

---

Electronic Thesis and Dissertation Repository

---

11-10-2020 1:00 PM

## Making Sense of Online Public Health Debates with Visual Analytics Systems

Anton Ninkov, *The University of Western Ontario*

Supervisor: Sedig, Kamran, *The University of Western Ontario*

A thesis submitted in partial fulfillment of the requirements for the Doctor of Philosophy degree in Library & Information Science

© Anton Ninkov 2020

Follow this and additional works at: <https://ir.lib.uwo.ca/etd>



Part of the [Cataloging and Metadata Commons](#), [Databases and Information Systems Commons](#), [Data Science Commons](#), [Graphics and Human Computer Interfaces Commons](#), and the [Health Sciences and Medical Librarianship Commons](#)

---

### Recommended Citation

Ninkov, Anton, "Making Sense of Online Public Health Debates with Visual Analytics Systems" (2020). *Electronic Thesis and Dissertation Repository*. 7429.  
<https://ir.lib.uwo.ca/etd/7429>

This Dissertation/Thesis is brought to you for free and open access by Scholarship@Western. It has been accepted for inclusion in Electronic Thesis and Dissertation Repository by an authorized administrator of Scholarship@Western. For more information, please contact [wlsadmin@uwo.ca](mailto:wlsadmin@uwo.ca).

## Abstract

Online debates occur frequently and on a wide variety of topics. Particularly, online debates about various public health topics (e.g., vaccines, statins, cannabis, dieting plans) are prevalent in today's society. These debates are important because of the real-world implications they can have on public health. Therefore, it is important for public health stakeholders (i.e., those with a vested interest in public health) and the general public to have the ability to make sense of these debates quickly and effectively. This dissertation investigates ways of enabling sense-making of these debates with the use of visual analytics systems (VASes). VASes are computational tools that integrate data analytics (e.g., webometrics or natural language processing), data visualization, and human-data interaction.

This dissertation consists of three stages. In the first stage, I describe the design and development of a novel VAS, called VINCENT (VISual aNalytiCs systEm for investigating the online vacciNe debaTe), for making sense of the online vaccine debate. VINCENT helps users to make sense of data (i.e., online presence, geographic location, sentiments, and focus) from a collection of vaccine focused websites. In the second stage, I discuss the results of a user study of VINCENT. Participants in the study were asked to complete a set of ten sense-making tasks that required investigating a provided set of websites. Based on the positive outcomes of the study, in stage three of the dissertation I generalize the findings from the first two stages and present a framework called ODIN (Online Debate entlTy aNalyzer). This framework consists of various attributes that are important to consider when analyzing online public health debates and provides methods of collecting and analyzing that data. Overall, this dissertation provides visual analytics researchers an in-depth analysis on the considerations and challenges for creating VASes to make sense of online public health debates.

**Keywords:** Visual Analytics, Online Public Health Debates, Vaccine Debate, Statin Debate, Cannabis Debate, Dieting Plan Debate, Webometrics, Natural Language Processing, Human-Data Interaction, Data Visualization

## Co-Authorship Statement

I was primarily responsible for the research presented in this dissertation, including conceptualization, study design, data collection, analysis, and manuscript preparation. The concepts presented in Chapters 3, 4, and 5 were written with the guidance and assistance of Dr. Sedig. Chapters 3 and 4 have been published as co-authored papers, and Chapter 5 is submitted for publication. Dr. Sedig has read this statement and is in agreement with it.

## Dedication

This work is dedicated to all the people of the world.

## Acknowledgments

There are many people who have provided various forms of guidance, assistance, mentorship, and encouragement to the completion of this project. I would like to take this opportunity to name and thank these people. If I forget anyone, as I'm sure will be the case, my deepest apologies and know that your contributions did not go unnoticed.

First, I would like to thank my supervisor, Dr. Kamran Sedig. When I started my PhD, I was a student with big ideas and many interests but needed direction and focus to harness them. I was interested in examining how interaction with information can help people learn and think about difficult concepts more clearly. After meeting Kamran, I knew that he would be able to help me develop my ideas in a meaningful and impactful way. I am very proud of the work we have done these past few years. I believe the ideas and concepts explored in this dissertation are steppingstones toward the development of tools that will enable people to think deeper and in ways we currently can't imagine. Thank you, Kamran, not only for guiding me during this research but for also helping me become a better and more complete person.

I would like to thank my thesis committee, Dr. Liwen Vaughan and Dr. Victoria Rubin, for their help in shaping my research through its various phases. Dr. Vaughan commenced my journey down the path of studying webometrics and guided me through publishing my first two research articles. Liwen, thank you for giving me some much-needed critical advice on how to conduct studies and publish research. Dr. Rubin introduced me to natural language processing, which has been an important component of my research. Thank you, Vicki, for all the resources you shared with me and suggestions you made to help me consider how to incorporate natural language processing in my research. As well, thank you for your insightful and considerate feedback throughout my dissertation. I would also like to thank all of the participants in my study. Without you, this research would not have been possible. I appreciate the time you gave, both in completing the sense-making and interview sessions.

I must thank all of my friends who made this journey possible. To my fellow FIMS graduate students (Yimin Chen, Nafiz Zaman Shuva, Nicole Dalmer, Jennifer Opoku, Jill Veenedaal, Zak Bronson & many others) and London based friends (Evan Nagel, Tanner Forester), who supported me in so many ways and made my time there so pleasant, thank you. To my fellow colleagues at the INSIGHT lab (Brent Davis, Jon Demelo, Amir Haghighati, Elaine Zibrowski,

Mozhgan Parsa, Parinaz Nasr), your support both technically and morally was greatly appreciated and I thank you. To my friends in Ottawa (Dave Lowe, Matt Richler, Caitlin McIntyler, Marc-André Mineau, Kimberley Do, Andrew Nguyen, Steve Closs) and from afar (Ishtar Laguna Monroy, Vinicius Paranhos, Andy Read, Jesse Morse, Robbie Irvine, James Joseph Hamme Jr.) thank you for always reminding me who I am and making my life more pleasant and fulfilling.

I would also like to thank Paul Benedetti, a colleague and former professor at FIMS, for all his help. You have been a mentor and great friend to me. Thank you for always being willing to listen and discuss anything with me. It was always great to work with you.

I would like to extend a special thanks to two dear friends of mine: a fellow LIS PhD graduate from FIMS, Dr. Claire Burrows, and her husband Dr. Tom Corkett, who for the better part of 2 years hosted me in their home in London while I worked and studied remotely from Toronto. Knowing that the door to their house was always open for me and that there would be dinner and loving animals ready for me made my complicated situation at the time much more bearable.

To my parents, Zoran and Marie, words can't express how grateful I am for everything you have done for me. You have supported me in every way imaginable my whole life and have believed in me even when I wasn't sure if I believed in myself. You share in all my successes forever. To my brother, Milan and his wife Hannah, thank you so much for your unwavering support. To my sister Gabrielle, thank you so much for always making me smile and keeping my dreams big. To Jean, thank you for always encouraging me, especially with computers. To my uncles, aunts, and cousins, thank you for helping me along the way and making me who I am today. I love you all.

Finally, I would like to thank my incredible wife, Marie-France, who I could not have done this without. Thank you as well to her whole family who have encouraged me throughout this process. Marie-France, since I met you, you have inspired me to be the best version of myself I can be. I recognize that I may not have always been at my best throughout the process of writing this dissertation. However, you gave me time, patience, encouragement, and support even when I had difficulty accepting them. For that, I am forever grateful. You deserve all the happiness in the world, and I hope that this shared achievement brings you some. And don't forget: "May the four winds blow you safely home".

## Contents

Abstract.....	ii
Co-Authorship Statement.....	iii
Dedication.....	iv
Acknowledgments.....	v
Contents .....	vii
List of Tables .....	xi
List of Figures.....	xii
Chapter 1 – Introduction.....	1
1.1. Research Purpose and Questions .....	3
1.2. Contributions.....	4
1.3. Motivation.....	6
1.4. Structure of the Dissertation .....	7
1.5. References.....	9
Chapter 2 – Background .....	12
2.1. Information Spaces, Sense-Making, and Distributed Cognition .....	12
2.2. Visual Analytics Systems .....	14
2.2.1. Data Analytics, Webometrics, and Natural Language Processing .....	16
2.2.2. Data Visualization.....	22
2.2.3. Human-Data Interaction .....	23
2.3 Design Frameworks .....	24
2.4. Online Public Health Debates.....	25
2.5. Summary .....	28
2.6. References.....	29
Chapter 3 - VINCENT: A visual analytics system for investigating the online vaccine debate ..	41
Abstract.....	42
3.1. Introduction.....	42
3.2. Background.....	43
3.2.1. Vaccine Debate .....	44
3.2.2. Visual Analytics Systems (VASes) .....	45
3.2.3. Webometrics .....	47
3.2.4. Natural Language Processing (NLP) .....	48
3.3. System Design .....	49
3.3.1. Analytics Engine.....	50

3.3.2. Data Visualizations .....	53
3.3.3. Human-Data Interactions .....	57
3.4. Summary and Conclusions .....	62
3.4.1. Limitations .....	64
3.4.2. Future research.....	65
3.6. References.....	66
Chapter 4 - The Online Vaccine Debate: Study of A Visual Analytics System.....	73
Abstract.....	74
4.1. Introduction.....	75
4.2. Background.....	76
4.2.1. Online Public Health Debates.....	77
4.2.2. Visual Analytics Systems (VASes) .....	78
4.3. Methodology .....	79
4.3.1. VINCENT: Treatment Instrument.....	80
4.3.2. Sense-Making Session .....	83
4.3.3. Interview Session.....	86
4.4. Results.....	86
4.4.1. Performance Results .....	87
4.4.2. Response to VINCENT .....	94
4.4.3. Usability of VINCENT.....	107
4.5. Discussion and Conclusions .....	112
4.5.1. Overall .....	112
4.5.2. Webometrics .....	113
4.5.3. NLP.....	115
4.5.4. Considerations for Developing VASes for Online Public Health Debates .....	116
4.5.5. Limitations .....	117
4.5.6. Future Research .....	118
4.6. References.....	120
Chapter 5 - Online public health debates: A framework-based approach to visual analytics ....	124
Abstract.....	125
5.1. Introduction.....	125
5.2. Background.....	128
5.2.1. Information Spaces & Sense-Making .....	128
5.2.2. VASes .....	129



5.2.3. Online Public Health Debates .....	131
5.3. ODIN.....	135
5.4. Online Public Health Debates – Four Case Studies .....	140
5.4.1. Vaccines .....	140
5.4.2. Cannabis.....	144
5.4.3. Statins.....	148
5.4.4. Dieting Plans.....	153
5.5. ODIN-Based Design of Visual Analytics Systems.....	159
5.5.1. Vaccines.....	160
5.5.2. Statins.....	163
5.5.3. Cannabis.....	165
5.5.4. Dieting Plans.....	168
5.6. Discussion .....	170
5.6.1. Limitations & Future Research.....	171
5.7. References.....	173
Chapter 6 – Conclusions .....	188
6.1. Dissertation Summary.....	188
6.2. Research Contributions and Conclusions .....	188
6.3. Discussion .....	190
6.4. Limitations .....	193
6.5. Future Research .....	196
6.6. References.....	200
Appendices.....	203
Appendix A - Set of Websites .....	203
Appendix B - Tasks .....	204
Appendix C - Defined Terms.....	208
Appendix D - No System Post-Tasks Questionnaire .....	209
Appendix E - System Post-Tasks Questionnaire .....	210
Appendix F - Interview Questions.....	213
Appendix G – Ethics Approval for Study.....	214
Appendix H – Flyer for Study .....	215
Appendix I – Poster for Study .....	216
Appendix J – In-Class Verbal Recruitment Script.....	217
Appendix K – Participation Number Form.....	218

Appendix L – Letter of Information and Consent for Exploration Session.....	219
Appendix M – Letter of Information and Consent for Interview Session .....	224
Curriculum Vitae .....	228

## List of Tables

Table 1 Overall Achievement Results .....	89
Table 2 Online Presence Tasks .....	90
Table 3 Online Presence Sub-Tasks .....	91
Table 4 Websites' Locations Tasks .....	92
Table 5 Word Frequency Tasks .....	93
Table 6 Text-Based Emotion Analysis Tasks.....	94
Table 7 Overall Easiness of and Confidence in Response to Tasks .....	96
Table 8 Online Presence Tasks Easiness .....	101
Table 9 Geographic Location Tasks Easiness .....	103
Table 10 Focus Tasks Easiness.....	105
Table 11 Website Text Emotion Tasks Easiness .....	106
Table 12 Attributes of ODEs .....	139

## List of Figures

Figure 1 Sense-Making Loop for VASes.....	16
Figure 2 VINCENT: A Visual Analytic System.....	50
Figure 3 MDS Similarity Map.....	54
Figure 4 Word Cloud for all Websites.....	55
Figure 5 Map of Website Locations.....	56
Figure 6 Emotion Bar Chart.....	57
Figure 7 Website Selection Interaction.....	58
Figure 8 VINCENT after Website Selection Interaction.....	58
Figure 9 Global Filtered Selection (North Eastern North America).....	59
Figure 10 Filtered Vaccine Selection (MMR).....	60
Figure 11 Hover to Expand Information Box.....	61
Figure 12 Navigate Map of Website Locations.....	62
Figure 13 Screenshot of VINCENT.....	81
Figure 14 Boxplot of Tasks Completed.....	88
Figure 15 Boxplot of Average Score.....	88
Figure 16 Box Plot of Easiness.....	94
Figure 17 Box Plot of Confidence.....	95
Figure 18 Sense-Making Loop.....	131
Figure 19 Vaxxter Website “News” Page.....	143
Figure 20 Voices for Vaccines’ Website Homepage.....	144
Figure 21 SAM Website “The Victims of Marijuana” Page.....	147
Figure 22 NORML Website Homepage.....	148
Figure 23 Dr. Joseph Mercola Website Article on Statins.....	152
Figure 24 Take Cholesterol to Heart Website Homepage.....	152

