components	Attributes
User attribute on social media	<ul style="list-style-type: none"> • Age • Sex • Geolocation • Income • Education • Occupation • Party affiliations • Region affiliations • Following • Followed
Inferred Intention	<ul style="list-style-type: none"> • Attitude

	<ul style="list-style-type: none"> ○ Acceptance ○ Non-commitment ○ Rejection ● Emotion <ul style="list-style-type: none"> ○ Positive ○ Negative ○ Neutral ○ Mixed ● Perception ● Ideology
Online Reaction Behavior	<ul style="list-style-type: none"> ● Share ● Like/dislike ● Comment ● Post ● Bookmark ● Nothing
Offline Reaction Behavior	<ul style="list-style-type: none"> ● Agreement ● Disagreement

Table 3-2 Attributes of each components in the SoMeIL conceptual framework

Next, we define "online reaction behaviour" as it occurs "after" a user views a social media post. We can measure collective online behaviours via the numbers of likes, shares, and others (Table 3-2). These attributes are not mutually exclusive because a person can have multiple reactions after viewing a post. Besides, we add an attribute called "nothing" to reflect that an individual may have no

reaction at all, or a reaction that is not captured by the social media platform. For example, the user may laugh so hard in reality but doesn't even "like" the post after viewing a hilarious post. The "nothing" attribute is theoretically the same as "non respondent bias" in survey research. Although there are other digital tracking tools to help infer viewers without any online reactions, researchers have been unable to directly access or retrieve such information since social media companies can decide what information can be available to researchers. We are especially interested in online behaviour, or its propagation patterns because it can be used to infer or confirm collective inferred intentions, as measured by emotions, attitudes, or perceptions. For instance, digital marketing research on Twitter has long estimated the number of users sharing similar opinions (i.e., acceptance) by the number of likes and retweets of a given tweet, whereas disagreements (i.e., rejection) can be reflected by the number of replies.⁴⁶ Dividing the latter by the former, if the resulting "Twitter ratio" is at least 0.5, it indicates positive or neutral responses, whereas below 0.5 suggests negative responses.⁴⁶ Therefore, by collecting and analysing the attributes within the online behaviours, scholars can better understand or estimate what inferred intentions of the 'quiet majority' users are since approximately 10% of users produce 90% of content on Twitter, for example.⁴⁷ Online behaviour can be used to infer people's behavioural intentions. For example, if someone tweeted that they would get COVID-19 vaccinated as soon as they became eligible, and the tweet resulted in 1,500 likes and 2,500 retweets, it was estimated approximately 3,501 pro-vaccine people. Nonetheless, the number can be an overestimation considering that reactions are not mutually exclusive, or an underestimate since Twitter users do not really represent a general population in a given region. Additionally, such estimations may not apply to other social media.

As explained in the theory synthesis, existing models have theorized that behaviours can be attributed to attitudes, perceptions, and emotions, but it has remained challenging to clearly distinguish them because they are interrelated and cannot be easily measured. Researchers have inferred associations among attitudes, perceptions, and emotions in various ways,^{11, 48} but we decide to group them together in our framework as "inferred intention" In our opinion, it is unnecessary to distinguish them since they can be used interchangeably or along with each other in different contexts. It becomes more important to infer potential behavioural intentions using attitudes, perceptions, emotions, or ideologies. we have adopted SJT to infer intentions (Table 3-2) because this makes it more feasible when using NLP techniques to analyse social media data, especially in

infodemiological studies. For example, when investigating public intentions toward COVID-19 vaccination, acceptance can be theoretically associated with pro-vaccine individuals, rejection probably suggests anti-vaccine people, and non-commitment might be regarded as a proxy for vaccine-hesitant people as evidenced by prior research.⁴⁹ Yet we acknowledge that there are limitations in this assumption, so we need to be careful in how we interpret data and ascribe intentions based on our categorization of individuals. To better understand public discourse on social media, a promising ML technique—stance detection—can be applied to infer whether or not people’s attitudes toward a give topic.^{11, 49, 50} For example, whether or not people support or oppose the COVID-19 vaccination. In addition to stance detection,⁵⁰ a common way to infer attitudes in existing infodemic studies involves topic modelling and sentiment analyses.^{11, 12, 13} Depending on models of sentiment analyses, emotions can be categorized at basic levels (i.e., positive, neutral, and negative) or more detailed levels (e.g., sad, anger, happy, joy, etc.).^{51, 52} However, according to our research experiences and other infodemic studies, sentiment analysis can still result in misclassifications regardless of levels.^{53, 54, 55} Therefore, our framework remains conventional to maintain emotions at basic levels with an additional level called “mixed” sentiment. The “mixed” attribute is added to address possible misclassifications in the “neutral” category resulting from sentiment analysis. When a tweet includes an approximately equivalent number of positive and negative words, it’s classified as “neutral” by the sentiment analysis. However, this doesn’t mean the tweet is really “neutral” because it can actually be “positive,” “negative,” or “mixed” overall, depending on its context.^{53, 54, 55} Misclassifications often occur in ironic or humorous tweets.^{53, 54, 55} The “mixed” feeling in our framework refers to an equal amount of positive and negative feelings expressed simultaneously in a tweet without being “positive” or “negative” overall. For instance, if someone tweets equal number of concerns and favours towards COVID-19 vaccines without explicit conclusions, this tweet can be regarded as “mixed” by humans, but it’s likely classified as “neutral” by sentiment analysis. Although we incorporate the ‘mixed’ attribute in our framework, we acknowledge that existing sentiment analyses have not been sophisticated or advanced enough to categorize such “mixed” feelings. In addition, even humans cannot interpret ‘mixed’ feelings consistently given external social-cultural factors, similar to humours are different in different cultures. Therefore, improvements are still needed.

For the “offline reaction behaviour” shown in Figure 3-1, although boundaries between our physical and digital worlds have become less distinguishable, it remains unclear whether or not people really react upon information received from social media. Some may have consistent online and offline reaction behaviours, another may have contradictory online and offline reaction behaviors, and others may only have either online or offline reaction behaviours. Even if individuals tweet or like a tweet indicating that they are willing to get vaccinated, it remains inconclusive unless they later share a selfie or their vaccination record on social media to prove that they, in fact, get COVID-19 vaccinated. In this case, their self-reported offline reaction behaviour matches their online reaction behaviour. Their self-reported offline reaction behaviour is also adherent to public health interventions. Therefore, one’s self-reported offline reaction behaviour can be inferred in two ways: one is whether an individual’s online and offline self-reported behaviours are consistent, and the other is whether their self-reported offline behaviour follows the public health interventions. The “offline reaction behaviour” in the COVID-19 vaccination example has been primarily self-reported if using only social media data. However, there are other data, such as administrative data, that can possibly provide directly measured “offline reaction behaviour” instead of self-reported data like social media or survey.

The initial, preliminary validation of part of the SoMeIL conceptual framework occurred in several ways. Firstly, co-authors in the study, except S-FT, have served as epidemiological subject matter experts to be consulted and contributed to the development of the SoMeIL framework. After three runs of in-depth discussions, the initial consensus of the SoMeIL framework were reached. Next, the framework was presented in the Society for Epidemiologic Research (SER) 2023 Annual meeting to collect expert feedback to further revise the framework. In addition, the framework was sent to S-FT’s connections in the WHO ‘s infodemic manager training to gather brief feedback. However, given different backgrounds and knowledge from the subject matter experts and ever-changing infodemics during the pandemic, the current SoMeIL conceptual framework have represented the minimal agreements among differ stakeholders in a very difficult time.

Our life has become more digitalized, and younger generations have tended to share their life on social media, but it has not made it easier to associate offline behaviour with online information due to privacy concerns and available social media data. Therefore, additional data resources or proxy

measures will be needed to investigate associations between online information on social media and offline behaviours at population levels. For instance, administrative or census data regarding COVID-19 vaccine administration can be linked with social media data to estimate or confirm vaccination coverage in a region. Similar approaches can be applied when there are outbreaks of infectious diseases.

3.5 Discussion

The SoMeIL conceptual framework (Figure 3-1) consists of five major constructs inspired from existing theories. Dashed boundaries indicate that online and offline environments have become less distinctive as information flows. Arrows represent potential associations among these components and how they influence or self-feed each other as the framework gives a sense of loop. Attributes of each construct can be inferred or measured via advanced NLP or ML techniques if data are available and of high quality. Although we have used NLP techniques to explain our conceptual framework throughout this paper based on our study published during the COVID-19 pandemic,⁵⁶ the proposed framework is not limited to quantitative infodemiological research only. That is, the proposed conceptual framework can be applied in qualitative research.

There are several limitations in the SoMeIL conceptual framework. Firstly, more evaluations need to be conducted since this is a new conceptual framework. Furthermore, given that the SoMeIL conceptual framework primarily focuses on social media, it is acknowledged that this proposed framework can only be useful in more digitalized populations, cultures, or nations. Besides, with new social media platforms popping up, data formats and types can change given different platform designs. Therefore, the SoMeIL framework may need to be revised to reflect and investigate non-textual data, such as videos and images. Although there are advanced NLP and ML techniques that can analyse videos and images, they have not been well adapted in current infodemiological or social listening studies. In addition, more expert reviews and inputs are needed after the pandemic. Lastly, each social media has different user characteristics. Therefore, social media data can be biased. Researchers will need to be careful when interpreting findings from different social media platforms even with our proposed conceptual framework.

Although existing health behaviours, communications, and latest infodemic theories have been used in infodemiological studies, these theories have not reflected well the distinctive nature of social

media in the current complex information ecosystems. Therefore, a novel conceptual framework—social media infodemic listening (SoMeIL) for public health behaviour—is proposed to help future infodemiological research. We acknowledge that the SoMeIL conceptual framework still needs validations for its efficacy, safety and usability. We anticipate the SoMeIL conceptual framework will be revised as more studies will be conducted in the future.

3.6 References

1. Wilhelm E, Ballalai I, Belanger M-E, Benjamin P, Bertrand-Ferrandis C, Bezbaruah S, et al. Measuring the burden of infodemics: Summary of the methods and results of the fifth WHO Infodemic Management Conference. *JMIR Infodemiology* [Internet]. 2023;3:e44207. Available from: <http://dx.doi.org/10.2196/44207>
2. Calleja N, AbdAllah A, Abad N, Ahmed N, Albarracin D, Altieri E, et al. A public health research agenda for managing infodemics: Methods and results of the first WHO infodemiology conference. *JMIR Infodemiology* [Internet]. 2021;1(1):e30979. Available from: <http://dx.doi.org/10.2196/30979>
3. Tangcharoensathien V, Calleja N, Nguyen T, Purnat T, D'Agostino M, Garcia-Saiso S, et al. Framework for managing the COVID-19 infodemic: Methods and results of an online, crowdsourced WHO technical consultation. *J Med Internet Res* [Internet]. 2020;22(6):e19659. Available from: <http://dx.doi.org/10.2196/19659>
4. Borges do Nascimento IJ, Beatriz Pizarro A, Almeida J, Azzopardi-Muscat N, André Gonçalves M, Björklund M, et al. Infodemics and health misinformation: a systematic review of reviews. *Bull World Health Organ* [Internet]. 2022;100(9):544–61. Available from: <http://dx.doi.org/10.2471/blt.21.287654>
5. Suarez-Lledo V, Alvarez-Galvez J. Prevalence of health misinformation on social media: Systematic review. *J Med Internet Res* [Internet]. 2021;23(1):e17187. Available from: <http://dx.doi.org/10.2196/17187>
6. Wang Y, McKee M, Torbica A, Stuckler D. Systematic literature review on the spread of health-related misinformation on social media. *Soc Sci Med* [Internet]. 2019;240(112552):112552. Available from: <http://dx.doi.org/10.1016/j.socscimed.2019.112552>
7. Ngai EWT, Tao SSC, Moon KKL. Social media research: Theories, constructs, and conceptual frameworks. *Int J Inf Manage* [Internet]. 2015;35(1):33–44. Available from: <http://dx.doi.org/10.1016/j.ijinfomgt.2014.09.004>
8. Aral S. *The Hype Machine: How social media disrupts our elections, our economy, and our health--and how we must adapt*. Currency; 2021.

9. Purnat TD, Vacca P, Czerniak C, Ball S, Burzo S, Zecchin T, et al. Infodemic signal detection during the COVID-19 pandemic: Development of a methodology for identifying potential information voids in online conversations. *JMIR Infodemiology* [Internet]. 2021;1(1):e30971. Available from: <http://dx.doi.org/10.2196/30971>
10. Purnat TD, Wilson H, Nguyen T, Briand S. EARS – A WHO platform for AI-supported real-time online Social Listening of COVID-19 conversations. In: *Studies in Health Technology and Informatics*. IOS Press; 2021.
11. ALDayel A, Magdy W. Stance detection on social media: State of the art and trends. *Inf Process Manag* [Internet]. 2021;58(4):102597. Available from: <http://dx.doi.org/10.1016/j.ipm.2021.102597>
12. Medhat W, Hassan A, Korashy H. Sentiment analysis algorithms and applications: A survey. *Ain Shams Eng J* [Internet]. 2014;5(4):1093–113. Available from: <http://dx.doi.org/10.1016/j.asej.2014.04.011>
13. Vayansky I, Kumar SAP. A review of topic modeling methods. *Inf Syst* [Internet]. 2020;94(101582):101582. Available from: <http://dx.doi.org/10.1016/j.is.2020.101582>
14. Boender TS, Schneider PH, Houareau C, Wehrli S, Purnat TD, Ishizumi A, et al. Establishing infodemic management in Germany: A framework for social listening and integrated analysis to report infodemic insights at the national public health institute. *JMIR Infodemiology* [Internet]. 2023 [cited 2023 Aug 2];3:e43646. Available from: <https://preprints.jmir.org/preprint/43646>
15. Schillinger D, Chittamuru D, Ramírez AS. From “infodemics” to health promotion: A novel framework for the role of social media in public health. *Am J Public Health* [Internet]. 2020;110(9):1393–6. Available from: <http://dx.doi.org/10.2105/ajph.2020.305746>
16. Nan X, Iles IA, Yang B, Ma Z. Public health messaging during the COVID-19 pandemic and beyond: Lessons from communication science. *Health Commun* [Internet]. 2022;37(1):1–19. Available from: <http://dx.doi.org/10.1080/10410236.2021.1994910>
17. Vraga EK, Jacobsen KH. Strategies for effective health communication during the Coronavirus pandemic and future emerging infectious disease events. *World Med Health Policy* [Internet]. 2020;12(3):233–41. Available from: <http://dx.doi.org/10.1002/wmh3.359>

18. Borsboom D, van der Maas HLJ, Dalege J, Kievit RA, Haig BD. Theory construction methodology: A practical framework for building theories in psychology. *Perspect Psychol Sci* [Internet]. 2021;16(4):756–66. Available from: <http://dx.doi.org/10.1177/1745691620969647>
19. Walker LO, Avant KC. *Strategies for theory construction in nursing*. 2018.
20. Limbu YB, Gautam RK, Pham L. The Health Belief Model applied to COVID-19 vaccine hesitancy: A systematic review. *Vaccines (Basel)* [Internet]. 2022;10(6):973. Available from: <http://dx.doi.org/10.3390/vaccines10060973>
21. U. S. Department Human Services, National Cancer Institute (Estats Units d'Amèrica), National Health. *Theory at a Glance: A guide for Health Promotion Practice*. North Charleston, SC: Createspace Independent Publishing Platform; 2012.
22. Yang ZJ, Aloe AM, Feeley TH. Risk information seeking and processing model: A meta-analysis: RISP meta-analysis. *J Commun* [Internet]. 2014;64(1):20–41. Available from: <http://dx.doi.org/10.1111/jcom.12071>
23. Yang JZ, Liu Z, Wong JCS. Information seeking and information sharing during the COVID-19 pandemic. *Commun Q* [Internet]. 2022;70(1):1–21. Available from: <http://dx.doi.org/10.1080/01463373.2021.1995772>
24. Malik A, Mahmood K, Islam T. Understanding the Facebook users' behavior towards COVID-19 information sharing by integrating the theory of planned behavior and gratifications. *Inf Dev* [Internet]. 2021;026666692110493. Available from: <http://dx.doi.org/10.1177/02666669211049383>
25. Scannell D, Desens L, Guadagno M, Tra Y, Acker E, Sheridan K, et al. COVID-19 vaccine discourse on twitter: A content analysis of persuasion techniques, sentiment and mis/disinformation. *J Health Commun* [Internet]. 2021;26(7):443–59. Available from: <http://dx.doi.org/10.1080/10810730.2021.1955050>
26. Beauchamp MR, Crawford KL, Jackson B. Social cognitive theory and physical activity: Mechanisms of behavior change, critique, and legacy. *Psychol Sport Exerc* [Internet]. 2019;42:110–7. Available from: <http://dx.doi.org/10.1016/j.psychsport.2018.11.009>

27. Carpenter CJ. A meta-analysis of the effectiveness of health belief model variables in predicting behavior. *Health Commun* [Internet]. 2010;25(8):661–9. Available from: <http://dx.doi.org/10.1080/10410236.2010.521906>
28. Gagne C, Godin G. The theory of planned behavior: Some measurement issues concerning belief-based variables. *J Appl Soc Psychol* [Internet]. 2000;30(10):2173–93. Available from: <http://dx.doi.org/10.1111/j.1559-1816.2000.tb02431.x>
29. J. Kitchen P, Kerr G, E. Schultz D, McColl R, Pals H. The elaboration likelihood model: review, critique and research agenda. *Eur J Mark* [Internet]. 2014;48(11/12):2033–50. Available from: <http://dx.doi.org/10.1108/ejm-12-2011-0776>
30. Marks DF. The COM-B system of behaviour change: Properties, problems and prospects. *Qeios* [Internet]. 2020; Available from: <http://dx.doi.org/10.32388/u5mttb.2>
31. Nigg CR, Geller KS, Motl RW, Horwath CC, Wertin KK, Dishman RK. A research agenda to examine the efficacy and relevance of the Transtheoretical Model for physical activity behavior. *Psychol Sport Exerc* [Internet]. 2011;12(1):7–12. Available from: <http://dx.doi.org/10.1016/j.psychsport.2010.04.004>
32. Schunk DH, DiBenedetto MK. Motivation and social cognitive theory. *Contemp Educ Psychol* [Internet]. 2020;60(101832):101832. Available from: <http://dx.doi.org/10.1016/j.cedpsych.2019.101832>
33. Azer J, Blasco-Arcas L, Harrigan P. #COVID-19: Forms and drivers of social media users' engagement behavior toward a global crisis. *J Bus Res* [Internet]. 2021;135:99–111. Available from: <http://dx.doi.org/10.1016/j.jbusres.2021.06.030>
34. Naeem M. Do social media platforms develop consumer panic buying during the fear of Covid-19 pandemic. *J Retail Consum Serv* [Internet]. 2021;58(102226):102226. Available from: <http://dx.doi.org/10.1016/j.jretconser.2020.102226>
35. Oh S-H, Lee SY, Han C. The effects of social media use on preventive behaviors during infectious disease outbreaks: The mediating role of self-relevant emotions and public risk perception. *Health Commun* [Internet]. 2021;36(8):972–81. Available from: <http://dx.doi.org/10.1080/10410236.2020.1724639>

36. Steinert S. Corona and value change. The role of social media and emotional contagion. *Ethics Inf Technol* [Internet]. 2021;23(S1):59–68. Available from: <http://dx.doi.org/10.1007/s10676-020-09545-z>
37. Dredze M, Paul MJ, Bergsma S, Tran H. Carmen: A Twitter Geolocation System with Applications to Public Health [Internet]. 2013. Available from: <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.309.6126>
38. Sloan L, Morgan J. Who tweets with their location? Understanding the relationship between demographic characteristics and the use of geoservices and geotagging on twitter. *PLoS One* [Internet]. 2015;10(11):e0142209. Available from: <http://dx.doi.org/10.1371/journal.pone.0142209>
39. Stock K. Mining location from social media: A systematic review. *Comput Environ Urban Syst* [Internet]. 2018;71:209–40. Available from: <http://dx.doi.org/10.1016/j.compenvurbsys.2018.05.007>
40. Jiang X, Su M-H, Hwang J, Lian R, Brauer M, Kim S, et al. Polarization over vaccination: Ideological differences in Twitter expression about COVID-19 vaccine favorability and specific hesitancy concerns. *Soc Media Soc* [Internet]. 2021;7(3):205630512110484. Available from: <http://dx.doi.org/10.1177/20563051211048413>
41. Schmidt AL, Zollo F, Scala A, Betsch C, Quattrociocchi W. Polarization of the vaccination debate on Facebook. *Vaccine* [Internet]. 2018;36(25):3606–12. Available from: <http://dx.doi.org/10.1016/j.vaccine.2018.05.040>
42. Mønsted B, Lehmann S. Characterizing polarization in online vaccine discourse—A large-scale study. *PLoS One* [Internet]. 2022;17(2):e0263746. Available from: <http://dx.doi.org/10.1371/journal.pone.0263746>
43. Rathje S, He JK, Roozenbeek J, Van Bavel JJ, van der Linden S. Social media behavior is associated with vaccine hesitancy. *PNAS Nexus* [Internet]. 2022;1(4). Available from: <http://dx.doi.org/10.1093/pnasnexus/pgac207>

44. Yuan X, Schuchard RJ, Crooks AT. Examining emergent communities and social bots within the polarized online vaccination debate in Twitter. *Soc Media Soc* [Internet]. 2019;5(3):205630511986546. Available from: <http://dx.doi.org/10.1177/2056305119865465>
45. Shakeri Hossein Abad Z, Butler GP, Thompson W, Lee J. Physical activity, sedentary behavior, and sleep on Twitter: Multicountry and fully labeled public data set for digital public health surveillance research. *JMIR Public Health Surveill* [Internet]. 2022;8(2):e32355. Available from: <http://dx.doi.org/10.2196/32355>
46. Salamander G. Your tweet got ratioed, what next? [Internet]. *eclincher*. Gil Salamander; 2022 [cited 2023 Aug 2]. Available from: <https://eclincher.com/your-tweet-got-ratioed-what-next>
47. Understanding users of social networks [Internet]. HBS Working Knowledge. 2009 [cited 2023 Aug 2]. Available from: <https://hbswk.hbs.edu/item/understanding-users-of-social-networks>
48. Bahamonde-Birke FJ, Kunert U, Link H, Ortuzar J de D. About attitudes and perceptions: Finding the proper way to consider latent variables in discrete choice models. *SSRN Electron J* [Internet]. 2015; Available from: <http://dx.doi.org/10.2139/ssrn.2603218>
49. Nyawa S, Tchuente D, Fosso-Wamba S. COVID-19 vaccine hesitancy: a social media analysis using deep learning. *Ann Oper Res* [Internet]. 2022; Available from: <http://dx.doi.org/10.1007/s10479-022-04792-3>
50. Cao R, Luo X, Xi Y, Qiao Y. Stance detection for online public opinion awareness: An overview. *Int J Intell Syst* [Internet]. 2022;37(12):11944–65. Available from: <http://dx.doi.org/10.1002/int.23071>
51. Nandwani P, Verma R. A review on sentiment analysis and emotion detection from text. *Soc Netw Anal Min* [Internet]. 2021;11(1). Available from: <http://dx.doi.org/10.1007/s13278-021-00776-6>
52. Zadra JR, Clore GL. Emotion and perception: the role of affective information: Emotion and perception. *Wiley Interdiscip Rev Cogn Sci* [Internet]. 2011;2(6):676–85. Available from: <http://dx.doi.org/10.1002/wcs.147>

53. Davis JJ, O'Flaherty S. Assessing the accuracy of automated Twitter sentiment coding. *Academy of Marketing Studies Journal*, suppl.Special Issue 2012;16:35-50.
54. He L, Yin T, Zheng K. They May Not Work! An evaluation of eleven sentiment analysis tools on seven social media datasets. *J Biomed Inform [Internet]*. 2022;132(104142):104142. Available from: <http://dx.doi.org/10.1016/j.jbi.2022.104142>
55. Lee S, Ma S, Meng J, Zhuang J, Peng T-Q. Detecting sentiment toward emerging infectious diseases on social media: A validity evaluation of dictionary-based sentiment analysis. *Int J Environ Res Public Health [Internet]*. 2022;19(11):6759. Available from: <http://dx.doi.org/10.3390/ijerph19116759>
56. Huang S-H, Tsao S-F, Chen H, Bin Noon G, Li L, Yang Y, et al. Topic modelling and sentiment analysis of tweets related to Freedom Convoy 2022 in Canada. *Int J Public Health [Internet]*. 2022;67. Available from: <http://dx.doi.org/10.3389/ijph.2022.1605241>

Chapter 4: A computer-assisted qualitative validation and demonstration of the SoMeIL conceptual framework

Status: Currently under journal review.

Authors: Shu-Feng Tsao, Helen Chen, Zahid A. Butt

The Theory Construction Methodology (TCM) serves as a structured framework for the development and evaluation of new theories. When a novel theory is introduced, it undergoes a rigorous process of assessment to gauge its adequacy and associations. This evaluative phase is pivotal in ensuring that the theory is not only comprehensive but also reflects the complexities of the subject matter accurately. If the assessments yield new insights or reveal differences from existing knowledge, researchers have the opportunity to revisit and revise the theory. This iterative process within TCM continues until the theory achieves stability and demonstrates its robustness in explaining the phenomena it seeks to address. In the context of this discussion, TCM is the adopted method for developing the Social Media Infodemic Listening (SoMeIL) for public health behaviours conceptual framework. Once the framework is proposed in a previous study (Chapter 3), the subsequent step is its validation. Validation is a critical phase in the research process since it verifies the SoMeIL framework's utility and applicability. This validation process can take two major types: qualitative and quantitative. Both approaches serve distinct purposes and offer unique insights into the SoMeIL conceptual framework.

Qualitative validation methods are rooted in the exploration of the depth and richness of the conceptual framework. They provide a nuanced understanding of the framework's relevance and its ability to capture the complexity of the subject matter. Several qualitative validation techniques are available, each offering a specific angle of scrutiny. Firstly, qualitative observation involves researchers using their five senses to collect data. This approach allows them to immerse themselves in the phenomenon under investigation, gaining firsthand insights. For instance, in the context of SoMeIL, researchers might observe how individuals engage with social media platforms and how this relates to their information literacy skills. Secondly, qualitative case studies involve an in-depth exploration of a complex phenomenon in the real world. Researchers identify various factors and variables interacting with each other, providing a holistic view of the subject. In the SoMeIL

framework, a case study could involve examining specific components of the conceptual framework where they interact and reveal unique insights. Next, interviews and focus groups offer opportunities to gather qualitative data directly from participants. Through open-ended questions and discussions, researchers can dive deeper into the social phenomena being studied. In the context of SoMeIL conceptual framework, interviews could be used to ask individuals or subject matter experts about their experiences and challenges in navigating social media for information consumption and their behaviours. Expert reviews are another valuable qualitative method for gathering insights and consensus from subject matter experts. In the case of SoMeIL, experts in health communication, behavioural sciences, social media, and related fields can review the framework and provide feedback, helping refine and validate it. Last but not least, grounded theory is a qualitative research method that involves developing or validating a framework by analyzing various qualitative data. Researchers identify patterns and themes within the data, allowing the framework to emerge organically. However, the grounded theory typically assumes that researchers approach the analysis without being heavily influenced by prior knowledge or experiences, which may not align with the context of this PhD research. Grounded theory necessitates an open-minded approach, free from preconceived notions. Since the researchers have acknowledged influences from prior knowledge and experiences, opting for grounded theory could lead to biased results. Instead, a case study approach is chosen to conduct a preliminary qualitative validation in this chapter. Case studies provide the flexibility to explore specific instances within the SoMeIL framework while accommodating the researchers' prior knowledge and experiences. This approach allows for a focused examination of how the conceptual framework applies to real-world scenarios, providing valuable insights into its practicality and relevance.

Although qualitative validation methods offer depth and context, quantitative validation methods bring a different dimension to the assessment of a conceptual framework. These methods focus on numerical data and statistical analysis to gauge the framework's reliability and generalizability. Quantitative validation methods may include surveys, experiments, or statistical analyses of data. These approaches will help quantify the relationships and associations proposed by the SoMeIL framework, providing empirical evidence of its effectiveness. Therefore, this chapter presents a quantitative validation of part of the SoMeIL conceptual framework. The partial components of the SoMeIL conceptual framework involve in this study include online reaction behaviours, emotion,

intention, and self-reported COVID-19 vaccination behaviour as the offline reaction behaviour. The health information in this study is the massive governmental vaccination campaigns that encourage people to take the first dose of the COVID-19 vaccines in Canada. This study aims to validate how online reaction behaviour and emotion can be associated with the vaccination intention, which is further associated with the self-reported offline reaction behaviour.

In summary, the validation of the partial SoMeIL conceptual framework is a multifaceted process that encompasses both qualitative and quantitative methods in this PhD research. The choice of a preliminary qualitative validation through a case study approach is made to accommodate the researchers' prior knowledge and experiences while providing a rich understanding of how the framework operates in real-world scenarios in this chapter. Subsequently, in Chapter 5, the quantitative validation will further strengthen the framework's credibility by using statistical analysis. Through both qualitative and quantitative validations, the SoMeIL framework will emerge as a robust tool for understanding and addressing the complex interactions of information circulating on social media and people's public health behaviours in today's digital age. The following is the abstract of this preliminary qualitative study, followed by its full-text manuscript.

4.1 Abstract

Background

As the number of health infodemic and social listening research has increased, different approaches have been used in such studies. The social media Infodemic listening (SoMeIL) for public health behaviours conceptual framework has been proposed as a systematic way or theoretical lens to conduct similar investigations in the future. Given the novelty of the SoMeIL conceptual framework, validations are needed. Therefore, this study aims to provide preliminary validation for partial components in the framework. Another objective is to demonstrate the application of the SoMeIL conceptual framework.

Methods

An existing clean Twitter dataset about the Canadian Freedom Convoy is used. It includes 560,140 unique English tweets from 15 January to 14 February 2022. The Latent Dirichlet Allocation (LDA)

topic modelling and qualitative thematic analysis are employed to infer attitudes of Twitter users discussing the convoy.

Findings

The LDA topic modelling has generated five themes. Four of them have appeared to be in favour of the convoy. A random sample of 500 tweets from each topic were used in the qualitative thematic analysis. The results of the qualitative thematic analysis have shown voices against the convoy within each topic generated from the LDA topic modelling. The major difference between those supporting and opposing the convoy is their ideology of freedom.

Interpretation

The study has preliminarily validated partial components of the SoMeIL conceptual framework, including inferring intentions by identifying public attitudes from public discourse regarding the convoy. It has also demonstrated how the SoMeIL conceptual framework can be applied by using both the LDA topic modelling and qualitative thematic analysis. This can help researchers to better understand public discourse on Twitter by leveraging strengths from both quantitative and qualitative methods.

Funding

This study was supported by the 2023-24 Ontario Graduate Scholarship.

Keywords: conceptual framework, topic modelling, qualitative thematic analysis, Twitter, validation

4.2 Introduction

Since the COVID-19 pandemic, social media data have been intensively studied for researchers to better understand public discourse, attitudes, opinions, perceptions, and so on. Equipped with advanced computing power and natural language processing (NLP) or artificial intelligence (AI) techniques, researchers can analyse a sheer volume of social media data efficiently.^{1,2,3} Topic modelling and sentiment analysis have been two common analyses in such social media research. Out of various topic modelling, the Latent Dirichlet Allocation (LDA) based topic modelling has been widely applied with different customizations.^{1,2,3} For example, LDA topic modelling accounted for 155 out of 193 articles, or 80% of the total articles included in a systematic review that surveyed

current topic modelling techniques used for social media data analysis.³ Similarly, qualitative researchers have used comparable techniques, such as thematic analysis and discourse analysis, to investigate public discourse on social media.⁴ However, qualitative social media research is generally limited to smaller data sizes than quantitative studies which collect and analyse millions, if not billions, of data. Therefore, researchers have investigated how to combine both topic modelling and qualitative analysis together for automated topic extractions and qualitative thematic analysis for in-depth explorations and interpretations.⁴⁻¹⁰ This integrated approach allows for a more comprehensive analysis of complex social media data since researchers can gain a richer understanding of social media content, including the identification of both latent topics and nuanced qualitative themes.⁵⁻¹⁰ The integration of topic modelling and qualitative analysis also increases research efficiency since LDA topic modelling automates the initial topic identification process, making it efficient for handling large datasets.⁶⁻⁹ This automation can save time and resources compared to manual coding of all data in conventional qualitative analysis.⁶⁻⁹ The efficiency is especially important for an infectious disease outbreak, such as COVID-19, and corresponding crisis communications in today's complex information ecosystem as proposed in the Social Media Infodemic Listening (SoMeIL) for public health behaviours conceptual framework.¹¹ Overall, by using both LDA topic modelling and qualitative thematic analysis, researchers can leverage the strengths of both quantitative efficiency and qualitative depth.

As demonstrated in the SoMeIL conceptual framework,¹¹ it is crucial to have a systematic approach to better understand public discourse on social media efficiently with assistance from the latest NLP or AI techniques. According to the SoMeIL framework,¹¹ when people encounter health information on social media, researchers can investigate how people react to the given information online or offline. Since the SoMeIL conceptual framework¹¹ is new, it is essential to validate the framework according to the theory construction methodology (TCM).¹² Therefore, this study aimed to preliminarily validate the SoMeIL conceptual framework with partial components, and to demonstrate the preliminary utility of the framework. Following prior studies,⁴⁻¹⁰ the LDA topic modelling and qualitative thematic analysis are used to infer Canada people's attitude toward the Canadian Freedom Convoy on Twitter.¹³ More specifically, the health information in this case is the tweets related to the 2022 Canadian Freedom Convoy.¹³ People's attitudes and any self-reported

online or offline reaction behaviours are inferred via the LDA topic modelling and qualitative thematic analysis.

4.3 Methods

4.3.1 Data Collection and Preprocessing

Data came from an existing cleaned Canada Freedom Convoy dataset from a different study.¹³ The clean dataset consists of 560,140 unique English tweets from 15 January to 14 February 2022, when the Freedom Convoy occurred in Canada in 2022.¹³ Several steps were involved to prepare these tweets for the LDA topic modelling. Firstly, all tweets were converted to lowercase, and Unicode strings were converted to ASCII. This conversion removed non-text elements, such as URLs, punctuation, special characters, emojis, emoticons, numbers within words, numbers in sentences, and additional spaces. Repeating character sequences longer than three characters, such as "hahaha," were condensed to three-letter sequences like "hah."¹³ Secondly, common stop words were removed using the spaCy English stop words dictionary.¹⁴ Then, the tweets were lemmatised using the WordNetLemmatizer¹⁵ and stemmed using the PorterStemmer¹⁵. Words were further grouped into phrases using the Phrases tool from Python's Gensim package.¹⁵ The next step involved tokenization using SKLearn's CountVectorizer¹⁶ to count the occurrences of each token in the dataset. Terms that appeared in more than 90% of the total dataset or words that appeared less than ten times were filtered out to eliminate both frequent and infrequent terms according to existing literature.¹³ During this manual review, synonyms of keywords that emerged across all topics were identified and treated as additional stop words. This was done to enhance the clarity and interpretation of the emerging topics during further topic modelling and data analysis.

4.3.2 Data Analysis

Unsupervised LDA topic modelling was applied using Python's SKLearn package¹⁷ to identify potential keywords for various topics. The topic modelling process was configured to search for up to 15 topics. Topic optimization was performed based on coherence scores. Once the number of the optimal topics generated from the LDA topic modelling is identified, a random sample of 500 tweets from each topic is selected for the following qualitative thematic analysis to better understand the overall context in the Freedom Convoy discourse. There are six steps in the qualitative thematic

analysis: (1) familiarised with the data, (2) generate initial codes, (3) identify themes, (4) review the identified themes, (5) define the identified themes, and finally (6) report the results.¹⁸⁻¹⁹

4.3.3 Ethical Approval

This study was approved by the University of Waterloo Office of Research Ethics (#43961).

4.3.4 Funding

This study was supported by the 2023-24 Ontario Graduate Scholarship.

4.4 Results

Table 4-1 shows the five topics generated from the LDA topic modelling along with the top 15 frequent words in each topic. These five topics are chosen based on the research question. That is, public discourse and inferred attitude towards the Freedom Convoy movement and researchers' interpretations with corresponding keywords from each topic. Typically, coherence scores tend to rise with an increase in the number of topics. Nevertheless, as the number of topics grows, they have become finer-grained clusters that correspond to minor events within the convoy movement. This trend leads to fewer topics effectively capturing the broader overarching themes of the movement. Therefore, the optimal number of topics is determined to be five in this case.