Figure 25 Take Cholesterol to Heart Website, Cholesterol Facts Page.....	153
Figure 26 Association for Size Diversity and Health Website, Healthy at Every Size Page ....	157
Figure 27 Vegan Website, Lifestyle Page.....	158
Figure 28 Vegan Website, Chicken Killing Foam Page .....	158
Figure 29 VINCENT, a VAS for Vaccines .....	160
Figure 30 VAS for Statins.....	163
Figure 31 VAS for Cannabis.....	165
Figure 32 VAS for Cannabis with Selections Made.....	166
Figure 33 VAS for Dieting Plans.....	168

## Chapter 1 – Introduction

Online debates on a wide variety of issues are a frequent occurrence in today's society. A common subset of online debates that merits scrutiny is online public health debates. Many people today are deeply concerned about their physical well-being, and as a result of this concern, they invest substantial time and effort into investigating various health-related topics. To learn about health topics, people rely on the Internet as a key source of information (S. Fox & Duggan, 2013). What they encounter in their searches can influence how they come to make decisions, what they believe, and how they act (S. Fox & Duggan, 2013; Miller & Bell, 2012). As a result of this influence, many websites and social media profiles that offer a plethora of information on such topics have emerged (Kitchens, Harle, & Li, 2014; Swar, Hameed, & Reyshav, 2017). Although some of this information is very informative, making sense of the large quantity and varying quality of information pertaining to a given debate can be a difficult and confusing endeavour (Seymour, Getman, Saraf, Zhang, & Kalenderian, 2015), especially when different pieces of information conflict with one another (Truumees et al., 2020; Yoon, Sohn, Choi, & Jung, 2017).

An example of a debated health-related topic of great interest to many is vaccines. Although vaccines in general have been described as one of the greatest medical achievements of all time (Chan, 2018), for a variety of reasons many people have fears about the possible negative effects that vaccines can have on their own or their families' health (Kata, 2010, 2012). This negative attitude towards the use of vaccines has become increasingly common in recent years, and people are spending more time educating themselves on this topic than they ever have before. As people explore vaccine-related information, they are likely to confront various sub-topics of the debate that are at odds with one another. For example, on the Internet, there is much conflicting information on whether vaccines cause autism in children (Jang, McKeever, McKeever, & Kim, 2019). The consequence of conflicting information is the emergence of online camps that have different positions and arguments. In this dissertation, I refer to the collection of all such positionings and entrenched camps on a topic as an "online public health debate."

A variety of other online public health debates also involve intensely held views that have social and other impacts on our world, such as the decisions on what people eat, the safety of statins or cannabis, as well as whether e-cigarettes are a healthier alternative to smoking (Kickbusch, 2009; Morphett, Herron, & Gartner, 2019; Velardo, 2015). The impact that such debates can have on public health decisions requires that these debates be investigated quickly and effectively. To date, such investigations have been completed by researchers drawing expertise from a range of disciplines. While such cross-disciplinary research is necessary, in order to complete thorough investigations of online public health debates, it would be helpful to remove the barriers encountered that make it difficult to make sense of these debates and the information contained within them. Furthermore, there is a need for tools that give public health stakeholders as well as the general public the ability to investigate online public health debates<sup>1</sup> (or online debates in general), an activity that requires the person undertaking the investigation to perform sense-making tasks.

In the context of the online vaccine debate, examples of sense-making tasks that the investigating person needs to complete include: 1) determining if there are more pro-vaccine or anti-vaccine websites; 2) comparing the online presence of various websites; 3) identifying geographic regions with greater concentrations of pro-vaccine or anti-vaccine websites; 4) comparing if two websites focus on similar issues related to the vaccine debate. Completing tasks such as these can be challenging. To help facilitate these tasks, computational tools can be useful, or even essential.

Computational tools can help to alleviate some of the difficulties encountered while making sense of online public health debates (Jonassen, 1995; Liu, Nersessian, & Stasko, 2008; Sedig, Klawe, & Westrom, 2001). Many different subsets of computational tools exist. Two such examples that have been particularly useful for stakeholders in public health are data analytics tools and interactive visualization tools (Ola & Sedig, 2014). While these computational tools are beneficial in addressing some of the challenges faced by public health stakeholders, they are

---

<sup>1</sup> When I refer to "investigating online public health debates" throughout this dissertation, I am referring only to making sense of the data and the structures of these debates, not determining the validity or merit of any specific debate position.

alone not enough to completely support cognitive activities that are based on large and complex information spaces (Keim, Mansmann, & Thomas, 2010). A subset of computational tools that integrates both data analytics and interactive visualization tools to make sense of complex information spaces such as online public health debates is visual analytics systems (VASes) (Sedig & Parsons, 2016).

VASes are made up of three components: data analytics, data visualizations, and human-data interaction (Pirulli & Card, 2005; Sedig & Parsons, 2016). VASes can enable the general populace to quickly examine and make sense of complex information like online public health debates. Beyond the general populace, other stakeholders such as public health practitioners and policy makers also need to make sense of these debates quickly and accurately. It is therefore important not only for such systems to help stakeholders to perform the necessary sense-making activities, but also for them to be human-centered and thus adapted to users' needs. In order to fulfill these criteria, generalizing the process and developing frameworks are important (Sedig, Naimi, & Haggerty, 2017).

## 1.1. Research Purpose and Questions

There is currently a great need for research about online public health debates. With the increasing presence of the Internet in many people's everyday life, never before has so much information been instantaneously available to the general public. This access to information can result in many positive outcomes, such as the ability for people to be aware of up to date information and news items. However, with so much misinformation and disinformation (C. Fox, 1983) being widely distributed online, especially in relation to health-focused topics, it is critical for researchers to develop ways to make sense of how this information manifests itself as online debates. VASes can be of particular use in this effort.

Research on visual analytics has experienced a period of great growth and interest in recent years, and it has had massive implications for connected research in areas such as social sciences, business, and health care. Machine learning and artificial intelligence tools for analyzing various types of information posted on the Internet are being created and improved all the time. Simultaneously, these tools are becoming more accessible to researchers via both open-access and paid-access. In parallel to this, tools for developing interactive data visualizations



such as (but not limited to) Tableau, D3.js, Google Analytics, or Microsoft Power BI are all making the development of VASes more accessible and these tools have become more prevalent among researchers in recent years. Merging these two rapidly developing areas could be a fruitful way to facilitate investigating information from online public health debates.

While there has been visual analytics research that has looked into online debates and the spread of misinformation online (Kwon et al., 2015; Steed, Drouhard, Beaver, Pyle, & Bogen, 2015), none has specifically examined the problems related to online public health debates using webometrics and natural language processing data analytics techniques. This is for two primary reasons. First, as mentioned, tools that make this type of research feasible have not always existed. Second, online public health debates are more prevalent today than previously, and there is concern about the growing number and intensity of online debates related to public health. This growth is likely connected to the rising popularity of alternative medicine and distrust in medical institutions (Coulter & Willis, 2004; Lanzarotta & Ramos, 2018; Machado, de Siqueira, & Gitahy, 2020).

Opportunity now exists for research on VASes that help people make sense of online debates. Furthermore, there is particularly a need to understand: 1) if VASes that facilitate sense making of online public health debates can be developed; 2) if so, how useful they are to people in making sense of these debates; and 3) how to generalize the process of creating these systems. With that in mind, the research questions that this dissertation addresses are as follows:

- *Is it feasible to integrate webometrics, natural language processing, data visualization, and human-data interaction into a VAS for making sense of online public health debates?*
- *Can VASes help to facilitate making sense of online public health debates, and if they can, how might they do so?*
- *What considerations must go into the development of VASes for making sense of online public health debates?*

## 1.2. Contributions

The key contributions of this dissertation are: 1) Design and implementation of VINCENT (Visual aNalytiCs systEM for investigating the online vacciNe debaTe), a novel VAS created to

facilitate the investigation of the online vaccine debate; and 2) ODIN (Online Debate entity analyzer), a framework developed for creating VASes for online public health debates based on what I learned from developing VINCENT. These contributions were developed over three stages. In the first stage, I designed and developed VINCENT. In the second stage, I user-tested VINCENT to see if it could help facilitate sense-making for the online vaccine debate. In the third stage, based on the promising findings of the user-testing of VINCENT, I created ODIN to guide the development of future VASes for making sense of online public health debates.

In the first stage, I was able to determine that it was feasible to integrate data analytics about information from a selection of websites with data visualization and human-data interaction in a seamless manner. VINCENT was designed to help users explore visualizations of data from a group of vaccine-focused websites. These websites vary in their position on vaccines, online presence, topics of focus about vaccines, geographic location, and sentiment towards the efficacy and morality of vaccines. While numerous VASes had been developed and studied previously, the novelty of VINCENT is that it integrates (specifically) webometrics (e.g., co-link analysis), natural language processing (e.g., text-based emotion analysis), data visualization, and human-data interaction. The result of this stage of the research showed that developing such a system was feasible and as such, I discuss some of the challenges of doing so.

In the second stage, I present the results of a user study of VINCENT to determine whether and how the system helped users to make sense of the online vaccine debate. The study compared people's ability to complete sense-making tasks about the online vaccine debate with and without VINCENT. The findings of this study showed that users of VINCENT performed better on the tasks, found it easier to perform the tasks, and had more confidence in their performance on the tasks than those who did not use the system. Additionally, based on user feedback from this stage, I identified a few issues associated with the system that should be taken into consideration when developing future VASes for online public health debates.

Based on the positive results from the first two stages, in the third stage I developed ODIN (Online Debate entity analyzer). ODIN is a framework for generalizing online public health debates and is based on a construct that I call Online Debate Entities (ODEs). ODEs are the various online sources of information pertaining to a given subject. That is, they have online

presences (e.g., websites or Twitter profiles) and belong to organizations and people who debate public health topics. The ODIN framework assists with the analysis of various attributes of ODEs which are needed to permit stakeholders to quickly make sense of online public health debates, or any online debates. In this framework, I identify and define seven online debate attributes (presence, shared presence, geographic location, age, registrant, focus, emotion). Using four examples of online public health debates (vaccines, cannabis, statins, and dieting plans), I demonstrate how ODIN can be used not only for systematizing analysis of ODEs, but also for helping design framework-based VASes that facilitate stakeholders' investigations of other online public health debates.

### 1.3. Motivation

I first started researching topics related to online public health debates in 2015. I began by working on a webometric study of the online vaccine debate. In that study, my coauthor and I found that the debate was very polarized and that many important distinctions could be made between the two contradicting positions of the debate (anti-vaccine and pro-vaccine) (Ninkov & Vaughan, 2017). After this study, we continued the work on the online vaccine debate with another project in which we demonstrated how Multi-Dimensional Scaling, based on co-link analyses, could be useful for identifying clusters of similar websites within each position of the debate (Vaughan & Ninkov, 2018).

These two research projects highlighted the complexity of the online vaccine debate and prompted me to further examine online public health debates, though from a different perspective. I perceived a need to go beyond investigating and reporting the structures and characteristics of online public health debates and to consider new ways of making the complex information from these debates accessible to the general public and stakeholders. This type of consideration was needed because the relevant data for making sense of online public health debates is not always easy to access or understand (Marshall & Bly, 2005). In other words, owing to a lack of information that they can access or meaningfully engage with, people who attempt to investigate these debates without any assistance can miss out on developing a clear picture about what is occurring. For example, if in their investigation people only come across websites that take a certain position in a debate, they may think that there is no debate or that one

position is much more representative than alternate views. It has been demonstrated that VASes are useful in helping people with sense-making efforts (Jonassen, 1995; Liu et al., 2008; Sedig et al., 2001).

The importance of VASes can be appreciated if one considers the concept of distributed cognition (Hutchins, 2006; Salomon, 1993), which proposes off-loading some of the cognitively taxing tasks people face onto computers so that they are freed to think about bigger issues. In other words, by developing VASes that make information from online public health debates accessible to users, it is possible to help them to make sense of the information by having the VAS analyze the data – though not interpret the results - for them. VASes that do this are valuable because most members of the general public and stakeholders are not experts in statistics or data analytics. The threshold to understand the data analysis techniques required to make sense of online public health debates creates a barrier between individuals and the information they need to fully understand the structures of these debates. This research was sparked by my interest in examining how VASes can be developed to eliminate these barriers and help to support those who want or need to make sense of online public health debates.

Furthermore, aside from developing VASes for specific instances of online public health debates, it is important to generalize the process of analyzing these debates to help facilitate research in this area of application of visual analytics. The goal of developing an adequate framework capable of generalizing the analysis requirements of online public health debates was to provide future researchers with a basis for understanding how these debates can be analyzed. Based on the generalizations of this framework, in context with the rest of this dissertation, there is great potential for future growth in this research area.

#### **1.4. Structure of the Dissertation**

The dissertation will be structured as follows. Chapter 2 provides an overview of the background of the topics related to this research. The chapter includes discussions on information spaces, sense-making, distributed cognition, visual analytics, design frameworks, and online public health debates. Chapter 3 is a discussion on the development of VINCENT. In this chapter, explanations are given for the design decisions made regarding the data analytics, data visualizations, and human-data interactions of the system. Chapter 4 discusses the results of a

user study of VINCENT. In this chapter, the experimental design of the user study is described, and the results of the study are analyzed and discussed. Chapter 5 builds on the work and findings from Chapters 3 and 4 to present a framework to guide the development of VASes for online public health debates. Chapter 6 is a reflection on the research conducted in this dissertation. In this chapter, the findings of this research are discussed, the contributions and limitations of the dissertation are stated, and recommendations for future research in the area are considered.

## 1.5. References

- Chan, M. (2018). *Ten years in public health 2007-2017: Report By Dr Margaret Chan Director-general World Health Organization*. World Health Organization.
- Coulter, I. D., & Willis, E. M. (2004). The rise and rise of complementary and alternative medicine: a sociological perspective. *Medical Journal of Australia*, 180(11), 587–589.
- Fox, C. J. (1983). *Information and misinformation: An investigation of the notions of information, misinformation, informing, and misinforming*. Westport, CT: Greenwood Press.
- Fox, S., & Duggan, M. (2013). Health online 2013. *Health*, 2013, 1-55.
- Hutchins, E. (2006). *Cognition in the wild*. Cambridge, MA: MIT Press.
- Jang, S. M., Mckeever, B. W., Mckeever, R., & Kim, J. K. (2019). From social media to mainstream news: The information flow of the vaccine-autism controversy in the US, Canada, and the UK. *Health Communication*, 34(1), 110–117.
- Jonassen, D. H. (1995). Computers as cognitive tools: Learning with technology, not from technology. *Journal of Computing in Higher Education*, 6(2), 40–73.
- Kata, A. (2010). A postmodern Pandora's box: Anti-vaccination misinformation on the Internet. *Vaccine*, 28(7), 1709–1716. <https://doi.org/10.1016/j.vaccine.2009.12.022>
- Kata, A. (2012). Anti-vaccine activists, Web 2.0, and the postmodern paradigm - An overview of tactics and tropes used online by the anti-vaccination movement. *Vaccine*, 30(25), 3778–3789. <https://doi.org/10.1016/j.vaccine.2011.11.112>
- Keim, D. A., Mansmann, F., & Thomas, J. (2010). Visual analytics: how much visualization and how much analytics? *ACM SIGKDD Explorations Newsletter*, 11(2), 5–8.
- Kickbusch, I. (2009). Health literacy: engaging in a political debate. *International Journal of Public Health*, 54(3), 131–132.
- Kitchens, B., Harle, C. A., & Li, S. (2014). Quality of health-related online search results. *Decision Support Systems*, 57, 454–462.
- Kwon, B. C., Kim, S.-H., Lee, S., Choo, J., Huh, J., & Yi, J. S. (2015). VisOHC: Designing visual analytics for online health communities. *IEEE Transactions on Visualization and Computer Graphics*, 22(1), 71–80.

- Lanzarotta, T., & Ramos, M. A. (2018). Mistrust in Medicine: The Rise and Fall of America's First Vaccine Institute. *American Journal of Public Health, 108*(6), 741–747.
- Liu, Z., Nersessian, N., & Stasko, J. (2008). Distributed cognition as a theoretical framework for information visualization. *IEEE Transactions on Visualization and Computer Graphics, 14*(6).
- Machado, D. F. T., de Siqueira, A. F., & Gitahy, L. (2020). Natural stings: alternative health services selling distrust about vaccines on YouTube.
- Marshall, C. C., & Bly, S. (2005). Saving and using encountered information: implications for electronic periodicals. In *Proceedings of the Sigchi conference on human factors in computing systems* (pp. 111–120). ACM.
- Miller, L. M. S., & Bell, R. A. (2012). Online health information seeking: the influence of age, information trustworthiness, and search challenges. *Journal of Aging and Health, 24*(3), 525–541.
- Morphett, K., Herron, L., & Gartner, C. (2019). Protectors or puritans? Responses to media articles about the health effects of e-cigarettes. *Addiction Research & Theory, 1–8*.
- Ninkov, A., & Vaughan, L. (2017). A webometric analysis of the online vaccination debate. *Journal of the Association for Information Science and Technology, 68*(5), 1285–1294. <https://doi.org/10.1002/asi.23758>
- Ola, O., & Sedig, K. (2014). The challenge of big data in public health: an opportunity for visual analytics. *Online Journal of Public Health Informatics, 5*(3), e223, 1–21. <https://doi.org/10.5210/ojphi.v5i3.4933>
- Pirolli, P., & Card, S. (2005). The sensemaking process and leverage points for analyst technology as identified through cognitive task analysis. In *Proceedings of international conference on intelligence analysis* (Vol. 5, pp. 2–4). McLean, VA, USA.
- Salomon, G. (1993). No distribution without individuals' cognition: A dynamic interactional view. *Distributed Cognitions: Psychological and Educational Considerations, 111–138*.
- Sedig, K., Klawe, M., & Westrom, M. (2001). Role of interface manipulation style and scaffolding on cognition and concept learning in learnware. *ACM Transactions on Computer-Human Interaction (TOCHI), 8*(1), 34–59.

- Sedig, K., Naimi, A., & Haggerty, N. (2017). Aligning Information Technologies With Evidencebased Health-care Activities: A Design And Evaluation Framework. *Human Technology, 13*(2).
- Sedig, K., & Parsons, P. (2016). *Design of visualizations for human-information interaction: A pattern-based framework*. San Rafael, California (1537 Fourth Street, San Rafael, CA 94901 USA): Morgan & Claypool.
- Seymour, B., Getman, R., Saraf, A., Zhang, L. H., & Kalenderian, E. (2015). When advocacy obscures accuracy online: digital pandemics of public health misinformation through an anti-fluoride case study. *American Journal of Public Health, 105*(3), 517–523.
- Steed, C. A., Drouhard, M., Beaver, J., Pyle, J., & Bogen, P. L. (2015). Matisse: A visual analytics system for exploring emotion trends in social media text streams. In *2015 IEEE International Conference on Big Data (Big Data)* (pp. 807–814). IEEE.
- Swar, B., Hameed, T., & Reychav, I. (2017). Information overload, psychological ill-being, and behavioral intention to continue online healthcare information search. *Computers in Human Behavior, 70*, 416–425.
- Truumees, D., Duncan, A., Mayer, E. K., Geck, M., Singh, D., & Truumees, E. (2020). Cross sectional analysis of scoliosis-specific information on the internet: potential for patient confusion and misinformation. *Spine Deformity, 1*–9.
- Vaughan, L., & Ninkov, A. (2018). A new approach to web co-link analysis. *Journal of the Association for Information Science and Technology, 69*(6), 820–831.
- Velardo, S. (2015). The nuances of health literacy, nutrition literacy, and food literacy. *Journal of Nutrition Education and Behavior, 47*(4), 385–389.
- Yoon, H., Sohn, M., Choi, M., & Jung, M. (2017). Conflicting online health information and rational decision making: implication for cancer survivors. *The Health Care Manager, 36*(2), 184–191.



## Chapter 2 – Background

There are several concepts that require additional context to understand the ideas presented in this dissertation. In this section, I provide this context. This section will be structured as follows: Section 2.1. discusses information spaces, sense-making, and distributed cognition; Section 2.2. provides an overview of visual analytics systems (VASes); Section 2.3. describes and provides examples of design frameworks; and Section 2.4. provides an overview of online public health debates.

### 2.1. Information Spaces, Sense-Making, and Distributed Cognition

Information spaces are bodies of information that are thought to have spatial characteristics (Sedig & Parsons, 2016). Compared to the related concept of “data”, which refers to information that has already been discerned and recorded, an information space is a useful concept for visual analytics research because it offers the freedom to conceptualize unstructured or abstract information. Information spaces are made up of information items (e.g., entities, properties and relationships) that exist at various levels of granularity. For online debates, information spaces are important because these debates are unstructured and the necessary data are not always readily or easily available (Moreno, Ozogul, & Reisslein, 2011; Snyder, 2014; J. C. Thomas, Diament, Martino, & Bellamy, 2012). The entities, properties, and/or relationships may not be clear to researchers, and the methods by which researchers can measure these characteristics may also not be apparent. Making sense of an information space such as an online debate is an example of a cognitive activity.

Cognitive activities are part of everyday life. Examples of everyday cognitive activities can include things such as: 1) preparing a schedule for the day; 2) evaluating the quality of a news item; 3) ranking your favourite musician’s albums. Cognitive activities that are information-intensive and involve intense human cognition can be further described as complex cognitive activities (Ericsson & Hastie, 1994; Funke, 2010) and have two distinct characteristics: 1) they require the use of complex psychological processes and 2) they exist in the presence of complex conditions (Knauff & Wolf, 2010). An example of a complex cognitive activity that falls within the primary focus of this dissertation, is making sense of online public health debates.

The concept of sense-making has been used in different ways by researchers from a variety of disciplines. One common conception of sense-making, especially in the areas of Communication and Library and Information Science, is that it is an approach for studying the making of sense in a communication situation (Dervin & Frenette, 2000). The approach was developed to assess how people make sense of “their intersections with institutions, media, messages, and situations” (Dervin, 1999). According to this understanding, sense-making involves studying how people move through time and space to bridge gaps in their understanding of a variety of things, and the goal of the researcher is to understand this process more deeply.

In this dissertation, however, I consider sense-making as a cognitive activity in which people gradually develop mental models of an information space about which they have insufficient knowledge (Klein, Moon, & Hoffman, 2006; Sedig & Parsons, 2013). Sense-making activities are formed from a set of tasks, some of which can include: scanning the information space, selecting relevance of items, and examining items in more detail (Pirolli & Card, 2005). Sense-making can require people to take complex information and uncover from it meaning that otherwise could go unnoticed. The problems that sense-making tasks try to solve are often ill-structured and open ended (Buchel & Sedig, 2016). Sense-making involves users establishing goals, discovering how an information space is structured, and determining what questions to ask as well as how the answers to those questions should be organized (Russell, Stefik, Pirolli, & Card, 1993). A challenge that individuals face when attempting to complete sense-making activities can be that relevant information is not always easy to access, stored in the proper format, or kept in the correct locations (Marshall & Bly, 2005). People may find the sense-making tasks that are required for complex information spaces difficult. To help them to deal with this difficulty, tools that distribute the cognitive load of sense-making are important.

Distributed cognition (Hutchins, 2006; Norman, 1993; Salomon, 1993) is a theory about human cognition that can help researchers to understand how to help people to perform sense-making activities on complex information spaces. Distributed cognition extends the reach of what is cognitive beyond just the individual and includes in the process things from the environment such as other individuals or cognitive tools (Hollan, Hutchins, & Kirsh, 2000). The theory views these external resources as an extension of the mind. People’s attempts at sense-making activities are more successful if done in collaboration with tools that help distribute cognition because

these activities involve coordinating use of knowledge structures emanating from the mind, environment, and other individuals (Keim, Kohlhammer, Ellis, & Mansmann, 2010). Sense-making can be cognitively taxing on individuals and this results in it being difficult or, in some cases, impossible to process the necessary data without the assistance of tools. Computational tools like VASes have been found to be beneficial in helping people to complete these types of activities because they distribute some of the cognitive load of the tasks undertaken (Jonassen, 1995; Z. Liu, Nersessian, & Stasko, 2008; Pohl, Smuc, & Mayr, 2012; Sedig, Klawe, & Westrom, 2001).

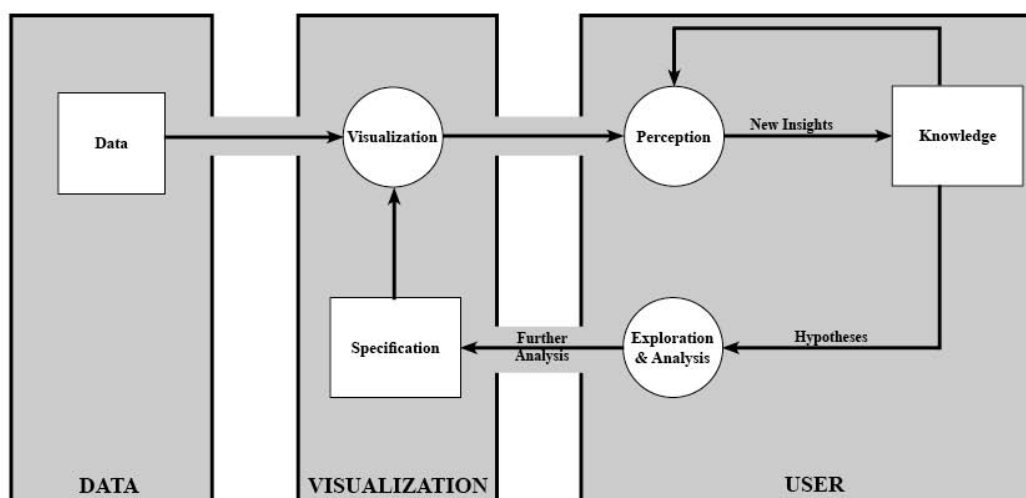
## 2.2. Visual Analytics Systems

People are often victims of information overload, which can result in them becoming lost in data and overwhelmed by what it means (Keim et al., 2008). With the rise of “big data” (Higgins et al., 2011; O’carroll, Cahn, Auston, & Selden, 1998; Ola, 2017), the need to develop visual analytic systems has never been more urgent. As a result of the growth and widespread adoption of the Internet in the general public’s daily life, people are bombarded with vast amounts and many formats of data, and there is an emerging need for information scientists to develop new ways of helping people to think about and understand the information.