Topic	Top 15 Words	Inferred Attitudes
1	support, covid, stand, cdnpoli, govern, arrest, driver, news, call, peopl, rally, terrorist, tyranni, time, today	Convoy supporters blame the media that negatively label the convoy.
2	trudeau, like, world, video, thank, peac, honkhonk, love, look, speak, movement, share, power, lie, flag	Convoy supporters blame the politicians who disagree with the convoy.
3	mandat, protest, report, live, end, start, vaccin_mandate, want, stop, govern, vaccin, country, way, american, ottawa	Convoy supporters argue that the vaccine mandate should be stopped.

4	ottawa, polic, day, come, break, ontario, weekend, week, protestor, head, kid, help, citi, thousand, actual	Convoy opponents dissatisfy with the police enforcement and blame the supporters bringing their kids.
5	peopl, right, medium, know, need, donat, think, thing, gofundm, go, want, let, fund, watch, organ	Convoy supporters call for donations.

Table 4-1: Latent Dirichlet allocation topic modeling results

After the five themes are identified from the LDA topic modelling, 500 tweets from each topic are randomly selected, resulting in a total of 2,500 tweets for the qualitative thematic analysis. Figure 4-1 shows the themes resulting from the qualitative thematic analysis. Except the fourth topic, all the other topics generated from the LDA topic modelling have generally showed supporting attitudes for the convoy. However, in these topics, there are tweets against such attitudes, but the LDA topic modelling have grouped them together given similar frequent words. The qualitative thematic analysis has been helpful to find different voices in each topic resulted from the LDA topic modelling.

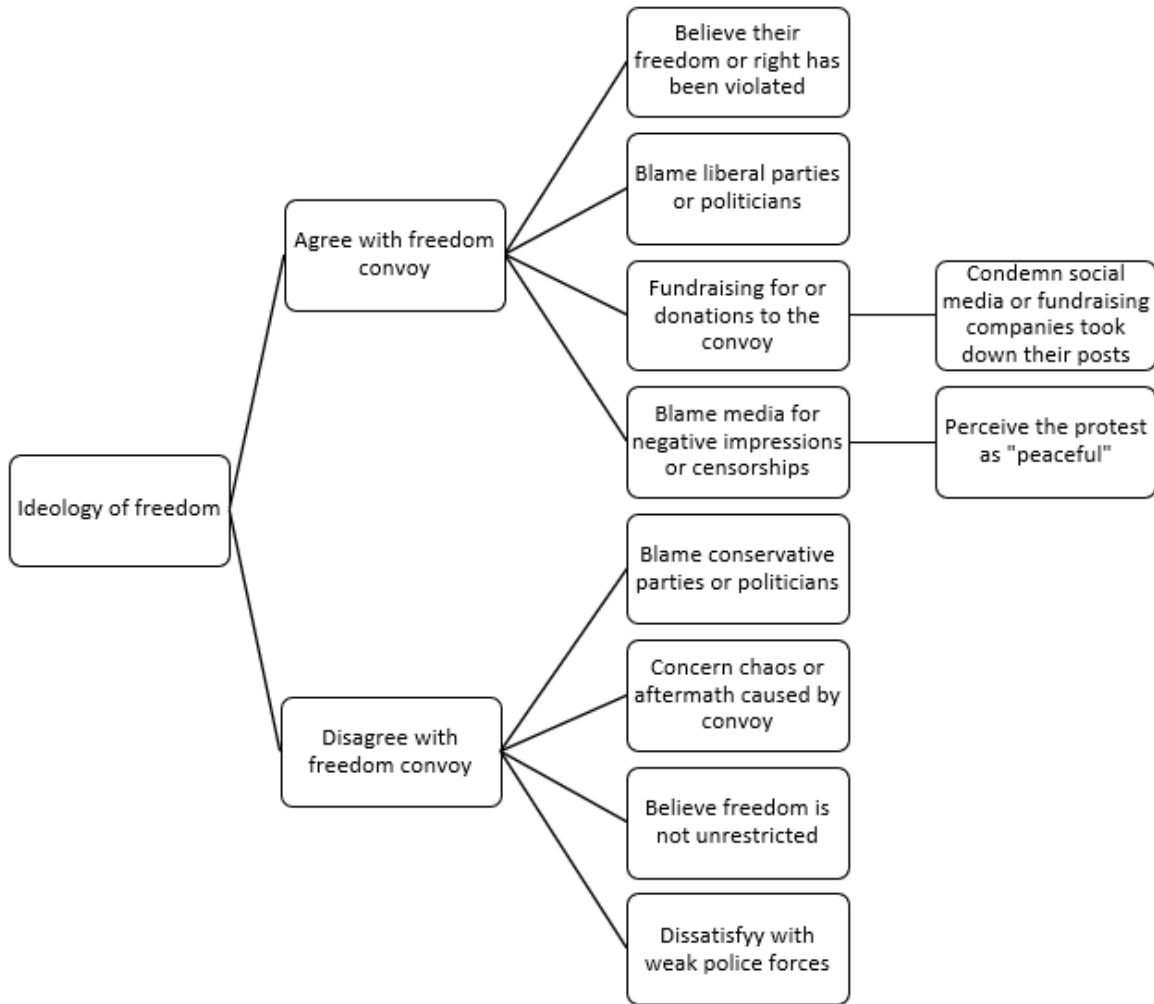


Figure 4-1: Qualitative thematic analysis results

The underlying difference between those who support and oppose the convoy is the ideology of freedom. That is, people have different interpretations of their fundamental right of freedom. The convoy supporters have strongly believed their right of freedom have been severely violated by the vaccine mandate implemented by the governments. In contrast, convoy opponents have argued that individual freedom should not be unlimited, so people should not put individuals before others.

4.5 Discussion

In the example of Canadian Freedom Convoy, the ideology of freedom has emerged as an overarching theme resulted from the qualitative thematic analysis, illustrating an explicit difference in freedom perceptions between the convoy supporters and opponents. The results of the qualitative thematic analysis actually appeared in each topic generated from the LDA topic modelling since a tweet can have more than one themes coded during the qualitative thematic analysis. In other words, themes emerged from the qualitative thematic analysis has provided an overall context across all topics generated from the LDA topic modelling. Among the convoy advocates, they strongly believed the vaccine mandates has fundamentally violated their fundamental right of freedom, as one tweet illustrated:

*...the Convoy stands for freedom, and a democracy, as opposed to what we are experiencing....
Borderline China! No more mandates!*

Similarly, some tweets articulated that they were not “anti-vaxxers” as news media labelled them. They said that they opposed “mandates” rather than against “vaccines,” and the whole point, the protest supporters argued, was that they should have freedom to choose or decide for themselves rather than being forced to follow whatever the government had them do, such as the “zero COVID-19 policy” imposed by China on its own people. In contrast, people against the convoy implied such protesters selfish and believed that the freedom was not unrestricted, as one tweet explained:

The thing is, the Charter has been tested over things like public health measures and workplace safety before, and it's not a right to do what you want, all the time, anywhere, with anyone.

Given different viewpoints, it is not uncommon that each side criticized politicians or political parties at the other side since the ideology of freedom is also shaped by political affiliations. The protest proponents vilified the liberal politicians or parties, especially pointing to the Canada prime minister, as the following tweet demonstrated:

When you allow a sick leader's ego to hold a country hostage you become a fascist sympathizer. My mother lived under Nazis occupation. She said the righteous ones are always the real fascists

The same logic applied to people against the convoy as they blamed the conservative politicians or parties that supported the protest: “*Shameful for the Conservatives to embrace these loons. And not*

surprising.” Furthermore, people who supported the convoy criticized news media for negative reports or coverages by, for example, arguing:

The media coverage of the #FreedomConvoy2022 has largely been nothing short of disgraceful disinformation. Having failed to ignore the convoy into oblivion, @XXX + the leftist press seized upon a dishonest and cliché narrative and set out to tell the story accordingly.

In other words, the convoy supporters defended their protest not as terrible as the media painted them, and the media used political propaganda against them by imposing disproportionately negative biases. Some went even further to tweet the convoy was just “peaceful” protesters as opposed to “violent lawbreakers,” “fringe minority,” “far-right extremists,” or “white supremacists” as the media called them. Given the perception of “peaceful” protest, people backing the convoy also condemned the state of emergency declared by the Ottawa Mayor to use increasing police forces to break down the protest. Compared to the “peaceful” convoy perceived by the supporters, the convoy opponents, consistent with the media, considered the protest “violent” and caused numerous harmful chaos or aftermaths. For instance, one tweet showed its disapproval:

So Canada's 'Freedom Convoy', opposing vax mandate for truckers,

- Harrassed homeless shelter soup kitchen demanding food & assaulted a homeless person,

- Defaced the Terry Fox Memorial Statue by draping it with an upside down CA

- Stood on The Tomb of The Unknown Soldier

WTAF?!

People against the convoy also criticized the police. However, unlike criticisms from the convoy supporters, people who disagreed with the protest perceived insufficient police enforcement to contain the protest before the state of emergency was declared. Although these people acknowledged that the truckers had the right to protest, the protesters should not just block Ottawa the way they wanted and disrupted residents’ daily lives. In addition, they condemned the police who showed supports for the convoy by not seriously enforcing the laws or donating to the convoy, just similar to the convoy proponents tweeted their supports via donations: “*In for \$20, wish i could kick in a lot more...a patriot donated 10k anonymously, God Bless you!*” In the meantime, there were fundraising

tweets calling for donations. When the fundraising post was suspended and donations were frozen, the convoy supporters unleashed their anger and tried to find other alternatives. They heavily criticized social media that censored their posts, in addition to the fundraising websites that took down their pages. On the contrary, the convoy opponents showed support for donation suspensions.

The example has showed that using both the NLP topic modelling and qualitative thematic analysis, researchers could better understand public discourse given its contexts on a social media in a relatively shorter period compared to only thematic analysis or other qualitative studies. Furthermore, the demonstration has shown how to apply the conceptual framework to better understanding, or “social media listening,” of an event and to infer behavioural intentions more efficiently. The convoy supporters could be assumed that they had been relatively resistant to the COVID-19 vaccinations than the convoy opponents. In other words, the protesters were less likely to get vaccinated than others as evidence by their perceived ideology of freedom.

This study has also preliminarily validated partial components of the SoMeIL conceptual framework,¹¹ including the inferred intention by identifying attitudes (i.e., agree or disagree with the convoy) and self-reported offline reaction behaviours posted on Twitter. For the convoy supporters, they have self-reported donations to the fundraising webpage in real life. On the contrary, people against the convoy have mostly complained the aftermath caused by the convoy on Twitter instead of any self-reported offline reaction behaviours in the real life, such as organizing a counterprotest toward the convoy.

There are several limitations of this study. Firstly, there could be biases in the chosen keywords used for data collection. As Table 4-1 have showed most topics were leaning towards support for the convoy, it is likely those not in favor of the convoy might not use similar keywords or hashtags. Another limitation is that only English tweets were included in the study. Therefore, the results could not be generalised other social media. Furthermore, there are different topic modelling techniques. Topics generated from other techniques could group topics differently than the LDA topic modelling used in this study. Similarly, other researchers might have identified and coded themes differently in the qualitative thematic analysis. Nonetheless, this study has demonstrated that using both LDA topic modelling and qualitative thematic analysis can better understand the overall public discussions

efficiently by leveraging strengths from both quantitative and qualitative methods. In addition, the study has provided preliminary evidence to partially validate the SoMeIL conceptual framework.¹¹

For future qualitative validation studies, it is expected that other components of the SoMeIL conceptual framework will be validated. For example, for health information, researchers can investigate how information formatting, quality, credibility, or moderation have influenced people's online or offline reaction behaviours. Researchers can also investigate how social media algorithms or platform designs have attracted different types of users and drive the engagements on social media. This can help health professionals or organizations to tailor communication strategies on different social media platforms given different platform interfaces and user profiles. Furthermore, future research can study how the SoMeIL conceptual framework can be applied to different countries or cultures with adjustments.

4.6 Conclusion

In summary, the study has preliminarily validated partial components of the SoMeIL conceptual framework. It has also demonstrated how to apply the SoMeIL conceptual framework by leveraging LDA topic modelling and qualitative thematic analysis. Five topics have resulted from the LDA topic modelling, but the qualitative thematic analysis has provided further contexts in these topics by inferring the attitudes and corresponding reactions.

4.7 References

1. Rana TA, School of Computer Sciences, Universiti Sains Malaysia (USM), Malaysia, Cheah Y-N, Letchmunan S, School of Computer Sciences, Universiti Sains Malaysia (USM), Malaysia, School of Computer Sciences, Universiti Sains Malaysia (USM), Malaysia. Topic modeling in sentiment analysis: A systematic review. *J ICT Res Appl* [Internet]. 2016;10(1):76–93. Available from: <http://dx.doi.org/10.5614/itbj.ict.res.appl.2016.10.1.6>
2. Sandhiya R, Boopika AM, Akshatha M, Swetha SV, Hariharan NM. A review of topic modeling and its application [Internet]. *Handbook of Intelligent Computing and Optimization for Sustainable Development*. Wiley; 2022. p. 305–22. Available from: <http://dx.doi.org/10.1002/9781119792642.ch15>
3. Laureate CDP, Buntine W, Linger H. A systematic review of the use of topic models for short text social media analysis. *Artif Intell Rev* [Internet]. 2023; Available from: <http://dx.doi.org/10.1007/s10462-023-10471-x>
4. Fu J, Li C, Zhou C, Li W, Lai J, Deng S, et al. Methods for analyzing the contents of social media for health care: Scoping review. *J Med Internet Res* [Internet]. 2023;25:e43349. Available from: <http://dx.doi.org/10.2196/43349>
5. Ancheta JR, Sy C, Maceda L, Oco N, Roxas R. Computer-assisted thematic analysis of Typhoon Fung-Wong tweets. In: *TENCON 2017 - 2017 IEEE Region 10 Conference*. IEEE; 2017.
6. Nikolenko SI, Koltcov S, Koltsova O. Topic modelling for qualitative studies. *J Inf Sci* [Internet]. 2017;43(1):88–102. Available from: <http://dx.doi.org/10.1177/0165551515617393>
7. Brookes G, McEnery T. The utility of topic modelling for discourse studies: A critical evaluation. *Discourse Stud* [Internet]. 2019;21(1):3–21. Available from: <http://dx.doi.org/10.1177/1461445618814032>
8. Jacobs T, Tschötschel R. Topic models meet discourse analysis: a quantitative tool for a qualitative approach. *Int J Soc Res Methodol* [Internet]. 2019;22(5):469–85. Available from: <http://dx.doi.org/10.1080/13645579.2019.1576317>

9. Gillies M, Murthy D, Brenton H, Olaniyan R. Theme and topic: How qualitative research and topic modeling can be brought together. 2022; Available from: <http://dx.doi.org/10.48550/ARXIV.2210.00707>
10. Golos AM, Guntuku SC, Piltch-Loeb R, Leininger LJ, Simanek AM, Kumar A, et al. Dear Pandemic: A topic modeling analysis of COVID-19 information needs among readers of an online science communication campaign. *PLoS One* [Internet]. 2023;18(3):e0281773. Available from: <http://dx.doi.org/10.1371/journal.pone.0281773>
11. Tsao S-F, Chen H, Meyer S, Butt ZA. Proposing a conceptual framework: social media listening for public health behavior. 2023; Available from: <http://dx.doi.org/10.48550/ARXIV.2308.02037>
12. Borsboom D, van der Maas HLJ, Dalege J, Kievit RA, Haig BD. Theory construction methodology: A practical framework for building theories in psychology. *Perspect Psychol Sci* [Internet]. 2021;16(4):756–66. Available from: <http://dx.doi.org/10.1177/1745691620969647>
13. Huang S-H, Tsao S-F, Chen H, Bin Noon G, Li L, Yang Y, et al. Topic modelling and sentiment analysis of tweets related to Freedom Convoy 2022 in Canada. *Int J Public Health* [Internet]. 2022;67. Available from: <http://dx.doi.org/10.3389/ijph.2022.1605241>
14. Honnibal M, Montani I. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. 2017.
15. Rehurek R, Sojka P. Gensim–python framework for vector space modelling. NLP Centre, Faculty of Informatics, Masaryk University, Brno, Czech Republic. 2011;3(2).
16. Maier D, Waldherr A, Miltner P, Wiedemann G, Niekler A, Keinert A, et al. Applying LDA topic modeling in communication research: Toward a valid and reliable methodology. *Commun Methods Meas* [Internet]. 2018;12(2–3):93–118. Available from: <http://dx.doi.org/10.1080/19312458.2018.1430754>
17. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R, Dubourg V, Vanderplas J. Scikit-learn: Machine learning in Python. *the Journal of machine Learning research*. 2011 Nov 1;12:2825-30.

18. Maguire M, Delahunt B. Doing a thematic analysis: A practical, step-by-step guide for learning and teaching scholars. *AISHE-J* [Internet]. 2017 [cited 2023 Sep 15];9(3). Available from: <http://ojs.aishe.org/index.php/aishe-j/article/view/335>
19. Kiger ME, Varpio L. Thematic analysis of qualitative data: AMEE Guide No. 131. *Med Teach* [Internet]. 2020;42(8):846–54. Available from: <http://dx.doi.org/10.1080/0142159x.2020.1755030>

Chapter 5: Validating part of the social media infodemic listening conceptual framework using structural equation modelling

Status: Currently under journal review.

Authors: Shu-Feng Tsao, Helen Chen, Zahid A. Butt After completing the preliminary qualitative validation presented in Chapter 4, this study now turns to a quantitative approach for validating partial components of the SoMeIL conceptual framework. In this chapter, structural equation modelling (SEM), a powerful statistical analysis technique is employed since it is widely used in social science, psychology, health, economics, and other related fields. SEM allows researchers to investigate complex relationships among variables, both observable (measured) and latent (unobservable or hypothetical), offering a robust tool for understanding complex associations. SEM represents an advanced statistical technique beyond typical regression analysis. Regression analysis is a special type of SEM, and it typically focuses on understanding the relationship between one dependent variable and one or more independent variables. Unlike regression models, SEM allows multiple dependent variables into the modelling simultaneously. This enables researchers to explore not only direct relationships but also indirect relationships, often referred to as paths, among various variables. Therefore, SEM is particularly well-suited for testing complex hypotheses.