One theory that can help with the conceptualization of VASes is general systems theory, which can be thought of as the science of systems (Laszlo & Clark, 1972). Systems theory views a system as composed of entities, properties, and relationships (Skyttner, 2005). VASes are complex, multi-level systems, consisting of systems within systems (Sedig & Parsons, 2016). These multi-level systems consist of super-systems, systems made up of other systems, and sub-systems (Skyttner, 2005). Applying this understanding of systems theory to VASes makes it possible to better understand how VASes work. When building and examining VASes, the interactions that the user has with the system can impact any of these levels. At higher levels, interactions with the super-system will change the broader display of VASes. At lower levels, the interaction sub-system will change specific components of the system. These interactions, regardless of level, are important to the functioning of VASes and necessary for making sense of the data being presented.

VASes can help people to deal with the problems associated with big data by integrating human insight with the powerful data analytics of computers (e.g., machine learning) with meaningful data visualization and human-data interaction (J. J. Thomas & Cook, 2005). The purpose of this integration is to help users of the system complete complex cognitive activities such as discovering patterns and knowledge in data, making decisions, or planning based on data. Visual analytic research is focused on developing systems that allow people to do whatever they need to do with the data. As Börner states: “Just like the microscope, invented many centuries ago, allowed people to view and measure matter like never before, (visual) analytics is the modern equivalent to the microscope” (Börner, 2015).

Sense-making requires people to rapidly compare and contrast information items without being hindered by the information’s format or location (Keel, 2007). VASes can be particularly useful in this effort. The sense-making loop displayed in Figure 1 structures the knowledge discovery process that a VAS user undertakes, and it clarifies how people use these systems to make sense of an information space (Keim et al., 2008; Van Wijk, 2005). In the sense-making loop, data is first analyzed and fed into the VAS. To make sense of the analyzed data, the user perceives it through a visualization. The user then gathers new insights from the visualizations of the data by performing interactions on the VAS. Based on the new knowledge they gain from these insights, users can further analyze the data and explore the information in ways that they would not have been capable of doing previously. The user will start to ask new questions, develop new hypotheses, and begin interacting with and adjusting the specifications of the VAS to see the data in new ways and conduct further analysis. The exact process (i.e., number of times one goes through the loop and the specifications) will vary depending on the tasks at hand. In other words, users control the visualizations they see and the way in which the data is analyzed based on the observations and questions they formulate while using the VAS. The sense-making loop is repeated as the user completes more tasks and generates further insights from the data.

**Figure 1***Sense-Making Loop for VASes<sup>2</sup>*

Having explained how people use VASes to make sense of information spaces, I will now discuss VASes' specific components. VASes are made up of three integrated components: data analytics, data visualizations, and human-data interaction. In the following sub-sections (2.2.1. and 2.2.2.), I will explain in further detail these components and describe some of the specific methods that can use.

### 2.2.1. Data Analytics, Webometrics, and Natural Language Processing

The data analytics component of visual analytics systems includes all tasks that deal with the collection, processing, and analysis of data. Data analytics is extremely broad and encompasses a wide variety of methods and approaches. Two specific areas of data analytics that are considered in this dissertation and when examining online public health debates are webometrics and natural language processing (NLP). I will discuss these two areas in the remainder of this section.

---

<sup>2</sup> The Sense-Making Loop has been based on (Keim et al., 2008; Van Wijk, 2005)

### 2.2.1.1. Webometrics

Webometrics, a research area of information science that comes from the related area of bibliometrics, can be described as the “quantitative study of web-related phenomena” (Thelwall, Vaughan, & Björneborn, 2005). Because so many people use and rely on the Internet, it is important to investigate the various metrics used for analyzing the data it generates. There are two types of webometrics analysis: evaluative and relational (Stuart, 2014; Thelwall, 2008).

Evaluative webometrics can include examining webpages and social media profiles for attributes such as (but not limited to) the number of external inlinks or followers they receive (i.e., online presence (Ninkov & Vaughan, 2017; Vaughan & Ninkov, 2018)), website geographic location (Ninkov & Vaughan, 2017; Stuart, 2014; Thelwall, 2004), age of the website (Jain & Gupta, 2016; Kefi & Perez, 2018; Kend & Goode, 2018; Zheng et al., 2019), and registrant of the website (Ninkov & Vaughan, 2017; Rothenfluh & Schulz, 2018). In the following paragraphs, I will discuss each of these attributes.

All websites or social media profiles have an online presence - that is, they all attract some level of attention from other online sources and/or web users. In the context of online debates, the more presence the sites or profiles have, the more popularity and/or influence they will have within the debate that they address. Presence can be quantified by various metrics. On the general web, inlinks (Ninkov & Vaughan, 2017; Thelwall, 2004; Thelwall, Sud, & Wilkinson, 2012), website traffic (Baka & Leyni, 2017; Brumshteyn & Vas'kovskii, 2017), and website rankings from online resources such as Alexa (alexa.com), MOZ (moz.com), or Majestic (majestic.com) are all powerful tools for measuring presence. On Twitter, metrics like followers (McCoy, Nelson, & Weigle, 2017; Triemstra, Poepelman, & Arora, 2018) or follower/following ratios (Anger & Kittl, 2011; Borgmann et al., 2016) have been used as indicators of online presence.

The geographic location of a website or social media profile describes where it is located in the world. This data can be collected through means of content analysis of the website or social media profile (Halavais, 2000; Holmberg & Thelwall, 2009; Ninkov & Vaughan, 2017). In the case of websites, services like WHOIS (<https://lookup.icann.org/>) that provide registration information for websites have been used to study geographic location of websites in the past

(Janc, 2016; Tsou et al., 2013). Location information regarding tweets is provided by the website itself (Stefanidis et al., 2017; Waseem & Hovy, 2016). It is worth noting here that WHOIS is no longer a viable service because of changes to privacy requirements as a result of General Data Protection Regulation of the European Union. However, these laws and services are frequently changing, and it is possible that new services that fill this need will emerge in the coming months or years.

The age of a website or social media profile corresponds to the length of time that it has been registered. For websites, the Internet Archive's Wayback Machine (<http://web.archive.org/>) is a resource that makes it possible to see how long a website has been active. On Twitter, the age of the profile is visible on the profile's homepage. The age of a website or social media profile can be used as a way to assess the evolution of a topic online (Jain & Gupta, 2016; Kefi & Perez, 2018; Kend & Goode, 2018; Zheng et al., 2019). For example, if there are many websites that take Position A on a particular topic and those sites are on average older than websites that take Position B on that same topic, this is an indication that Position B is newer to the debate or has only become popular more recently.

The registrant is the person or organization that owns a website or social media profile. The registrant is not always easy to determine based on the site's or profile's name and surface-level appearance; this makes uncovering the registrant an important task. Content analyses along with a registrant classification system can be used to collect this data (Ninkov & Vaughan, 2017; Rothenfluh & Schulz, 2018). In addition, services like WHOIS (subject to the limitations already discussed previously in this section) have been used by researchers to obtain information on the registration of a website.

Relational webometrics are different from evaluative webometrics and focus on "providing an overview of the relationships between different actors" (Stuart, 2014). Co-occurrence measurements to indicate similarity are important for relational analysis in webometrics (Stuart, 2014; Thelwall, 2004, 2009). The concept behind this method is that the more entities share occurrences (e.g., inlinks), the more likely they are to be similar in some way (Thelwall, 2008). This method can apply to webometrics in the study of co-links to help analyze similarity in terms of shared online presence between websites (Ortega & Aguillo, 2008; Vaughan & You, 2006,

2008, 2010). To represent and examine co-link data, numerous studies have been conducted with multi-dimensional scaling (MDS) to analyze business (Vaughan & You, 2006), university (Thelwall & Wilkinson, 2004), government (Holmberg, 2009), and political domains (Kim, Barnett, & Park, 2010; Romero-Frías & Vaughan, 2010). All these studies found that using MDS to analyze co-links generated worthwhile insights into the data.

Hyperlinks can be particularly useful to facilitate relational webometric research by providing “an overview of the relationships between different actors” (Stuart, 2014). Co-links, or hyperlinks shared by two websites, can be used as an indicator of relatedness and similarity between two websites (Stuart, 2014; Thelwall, 2009). The more co-links two websites have, the more likely they are to be related to each other. The origins of co-link analysis come from co-citation analysis, a bibliometric research method where relations between articles and authors are evaluated by looking at the citations they share (Leydesdorff & Vaughan, 2006; Marshakova, 1973; Small, 1973).

Normalization of co-link data can also be an important step in conducting co-link analysis. While the number of co-links is an important indicator of relatedness or similarity, the number of total inlinks a site has can also be an important consideration. To this end, the Jaccard Index has been used in the past (Leydesdorff & Vaughan, 2006; Small, 1973; Vaughan & Ninkov, 2018). This normalization method is explained in the following equation, as outlined in (Small, 1973; Vaughan & You, 2006):

$$\text{NormalizedColinkCount} = n(A \cap B) / n(A \cup B)$$

This equation can be understood as follows: A is the set of web pages that link to site X, B is the set of web pages that link to site Y,  $n(A \cap B)$  is the number of pages which link to both site X and Y (raw co-link count), and  $n(A \cup B)$  is the number of pages which link to either site X or Y (Vaughan & You, 2006).

MDS uses proximity to indicate how similar or different two objects are (Kruskal & Wish, 1978). By using MDS to map out relationships between websites with co-link analysis, people can identify the hidden structures in data, ascertain groupings, and assess the relatedness of websites, which in turn can help make navigating the data easier and more comprehensible



(Kruskal & Wish, 1978). MDS allows the generation of visual representations of what otherwise would be an invisible space, and it has the potential to reveal information about the world that could go unnoticed to viewers. MDS mapping can be based on two different proximity measures: similarity or dissimilarity. Co-link analysis is a representation of similarity (shared inlinks). Using the PROXSCAL option in SPSS for MDS, websites can be plotted appropriately according to their similarity measures (Leydesdorff & Vaughan, 2006).

Co-link analysis has been used in webometric research to study a variety of domains. For example, this data has been used to “pair business websites as a measure of similarity between two companies” when comparing the global market to the Chinese market (Vaughan & You, 2006). Using MDS mapping, the various sectors of the telecommunication industry have been visualized for both the global market and Chinese market. The researchers combined these findings with measurements of traditional inlink counts to determine the most competitive companies for each grouping and make observations. Again using the telecommunication industry as an example, the researchers further explored the benefit of combining co-link data with content data (keywords) (Vaughan & You, 2008) and the feasibility of using co-word data instead of co-link data (Vaughan & You, 2010) to uncover business information. Similarly, co-links and MDS have also been used to study the Nordic academic web space (Ortega & Aguillo, 2008). These researchers successfully presented the web space of Nordic academia and demonstrate the benefits of using an asymmetric matrix, which removes non-academic (relevant) sites from the co-link data (compared to a symmetric matrix which includes all inlinks).

#### *2.2.1.2. Natural Language Processing*

NLP is a vast area of research that focuses on using computational methods to understand and produce human language content (Hirschberg & Manning, 2015). NLP can encompass many different focuses, two of which that are involved in this dissertation: text-based emotion detection and word/phrase frequency (B. Liu, 2015). Text-based emotion detection has been and continues to be examined in NLP research (B. Liu, 2015; Ptaszynski, Masui, Rzepka, & Araki, 2014; Rubin, Stanton, & Liddy, 2004; Tokuhisa, Inui, & Matsumoto, 2008). There have been many different approaches to examine text based emotion (B. Liu, 2015). For example, Watson and Tellegen’s Circumplex Theory of Affect has been applied as a potential model for emotion

classification using an eight-fold categorization of emotions based on the axis of positive affect and negative affect (Rubin et al., 2004). Parrott's model of emotion classification includes primary, secondary, and tertiary emotions, which are useful because basic emotions are often not fine grained enough to label a text containing emotion (B. Liu, 2015). The Emotive Expression Dictionary by Nakamura is a Japanese-specific emotion dictionary that proposes ten emotion types that are most common and appropriate to Japanese language and culture (Ptaszynski et al., 2014). These are just a few of the many different types of emotion classification labeling systems that exist. Tokuhisa et al. studied emotion classification in conjunction with sentiment analysis (Tokuhisa et al., 2008). The researchers used sentiment polarity as a higher-level classification and then incorporated emotion classification as a way to achieve a more fine-grained analysis of the sentiment shared, and they built their emotion classification system using a k Nearest Neighbor algorithm, a supervised machine-learning system. They created an emotion-provoking corpus that was a lexicon of emotion words corresponding to ten emotions (happiness, pleasantness, relief, fear, sadness, disappointment, unpleasantness, loneliness, anxiety, and anger) (Tokuhisa et al., 2008). They found that by combining sentiment analysis with emotion classification, more powerful results that improved sentiment polarity identification could be achieved.