A distinctive feature of SEM is its incorporation of latent variables, which are theoretical constructs that cannot be directly observed or measured but are inferred from a set of measured indicators. For instance, the concept of "intention" in the SoMeIL conceptual framework. Intentions represent a latent variable because they are not directly measurable but can be inferred from observable variables retrieved from social media, such as the number of likes. In the SoMeIL framework, latent variables such as attitudes, perceptions, emotions, and behavioral intentions are central to understanding how health information people receive on social media can eventually associated with their behaviours online and offline. To provide context for this study, it's crucial to revisit the SoMeIL conceptual framework. The framework outlines a complex view of the factors influencing an individual's ability to critically assess and engage with information on social media platforms. At its core, the framework includes latent variables such as attitudes towards given health information, perceptions of information credibility or quality, emotional responses to the health information, and finally,

behavioural intentions, the most critical factor that is assumed to influence people's behaviours in real life. This study's primary objective is to validate the inference of behavioral intentions from other observable variables extracted from Twitter data. Behavioral intentions, which are central to understanding how individuals react with information on social media, represent a latent variable in the framework. In this context, SEM serves as an appropriate method to quantitatively assess and validate certain components of the SoMeIL conceptual framework.

In SEM, the analysis begins by specifying a measurement model. The measurement model illustrates how observed variables are linked to their respective latent variables. This step typically involves confirmatory factor analysis (CFA), which quantifies the strength of the relationship between each observed variable and its associated latent variable. In this study, this means quantifying the relationship between measured variables (i.e., online reaction behaviours), such as the number of likes and the number of shares retrieved from Twitter, and the latent variables of behavioral intentions or engagements. Following the measurement model, researchers proceed to construct the structural model. This structural model investigates the complex relationships among the latent variables and observed variables. Researchers specify their hypothesized paths or associations between these variables and estimate the strengths of these relationships. In essence, the structural model allows researchers to explore the complex interplay between online reaction behaviours, behavioral intentions, and the self-reported offline reaction behaviours within the SoMeIL framework. Once the SEM model is established, researchers rely on a range of fit indices to assess how well the model aligns with the observed data. Commonly used fit indices include but not limited to the chi-squared statistic, the Root Mean Square Error of Approximation (RMSEA), and the Comparative Fit Index (CFI). These indices provide insights into the goodness of model fit and provide a measure of how well the model represents the relationships inherent in the data. A good model fit indicates that the proposed SEM model is a suitable representation of the observed data, increasing the SoMeIL framework's validity. In cases where the initial hypothesized SEM model does not align well with the observed data, researchers have the flexibility to make necessary adjustments. These adjustments may involve adding or removing paths between variables, permitting error correlations between variables, or altering the overall model structure. Such modifications aim to enhance the model's fit and ensure that it faithfully captures the underlying relationships among the variables according to the conceptual framework.

SEM has been applied across a wide spectrum of research fields. In psychology and health sciences, researchers frequently employ SEM to study complex and unobservable variables, often referred to as latent variables. These can include constructs like personality traits, intentions, knowledge, or mental health states. SEM provides a valuable tool for understanding the intricate interplay of these variables and their impact on human behaviours. Furthermore, SEM is a versatile tool in the social sciences. It enables researchers to investigate the impact of various factors on social phenomena. For instance, in education research, SEM can be employed to examine how student achievement is influenced by a combination of factors, including socioeconomic status, teaching methods, and parental involvement. Overall, SEM stands as a powerful and flexible statistical approach for exploring complex relationships among variables in diverse research domains. This chapter has provided an overview of SEM and its role in quantitatively validating some components of the SoMeIL conceptual framework. With SEM's assistance, this study aims to shed light on the behavioral intentions underpinning information engagement on social media platforms. By employing this robust analytical technique like SEM, this study aims to preliminarily validate certain components of the SoMeIL framework, contributing to a deeper understanding of how individuals' behaviours can be influenced by health information circulating on social media in the digital age. This chapter begins with the study's abstract, followed by the full-text manuscript.

5.1 Abstract

Background

Existing literature has shown various factors promoting or hindering people's intentions for COVID-19 vaccination, and studies using structural equation modelling (SEM) has been a common approach for such research to validate associations. We have proposed a conceptual framework, called social media listening (SoMeIL) for public health behaviours, hypothesising parameters retrieved from social media platforms can be used to infer people's intentions for the vaccination behaviours. Therefore, this study aimed to preliminarily validate some components of the SoMeIL conceptual framework using SEM and Twitter data. It also examined the feasibility of using Twitter data in SEM research.

Methods

A total of 2,420 of English language tweets in Toronto or Ottawa, Ontario, Canada, were collected from March 8 to June 30, 2021. Confirmatory factor analysis and SEM was applied to validate our proposed conceptual framework.

Findings

The results showed that sentiment scores, the log-numbers of a tweet's favourites and retweets, and the log-numbers of a user's favourites, followers, and public lists had significant direct associations with the self-reported COVID-19 vaccination intention. The sentiment score of a tweet had the strongest relationship, whereas the number of followers for a user had the weakest relationship with the intentions of COVID-19 vaccine uptake.

Interpretation

The findings have preliminarily validated some components of the SoMeIL conceptual framework by testing associations between the self-reported COVID-19 vaccination intention and sentiment scores, the log-numbers of a tweet's favourites, a tweet's retweets, a user's favourites, a user's followers, and a user's public lists. This study also demonstrated the feasibility of using Twitter data in SEM research. More importantly, it preliminarily validated that, in the SoMeIL framework, these six components as online reaction behaviours could be used to infer the self-reported COVID-19 vaccination intention.

Funding

This study was supported by the 2023-24 Ontario Graduate Scholarship.

Keywords: structural equation modelling, Twitter, COVID-19; vaccine intention

5.2 Introduction

Throughout the COVID-19 pandemic, social media has played a substantial role in shaping public perceptions and attitudes toward COVID-19 vaccination.¹⁻² Given extreme interventions like lockdowns to contain the COVID-19 transmission before vaccines were available, people have increasingly connected and relied on digital channels, such as social media, to receive information related to COVID-19. Although social media platforms can be useful tools for disseminating accurate and helpful information, they have also fueled vaccine hesitancy in various ways.¹⁻⁶ The spread of

misinformation about COVID-19 vaccines has been breeding grounds for vaccine hesitancy, given conspiracy theories and other misleading information regarding vaccine safety and efficacy, polarization, and emotions, can easily go viral and create doubts among users.¹⁻⁶ Besides typical online questionnaires or qualitative analysis, researchers have applied machine learning (ML) or artificial intelligence (AI) techniques to investigate and better understand public discourse and sentiments, inferring people's COVID-19 vaccine intention.⁷⁻⁹ The World Health Organization has coined "social listening" to describe such activities and deployed its Early AI-supported Response with Social Listening (EARS) platform during the pandemic.¹⁰

Social listening studies have adopted existing theories from health behaviours, communication, and behavioural sciences.⁷⁻¹⁰ However, there are limitations in these models, and they don't really reflect the current complex information ecosystems. Literature has shown various social listening studies investigating impacts of exposure to information circulating on social media platforms on their COVID-19 vaccination intentions or behaviours.¹¹ Such research has generally been done using surveys and statistical analyses, such as structural equation modelling (SEM), to identify any associations.¹²⁻¹⁷ In other words, SEM has been widely used to investigate factors influencing people's intentions or attitudes toward the COVID-19 vaccination with different theories and variables.¹²⁻¹⁷ For example, several studies have adopted health behavioural theories, such as health belief model (HBM), theory of planned behaviour (TPB), and extended parallel process model (EPPM), to investigate factors that encourage or discourage the COVID-19 vaccine uptake.¹²⁻¹⁷ In general, respondents were more likely to become vaccinated if they perceived a higher risk of being infected with the COVID-19 virus, perceived greater benefits of vaccines, and subjective norm.¹²⁻¹⁷ Online survey has been primarily used in the SEM research given its advantages, such as cost effectiveness, easy administration, global outreach, and efficiency.¹⁸ However, survey research has some limitations, including nonrespondent bias, recall bias, assumed honesty, respondents' misunderstanding or misinterpretation of questions, and others.¹⁸ Although social media data has been used in numerous COVID-19 social listening studies,^{7, 10} it is rare to find SEM studies using social media data. Compared to SEM, although researchers can apply ML or AI techniques to analyse large amount of social media data, such studies don't really show any statistical relationships like SEM.

Accordingly, a new conceptual framework—social media infodemic listening (SoMeIL) for public health behaviour—has been proposed to address the multifaceted health infodemics on social media somewhere else.¹⁹ The proposed conceptual framework has theorised that people’s online reaction behaviours can indicate their intentions to uptake the COVID-19 vaccines, for example.¹⁹ In other words, parameters derived from social media platforms, such as the numbers of likes and shares to a given post, can be used as a proxy to infer people’s self-reported intentions for the COVID-19 vaccination behaviours in real life. Given our special interest in social media and its critical role in health infodemics and thus people’s behaviours, it is important to directly use social media data to validate such associations. SEM has been commonly used to validate conceptual frameworks where latent variables involve with survey data,¹²⁻¹⁷ but social media data has not been directly and extensively used in SEM analysis. Although many studies have investigated how social media has influenced people’s intentions to become COVID-19 vaccinated, they primarily rely on questionnaires to collect data,¹²⁻¹⁷ with few studies have requested participants to provide their social media posts. Using social media data is conceptually similar to typical SEM research with online surveys since social media data share the same benefits while some limitations are mitigated. Ideally, researchers can retrieve as many relevant parameters and data as social media platform’s application programming interface (API) allows. Thus, the sample size of social media data is generally not an issue. Social media data may have similar nonrespondent bias from inactive users or users not on a given social media platform, but the bias may be mitigated by the numbers of likes, shares, or other parameters. Furthermore, since researchers do not need to design the questions, there is no need to assume respondents’ honesty and worry about respondents’ misunderstanding or misinterpretation of the questions.

Therefore, this study aimed to validate online reaction behaviours, intentions, and self-reported offline reaction behaviours in the SoMeIL conceptual framework with SEM using Twitter data. It also demonstrated the feasibility of using Twitter data in SEM research. Figure 4-1 shows the proposed SEM model derived from part of the SoMeIL conceptual framework and corresponding hypotheses. Directly measured variables are represented by rectangles and latent variables by circles. The health information in this case is the massive vaccination campaign that encourage people in Canada to take the first dose of COVID-19 vaccines. Online reaction behaviours include sentiment scores (i.e., emotion in the framework), the log-number of favourites, the log-number of retweets, the log-number

of user favourites, the log-number of user followers, the log-number of user friends, and the log-number of times when user listed in a tweet. The offline reaction behaviour is the self-reported vaccination or not in a tweet. We theorized positive associations in all hypotheses as below:

- H₁: There is a significant relationship between a tweet's sentiment score and tweet engagement. That is, emotion expressed in a tweet is statistically associated with the tweet engagement, positively or negatively.
- H₂: There is a significant relationship between the log-number of a tweet's favourites and tweet engagement. That is, the log-number of likes (i.e., favourites) a tweet receives is statistically associated with the tweet engagement, positively or negatively.
- H₃: There is a significant relationship between the log-number of a tweet's retweets and tweet engagement. That is, the log-number of sharing (i.e., retweets) a tweet has is statistically associated with the tweet engagement, positively or negatively.
- H₄: There is a significant relationship between tweet engagement and the COVID-19 vaccination. That is, the tweet engagement is statistically associated with self-reported offline vaccination behaviour (i.e., vaccinated), positively or negatively.
- H₅: There is a significant relationship between the log-number of a user's favourites and user engagement. That is, the log-number of favourites (i.e., likes) a user receives is statistically associated with the user engagement, positively or negatively.
- H₆: There is a significant relationship between the log-number of a user's followers and user engagement. That is, the log-number of followers a user has is statistically associated with the user engagement, positively or negatively.
- H₇: There is a significant relationship between the log-number of a user's friends and user engagement. That is, the log-number of friends a user has is statistically associated with the user engagement, positively or negatively.
- H₈: There is a significant relationship between the log-number of a user's public lists and user engagement. That is, the log-number of times a user is mentioned in another tweet

(i.e., public lists) is statistically associated with the user engagement, positively or negatively.

- H₉: There is a significant relationship between user engagement and COVID-19 vaccination. That is, the tweet engagement is statistically associated with self-reported offline vaccination behaviour (i.e., vaccinated), positively or negatively.

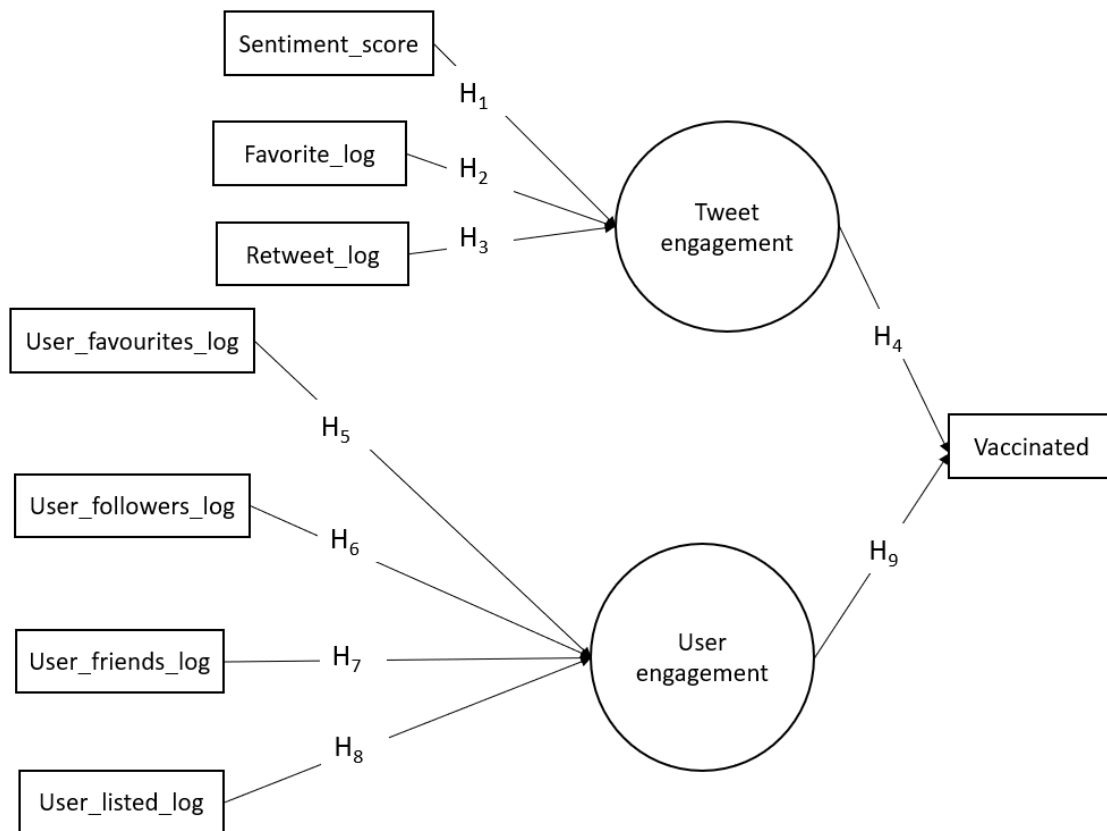


Figure 5-1: Proposed conceptual framework

5.3 Methods

5.3.1 Data collection

This study utilized a cross-sectional design since we were interested to know the COVID-19 vaccination behaviours in Ontario adults in Toronto and Ottawa when the first dose of COVID-19 vaccines became available via online appointments. English tweets related to the COVID-19

pandemic from March 8 to June 30, 2021, were retrieved via Twitter’s Academic API using keywords and hashtags listed in Table 5-1. This resulted in approximately two billion tweets. Next, tweets were narrowed down if “Toronto” or “Ottawa” was included in tweets or identified in users’ locations to gather as many tweets as possible in Toronto or Ottawa, Ontario, Canada. This was done to cope with missing geolocations indicated in literature,²⁰ resulting in approximately four million tweets. The following are definitions of the parameters included for the data preprocessing and analysis:

- Text: tweet.²¹
- Favorite_count: how many times this tweet has been liked by Twitter users.²¹ It is used as part of the online reaction behaviours described in the SoMeIL conceptual framework.
- Retweet_count: number of times this tweet has been retweeted,²¹ which means sharing. It is used as part of the online reaction behaviours described in the SoMeIL conceptual framework.
- User_favourites_count: number of followers this account currently has.²² It is used as part of the online reaction behaviours described in the SoMeIL conceptual framework.
- User_followers_count: number of followers this account currently has.²² It is used as part of the online reaction behaviours described in the SoMeIL conceptual framework.
- User_friends_count: number of users this account is following (i.e., their “followings”).²² It is used as part of the online reaction behaviours described in the SoMeIL conceptual framework.
- User_listed_count: number of public lists that this user is a member of.²² It is used as part of the online reaction behaviours described in the SoMeIL conceptual framework.

Keywords (using OR in queries)	Hashtags (using OR in queries)
corona, rona, coronavirus, covid, covid19, covid-19, sarscov2, sars cov2, sars cov 2, covid_19, ncov, ncov2019, 2019-ncov, pandemic, 2019ncov, covidots, 121hinese virus, wuhan virus, kung flu, jab, vax, vaccinated, immunization(s), herd immunity,	#corona, #rona, #coronavirus, #covid, #covid19, #covid-19, #sarscov2, #covid_19, #ncov, #ncov2019, #2019-ncov, #pandemic #2019ncov, #covidots, #chinesevirus, #wuhanvirus, #kungflu, #coronaupdate, #coronawarriors, #vaccine, #vaccines,

vaccine, vaccines, corona vaccine, corona vaccines	#coronavaccine, #coronavaccines, #herdimmunity
--	--

Table 5-1: Keywords and hashtags for data collection

To prepare for the following sentiment analysis, Twitter handles (i.e., @username), uniform resource locator links (URLs), punctuations, stopwords, and retweets were removed in accordance with existing studies.²³⁻²⁴ Then words in a tweet were converted to their most general form.²³⁻²⁴ This was done using the Natural Language Toolkit (NLTK) package version 3.8.1.²⁵

5.3.2 Measures

Except for texts, the other independent variables were added one and then transformed using natural logarithm given the presence of zeros, since the natural logarithm of one remains zero. Next, to prepare for the sentiment score shown in Figure 4-1, the Valence Aware Dictionary and sEntiment Reasoner (VADER) was used to calculate the sentiment compound score, resulting in a continuous value normalized between -1 and 1 for each tweet in the subset.²⁶ The sentiment score was considered as “emotion” in the SoMeIL conceptual framework.

To prepare for the dependent variable “vaccinated” shown in Figure 4-1, a subset of the four million tweets was created by retrieving tweets that included “appoint,” “jab,” “shot,” and “vaccin.” We manually reviewed and labelled tweets as “1” if users explicitly self-reported that they were looking or waiting for a vaccine appointment, or if they were vaccinated with COVID-19 vaccine already. Tweets were labelled as “0” if users explicitly self-reported that they were hesitant or against the COVID-19 vaccines. Other tweets were excluded if they didn’t have any explicit expressions about the COVID-19 vaccination or if they were news, although they were still relevant to the overall pandemic and vaccine roll out in Canada. The subset ended up with 2,420 English tweets with unique 2,420 users as it was comparable with sample sizes from survey respondents in the existing SEM literature.¹²⁻¹⁷

5.3.3 Statistical Analysis

Descriptive statistics, such as means or frequencies, standard deviations, and Spearman correlations, were used to describe the measures in the proposed model (Table 5-2), except for the latent variables.

Spearman correlations were calculated to account for outliers and non-normal distributions in some measured variables even after the data transformation via natural logarithm (Appendix A). Then the confirmatory factor analysis (CFA) with the diagonally weighted least squares (DWLS), also known as robust WLS, was used to test the “fit” of the observed variables for each latent variable. The robust WLS was specified because some measured variables still violated the normal distribution assumption after the data transformation.²⁷⁻²⁹ For each CFA model, variables were removed until fit indices, including chi-square, comparative fit index (CFI), goodness of fit (GFI), adjusted goodness of fit (AGFI), Tucker-Lewis index (TLI), and root mean square error of approximation (RMSEA), were acceptable. For CFI, GFI, AGFI and TLI, ≥ 0.90 is generally considered acceptable, and ≥ 0.95 is good. RMSEA ≤ 0.08 is recommended.³⁰⁻³¹ After CFA, SEM was performed to test the proposed model (Model 1) in Figure 4-1 with the DWLS and the same recommended criteria for the fit indices. Model 1 would be optimized if the model fit indices suggested better models could be found according to the proposed conceptual framework and correlation matrix. All the data analyses were done in Jupyter Notebook available in Anaconda version 4.3.3, with the semopy package used for CFA and SEM.³²⁻³³

5.3.4 Ethical approval

This study was approved by the University of Waterloo Office of Research Ethics (#43961).

5.3.5 Funding

This study was supported by the 2023-24 Ontario Graduate Scholarship.

5.4 Results

Descriptive statistics and correlations of measured variables are shown in Table 5-2 and Table 5-3, respectively. Results from CFA are shown in Table 5-4. The latent variable “tweet_engagement” was saturated, and the latent variable “user_engagement” had good fit indices except RMSEA, which was larger than the recommended 0.08. When both latent variables were combined in the full measurement model, CFA revealed borderline fit indices since they were close to the acceptable cut offs, and the RMSEA of the full measure model also decreased slightly.

Measures	Mean (SD) or n (%)	Min	Max
Vaccinated (1)	2,173 (89.79%)	-	-
Sentiment_score	0.37 (0.57)	-0.95	0.98
Favorite_log	0.53 (1.13)	0	7.74
Retweet_log	3.35 (2.63)	0	9.46
User_favourites_log	9.84 (2.05)	0	14.15
User_followers_log	6.64 (1.90)	0	13.16
User_friends_log	6.80 (1.38)	0	12.13
User_listed_log	2.21 (1.92)	0	8.47

Table 5-2: Descriptive statistics of measured variables

	1	2	3	4	5	6	7	8
1. Vaccinated	-							
2. Sentiment_score	0.30***	-						
3. Favorite_log	0.03	-0.11***	-					
4. Retweet_log	0.18***	0.40***	-0.49***	-				
5. User_favourites_log	0.01	-0.06**	-0.15***	0.20***	-			
6. User_followers_log	0.05*	-0.11***	0.18***	-0.09***	0.39***	-		
7. User_friends_log	0.01	-0.05*	0.03	0.003	0.44***	0.73***	-	
8. User_listed_log	0.17***	-0.05*	0.18***	-0.10***	0.15***	0.73***	0.49***	-

*p<0.05 ** p<0.01 ***p<0.001

Table 5-3: Spearman correlations for measured variables

	Tweet_engagement	User_engagement	Complete measurement model
Degrees of freedom	0	2	13
chi-square p-value⁺	-	<0.05	<0.05

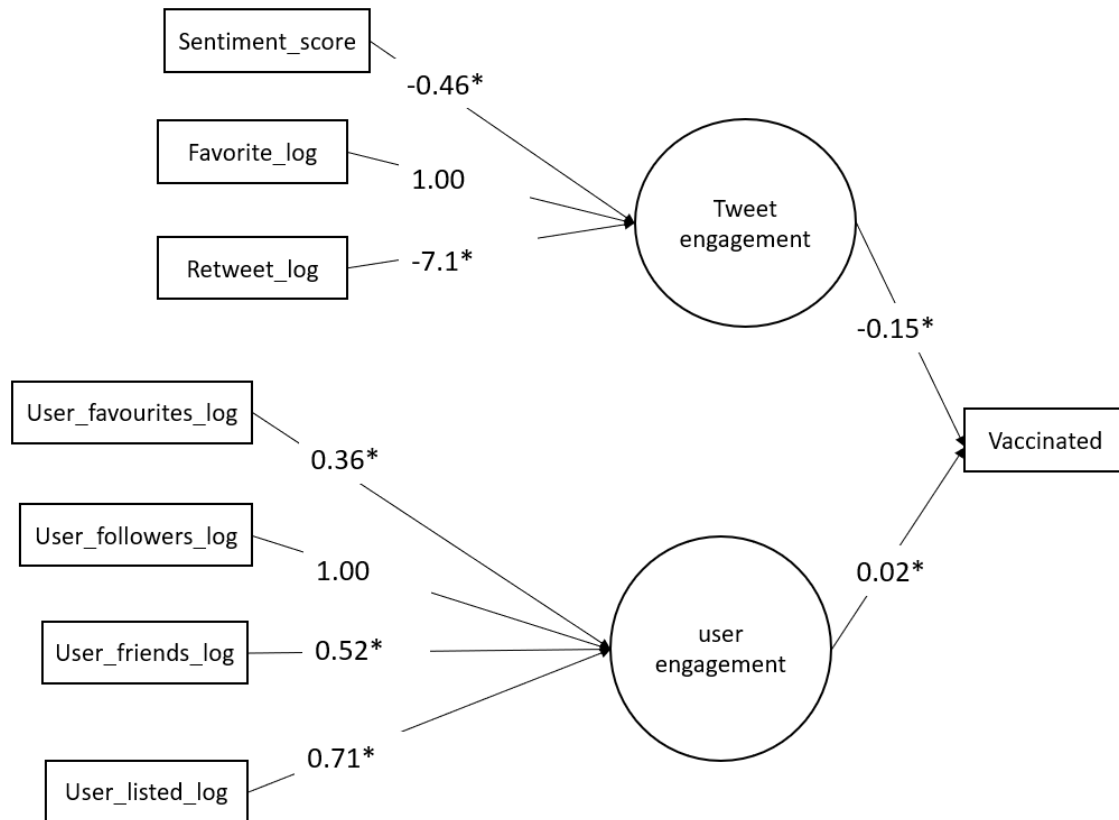
CFI	0.98	0.97	0.88
GFI	0.98	0.97	0.88
AGFI	-	0.90	0.80
NFI	0.98	0.97	0.88
TLI	-	0.91	0.80
RMSEA	∞	0.12	0.12

Table 5-4: Fit statistics for each latent variable and full measurement model

Given the borderline CFA results using DWLS and Twitter data instead of typical surveys, we still decided to test Model 1 using SEM. Figure 5-2 demonstrates Model 1, and its model fit indices were shown in Table 5-5. As Figure 5-2 illustrates, two hypotheses, H₂ and H₆, were not supported because they did not have a statistically significant association. Instead, SEM suggested that the log-number of a tweet's favourites and the log-number of a user's followers were fixed in the model as references:

- H₁: There is a significant relationship between a tweet's sentiment score and tweet engagement (p<0.05).
- H₂: There is a significant relationship between the log-number of a tweet's favourites and tweet engagement (p-value was not provided).
- H₃: There is a significant relationship between the log-number of a tweet's retweets and tweet engagement (p<0.05).
- H₄: There is a significant relationship between tweet engagement and the COVID-19 vaccination (p<0.05).
- H₅: There is a significant relationship between the log-number of a user's favourites and user engagement (p<0.05).
- H₆: There is a significant relationship between the log-number of a user's followers and user engagement (p-value was not provided).
- H₇: There is a significant relationship between the log-number of a user's friends and user engagement (p<0.05).

- H₈: There is a significant relationship between the log-number of a user's public lists and user engagement ($p < 0.05$).
- H₉: There is a significant relationship between user engagement and COVID-19 vaccination ($p < 0.05$).



* $p < 0.05$

Figure 5-2: Model 1

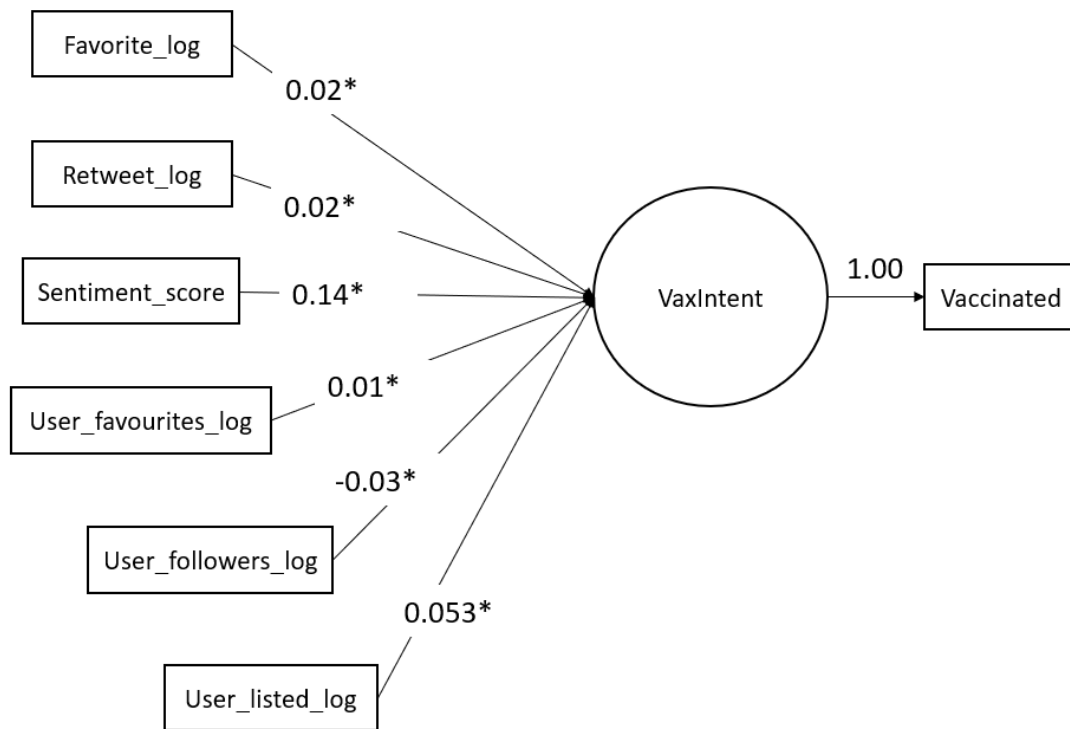
	Model 1	Model 2
Degrees of freedom	18	20
chi-square p-value ⁺	$P < 0.05$	1.0000
CFI	0.8321	1.0079

GFI	0.8287	1.0000
AGFI	0.7335	1.0000
NFI	0.8267	1.0000
TLI	0.7388	1.0106
RMSEA	0.1219	0.0000

+: The chi-squared p-value is not recommended to be considered regardless of SEM models because it was heavily influenced by the sample size.

Table 5-5: Model fit indices for two SEM models

However, the fit indices of Model 1 in Table 5-5 indicated that the model could be optimized. According to the results from CFA and SEM for Model 1, it was hypothesized that instead of two latent variables, one latent variable might be better. Figure 5-3 showed the final SEM model (Model 2) found after model revisions based on the proposed conceptual framework. The model indices of Model 2 are also included in Table 5-5. According to Model 2, the log-number of a user's friends were removed, and the remaining variables had statistically significant relationships with the latent variable. Nonetheless, it was not straightforward to interpret the estimated coefficients and standard errors when the variables were transformed with natural logarithm. Therefore, Table 5-6 and Table 5-7 showed coefficients and their standard errors for each variable in Model 1 and Model 2, respectively, after the estimates were converted back.



*p<0.05

Figure 5-3: Model 2

Variables	Coefficient ± Standard error
Sentiment score → tweet engagement	0.634 ± 1.04
Favourite → tweet engagement	2.72
Retweet → tweet engagement	0.001 ± 3.24
Tweet engagement → vaccinated	0.86 ± 1.02
User favourites → user engagement	1.43 ± 1.02
User followers → user engagement	2.72
User friends → user engagement	1.67 ± 1.01
User listed → user engagement	2.03 ± 1.02

User engagement → vaccinated	1.02 ± 1.00
------------------------------	-------------

Table 5-6: Model 1 estimates after conversion

Variables	Coefficient ± Standard error
Sentiment score → VaxIntent	1.15 ± 1.01
Favourite → VaxIntent	1.02 ± 1.00
Retweet → VaxIntent	1.02 ± 1.00
User favourites → VaxIntent	1.0 ± 1.00
User followers → VaxIntent	0.97 ± 1.0
User listed → VaxIntent	1.05 ± 1.0

Table 5--5-7: Model 2 estimates after conversion

5.5 Discussion

The current study was conducted to preliminarily evaluate the online reaction behaviours, emotion, intention, and self-reported offline behaviour proposed in the SoMeIL conceptual framework.¹⁹ According to the SoMeIL conceptual framework,¹⁹ sentiment scores as emotion, the log-numbers of a tweet's favourites, retweets, a user's favourites, a user's followers, and a user's public lists as online reaction behaviours were investigated using SEM to assess their relationships with the self-reported COVID-19 vaccination as the offline reaction behaviour with a total of 2,420 English tweets. Shown in Table 5-6, most variables in Model 1 had positive associations, but the relationships between a tweet's sentiment score and tweet engagement, between the number of a tweet's retweet and tweet engagement, and between tweet engagement and vaccinated could vary. Similarly, in Model 2, the association between the number of a user's followers and the COVID-19 vaccination intention could be positive or negative (Table 5-7), whereas other variables in Model 2 had positive associations with the latent variable. However, Model 2 was the most optimal model according to its fit indices shown in Table 5-5, and all the variables in Model 2 had statistically significant relationships despite one unstable variable. According to Model 2 (Table 5-7), sentiment score showed the strongest positive relationship with the COVID-19 vaccination intention, followed by the number of public lists a user

belongs to. The number of a user's followers had the weakest association with the COVID-19 vaccination intention.

Overall, Model 2 has provided preliminary results validating the partial components within the SoMeIL conceptual framework given significant associations. That is, variables derived from social media platforms, such as Twitter, could be used to infer people's intentions of COVID-19 vaccine uptakes, which was the latent variable. Sentiment scores, which was calculated via the VADER sentiment analysis,²⁶ represented emotions and showed a significant association consistent with existing literature.^{2, 12-17} In other words, Twitter users who expressed generally positive sentiments toward the COVID-19 vaccines were more likely to become vaccinated against the pandemic.^{2, 12-17} Other variables also showed similar relationships. The more favourites and retweets a tweet received, or the more favourites, followers, public lists a user received, the more likely the user would accept the first dose of COVID-19 vaccines, although the number of a user's followers could have a negative effect at some cases.

Surprisingly, it appeared that outliers had little impacts on SEM since Model 2 met all the recommended criteria of the fit indices. In fact, when outliers were removed or replaced with medians, none of SEM models experimented converged. This remained unchanged even after different combinations of the measured variables were tested. For example, "favorite_log" was excluded because it became useless after removing its outliers or replacing outliers with its median, which was zero. This made the variable literally include only zeros since non-zero values were outliers. Even after "favorite_log" was excluded, other SEM models still failed to converge. Therefore, we hypothesized that without the "favorite_log" variable, the remaining data did not fit SEM models well.³⁴⁻³⁵ Therefore, although the assumption of no outliers in SEM was violated in the current study, the outliers actually included important information that should not be dropped from the modelling. Given the nature of social media data, the outliers could be legitimate since some tweets could receive more likes or shares, or some users could have more followers or receive from likes from others.

In addition to the preliminary validation of the partial components within the SoMeIL conceptual framework, this study might be the first one to use solely Twitter data in SEM research. The findings have shown a promising aspect to employ Twitter data in SEM research with proper theoretical

frameworks, but there were some limitations. Firstly, the generalizability of this study was limited since it did not represent Twitter users not included in the data and non-Twitter users. Furthermore, SEM was conducted in a cross-sectional manner, so it offered just a snapshot of the entire pandemic. In the future, longitudinal SEM could be experimented. However, unlike surveys, researchers would have no control over the frequency of people's tweeting behaviours. Some very active users might tweet daily, whereas others might tweet once in a while. It would require a lot of effort to find enough users who have similar tweeting frequencies to conduct a longitudinal SEM study, although it is not impossible. Data quality was another major limitation. For example, users' demographic information was primarily not available to researchers unless users self-identified their demographics on their Twitter profiles. It would require extensive manual identification or complex ML or AI techniques to retrieve or infer the users' complete demographic characteristics from Twitter data.³⁶⁻³⁷ However, this could potentially lead to even fewer representative samples since the majority of Twitter users have not included any descriptions about their demographics. Additionally, there are other methods to calculate sentiments [8-9], although the VADER sentiment analysis has been commonly used.²⁶ The data transformation via natural logarithm also limited the data quality due to information loss. In general, log transformations were not recommended for count data despite its common usage in linear models like regressions and SEM.³⁴⁻³⁵ Instead, it has been recommended to model count data with Poisson or negative binomial distributions.³⁵ Nonetheless, the Poisson or negative binomial distributions have not been made available in open-source SEM packages, such as semopy package.³²⁻
³³ However, we mitigated the concern using DWLS since it was used to address data that did not meet the normal distribution assumption.^{27-29, 34-35} Last but not least, ecological fallacy was a disadvantage in the SEM study. In other words, the findings should not be interpreted at the individual level.