One resource in particular that makes it possible for researchers to quickly automate text-based emotion analysis is IBM's Natural Language Understanding (NLU) API (Vergara, El-Khouly, El Tantawi, Marla, & Lak, 2017). The NLU API (formerly referred to as the AlchemyAPI) has been used by many researchers to study topics, sentiments, and emotions in a variety of texts (Meehan, Lunney, Curran, & McCaughey, 2013; Palomino, Taylor, Göker, Isaacs, & Warber, 2016; Rizzo & Troncy, 2011; Saif, He, & Alani, 2012). The NLU API allows researchers to either input text directly or pull text from URLs of webpages and return a number of different NLP analyses, one of which is emotion analysis. Furthermore, not only can the NLU API detect emotion on the entirety of a text/webpage, but can also return emotion scores for specified target words/phrases (Vergara et al., 2017). In addition, the NLU API has been examined in numerous studies for accuracy of the result as compared to a "gold standard" (Dale, 2018; Dolianiti et al., 2019). While the tool has been demonstrated as competent in a variety of settings, other tools that have outperformed it have been adopted in recent years – for example, BioBERT - especially in subject areas like those related to public health (Lee et al., 2019).

The study of word/phrase frequency in text has been examined and used in NLP research (Healey & Ramaswamy, 2011; Katsuki, Mackey, & Cuomo, 2015; McAuley, Leskovec, & Jurafsky, 2012). One of the important concerns in word/phrase frequency analysis is how to manage meaningless or unimportant words. In English, as in any language, there are many words that are repeated frequently that are not necessarily the key point of interest to a reader. Some of the more obvious examples of these words are “the”, “and”, and “of.” Other types of undesirable words can exist depending on the domain of interest (e.g., dates or numbers). To deal with this issue, the technique of filtering for a list of stop words has been used, and preliminary lists of these words have been created that allow researchers to automatically exclude words that are not of interest (Bird, Klein, & Loper, 2009).

It is also important for word/phrase frequency analyses to consider the length of the phrases that are of interest. This concept is known in NLP as the n-gram (Suen, 1979). If a researcher is examining a set of texts to understand the occurrences of frequent word phrases, an n-gram corpus would present the number of times each phrase with a word count of “n” appears. If “n” is 1, then single words such as “vaccine” or “statins” are of interest. If “n” is 3, then phrases of 3 words – for instance, “vaccines are good” or “vaccines cause autism” – are of interest. Various studies have used different n-grams to examine word frequency in a variety of ways. For example, one study successfully pulled key phrases from a scientific document using an “n” value between 1 and 4 (Kumar & Srinathan, 2008). Furthermore, this way of analyzing word combinations has been merged with other NLP methods (e.g., part of speech analysis) in various topic-extraction efforts (Hu & Liu, 2004; B. Liu, 2015; Pang & Lee, 2008).

### **2.2.2. Data Visualization**

Data visualizations are the second component of VASes and are the visual representations of the information derived from the analytics engine. Visualizations extend the capabilities of individuals to complete tasks by allowing them to analyze data in ways that would otherwise be difficult or impossible (Sedig, Parsons, & Babanski, 2012; Shneiderman, Plaisant, & Hesse, 2013). Visualizations are produced using various techniques, or in other words from “methods or templates that can be used in specific visualization design context” (Sedig & Parsons, 2016). Visualization techniques include visual marks, visual structures, and visual variables. It is

important to generalize these visualization techniques so as to better understand the nature of visualizations. The methods of implementing visualizations can vary depending on the domain (e.g., statistics, mathematics, business, or biology). However, the generalizations from these techniques remain the same for all.

Visual marks are the fundamental building blocks of visualizations. Visual marks are classified by the number of dimensions they use. This means that there are four possible type of marks: points (zero dimensions), lines (one dimension), areas (two dimensions), and volumes (three dimensions) (Börner, 2015; Card, Mackinlay, & Shneiderman, 1999; Munzner, 2015; Sedig & Parsons, 2016). Visual structures are made up of two or more visual marks that suggest a form of organization for the information and they can be abstract or concrete. Abstract structures are broader concepts and can include some variability. An example of a concrete visual structure could be a two-dimensional Cartesian coordinate system (Sedig & Parsons, 2016). Visual variables are properties of marks or structures that encode more information into a visualization. Through these various properties, more information can be encoded into a visualization. For example, visual variables can include things like the colour, size, shape, or texture of visual marks or structures (Sedig & Parsons, 2016).

### **2.2.3. Human-Data Interaction**

Human-data interaction is the third and final component of VASes. Interaction with a VASes allows the user to control the data they see and the way it is processed. VASes support users by allowing them to distribute the cognitive workload of analyzing data by interacting with it (Z. Liu et al., 2008; Salomon, 1993; Sedig & Parsons, 2016). Interaction can be thought of in terms of three components: actions that the user performs on a VAS, the resulting changes and reactions in the visualizations of the system, and the user's perception of the changes in the system (Sedig & Parsons, 2013). Each interaction supports different epistemic actions on information by the user. By performing actions on visualizations, the user is able to change and effect the data analytics of the system. Numerous interaction patterns have been identified; a selection of these include filtering, scoping, and drilling of data (Sedig & Parsons, 2013). Filtering is when a user can show or hide a subset of information in a visualization based on a criterion. Scoping is when a user dynamically works forward and backwards to review

compositional development and growth. Drilling is when a user is able to bring out and display interior deep information.

In response to the action of the user, the reaction in the VAS is visible through the changes to the visualizations of the system. Reactions that are not perceptible to the user within the VAS also occur (Sedig et al., 2012). These can include the ways in which the data is stored or the ways in which the data will be analyzed in future actions with system. Finally, the user perceives changes to the VAS to complete the interaction process. When the different actions, reactions, and perceptions combine, they allow a dialogue between the user and the system that facilitates complex cognitive activities. This is critical for activities like making sense of online public health debates, because the ways in which the data can be used are not necessarily fully understood by the developers of these systems. In other words, enabled by VASes, the user can interact and have a dialogue with the data that allows them to ask questions that they would not have known to ask when they began their investigation.

### 2.3 Design Frameworks

In the development of VASes, design frameworks can be an important resource for designers to achieve the best possible outcomes. In the most general sense, frameworks can be thought of as tools that can help guide the design process and organize the concepts (Sedig & Parsons, 2016). In the design process, it has been shown that four forms of support have the potential to be helpful for designers (Stolterman, 2008). These forms of support include: 1) simple tools or techniques (e.g., prototypes of systems); 2) frameworks that support decision-making (e.g., design patterns for VASes); 3) individual concepts that are open to interpretation in terms of how they can be implemented (e.g., affordance); and 4) high-level theoretical approaches that expand design thinking (e.g., human-centered design). With this in mind, we can see that design frameworks that support decision-making in the development of VASes are useful because they are malleable in the sense that designers can apply and adapt the framework as needed in the context of their situation. In another sense, design frameworks are important tools for designers of VASes because they provide them with support structures for thinking about the higher-level issues, thus freeing them to think more deeply about the specific context of the information space with which they are dealing (Sedig & Parsons, 2016). Moreover, design frameworks are

important because designers of VASes, just like the general public, tend to hold onto ideas even when faced with the shortcomings of these ideas or when made aware of better alternate solutions (Cross, 2004; Ullman, Dietterich, & Stauffer, 1988). As a result of this tendency for bias, developing design frameworks to support systematic thinking about the design process can help designers of VASes to achieve better results because they eliminate some of the potential avenues for that bias to manifest itself. There are numerous examples of design frameworks that have been developed for visual analytics (Ainsworth, 2006; Andrienko & Andrienko, 2013; Card et al., 1999; Hegarty, 2011; Sedig & Parsons, 2013; Sedig & Parsons, 2016). In the next two paragraphs, I will briefly describe two of these frameworks which serve as examples of how and why frameworks are important.

An important example of one of these design frameworks, for both static and interactive visualizations, approaches the design process from the highest level and is, therefore, the most generally applicable to visual analytics (Sedig & Parsons, 2016). This framework is pattern-based and breaks the design process into four stages: 1) information space; 2) visualization patterns and blending of patterns; 3) visualization techniques; 4) concrete encoding and interactions. This framework is useful for designers of VASes because it provides a roadmap for the design process to follow and generalizes the characteristics of visualizations and interactions so that they can be incorporated and blended together into a VAS.

Another example of a framework that has been developed for visual analytics that is less general than the one previously mentioned and more directed towards a specific applications for visual analytics, proposes an approach to externalizing spatio-temporal patterns (Andrienko & Andrienko, 2013). This framework consists of several components, including: 1) cartographic map display; 2) time series display; 3) interactive tools for clustering based on one or more of existing clustering methods; and 4) an interactive visual interface. This framework proposes several resources and tools that can be used to facilitate these components and provides practical examples of how the framework can be applied (Andrienko & Andrienko, 2013).

## **2.4. Online Public Health Debates**

Online debates are topics that are widely discussed on the Internet. They encompass two or more differing points of view. There are many well-documented examples of topics that are debated

online, including: gun control (Sridhar, Foulds, Huang, Getoor, & Walker, 2015; Walker, Tree, Anand, Abbott, & King, 2012), vaccines (Kata, 2010, 2012; Ninkov & Vaughan, 2017), abortion (Bloch, 2007; Hill, 2017), and climate change (Collins & Nerlich, 2015; Howarth & Sharman, 2015). Online debates are important for researchers to investigate not only because they are complex information spaces that are not easily understood, but also because they have been shown to influence the way in which people perceive an issue and have an impact on their real-world decisions (Kickbusch, 2009; Morphett, Herron, & Gartner, 2019; Velardo, 2015). People often rely on the Internet to gather the information they need to help them to form their opinions. An important area in which this manifests itself is health-related topics.

In this dissertation, I distinguish between what I have labelled “macro-level” and “micro-level” online debates. A macro-level online debate (Getman et al., 2018; Ninkov & Vaughan, 2017) is a topic that has various prominent websites and social media profiles (macro-level online debate entities) that provide mutually contradictory information. Macro-level online debate entities don’t necessarily engage with one another – it is their collective existence and inherent contradictions that makes up to the debate. A micro-level online debate (Herring et al., 2002; Nicholson & Leask, 2012; Oraby et al., 2017), on the other hand, is an instance of when a topic is discussed between specific social media profiles (micro-level online debate entities) on a website or web forum. Micro-level online debate entities are active and involve participants directly engaging with one another online to discuss and combat the topic of debate – evoking themes or based on evidence gathered from their understanding of the macro-level online debate. In this dissertation, I focus exclusively on macro-level debates and refer to this level of debate when I use the term “debates.”

Online public health debates are information spaces that are made up of various information items, including websites or social media profiles and their corresponding attributes (mentioned in Section 2.2.1.). The websites and social media profiles that participate in an online public health debate can have contrasting views and opinions from one another. This has been well documented for several online public health debate topics, including: vaccines (Getman et al., 2018; Nicholson & Leask, 2012; Seeman & Rizo, 2009), cannabis (Bilgri, 2016; Hasan & Ng, 2014), statins (Huesch, 2017; Navar, 2019), and dieting plans (Jauho, 2016; Mazzi, 2018).



The vaccine debate is one of - if not the most - widely debated of the online public health debates. The anti-vaccination movement may appear to be a new and emerging phenomenon in light of news coverage of outbreaks of various diseases that are preventable through vaccines – for example, measles (Davenport, 2017; Sun, 2017) or whooping cough (Sun, 2017; Yourex-West, 2017). However, anti-vaccination views and sentiments are not a recent development. Since the discovery of the first vaccines, vaccination has garnered much attention, both positive and negative. From the beginning, some have felt that the practice of vaccination violates personal freedoms and/or is ineffective, or even “unchristian” (Durbach, 2000). However, vaccines have had a tremendous impact on global health and have been described as “one of the greatest achievements of biomedical science and public health” by the Centers for Disease Control (Fox & Rainie, 2000).

There are many reasons for the persistence of anti-vaccine views online in spite of the medical community’s overwhelming support for immunization. Increasingly polarized political views (especially in the United States) have generated an environment in which the rejection of scientific facts has become more prevalent and accepted (Lewandowsky & Oberauer, 2016). This erosion of trust in science among segments of the population may contribute to this increased polarization. Even the 45<sup>th</sup> President of the United States of America, Donald Trump, tweeted in 2014: “Healthy young child goes to doctor, gets pumped with massive shot of many vaccines, doesn't feel good and changes - AUTISM. Many such cases!” (RealDonaldTrump, 2014). Trump has voiced further support for these anti-vaccine leaning views over his time as president in his public speeches and statements. In 2020, the anti-vaccine views persisted in the midst of the COVID-19 pandemic, which impacted so many aspects of life. Anti-vaccine sentiments have emerged as a result of efforts to develop a vaccine that might be useful as a way of mitigating the current global health crisis (Neumann-Böhme et al., 2020). Concerns about the potential COVID-19 vaccine range from the belief that it is going to be rushed and as such will not be safe to the idea that the disease itself is a hoax and the vaccine is a conspiracy.

Furthermore the rise in accessibility to, and widespread use of, the Internet has also played a role in amplifying the anti-vaccination movement (Kata, 2010, 2012). Kata states: “The connective power of the Internet brings together those previously considered on the fringe. Members of marginalized groups (e.g. Holocaust deniers, 9/11 ‘Truthers’, AIDS deniers) can easily and



uncritically interact with like-minded individuals online... Anti vaccine groups have harnessed postmodern ideologies and by combining them with Web 2.0 and social media, are able to effectively spread their messages” (Kata, 2012). The anti-vaccination movement relies on the Internet to spread its message using many tropes and tactics that promote their view that vaccines are or can be dangerous. For example, sentiments such as “I’m not against vaccines, I’m just for safe-vaccines” or “science has been wrong before” are often repeated in defense of anti-vaccination views (Kata, 2012).