Despite these limitations, this study has confirmed that not only Twitter data could be useful for SEM research, but it also partially and preliminarily validated the SoMeIL conceptual framework.¹⁹ That is, parameters retrieved from social media platforms like Twitter as online reaction behaviours could be used to infer people's self-reported intentions, which could be used as a proxy for users' vaccination behaviours in real life. For future research, we plan to apply ML or AI techniques, such as support vector machines (SVM) or Bidirectional Encoder Representations from Transformers (BERT) models, to correctly classify the self-reported offline reaction behaviours, so the data sample can be scaled up. Alternatively, instead of the self-reported offline reaction behaviours derived from

social media data, other data regarding the offline reaction behaviours can be collected and analysed to further validate the SoMeIL framework. Another is to collect different social media data, such as videos and images, to study how the SEM approach and the SoMeIL conceptual framework can be applied. For example, like typical SEM research, future studies can design a questionnaire to collect participants' demographics and request them to voluntarily give social media posts to researchers, so scholars can investigate how participants' online and offline reaction behaviours are associated with their demographics. Yet we acknowledge that collecting social media has become more and more difficult for researchers since social media platforms have started to restrict their API access for everyone else.

5.6 Conclusion

In conclusion, this study has provided preliminary validations to the proposed conceptual framework. The results showed that the six variables retrieved from Twitter had statistically significant relationships with the latent variable, which could be used as a proxy for people's self-reported COVID-19 vaccination uptake. This study also demonstrated that it was feasible to use Twitter data in SEM research. Further studies are needed to examine other SEM approaches and other social media to provide more validation to the proposed conceptual framework.

5.7 Declaration of interests

The authors have no conflicts of interest.

5.8 Data sharing statement

The subset and codes used in this study are available on [GitHub](#). However, tweets and other identifiers were not included to protect users' privacy according to Twitter's terms and conditions.

5.9 References

1. Cascini F, Pantovic A, Al-Ajlouni YA, Failla G, Puleo V, Melnyk A, et al. Social media and attitudes towards a COVID-19 vaccination: A systematic review of the literature. *EClinicalMedicine* [Internet]. 2022;48(101454):101454. Available from: <http://dx.doi.org/10.1016/j.eclinm.2022.101454>
2. Romate J, Rajkumar E, Gopi A, Abraham J, Rages J, Lakshmi R, et al. What contributes to COVID-19 vaccine hesitancy? A systematic review of the psychological factors associated with COVID-19 vaccine hesitancy. *Vaccines (Basel)* [Internet]. 2022 [cited 2023 Jul 22];10(11):1777. Available from: <https://www.mdpi.com/2076-393X/10/11/1777>
3. Skafle I, Nordahl-Hansen A, Quintana DS, Wynn R, Gabarron E. Misinformation about COVID-19 vaccines on social media: Rapid review. *J Med Internet Res* [Internet]. 2022;24(8):e37367. Available from: <http://dx.doi.org/10.2196/37367>
4. Zhao S, Hu S, Zhou X, Song S, Wang Q, Zheng H, et al. The prevalence, features, influencing factors, and solutions for COVID-19 vaccine misinformation: Systematic review. *JMIR Public Health Surveill* [Internet]. 2023;9:e40201. Available from: <http://dx.doi.org/10.2196/40201>
5. Thorakkattil SA, Abdulsalim S, Karattuthodi MS, Unnikrishnan MK, Rashid M, Thunga G. COVID-19 vaccine hesitancy: The perils of peddling science by social media and the lay press. *Vaccines (Basel)* [Internet]. 2022;10(7):1059. Available from: <http://dx.doi.org/10.3390/vaccines10071059>
6. Lieneck C, Heinemann K, Patel J, Huynh H, Leafblad A, Moreno E, et al. Facilitators and barriers of COVID-19 vaccine promotion on social media in the United States: A systematic review. *Healthcare (Basel)* [Internet]. 2022 [cited 2023 Jul 22];10(2):321. Available from: <https://www.mdpi.com/2227-9032/10/2/321>
7. Butt MJ, Malik AK, Qamar N, Yar S, Malik AJ, Rauf U. A survey on COVID-19 data analysis using AI, IoT, and social media. *Sensors (Basel)* [Internet]. 2023;23(12):5543. Available from: <http://dx.doi.org/10.3390/s23125543>
8. Alamoodi AH, Zaidan BB, Al-Masawa M, Taresh SM, Noman S, Ahmaro IYY, et al. Multi-perspectives systematic review on the applications of sentiment analysis for vaccine hesitancy.

Comput Biol Med [Internet]. 2021;139(104957):104957. Available from:
<http://dx.doi.org/10.1016/j.combiomed.2021.104957>

9. Umair A, Masciari E, Habib Ullah MH. Sentimental analysis applications and approaches during COVID-19: A survey. In: 25th International Database Engineering & Applications Symposium. New York, NY, USA: ACM; 2021. Doi: 10.1145/3472163.3472274.
10. Purnat TD, Wilson H, Nguyen T, Briand S. EARS – A WHO platform for AI-supported real-time online Social Listening of COVID-19 conversations. In: Studies in Health Technology and Informatics. IOS Press; 2021. Doi: 10.3233/SHTI210330.
11. Heyerdahl LW, Lana B, Giles-Vernick T. Rethinking the infodemic: Social media and offline action in the COVID-19 pandemic. In: Economics, Law, and Institutions in Asia Pacific. Singapore: Springer Singapore; 2022. p. 73–82.
12. Chu H, Liu S. Integrating health behavior theories to predict American’s intention to receive a COVID-19 vaccine. Patient Educ Couns [Internet]. 2021;104(8):1878–86. Available from:
<http://dx.doi.org/10.1016/j.pec.2021.02.031>
13. Fan C-W, Chen I-H, Ko N-Y, Yen C-F, Lin C-Y, Griffiths MD, et al. Extended theory of planned behavior in explaining the intention to COVID-19 vaccination uptake among mainland Chinese university students: an online survey study. Hum Vaccin Immunother [Internet]. 2021;17(10):3413–20. Available from: <http://dx.doi.org/10.1080/21645515.2021.1933687>
14. Irfan M, Shahid AL, Ahmad M, Iqbal W, Elavarasan RM, Ren S, et al. Assessment of public intention to get vaccination against COVID -19: Evidence from a developing country. J Eval Clin Pract [Internet]. 2022;28(1):63–73. Available from: <http://dx.doi.org/10.1111/jep.13611>
15. Mir HH, Parveen S, Mullick NH, Nabi S. Using structural equation modeling to predict Indian people’s attitudes and intentions towards COVID-19 vaccination. Diabetes Metab Syndr [Internet]. 2021;15(3):1017–22. Available from: <http://dx.doi.org/10.1016/j.dsx.2021.05.006>
16. Bui HN, Duong CD, Nguyen VQ, Vu NX, Ha ST, Le TT, et al. Utilizing the theory of planned behavior to predict COVID-19 vaccination intention: A structural equational modeling approach. Heliyon [Internet]. 2023;9(6):e17418. Available from:
<http://dx.doi.org/10.1016/j.heliyon.2023.e17418>

17. Drażkowski D, Trepanowski R. Reactance and perceived disease severity as determinants of COVID-19 vaccination intention: an application of the theory of planned behavior. *Psychol Health Med* [Internet]. 2022;27(10):2171–8. Available from: <http://dx.doi.org/10.1080/13548506.2021.2014060>
18. Evans JR, Mathur A. The value of online surveys: a look back and a look ahead. *Internet Res* [Internet]. 2018;28(4):854–87. Available from: <http://dx.doi.org/10.1108/intr-03-2018-0089>
19. Tsao S-F, Chen H, Meyer S, Butt ZA. Proposing a conceptual framework: social media listening for public health behavior. 2023; Available from: <http://dx.doi.org/10.48550/ARXIV.2308.02037>.
20. Huang B, Carley KM. A large-scale empirical study of geotagging behavior on Twitter. In: *Proceedings of the 2019 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining*. New York, NY, USA: ACM; 2019. Doi: 10.1145/3341161.3342870.
21. Tweet object [Internet]. Twitter. [cited 2023 Jul 23]. Available from: <https://developer.twitter.com/en/docs/twitter-api/v1/data-dictionary/object-model/tweet>
22. User object [Internet]. Twitter. [cited 2023 Jul 23]. Available from: <https://developer.twitter.com/en/docs/twitter-api/v1/data-dictionary/object-model/user>
23. Liu S, Liu J. Public attitudes toward COVID-19 vaccines on English-language Twitter: A sentiment analysis. *Vaccine* [Internet]. 2021;39(39):5499–505. Available from: <http://dx.doi.org/10.1016/j.vaccine.2021.08.058>
24. Reshi AA, Rustam F, Aljedaani W, Shafi S, Alhossan A, Alrabiah Z, et al. COVID-19 vaccination-related sentiments analysis: A case study using worldwide Twitter dataset. *Healthcare (Basel)* [Internet]. 2022;10(3):411. Available from: <http://dx.doi.org/10.3390/healthcare10030411>
25. Bird S, Klein E, Loper E. *Natural language processing with python: Analyzing text with the natural language toolkit*. O'Reilly Media; 2009.
26. Hutto C, Gilbert E. VADER: A parsimonious rule-based model for sentiment analysis of social media text. *Proceedings of the International AAAI Conference on Web and Social Media* [Internet]. 2014;8(1):216–25. Available from: <http://dx.doi.org/10.1609/icwsm.v8i1.14550>

27. Whittaker TA, Schumacker RE. A beginner's guide to structural equation modeling. 5th ed. London, England: Routledge; 2022.
28. Bowen NK, Guo S. Evaluating and improving CFA and general structural models. In: Structural Equation Modeling. Oxford University Press; 2011. p. 135–66.
29. Kupek E. Beyond logistic regression: structural equations modelling for binary variables and its application to investigating unobserved confounders. BMC Med Res Methodol [Internet]. 2006;6(1). Available from: <http://dx.doi.org/10.1186/1471-2288-6-13>
30. Hooper D, Coughlan J, Mullen MR. Structural equation modelling: Guidelines for determining model fit. Electron J Bus Res Meth [Internet]. 2008 [cited 2023 Jul 29];6(1):53-60-pp53-60. Available from: <https://academic-publishing.org/index.php/ejbrm/article/view/1224>
31. Peugh J, Feldon DF. “how well does your structural equation model fit your data?”: Is marcoulides and Yuan's equivalence test the answer? CBE Life Sci Educ [Internet]. 2020;19(3):es5. Available from: <http://dx.doi.org/10.1187/cbe.20-01-0016>
32. Igolkina AA, Meshcheryakov G. Semopy: A python package for structural equation modeling. Struct Equ Modeling [Internet]. 2020;27(6):952–63. Available from: <http://dx.doi.org/10.1080/10705511.2019.1704289>
33. Meshcheryakov G, Igolkina AA, Samsonova MG. Semopy 2: A structural Equation Modeling package with random effects in Python. 2021; Available from: <http://dx.doi.org/10.48550/ARXIV.2106.01140>
34. Xia Y, Yang Y. RMSEA, CFI, and TLI in structural equation modeling with ordered categorical data: The story they tell depends on the estimation methods. Behav Res Methods [Internet]. 2019;51(1):409–28. Available from: <http://dx.doi.org/10.3758/s13428-018-1055-2>
35. Green JA. Too many zeros and/or highly skewed? A tutorial on modelling health behaviour as count data with Poisson and negative binomial regression. Health Psychol Behav Med [Internet]. 2021;9(1):436–55. Available from: <http://dx.doi.org/10.1080/21642850.2021.1920416>

36. Golder S, Stevens R, O'Connor K, James R, Gonzalez-Hernandez G. Methods to establish race or ethnicity of Twitter users: Scoping review. *J Med Internet Res* [Internet]. 2022;24(4):e35788. Available from: <http://dx.doi.org/10.2196/35788>
37. Cesare N, Grant C, Nsoesie EO. Detection of user demographics on social media: A review of methods and recommendations for best practices. *arXiv preprint arXiv:1702.01807*. 2017 Feb 6:1-25.

Chapter 6 Conclusion

6.1 Summary of key findings

The scoping review in Chapter 2 was the first study that investigated a wide variety of literature about roles social media has played in the early pandemic. We reviewed articles published in peer-reviewed journals that analyzed social media data for different research questions. Six themes were found: (1) surveying public attitudes, (2) identifying infodemics, (3) assessing mental health, (4) detecting or predicting COVID-19 cases, (5) analyzing government responses to the pandemic, and (6) evaluating quality of health information in prevention education videos. At that point in time, findings identified limited research using ML or AI techniques to analyze social media data or developing real-time surveillance for COVID-19 with social media data included. However, the gaps have been bridged as the pandemic continued. Nonetheless, the scoping review led us to recognize different ways that social media data have been utilized, and the “public attitudes” and “infodemics” themes accounted for 60 out of 81 reviewed articles. Despite the abundant literature on these two themes, existing theories or theoretical constructs were rarely applied with topic modeling and sentiment analysis. Even if the existing theories have been applied, they have not well-reflected the role of the complex information ecosystem in the modern society where social media has intertwined both online and offline information channels. Therefore, the scoping review was important to identify further knowledge gap and theoretical needs to investigate how social media data can be systematically used in social listening practices.

The second study—the conceptual paper—in Chapter 3 contributed to a novel conceptual framework especially designed for social media infodemic listening (SoMeIL) for public health behaviors. The COVID-19 vaccination behavior was used as an example throughout the paper. We first reviewed and identified gaps in existing theories or theoretical constructs. Findings suggested that theories investigated how people’s behaviors on social media, such as pressing liking or sharing a post, could be used as proxies to be associated with their behaviors in real life. The results also implied an implicit assumption behind some theories: people would make a rational decision and thus behave accordingly after they were informed by the scientific information from health agencies or professionals. Other theories or theoretical constructs, however, indicated that emotions played an

important role in people's behaviors. Therefore, the SoMeIL conceptual framework was developed to address these limitations and gaps.

After the SoMeIL conceptual framework is developed, it is necessary to validate the framework according to TCM. The third study presented in Chapter 4 was a qualitative preliminary validation using the Canadian Freedom Convoy as an example by employing LDA topic modelling and qualitative thematic analysis. This study also demonstrated how the proposed conceptual framework could be applied to infer people's attitudes and what factors have led to different attitudes as well as reaction behaviours. This preliminary validation study has qualitatively investigated the attitudes and self-reported offline reaction behaviours within the SoMeIL conceptual framework. After the preliminary qualitative validation was completed, the next study in Chapter 5 offered preliminary quantitative validation.

Study IV in Chapter 5 validated the SoMeIL conceptual framework preliminarily using SEM by testing associations between the latent variable and observed variables retrieved from Twitter. The latent variable in this study was the self-reported intention to become COVID-19 vaccinated. Observed variables directly retrieved from Twitter included the log-numbers of a tweet's likes, a tweet's retweets, a user's favorites, a user's followers, a user's friends, and a user's public lists. All of these observed variables were regarded as "online reaction behaviours" as described in the SoMeIL conceptual framework. Sentiment scores of each tweet, another observed variable, were calculated using the VADER sentiment analysis. The sentiment score was considered as "emotion" based on the SoMeIL conceptual framework. The findings suggested statistically significant associations between the latent variable and the observed variables except the log-number of a user's friends. Therefore, the results suggested preliminary evidence that the latent variable, inferred by other observed variables, could be used as a proxy associated with people's self-reported vaccination behavior in real life. This study validated that some constructs, including online reaction behaviours, emotion, intention, and self-reported offline reaction behaviour in the SoMeIL conceptual framework were workable. It also demonstrated that analyzing Twitter data in SEM was feasible since surveys have been primarily used in SEM research.

6.2 Study limitations and strengths

There are some limitations in this dissertation. In the scoping review, no grey literature was included, and only peer-reviewed articles published in English were included. Besides, the scoping review included studies published in the early pandemic, so the identified knowledge gaps have already been filled as the pandemic continued. In the second study (i.e., the conceptual paper), only theories that met the inclusion criteria were included. Therefore, other useful theories were excluded from the analysis. This may lead to biases toward theories focusing on behavioral changes. In both Study III and Study IV, only English tweets related to vaccination events in Canada were included, so the findings were not representative to the general public in Canada and other countries. The quality of Twitter data has several issues. Firstly, geolocations can be primarily missing since most Twitter users choose not to self-disclose their real geolocations.^{36, 37, 38} Other users may include a vague geolocation, such as “on earth,” that does not provide any useful information given the research contexts.^{39, 40, 41} In addition, some locations in Canada have same names as other locations in different countries. For example, “Ontario, CA” can be Ontario in Canada or Ontario, California in the United States. Another common issue is informal languages or acronyms used in tweet. For instance, “va\$\$” has been commonly used by people against vaccination as the term is a combination of “vaccine” and dollar signs representing pharmaceutical companies profiting from the massive vaccination. Such informal languages can make data collection and cleaning biased if researchers do not know them before retrieving tweets from the API. This can also make data cleaning difficult as researchers tend to remove special characters, despite their meaning, when systematically preparing large numbers of collected tweets for NLP analyses, such as topic modeling and sentiment analysis.

Nonetheless, the strengths of this dissertation included the recognition and documentation of a great myriad of utilizations of social media data in studies related to the COVID-19 pandemic. Furthermore, the SoMeIL conceptual framework and its preliminary partial validations and application demonstrated that social media data could be used as proxies to infer self-reported people’s vaccination behaviors in real life in addition to other methods.

6.3 Directions for future research

Drawing from findings presented in this dissertation, the following are some recommendations for future applications or investigations regarding the proposed SoMeIL conceptual framework using qualitative, quantitative, or mix-methods research:

- The SoMeIL conceptual framework needs more validations using different social media data besides Twitter. When different social media data are combined, a relatively representative sample may be created. Furthermore, this PhD research have primarily textual data from Twitter. Other types of social media data, such as videos and images, can also be used in future validations. Additionally, other behaviors related to the public health, such as social distancing, handwashing, and uptakes of other vaccines, can also be used as different examples to validate the proposed conceptual framework.
- There are other SEM models, such as longitudinal SEM, can be applied to validate the proposed conceptual framework using social media data or surveys. When longitudinal analysis is applied, it is possible to use the number of daily COVID-19 cases or vaccine administrations as the dependent variable instead of the current self-reported vaccination behaviours used in the study. It is also important to examine the feasibility of analyzing various social media data in different SEM approaches. Multiple group SEM models can also be applied to investigate if associations among vaccine advocates will be different from that among vaccine hesitant or opponents.
- Although the VADER sentiment analysis generates standardized sentiment scores from -1 to +1, the VADER sentiment scores do not provide further insights other than positive, neutral, and negative sentiments when the scores are interpreted. The study also did not include non-textual data when the VADER sentiment analysis was applied. Therefore, other sentiment analysis or ML/AI techniques can be employed to calculate the sentiments with more comprehensive data or interpretable results.
- Other components of the SoMeIL conceptual framework also need to be validated. For example, for health information, future studies can investigate the quality and credibility of health information and how they influence people's online or offline reaction behaviours.