The polarity of the online vaccine debate has created a clear and observable divide (Ninkov & Vaughan, 2017) that could be having harmful effects on the health of the general population. It has been suggested by Brunson & Sobo that “providers and policymakers must begin to recognize the jagged, context-dependent, equifinal nature of how parents sort through vaccination-related information or account for their vaccination decisions in order to reverse declining vaccination rates” (Brunson & Sobo, 2017). There is a clear need for better methods of understanding online public health debates, such as those on vaccines.

## 2.5. Summary

This chapter has presented a summary of the important concepts related to this research. It has included a discussion of information spaces, sense-making, and distributed cognition; an overview of visual analytics systems (VASes); a description and examples of design frameworks; and an overview of online public health debates. With this provided background of these areas, it should help the reader contextualize and understand the approach to the research of this dissertation. This research examines how VASes can help people (both the general public and stakeholders) with making sense of online public health debates and ways to generalize the development of these systems. Online public health debates are complex information spaces that are difficult for people to make sense of because the quantity of data related to the debate and the difficulty of evaluating it in its various formats and locations. In the rest of this dissertation, we examine this problem and propose methods of making online public health debates easier for people to make sense of by using VASes.

## 2.6. References

- Ainsworth, S. (2006). DeFT: A conceptual framework for considering learning with multiple representations. *Learning and Instruction, 16*(3), 183–198.
- Andrienko, N., & Andrienko, G. (2013). A visual analytics framework for spatio-temporal analysis and modelling. *Data Mining and Knowledge Discovery, 27*(1), 55–83.
- Anger, I., & Kittl, C. (2011). Measuring influence on Twitter. In *Proceedings of the 11th international conference on knowledge management and knowledge technologies* (pp. 1–4).
- Baka, A. B. A., & Leyni, N. (2017). Webometric study of world class universities websites. *Qualitative and Quantitative Methods in Libraries, 105–115*.
- Bilgrei, O. R. (2016). From “herbal highs” to the “heroin of cannabis”: Exploring the evolving discourse on synthetic cannabinoid use in a Norwegian Internet drug forum. *International Journal of Drug Policy, 29*, 1–8.
- Bird, S., Klein, E., & Loper, E. (2009). *Natural language processing with Python: analyzing text with the natural language toolkit*. “O’Reilly Media, Inc.”
- Bloch, J. P. (2007). Cyber wars: catholics for a free choice and the online abortion debate. *Review of Religious Research, 165–186*.
- Borgmann, H., Woelm, J.-H., Merseburger, A., Nestler, T., Salem, J., Brandt, M. P., ... Loeb, S. (2016). Qualitative Twitter analysis of participants, tweet strategies, and tweet content at a major urologic conference. *Canadian Urological Association Journal, 10*(1–2), 39.
- Börner, K. (2015). *Atlas of Knowledge: Anyone Can Map*. The MIT Press.
- BrumshTEyn, Y. M., & Vas’kovskii, E. Y. (2017). Analysis of the webometric indicators of the main websites that aggregate multithematic scientific information. *Automatic Documentation and Mathematical Linguistics, 51*(6), 250–265.
- Brunson, E. K., & Sobo, E. J. (2017). Framing Childhood Vaccination in the United States: Getting Past Polarization in the Public Discourse. *Human Organization, 76*(1), 38–47.
- Buchel, O., & Sedig, K. (2016). From data-centered to activity-centered geospatial visualizations. *Geospatial Research: Concepts, Methodologies, Tools, and Applications: Concepts, Methodologies, Tools, and Applications, 1*, 246.
- Card, S. K., Mackinlay, J. D., & Shneiderman, B. (1999). Readings in Information Visualization: Using Vision to Think. In *Information Display* (Vol. 1st, p. 686). Retrieved from <http://portal.acm.org/citation.cfm?id=300679>

- Collins, L., & Nerlich, B. (2015). Examining user comments for deliberative democracy: A corpus-driven analysis of the climate change debate online. *Environmental Communication*, 9(2), 189–207.
- Cross, N. (2004). Expertise in design: an overview. *Design Studies*, 25(5), 427–441.
- Dale, R. (2018). Text analytics apis, part 1: The bigger players. *Natural Language Engineering*, 24(2), 317–324.
- Davenport, R. (2017, February 23). Number of measles cases in Nova Scotia rises to 7. *CBC*. Retrieved from <http://www.cbc.ca/news/canada/nova-scotia/measles-cases-halifax-public-health-1.3996118>
- Dervin, B. (1999). Chaos, order and sense-making: A proposed theory for information design. *Information Design*, 35–57.
- Dervin, B., & Frenette, M. (2000). Sense-Making Methodology: Communicating Communicatively. *Public Communication Campaigns*, 69.
- Dolianiti, F. S., Iakovakis, D., Dias, S. B., Hadjileontiadou, S. J., Diniz, J. A., Natsiou, G., Hadjileontiadis, L. J. (2019). Sentiment analysis on educational datasets: a comparative evaluation of commercial tools. *Educational Journal of the University of Patras UNESCO Chair*. <https://doi.org/10.26220/UNE.2987>
- Durbach, N. (2000). ‘They might as well brand us’: working-class resistance to compulsory vaccination in Victorian England. *Social History of Medicine*, 13(1), 45–63.
- Ericsson, K. A., & Hastie, R. (1994). Contemporary approaches to the study of thinking and problem solving. In *Thinking and problem solving* (pp. 37–79). Elsevier.
- Fox, S., & Rainie, L. (2000). The online health care revolution. Pew Internet & American life project. URL (Consulted November 2018): <https://www.pewinternet.org/2000/11/26/the-online-health-care-revolution/>.
- Funke, J. (2010). Complex problem solving: A case for complex cognition? *Cognitive Processing*, 11(2), 133–142.
- Getman, R., Helmi, M., Roberts, H., Yansane, A., Cutler, D., & Seymour, B. (2018). Vaccine hesitancy and online information: The influence of digital networks. *Health Education & Behavior*, 45(4), 599–606.
- Halavais, A. (2000). National borders on the world wide web. *New Media & Society*, 2(1), 7–28.

- Hasan, K. S., & Ng, V. (2014). Why are you taking this stance? Identifying and classifying reasons in ideological debates. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)* (pp. 751–762).
- Healey, C., & Ramaswamy, S. (2011). Visualizing twitter sentiment. *Sentiment Viz*. Retrieved from [https://www.csc2.ncsu.edu/faculty/healey/tweet\\_viz/tweet\\_app/](https://www.csc2.ncsu.edu/faculty/healey/tweet_viz/tweet_app/)
- Hegarty, M. (2011). The cognitive science of visual-spatial displays: Implications for design. *Topics in Cognitive Science*, 3(3), 446–474.
- Herring, S., Job-Sluder, K., Scheckler, R., & Barab, S. (2002). Searching for safety online: Managing "trolling" in a feminist forum. *The Information Society*, 18(5), 371–384.
- Higgins, J. W., Strange, K., Scarr, J., Pennock, M., Barr, V., Yew, A., ... Terpstra, J. (2011). "It's a feel. That's what a lot of our evidence would consist of": public health practitioners' perspectives on evidence. *Evaluation & the Health Professions*, 34(3), 278–296.
- Hill, R. L. (2017). The political potential of numbers: data visualisation in the abortion debate. *Women, Gender & Research*, 26(1), 83–96.
- Hirschberg, J., & Manning, C. D. (2015). Advances in natural language processing. *Sciencenat*, 349(6245), 261–266. <https://doi.org/10.1126/science.aaa8685>
- Hollan, J., Hutchins, E., & Kirsh, D. (2000). Distributed cognition: toward a new foundation for human-computer interaction research. *ACM Transactions on Computer-Human Interaction (TOCHI)*, 7(2), 174–196.
- Holmberg, K. (2009). *Webometric network analysis: Mapping cooperation and geopolitical connections between local government administration on the web*. Åbo Akademis förlag-Åbo Akademi University Press.
- Holmberg, K., & Thelwall, M. (2009). Local government web sites in Finland: A geographic and webometric analysis. *Scientometrics*, 79(1), 157–169. <https://doi.org/10.1007/s11192-009-0410-6>
- Howarth, C. C., & Sharman, A. G. (2015). Labeling opinions in the climate debate: A critical review. *Wiley Interdisciplinary Reviews: Climate Change*, 6(2), 239–254.
- Hu, M., & Liu, B. (2004). Mining Opinion Features in Customer Reviews. *19th National Conference on Artificial Intelligence*, 755–760. <https://doi.org/10.1145/1014052.1014073>

- Huesch, M. D. (2017). Commercial online social network data and statin side-effect surveillance: a pilot observational study of aggregate mentions on facebook. *Drug Safety*, 40(12), 1199–1204.
- Hutchins, E. (2006). *Cognition in the wild*. Cambridge, MA: MIT Press.
- Jain, A. K., & Gupta, B. B. (2016). Comparative analysis of features based machine learning approaches for phishing detection. In *2016 3rd international conference on computing for sustainable global development (INDIACom)* (pp. 2125–2130). IEEE.
- Janc, K. (2016). A global approach to the spatial diversity and dynamics of internet domains. *Geographical Review*, 106(4), 567–587.
- Jauho, M. (2016). The social construction of competence: Conceptions of science and expertise among proponents of the low-carbohydrate high-fat diet in Finland. *Public Understanding of Science*, 25(3), 332–345.
- Jonassen, D. H. (1995). Computers as cognitive tools: Learning with technology, not from technology. *Journal of Computing in Higher Education*, 6(2), 40–73.
- Kata, A. (2010). A postmodern Pandora's box: Anti-vaccination misinformation on the Internet. *Vaccine*, 28(7), 1709–1716. <https://doi.org/10.1016/j.vaccine.2009.12.022>
- Kata, A. (2012). Anti-vaccine activists, Web 2.0, and the postmodern paradigm - An overview of tactics and tropes used online by the anti-vaccination movement. *Vaccine*, 30(25), 3778–3789. <https://doi.org/10.1016/j.vaccine.2011.11.112>
- Katsuki, T., Mackey, T. K., & Cuomo, R. (2015). Establishing a link between prescription drug abuse and illicit online pharmacies: analysis of Twitter data. *Journal of Medical Internet Research*, 17(12), e280.
- Keel, P. E. (2007). EWall: A visual analytics environment for collaborative sense-making. *Information Visualization*, 6(1), 48–63.
- Kefi, H., & Perez, C. (2018). Dark Side of Online Social Networks: Technical, Managerial, and Behavioral Perspectives. *Encyclopedia of Social Network Analysis and Mining*, 535-556. doi:10.1007/978-1-4939-7131-2\_110217
- Keim, D., Andrienko, G., Fekete, J. D., Görg, C., Kohlhammer, J., & Melançon, G. (2008). Visual analytics: Definition, process, and challenges. *Lecture Notes in Computer Science Information Visualization*, 154-175. doi:10.1007/978-3-540-70956-5\_7

- Keim, D., Kohlhammer, J., Ellis, G., & Mansmann, F. (2010). *Mastering the information age solving problems with visual analytics*. Eurographics Association.
- Kend, M., & Goode, S. (2018). The Effect of Website Age on Reported Cash Flows. <http://hdl.handle.net/1885/149045>
- Kickbusch, I. (2009). Health literacy: engaging in a political debate. *International Journal of Public Health*, 54(3), 131–132.
- Kim, J. H., Barnett, G. A., & Park, H. W. (2010). A hyperlink and issue network analysis of the United States Senate: A rediscovery of the web as a relational and topical medium. *Journal of the Association for Information Science and Technology*, 61(8), 1598–1611.
- Klein, G., Moon, B., & Hoffman, R. R. (2006). Making sense of sensemaking 1: Alternative perspectives. *IEEE Intelligent Systems*, 21(4), 70–73.
- Knauff, M., & Wolf, A. G. (2010). Complex cognition: the science of human reasoning, problem-solving, and decision-making. *Cognitive Processing*, 11(2), 99–102. <https://doi.org/10.1007/s10339-010-0362-z>
- Kruskal, J. B., & Wish, M. (1978). *Multidimensional scaling*. Newbury Park: SAGE.
- Kumar, N., & Srinathan, K. (2008). Automatic keyphrase extraction from scientific documents using N-gram filtration technique. In *Proceedings of the eighth ACM symposium on Document engineering* (pp. 199–208).
- Laszlo, E., & Clark, J. W. (1972). *Introduction to systems philosophy*. Gordon and Breach New York.
- Lee, J., Yoon, W., Kim, S., Kim, D., Kim, S., So, C. H., & Kang, J. (2019). Biobert: pre-trained biomedical language representation model for biomedical text mining. *ArXiv Preprint ArXiv:1901.08746*.
- Lewandowsky, S., & Oberauer, K. (2016). Motivated rejection of science. *Current Directions in Psychological Science*, 25(4), 217–222.
- Leydesdorff, L., & Vaughan, L. (2006). Co-occurrence matrices and their applications in information science: extending ACA to the web environment. *Journal of the American Society for Information Science and Technology*, 57(12), 1616–1628.
- Liu, B. (2015). *Sentiment analysis: Mining opinions, sentiments, and emotions*. Cambridge University Press. <https://doi.org/10.1017/CBO9781139084789>

- Liu, Z., Nersessian, N., & Stasko, J. (2008). Distributed cognition as a theoretical framework for information visualization. *IEEE Transactions on Visualization and Computer Graphics*, 14(6).
- Marshakova, I. V. (1973). Bibliographic coupling system based on references. *Nauchno-Tekhnicheskaya Informatsiya Seriya, Ser, 2*, 3–8.
- Marshall, C. C., & Bly, S. (2005). Saving and using encountered information: implications for electronic periodicals. In *Proceedings of the Sigchi conference on human factors in computing systems* (pp. 111–120). ACM.
- Mazzi, D. (2018). “The diet is not suitable for all...”: On the British and Irish web-based discourse on the Ketogenic Diet. *Lingue Culture Mediazioni-Languages Cultures Mediation (LCM Journal)*, 5(1), 37–56.
- McAuley, J., Leskovec, J., & Jurafsky, D. (2012). Learning attitudes and attributes from multi-aspect reviews. In *Data Mining (ICDM), 2012 IEEE 12th International Conference on* (pp. 1020–1025). IEEE.
- McCoy, C. G., Nelson, M. L., & Weigle, M. C. (2017). University Twitter engagement: using Twitter followers to rank universities. *ArXiv Preprint ArXiv:1708.05790*.
- Meehan, K., Lunney, T., Curran, K., & McCaughey, A. (2013). Context-aware intelligent recommendation system for tourism. In *2013 IEEE International Conference on Pervasive Computing and Communications Workshops, PerCom Workshops 2013* (pp. 328–331). <https://doi.org/10.1109/PerComW.2013.6529508>
- Moreno, R., Ozogul, G., & Reisslein, M. (2011). Teaching with concrete and abstract visual representations: Effects on students’ problem solving, problem representations, and learning perceptions. *Journal of Educational Psychology*, 103(1), 32.
- Morphett, K., Herron, L., & Gartner, C. (2019). Protectors or puritans? Responses to media articles about the health effects of e-cigarettes. *Addiction Research & Theory*, 1–8.
- Munzner, T. (2015). *Visualization analysis & design*. Boca Raton: CRC Press.
- Navar, A. M. (2019). Fear-based medical misinformation and disease prevention: from vaccines to statins. *JAMA Cardiology*, 4(8), 723–724.



- Neumann-Böhme, S., Varghese, N. E., Sabat, I., Barros, P. P., Brouwer, W., Exel, J. V., . . . Stargardt, T. (2020). Once we have it, will we use it? A European survey on willingness to be vaccinated against COVID-19. *The European Journal of Health Economics*, 21(7), 977-982. doi:10.1007/s10198-020-01208-6
- Nicholson, M. S., & Leask, J. (2012). Lessons from an online debate about measles–mumps–rubella (MMR) immunization. *Vaccine*, 30(25), 3806–3812.
- Ninkov, A., & Vaughan, L. (2017). A webometric analysis of the online vaccination debate. *Journal of the Association for Information Science and Technology*, 68(5), 1285–1294. <https://doi.org/10.1002/asi.23758>
- Norman, D. A. (1993). *Things that make us smart: Defending human attributes in the age of the machine*. Cambridge, Mass: Perseus.
- O’carroll, P. W., Cahn, M. A., Auston, I., & Selden, C. R. (1998). Information needs in public health and health policy: results of recent studies. *Journal of Urban Health*, 75(4), 785–793.
- Ola, Oluwakemi, "The Design of Interactive Visualizations and Analytics for Public Health Data" (2017). Electronic Thesis and Dissertation Repository. 4953.
- Oraby, S., Reed, L., Compton, R., Riloff, E., Walker, M., & Whittaker, S. (2017). And that’s a fact: Distinguishing factual and emotional argumentation in online dialogue. *ArXiv Preprint ArXiv:1709.05295*.
- Ortega, J. L., & Aguillo, I. F. (2008). Visualization of the Nordic academic web: Link analysis using social network tools. *Information Processing & Management*, 44(4), 1624–1633.
- Palomino, M., Taylor, T., Göker, A., Isaacs, J., & Warber, S. (2016). The Online Dissemination of Nature–Health Concepts: Lessons from Sentiment Analysis of Social Media Relating to “Nature-Deficit Disorder.” *International Journal of Environmental Research and Public Health*, 13(1), 142.
- Pang, B., & Lee, L. (2008). Opinion Mining and Sentiment Analysis. *Foundations and Trends® in Information Retrieval*, 2(1–2), 1–135. <https://doi.org/10.1561/15000000011>
- Pirolli, P., & Card, S. (2005). The sensemaking process and leverage points for analyst technology as identified through cognitive task analysis. In *Proceedings of international conference on intelligence analysis* (Vol. 5, pp. 2–4). McLean, VA, USA.



- Pohl, M., Smuc, M., & Mayr, E. (2012). The user puzzle—explaining the interaction with visual analytics systems. *IEEE Transactions on Visualization and Computer Graphics*, 18(12), 2908–2916.
- Ptaszynski, M., Masui, F., Rzepka, R., & Araki, K. (2014). Emotive or Non-emotive: That is The Question. *ACL 2014*, 59.
- RealDonaldTrump. (2014). Healthy young child goes to doctor, gets pumped with massive shot of many vaccines, doesn't feel good and changes - AUTISM. Many such cases! Retrieved from <https://twitter.com/realdonaldtrump/status/449525268529815552?lang=en>
- Rizzo, G., & Troncy, R. (2011). Nerd: evaluating named entity recognition tools in the web of data. In *10th International Semantic Web Conference (ISWC'11), Demo Session, Bonn, Germany* (pp. 1–4).
- Romero-Frías, E., & Vaughan, L. (2010). European political trends viewed through patterns of Web linking. *Journal of the Association for Information Science and Technology*, 61(10), 2109–2121.
- Rothenfluh, F., & Schulz, P. J. (2018). Content, quality, and assessment tools of physician-rating websites in 12 countries: quantitative analysis. *Journal of Medical Internet Research*, 20(6), e212.
- Rubin, V. L., Stanton, J. M., & Liddy, E. D. (2004). Discerning emotions in texts. In *The AAAI Symposium on Exploring Attitude and Affect in Text (AAAI-EAAT)*.
- Russell, D. M., Stefik, M. J., Pirolli, P., & Card, S. K. (1993). The cost structure of sensemaking. In *Proceedings of the INTERACT'93 and CHI'93 conference on Human factors in computing systems* (pp. 269–276). ACM.
- Saif, H., He, Y., & Alani, H. (2012). Semantic sentiment analysis of twitter. *The Semantic Web—ISWC 2012*, 508–524.
- Salomon, G. (1993). No distribution without individuals' cognition: A dynamic interactional view. *Distributed Cognitions: Psychological and Educational Considerations*, 111–138.
- Sedig, K., Klawe, M., & Westrom, M. (2001). Role of interface manipulation style and scaffolding on cognition and concept learning in learnware. *ACM Transactions on Computer-Human Interaction (TOCHI)*, 8(1), 34–59.

- Sedig, K., & Parsons, P. (2013). Interaction design for complex cognitive activities with visual representations: A pattern-based approach. *AIS Transactions on Human-Computer Interaction*, 5(2), 84–113.
- Sedig, K., & Parsons, P. (2016). *Design of Visualizations for Human-Information Interaction: A Pattern-Based Framework. Synthesis Lectures on Visualization* (Vol. 4).  
<https://doi.org/10.2200/S00685ED1V01Y201512VIS005>
- Sedig, K., Parsons, P., & Babanski, A. (2012). Towards a Characterization of Interactivity in Visual Analytics. *Journal of Multimedia Processing and Technologies, Special Issue on Theory and Application of Visual Analytics*, 3(1), 12–28.  
<https://doi.org/10.1145/0000000.0000000>
- Seeman, N., & Rizo, C. (2009). Assessing and responding in real time to online anti-vaccine sentiment during a flu pandemic. *Healthcare Quarterly (Toronto, Ont.)*, 13, 8–15.
- Shneiderman, B., Plaisant, C., & Hesse, B. W. (2013). Improving healthcare with interactive visualization. *Computer*, 46(5), 58–66.
- Skyttner, L. (2005). *General systems theory: problems, perspectives, practice*. World scientific. ISBN 978-981-256-389-7.
- Small, H. (1973). Co-citation in the scientific literature: A new measure of the relationship between two documents. *Journal of the Association for Information Science and Technology*, 24(4), 265–269.
- Snyder, J. (2014). Visual representation of information as communicative practice. *Journal of the Association for Information Science and Technology*, 65(11), 2233–2247.
- Sridhar, D., Foulds, J., Huang, B., Getoor, L., & Walker, M. (2015). Joint models of disagreement and stance in online debate. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)* (pp. 116–125).
- Stefanidis, A., Vraga, E., Lamprianidis, G., Radzikowski, J., Delamater, P. L., Jacobsen, K. H., Crooks, A. (2017). Zika in Twitter: temporal variations of locations, actors, and concepts. *JMIR Public Health and Surveillance*, 3(2), e22.
- Stolterman, E. (2008). The nature of design practice and implications for interaction design research. *International Journal of Design*, 2(1).

- Stuart, D. (2014). *Web metrics for library and information professionals*. London. Facet.
- Suen, C. Y. (1979). N-gram statistics for natural language understanding and text processing. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, (2), 164–172.
- Sun, L. (2017, February 20). Trump energizes the anti vaccination movement in Texas. Washington Post. Retrieved from [https://www.washingtonpost.com/national/health-science/trump-energizes-the-anti-vaccine-movement-in-texas/2017/02/20/795bd3ae-ef08-11e6-b4ff-ac2cf509efe5\\_story.html](https://www.washingtonpost.com/national/health-science/trump-energizes-the-anti-vaccine-movement-in-texas/2017/02/20/795bd3ae-ef08-11e6-b4ff-ac2cf509efe5_story.html)
- Thelwall, M. (2004). *Link Analysis: An Information Science Approach*. Emerald Group Publishing Limited. Retrieved from <http://linkanalysis.wlv.ac.uk/index.html>
- Thelwall, M. (2008). Bibliometrics to webometrics. *Journal of Information Science*, 34(4), 605–621. <https://doi.org/10.1177/0165551507087238>
- Thelwall, M. (2009). Introduction to webometrics: Quantitative web research for the social sciences. *Synthesis Lectures on Information Concepts, Retrieval, and Services*, 1(1), 1–116.
- Thelwall, M., Sud, P., & Wilkinson, D. (2012). Link and co-inlink network diagrams with URL citations or title mentions. *Journal of the American Society for Information Science and Technology*, 63(4), 805–816.
- Thelwall, M., Vaughan, L., & Björneborn, L. (2005). Webometrics. *ARIST*, 39(1), 81–135.
- Thelwall, M., & Wilkinson, D. (2004). Finding similar academic Web sites with links, bibliometric couplings and colinks. *Information Processing & Management*, 40(3), 515–526.
- Thomas, J. C., Diament, J., Martino, J., & Bellamy, R. K. E. (2012). Using the “Physics” of notations to analyze a visual representation of business decision modeling. In *2012 IEEE Symposium on Visual Languages and Human-Centric Computing (VL/HCC)* (pp. 41–44). IEEE.
- Thomas, J. J., & Cook, K. a. (2005). Illuminating the path: The research and development agenda for visual analytics. *IEEE Computer Society*, 54(2), 184. <https://doi.org/10.3389/fmicb.2011.00006>
- Tokuhisa, R., Inui, K., & Matsumoto, Y. (2008). Emotion classification using massive examples extracted from the web. In *Proceedings of the 22nd International Conference on Computational Linguistics-Volume 1* (pp. 881–888). Association for Computational Linguistics.

- Triemstra, J. D., Poeppelman, R. S., & Arora, V. M. (2018). Correlations Between Hospitals' Social Media Presence and Reputation Score and Ranking: Cross-Sectional Analysis. *Journal of Medical Internet Research*, 20(11), e289.
- Tsou, M.-H., Yang, J.-A., Lusher, D., Han, S., Spitzberg, B., Gawron, J. M., ... An, L. (2013). Mapping social activities and concepts with social media (Twitter) and web search engines (Yahoo and Bing): a case study in 2012 US Presidential Election. *Cartography and Geographic Information Science*, 40(4), 337–348.
- Ullman, D. G., Dietterich, T. G., & Stauffer, L. A. (1988). A model of the mechanical design process based on empirical data. *Ai Edam*, 2(1), 33–52.
- Van Wijk, J. J. (2005). The value of visualization. In *VIS 05. IEEE Visualization, 2005*. (pp. 79–86). IEEE.
- Vaughan, L., & Ninkov, A. (2018). A new approach to web co-link analysis. *Journal of the Association for Information Science and Technology*, 69(6), 820–831.
- Vaughan, L., & You, J. (2006). Comparing business competition positions based on Web co-link data: The global market vs. the Chinese market. In *Scientometrics* (Vol. 68, pp. 611–628). <https://doi.org/10.1007/s11192-006-0133-x>
- Vaughan, L., & You, J. (2008). Content assisted web co-link analysis for competitive intelligence. *Scientometrics*, 77(3), 433–444. <https://doi.org/10.1007/s11192-007-1999-y>
- Vaughan, L., & You, J. (2010). Word co-occurrences on Webpages as a measure of the relatedness of organizations: A new Webometrics concept. *Journal of Informetrics*, 4(4), 483–491.
- Velardo, S. (2015). The nuances of health literacy, nutrition literacy, and food literacy. *Journal of Nutrition Education and Behavior*, 47(4), 385–389.
- Vergara, S., El-Khouly, M., El Tantawi, M., Marla, S., & Lak, S. (2017). Building Cognitive Applications with IBM Watson Services: Volume 7 Natural Language Understanding. In *Tech. rep.* (p. 98). IBM Corporation.
- Walker, M. A., Tree, J. E. F., Anand, P., Abbott, R., & King, J. (2012). A Corpus for Research on Deliberation and Debate. In *LREC* (Vol. 12, pp. 812-817).
- Waseem, Z., & Hovy, D. (2016). Hateful symbols or hateful people? predictive features for hate speech detection on twitter. In *Proceedings of the NAACL student research workshop* (pp. 88–93).