The impacts of social media algorithms and platform designs on people's understanding and behaviours also have not been fully fleshed out.

- Furthermore, more studies are needed to investigate the variability of the SoMeIL conceptual framework between one country or culture to another. The SoMeIL conceptual framework was developed primarily based on social media phenomena in Canada and countries with similar culture. Therefore, more evidence is needed to see how the framework can be applied to a different culture or country with adjustments.

This dissertation presented a new conceptual framework with preliminary application and validation. In other words, it is just the beginning. The proposed conceptual framework will be revised as more evidence is found. However, the proposed conceptual framework has provided a systematic and theoretical foundation for future social media listening for health organizations and public health policy makers. As social media has been inevitably integrated into people's daily life, especially among younger generations, it is essential to not only use social media for disseminating health information, but also investigate how social media has impacted people's behaviors in real life by inferring from public discourse and behaviors on social media.

6.4 Implications for public health practices and policies

Although health misinformation and vaccine hesitancy have long existed in human societies long before social media are created,⁴¹ social media have provided a channel allowing health information, in spite of its quality, to be distributed immediately to massive audiences without geographical restrictions given ubiquitous Internet and smartphones. Hospitals have communicated with their patients and provided services via social media before the pandemic.^{42, 43, 44} Health authorities or government agencies also need to spread their messages on social media to urge their residents taking preventive actions in addition to other communication channels.⁴⁵ As social media have integrated into people's daily lives worldwide, its dominance will make health infodemics have greater impacts on people. As a result, besides other conventional channels like surveys or word of mouth, it is crucial to "listen to" public discourse on different social media platforms and address emerging confusions, questions, and even misinformation in a timely manner.

Furthermore, as the SoMeIL conceptual framework illustrates, some indicators on social media, such as the numbers of likes and shares, can be used to infer people's vaccination behaviours in real life. Therefore, it can be adopted to forecast the vaccination coverage in the future for vaccine-preventable diseases. This can also help tailor communication strategies and address specific issues based on social media users' discussions and online behaviors to effectively reach different groups.⁴⁶ ⁴⁷ The proposed conceptual framework can be extended to other areas, such as symptom reports or behavioral patterns, to aid in public health decision-making and resource allocations.⁴⁸ By integrating social media into pandemic preparedness, health organizations and government authorities can harness its potential as a powerful tool to engage with the public, tackle health misinformation, and effectively respond to crises, ultimately helping to mitigate the impact of future pandemics.^{46, 47} Similar to the WHO's EARS platform,¹⁰ the proposed SoMeIL conceptual framework can be implemented as a way to provide real-time monitoring and surveillance. Literature has shown that social media can be used for early detection of emerging health threats and to track misinformation trends.^{7, 10, 11} Social media data can also complement traditional surveillance methods and help public health authorities respond quickly to potential outbreaks.^{7, 10, 11, 48} However, such effort has yet to be scaled up after the COVID-19 pandemic.

Overall, the proposed SoMeIL conceptual framework has provided a preliminary yet quantifiable way for social listening. It is recommended that future pandemic preparedness recognizes the significant roles that social media plays in shaping public perception, disseminating information, and influencing behaviors during a health crisis. Incorporating social media into pandemic preparedness strategies besides others can enhance communication, information sharing, and response efforts.

References

1. Cascini F, Pantovic A, Al-Ajlouni YA, Failla G, Puleo V, Melnyk A, et al. Social media and attitudes towards a COVID-19 vaccination: A systematic review of the literature. *EClinicalMedicine* [Internet]. 2022;48(101454):101454. Available from: <http://dx.doi.org/10.1016/j.eclinm.2022.101454>
2. Romate J, Rajkumar E, Gopi A, Abraham J, Rages J, Lakshmi R, et al. What contributes to COVID-19 vaccine hesitancy? A systematic review of the psychological factors associated with COVID-19 vaccine hesitancy. *Vaccines (Basel)* [Internet]. 2022 [cited 2023 Jul 22];10(11):1777. Available from: <https://www.mdpi.com/2076-393X/10/11/1777>
3. Skafle I, Nordahl-Hansen A, Quintana DS, Wynn R, Gabarron E. Misinformation about COVID-19 vaccines on social media: Rapid review. *J Med Internet Res* [Internet]. 2022;24(8):e37367. Available from: <http://dx.doi.org/10.2196/37367>
4. Zhao S, Hu S, Zhou X, Song S, Wang Q, Zheng H, et al. The prevalence, features, influencing factors, and solutions for COVID-19 vaccine misinformation: Systematic review. *JMIR Public Health Surveill* [Internet]. 2023;9:e40201. Available from: <http://dx.doi.org/10.2196/40201>
5. Thorakkattil SA, Abdulsalim S, Karattuthodi MS, Unnikrishnan MK, Rashid M, Thunga G. COVID-19 vaccine hesitancy: The perils of peddling science by social media and the lay press. *Vaccines (Basel)* [Internet]. 2022;10(7):1059. Available from: <http://dx.doi.org/10.3390/vaccines10071059>
6. Lieneck C, Heinemann K, Patel J, Huynh H, Leafblad A, Moreno E, et al. Facilitators and barriers of COVID-19 vaccine promotion on social media in the United States: A systematic review. *Healthcare (Basel)* [Internet]. 2022 [cited 2023 Jul 22];10(2):321. Available from: <https://www.mdpi.com/2227-9032/10/2/321>
7. Butt MJ, Malik AK, Qamar N, Yar S, Malik AJ, Rauf U. A survey on COVID-19 data analysis using AI, IoT, and social media. *Sensors (Basel)* [Internet]. 2023;23(12):5543. Available from: <http://dx.doi.org/10.3390/s23125543>

8. Alamoodi AH, Zaidan BB, Al-Masawa M, Taresh SM, Noman S, Ahmaro IYY, et al. Multi-perspectives systematic review on the applications of sentiment analysis for vaccine hesitancy. *Comput Biol Med* [Internet]. 2021;139(104957):104957. Available from: <http://dx.doi.org/10.1016/j.combiomed.2021.104957>
9. Umair A, Masciari E, Habib Ullah MH. Sentimental analysis applications and approaches during COVID-19: A survey. In: 25th International Database Engineering & Applications Symposium. New York, NY, USA: ACM; 2021. Doi: 10.1145/3472163.3472274.
10. Purnat TD, Wilson H, Nguyen T, Briand S. EARS – A WHO platform for AI-supported real-time online Social Listening of COVID-19 conversations. In: *Studies in Health Technology and Informatics*. IOS Press; 2021. Doi: 10.3233/SHTI210330.
11. Heyerdahl LW, Lana B, Giles-Vernick T. Rethinking the infodemic: Social media and offline action in the COVID-19 pandemic. In: *Economics, Law, and Institutions in Asia Pacific*. Singapore: Springer Singapore; 2022. p. 73–82.
12. Chu H, Liu S. Integrating health behavior theories to predict American’s intention to receive a COVID-19 vaccine. *Patient Educ Couns* [Internet]. 2021;104(8):1878–86. Available from: <http://dx.doi.org/10.1016/j.pec.2021.02.031>
13. Fan C-W, Chen I-H, Ko N-Y, Yen C-F, Lin C-Y, Griffiths MD, et al. Extended theory of planned behavior in explaining the intention to COVID-19 vaccination uptake among mainland Chinese university students: an online survey study. *Hum Vaccin Immunother* [Internet]. 2021;17(10):3413–20. Available from: <http://dx.doi.org/10.1080/21645515.2021.1933687>
14. Irfan M, Shahid AL, Ahmad M, Iqbal W, Elavarasan RM, Ren S, et al. Assessment of public intention to get vaccination against COVID -19: Evidence from a developing country. *J Eval Clin Pract* [Internet]. 2022;28(1):63–73. Available from: <http://dx.doi.org/10.1111/jep.13611>
15. Al-Dmour H, Masa’deh R, Salman A, Abuhashesh M, Al-Dmour R. Influence of social media platforms on public health protection against the COVID-19 pandemic via the mediating effects of public health awareness and behavioral changes: Integrated model. *J Med Internet Res* [Internet]. 2020;22(8):e19996. Available from: <http://dx.doi.org/10.2196/19996>

16. Cheng C, Espanha R. The impact of COVID-19-related information scanning via social media on Chinese intentions regarding coronavirus vaccinations. *Front Commun* [Internet]. 2023;7. Available from: <http://dx.doi.org/10.3389/fcomm.2022.1094850>
17. Beaudoin CE. Do social media matter? The effects of information seeking on COVID-19 psychological and behavioral processes. *Telemat Inform* [Internet]. 2023;83(102027):102027. Available from: <http://dx.doi.org/10.1016/j.tele.2023.102027>
18. Evans JR, Mathur A. The value of online surveys: a look back and a look ahead. *Internet Res* [Internet]. 2018;28(4):854–87. Available from: <http://dx.doi.org/10.1108/intr-03-2018-0089>
19. Tsao S-F, Chen H, Meyer S, Butt ZA. Proposing a conceptual framework: social media listening for public health behavior. 2023; Available from: <http://dx.doi.org/10.48550/ARXIV.2308.02037>
20. Huang B, Carley KM. A large-scale empirical study of geotagging behavior on Twitter. In: *Proceedings of the 2019 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining*. New York, NY, USA: ACM; 2019. Doi: 10.1145/3341161.3342870.
21. Tweet object [Internet]. Twitter. [cited 2023 Jul 23]. Available from: <https://developer.twitter.com/en/docs/twitter-api/v1/data-dictionary/object-model/tweet>
22. User object [Internet]. Twitter. [cited 2023 Jul 23]. Available from: <https://developer.twitter.com/en/docs/twitter-api/v1/data-dictionary/object-model/user>
23. Liu S, Liu J. Public attitudes toward COVID-19 vaccines on English-language Twitter: A sentiment analysis. *Vaccine* [Internet]. 2021;39(39):5499–505. Available from: <http://dx.doi.org/10.1016/j.vaccine.2021.08.058>
24. Reshi AA, Rustam F, Aljedaani W, Shafi S, Alhossan A, Alrabiah Z, et al. COVID-19 vaccination-related sentiments analysis: A case study using worldwide Twitter dataset. *Healthcare (Basel)* [Internet]. 2022;10(3):411. Available from: <http://dx.doi.org/10.3390/healthcare10030411>
25. Bird S, Klein E, Loper E. *Natural language processing with python: Analyzing text with the natural language toolkit*. O'Reilly Media; 2009.

26. Hutto C, Gilbert E. VADER: A parsimonious rule-based model for sentiment analysis of social media text. *Proceedings of the International AAAI Conference on Web and Social Media* [Internet]. 2014;8(1):216–25. Available from: <http://dx.doi.org/10.1609/icwsm.v8i1.14550>
27. Whittaker TA, Schumacker RE. *A beginner's guide to structural equation modeling*. 5th ed. London, England: Routledge; 2022.
28. Bowen NK, Guo S. Evaluating and improving CFA and general structural models. In: *Structural Equation Modeling*. Oxford University Press; 2011. p. 135–66.
29. Kupek E. Beyond logistic regression: structural equations modelling for binary variables and its application to investigating unobserved confounders. *BMC Med Res Methodol* [Internet]. 2006;6(1). Available from: <http://dx.doi.org/10.1186/1471-2288-6-13>
30. Hooper D, Coughlan J, Mullen MR. Structural equation modelling: Guidelines for determining model fit. *Electron J Bus Res Meth* [Internet]. 2008 [cited 2023 Jul 29];6(1):53-60-pp53-60. Available from: <https://academic-publishing.org/index.php/ejbrm/article/view/1224>
31. Peugh J, Feldon DF. “how well does your structural equation model fit your data?”: Is marcoulides and Yuan's equivalence test the answer? *CBE Life Sci Educ* [Internet]. 2020;19(3):es5. Available from: <http://dx.doi.org/10.1187/cbe.20-01-0016>
32. Igolkina AA, Meshcheryakov G. Semopy: A python package for structural equation modeling. *Struct Equ Modeling* [Internet]. 2020;27(6):952–63. Available from: <http://dx.doi.org/10.1080/10705511.2019.1704289>
33. Meshcheryakov G, Igolkina AA, Samsonova MG. Semopy 2: A structural Equation Modeling package with random effects in Python. 2021; Available from: <http://dx.doi.org/10.48550/ARXIV.2106.01140>
34. Xia Y, Yang Y. RMSEA, CFI, and TLI in structural equation modeling with ordered categorical data: The story they tell depends on the estimation methods. *Behav Res Methods* [Internet]. 2019;51(1):409–28. Available from: <http://dx.doi.org/10.3758/s13428-018-1055-2>

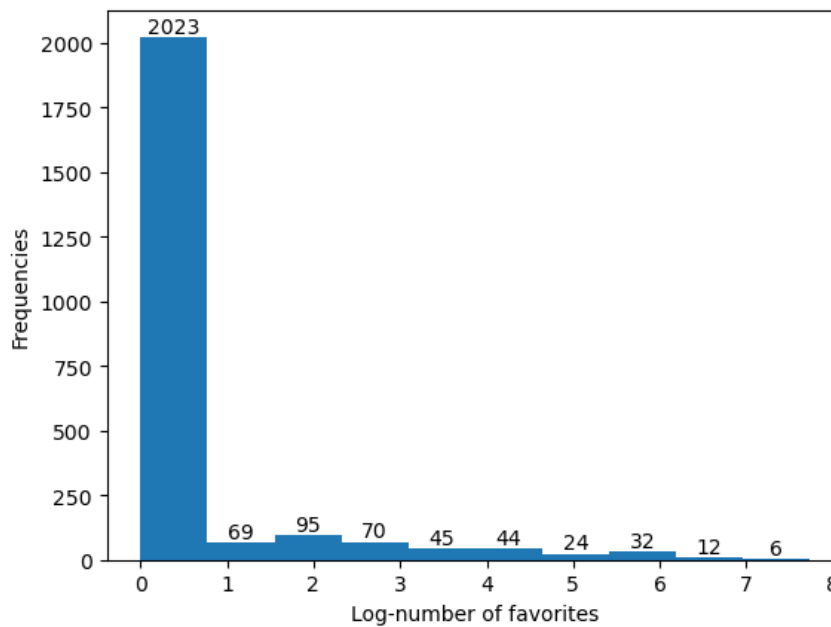
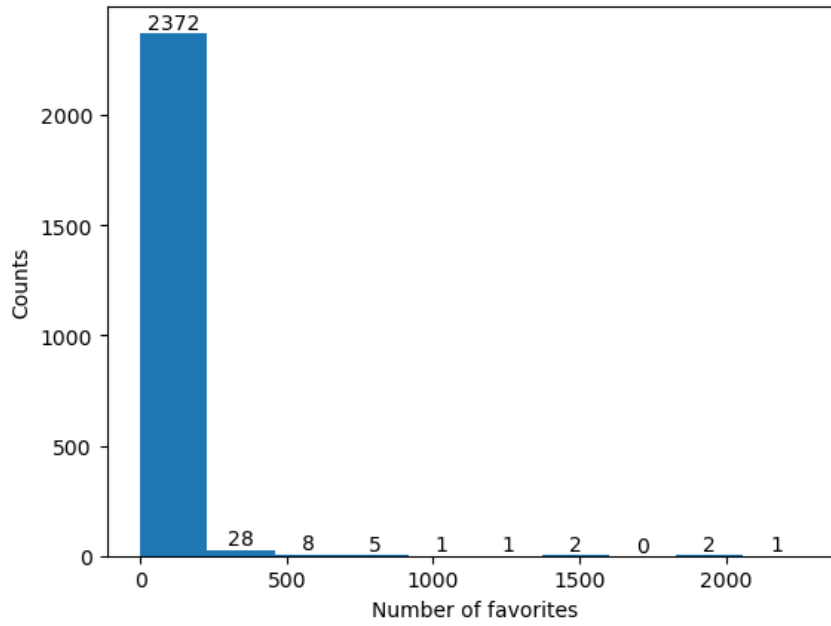
35. Egger R, Yu J. A topic modeling comparison between LDA, NMF, Top2Vec, and BERTopic to demystify Twitter posts. *Front Sociol* [Internet]. 2022;7. Available from: <http://dx.doi.org/10.3389/fsoc.2022.886498>
36. Vaismoradi M, Jones J, Turunen H, Snelgrove S. Theme development in qualitative content analysis and thematic analysis. *J Nurs Educ Pract* [Internet]. 2016;6(5). Available from: <http://dx.doi.org/10.5430/jnep.v6n5p100>
37. Kiger ME, Varpio L. Thematic analysis of qualitative data: AMEE Guide No. 131. *Med Teach* [Internet]. 2020;42(8):846–54. Available from: <http://dx.doi.org/10.1080/0142159x.2020.1755030>
38. Green JA. Too many zeros and/or highly skewed? A tutorial on modelling health behaviour as count data with Poisson and negative binomial regression. *Health Psychol Behav Med* [Internet]. 2021;9(1):436–55. Available from: <http://dx.doi.org/10.1080/21642850.2021.1920416>
39. Golder S, Stevens R, O'Connor K, James R, Gonzalez-Hernandez G. Methods to establish race or ethnicity of Twitter users: Scoping review. *J Med Internet Res* [Internet]. 2022;24(4):e35788. Available from: <http://dx.doi.org/10.2196/35788>
40. Cesare N, Grant C, Nsoesie EO. Detection of user demographics on social media: A review of methods and recommendations for best practices. *arXiv preprint arXiv:1702.01807*. 2017 Feb 6:1-25.
41. Nuwarda RF, Ramzan I, Weekes L, Kayser V. Vaccine hesitancy: Contemporary issues and historical background. *Vaccines (Basel)* [Internet]. 2022;10(10):1595. Available from: <http://dx.doi.org/10.3390/vaccines10101595>
42. Van de Belt TH, Berben SAA, Samsom M, Engelen LJ, Schoonhoven L. Use of social media by western European hospitals: Longitudinal study. *J Med Internet Res* [Internet]. 2012;14(3):e61. Available from: <http://dx.doi.org/10.2196/jmir.1992>
43. Griffis HM, Kilaru AS, Werner RM, Asch DA, Hershey JC, Hill S, et al. Use of social media across US hospitals: Descriptive analysis of adoption and utilization. *J Med Internet Res* [Internet]. 2014;16(11):e264. Available from: <http://dx.doi.org/10.2196/jmir.3758>

44. Richter JP, Muhlestein DB, Wilks CEA. Social media: how hospitals use it, and opportunities for future use. *J Healthc Manag* [Internet]. 2014 [cited 2023 Aug 2];59(6):447–60. Available from: https://journals.lww.com/jhmonline/fulltext/2014/11000/social_media__how_hospitals_use_it,_and.11.aspx
45. Moorhead SA, Hazlett DE, Harrison L, Carroll JK, Irwin A, Hoving C. A new dimension of health care: Systematic review of the uses, benefits, and limitations of social media for health communication. *J Med Internet Res* [Internet]. 2013;15(4):e85. Available from: <http://dx.doi.org/10.2196/jmir.1933>
46. Berg SH, O'Hara JK, Shortt MT, Thune H, Brønneck KK, Lungu DA, et al. Health authorities' health risk communication with the public during pandemics: a rapid scoping review. *BMC Public Health* [Internet]. 2021;21(1). Available from: <http://dx.doi.org/10.1186/s12889-021-11468-3>
47. Vraga EK, Jacobsen KH. Strategies for effective health communication during the Coronavirus pandemic and future emerging infectious disease events. *World Med Health Policy* [Internet]. 2020;12(3):233–41. Available from: <http://dx.doi.org/10.1002/wmh3.359>
48. Yang Y, Tsao S-F, Basri MA, Chen HH, Butt ZA. Digital disease surveillance for emerging Infectious Diseases: An early warning system using the internet and social media data for COVID-19 forecasting in Canada. In: *Caring is Sharing – Exploiting the Value in Data for Health and Innovation*. IOS Press; 2023. Available from: <https://doi.org/10.3233/shti230290>

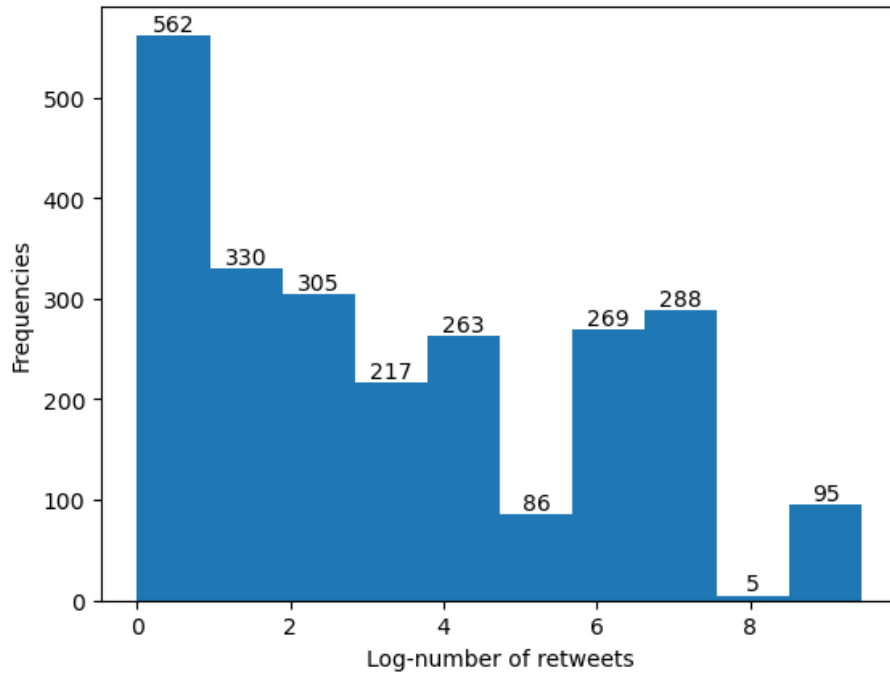
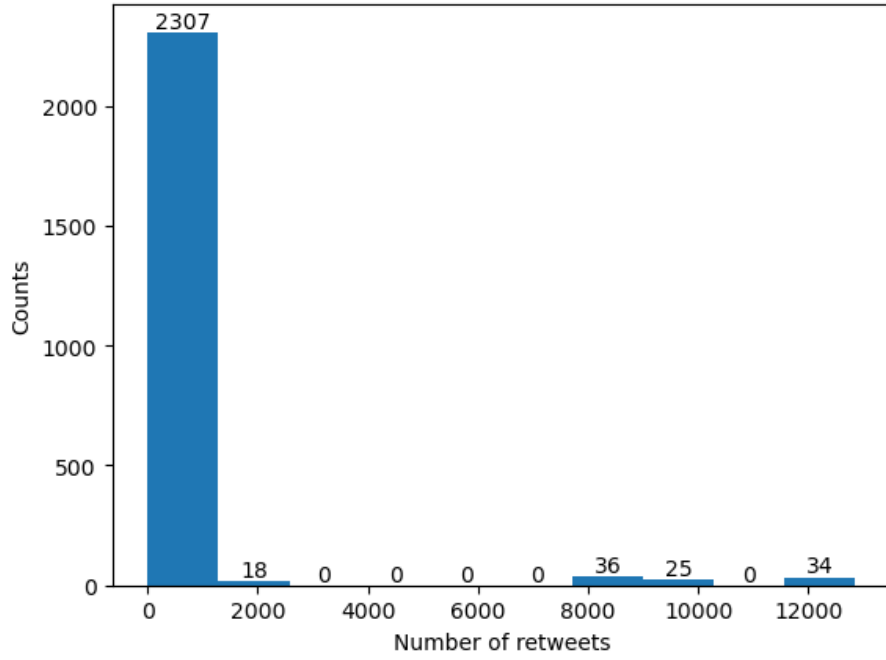
Appendix A: Chapter 3 Supplementary Materials

Figures of variables that were transformed via natural logarithm:

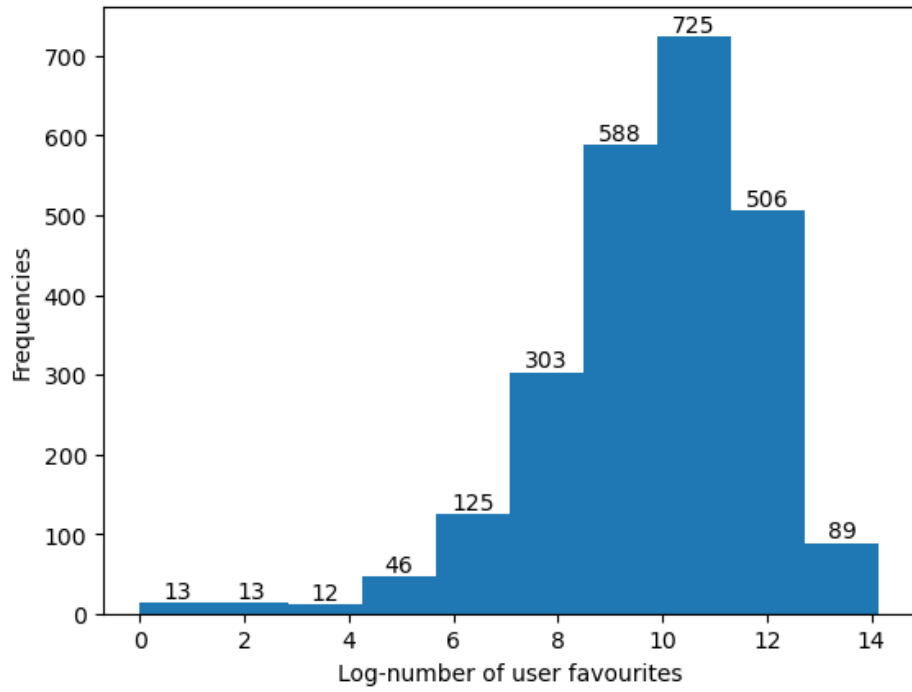
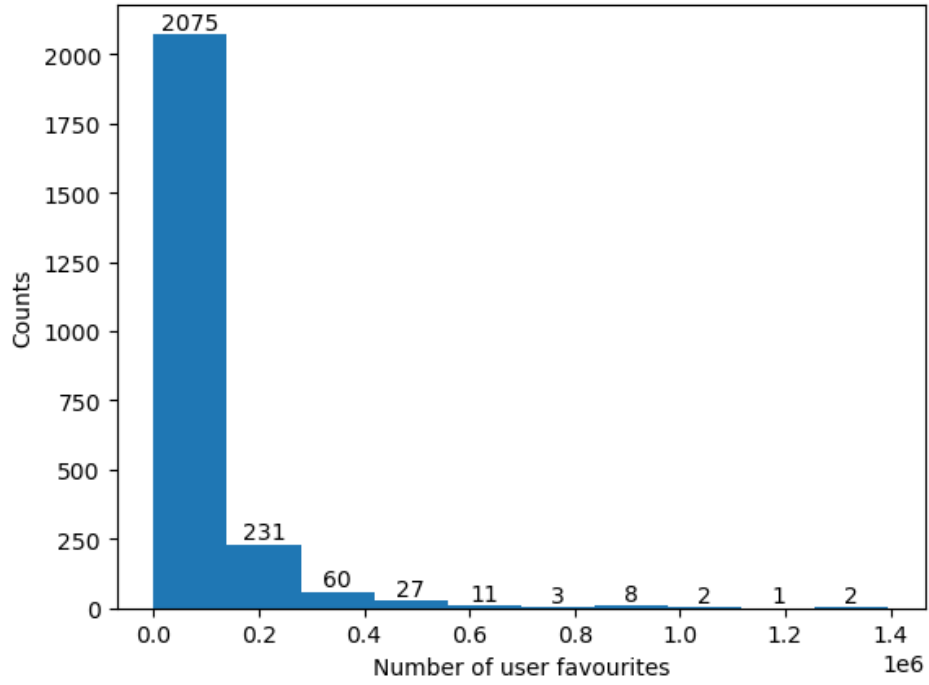
- Favorite count vs favorite log



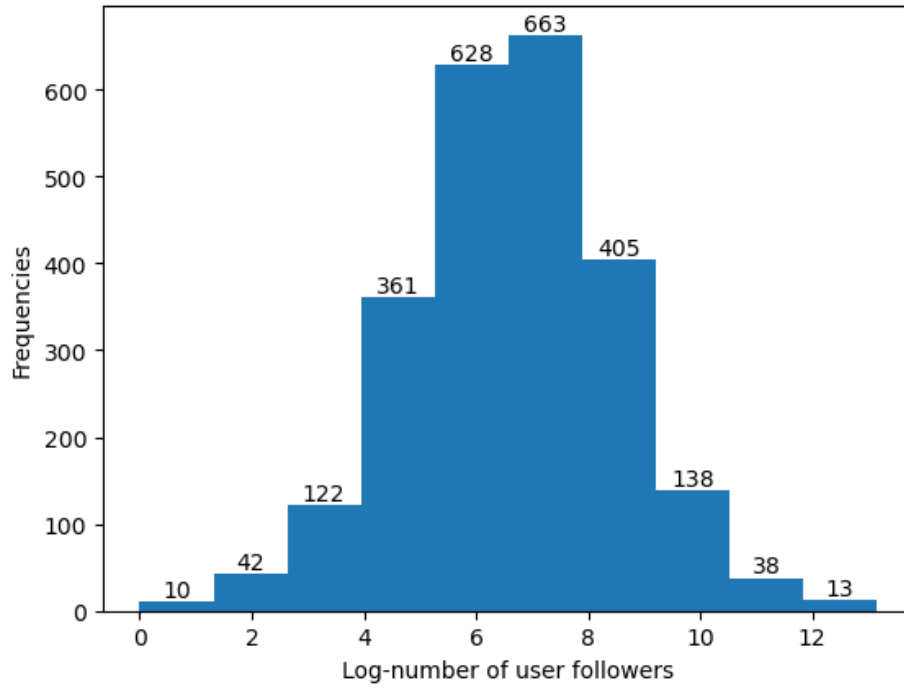
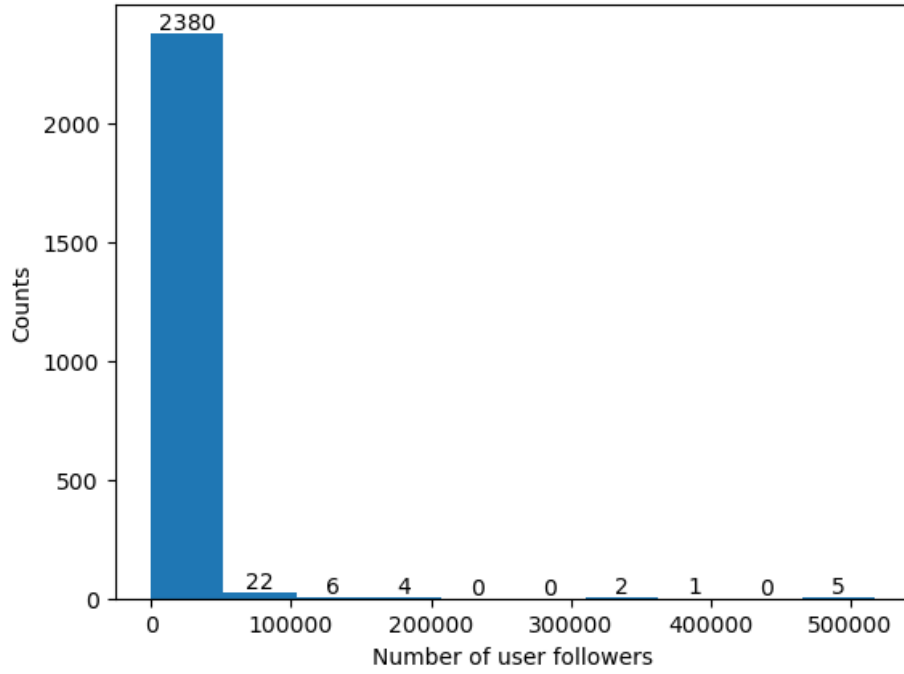
- Retweet count vs retweet log



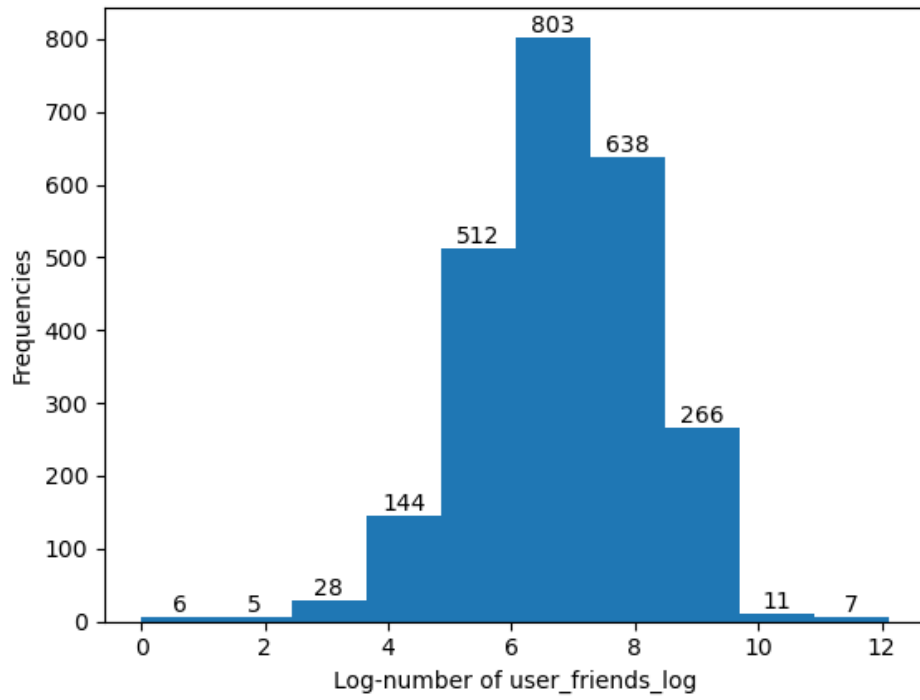
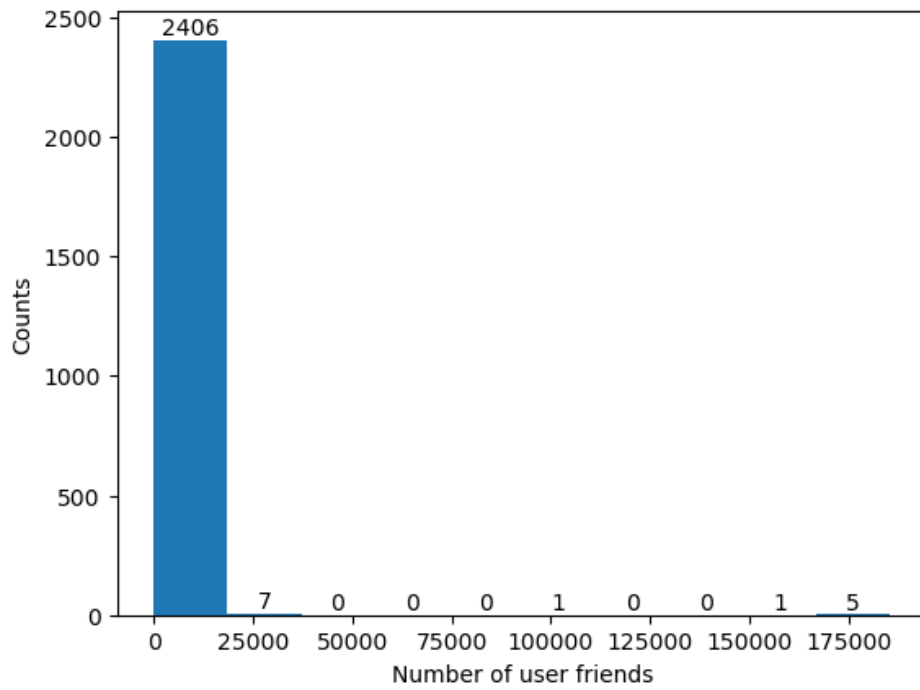
- User favourite count vs user favourite log



- User followers count vs user followers log



- User friends count vs user friends log



- User listed count vs user listed log