Yourex-West, H. (2017, June 7). Whooping cough outbreak declared across part of southern Alberta. *Global News*. Retrieved from <https://globalnews.ca/news/3509602/whooping-cough-outbreak-declared-across-part-of-southern-alberta/>

Zheng, H., Aung, H. H., Erdt, M., Peng, T., Sesagiri Raamkumar, A., & Theng, Y. (2019). Social media presence of scholarly journals. *Journal of the Association for Information Science and Technology*, 70(3), 256–270.

## Chapter 3 - VINCENT: A visual analytics system for investigating the online vaccine debate<sup>3</sup>

Anton Ninkov  
Western University  
Faculty of Information and Media Studies

Dr. Kamran Sedig  
Western University  
Faculty of Information and Media Studies & Department of Computer Science

---

<sup>3</sup> A version of this chapter has been published in:

Ninkov, A., & Sedig, K. (2019). VINCENT: A visual analytics system for investigating the online vaccine debate. *Online J. Public Health Inform.* 11(2). doi:10.5210/ojphi.v11i2.10114.

## Abstract

This paper reports and describes VINCENT, a visual analytics system that is designed to help public health stakeholders (i.e., users) make sense of data from websites involved in the online debate about vaccines. VINCENT allows users to explore visualizations of data from a group of 37 vaccine-focused websites. These websites differ in their position on vaccines, topics of focus about vaccines, geographic location, and sentiment towards the efficacy and morality of vaccines, specific and general ones. By integrating webometrics, natural language processing of website text, data visualization, and human-data interaction, VINCENT helps users explore complex data that would be difficult to understand, and, if at all possible, to analyze without the aid of computational tools.

The objectives of this paper are to explore A) the feasibility of developing a visual analytics system that integrates webometrics, natural language processing of website text, data visualization, and human-data interaction in a seamless manner; B) how a visual analytics system can help with the investigation of the online vaccine debate; and C) what needs to be taken into consideration when developing such a system. This paper demonstrates that visual analytics systems can integrate different computational techniques; that such systems can help with the exploration of online public health debates that are distributed across a set of websites; and that care should go into the design of the different components of such systems.

### 3.1. Introduction

As the use of the Internet expands, people engage in social discourse and debate in different areas of interest, generating a great deal of online data. One broad area of interest generating such online information is public health. Public health data is often large, complex, and difficult, if at all possible, to analyze without the aid of computational tools. Public health informatics is a research area that focuses on “the systematic application of information, computer science, and technology to public health practice, research, and learning” (O’Carroll, 2003). Visual analytics systems (VASes) can be of great utility in public health informatics (Ola & Sedig, 2014). VASes are computational tools that combine data visualization, human-data interaction, and data analytics. They allow users to interactively control data visualizations to change how data is

analyzed and presented to them. VASes make it possible for users to quickly make sense of online data that would otherwise be impossible or take more time and effort to accomplish.

In this paper, we report and describe a VAS designed to help public health stakeholders (users) make sense of data from websites involved in the online debate about vaccines. The VAS, VINCENT (VISual aNalytiCs systEm for investigating the online vacciNe debaTe), allows users to explore visualizations of data from a group of 37 vaccine-focused websites (listed in Appendix A). These websites range in their position on vaccines, topics of focus about vaccines, geographic location, and sentiment towards the efficacy and morality of vaccines, specific and general ones. While numerous VASes have been developed and studied previously, VINCENT is novel in that it integrates webometrics (i.e., co-link analysis), natural language processing (i.e., text-based emotion analysis), data visualization, and human-data interaction.

The research questions this paper examines are as follows:

- Is it feasible to integrate webometrics, natural language processing of website text, data visualization, and human-data interaction in a seamless manner to develop a VAS?
- Can such a VAS help with the investigation of the online vaccine debate?
- What are some of the considerations that need to go into developing such a system?

The remainder of this paper is organized as follows. Section 3.2. provides a conceptual and terminological background--i.e. vaccine debate, visual analytics systems, webometrics, and natural language processing. Section 3.3. describes the development of VINCENT and includes an in-depth discussion of the various components of the VAS. Section 3.4. provides a summary and conclusions.

## 3.2. Background

This section provides a conceptual and terminological background for this paper. We will first describe the issue that VINCENT aims to clarify--i.e. the vaccine debate. Next, we will review visual analytics. Finally, we will discuss the data analytics methods (webometrics and natural language processing) that are used in this research.



### 3.2.1. Vaccine Debate

In light of increased recent news coverage of outbreaks of diseases such as measles and whooping cough, the anti-vaccination movement appears to be a new and emerging phenomenon (Abbott, 2019; Oliviero, 2018; Otterman, 2019). The World Health Organization has listed the rise of the anti-vaccination campaign as a top ten health emergency in 2019 (Who.int, 2019). However, anti-vaccination views and sentiments are not a recent development. Since Edward Jenner's discovery of the smallpox vaccine, vaccination has garnered much attention both positive and negative. From the beginning, some have felt that the practice of vaccination is ineffective, violates personal freedoms, and is "unchristian" (Durbach, 2000). However, the Centers for Disease Control reports that vaccines have had a positive impact on global health and are "one of the greatest achievements of biomedical science and public health" (Fox & Rainie, 2000).

Despite the medical community's unified support of immunization, there are many reasons for the persistence of anti-vaccine views. There is some suggestion that increasingly polarized political views (especially in the United States) have generated an environment in which the rejection of scientific facts has become more prevalent and accepted (Lewandowsky & Oberauer, 2016). This erosion of trust in scientific findings among segments of the population may also contribute to this increased polarization. Additionally, the rise in accessibility to, and widespread use of, the Internet has played a role in amplifying the voice of the anti-vaccination movement (Kata, 2010, 2012). (Kata, 2012) states, "The connective power of the Internet brings together those previously considered on the fringe. Members of marginalized groups (e.g. Holocaust deniers, 9/11 'Truthers', AIDS deniers) can easily and uncritically interact with like-minded individuals online ... anti vaccine groups have harnessed postmodern ideologies and by combining them with Web 2.0 and social media, are able to effectively spread their messages." Hence, the Internet plays an important role in the anti-vaccination movement, helping spread their message and promoting their views on vaccination dangers.

The polarity of the vaccine debate is creating a clear divide and this has been revealed through both qualitative classification of inlinks (Ninkov & Vaughan, 2017) and quantitative co-link analysis (Vaughan & Ninkov, 2018). The divide is having harmful effects on the health of the

general population. “Providers and policymakers must begin to recognize the jagged, context-dependent, equifinal nature of how parents sort through vaccination-related information or account for their vaccination decisions in order to reverse declining vaccination rates” (Brunson & Sobo, 2017). Some of the themes of the discussion that have developed in this polarized debate include those related to autism and vaccines, evil government conspiracies, and technological developments (Mitra, Counts, & Pennebaker, 2016). A more automated approach that would allow an analysis of such online discussions and information could help illuminate this public health problem.

### 3.2.2. Visual Analytics Systems (VASes)

In today's environment of big data, people are often victims of information overload. They can get lost in and overwhelmed by the voluminous data and its meaning that they encounter (Keim et al., 2008). By combining human insight with powerful data analytics and integrated data visualizations and human-data interaction, VASes can help alleviate this problem. VASes can enable potential stakeholders to make sense of data. “Just like the microscope, invented many centuries ago, allowed people to view and measure matter like never before, (visual) analytics is the modern equivalent to the microscope” (Börner, 2015).

VASes are composed of three integrated components: an analytics engine, data visualizations, and human-data interactions (Sedig & Parsons, 2016; Sedig, Parsons, & Babanski, 2012). The analytics engine pre-processes and stores data (e.g., data cleaning & fusion), transforms it (e.g., normalization), and analyzes it (e.g., multi-dimensional scaling, emotion analysis) (Han, Pei, & Kamber, 2011). Examples of data analytics techniques that can be integrated into the analytics engine are webometrics and natural language processing (NLP). Data visualizations in a VAS can be visual representations of the information derived from the analytics engine. Visualizations extend the capabilities of individuals to complete tasks by allowing them to analyze data in ways that would be difficult or impossible to do otherwise (Sedig et al., 2012; Shneiderman, Plaisant, & Hesse, 2013). For instance, a scatterplot can be used to visually represent coordinates of entities, and this, in turn, helps the user determine quickly the proximity between data points. Human-data interaction is used in VASes to allow the user to control the data they see and the way the data is processed. Interaction in VASes supports users through distributing the workload

between the user and the system during their exploration and analysis of the data (Z. Liu, Nersessian, & Stasko, 2008; Salomon, 1993; Sedig & Parsons, 2016). Some examples of the numerous human-data interactions that can be incorporated into VASes include filtering, scoping, and drilling of data (Sedig & Parsons, 2013), with each interaction supporting different epistemic actions on information by the user.

One of the theories that can help with the conceptualization of VASes is general systems theory. Systems theory views a system as composed of entities, properties, and relationships (Skyttner, 2005). VASes are complex, multi-level systems, consisting of systems within systems (Sedig & Parsons, 2016). These multi-level systems consist of super-systems, systems made up of other systems, and sub-systems, together making up a super-system (Skyttner, 2005). With this understanding of systems theory, we can see how VASes work. When building and examining VASes, the interactions of the user with the system can have an impact on any of these levels. At the highest level, super-system interactions will change the overall display of the VAS. At lower levels, the interaction sub-system will change specific components of the system. These interactions, regardless of level, are important to the functioning of the VAS and necessary for making sense of the data being presented.

There are several resources available to assist in developing VASes. Two of the most widely used VAS resources include the open source D3.js JavaScript library (Bostock, Ogievetsky, & Heer, 2011) and Tableau software (Nair, Shetty, & Shetty, 2016). The advantage of D3.js is the almost limitless customization capabilities it offers, as it is bound only by programming constraints, and the fact that it is open source. However, the time, effort, and programming skills required by developers to create systems is greater for D3.js than other solutions, as there are fewer templates and starting points to work with. Tableau, on the other hand, is a proprietary data visualization software that provides users with the ability to develop interactive data visualizations with only minimal coding effort. One feature that makes Tableau particularly appealing is that there are several templates available to users to build their own interactive visualizations. As well, Tableau allows users to create dashboards easily, which place multiple interactive visualizations together in one system that automatically connects data together. While both D3.js and Tableau can be useful solutions for developing visual analytics, Tableau has been

used in this research because of its ability to create a functioning and useful visual analytics system while at the same time reducing the programming workload.

VASes incorporate one or more data analysis techniques including (but not limited to) supervised learning (i.e. decision trees or SVM), or cluster analysis (Keim et al., 2008). Previous VAS research has incorporated similar data analysis techniques to those used in VINCENT. For example, researchers have investigated how incorporating multi-dimensional scaling of co-occurrence data (discussed in Section 3.2.3.) in VASes help users investigate entities and identify clusters in a variety of data sets (Cao, Gotz, Sun, & Qu, 2011; Hund et al., 2016). As well, researchers have utilized emotion analysis (discussed in Section 3.2.4.) in VASes that help users investigate online text from both social media and the general web regarding a variety of topics (Beigi, Hu, Maciejewski, & Liu, 2016; Cho, Wesslen, Volkova, Ribarsky, & Dou, 2017; Pathak, Henry, & Volkova, 2017). Both these data analysis techniques have been implemented in VAS research independently of each other, however there have been no published studies examining the integration of the two techniques in a single VAS, as proposed in VINCENT.

### 3.2.3. Webometrics

Webometrics is the “quantitative study of web-related phenomena” (Thelwall, Vaughan, & Björneborn, 2005). With the ever-increasing adoption of the Internet, the various metrics used for analyzing its data, such as hyperlinks, become important to investigate. Two types of webometrics research methods exist: evaluative and relational (Stuart, 2014; Thelwall, 2008).

Evaluative webometrics can include examining webpages for properties such as (but not limited to) the number of external inlinks they receive (links directed to a website from another website) and the website location (Ninkov & Vaughan, 2017; Stuart, 2014; Thelwall, 2004). Examining the number of inlinks a website receives has been shown to be an indicator of performance in a variety of measures for organizations (Thelwall, 2001; Thelwall et al., 2005; Thelwall & Zuccala, 2008; Vaughan & Wu, 2004). Additionally, geographic location has demonstrated to be a valuable resource in conducting evaluative webometrics research (Holmberg & Thelwall, 2009; Ortega & Aguillo, 2008).

Relational webometrics focuses on “providing an overview of the relationships between different actors” (Stuart, 2014). Co-occurrence measurements to indicate similarity are important for relational analysis in webometrics (Stuart, 2014; Thelwall, 2004, 2009). The concept behind this method is that the more entities share occurrences, the more likely they are to be similar in some way (Thelwall, 2008). This method can apply to webometrics in the study of co-links to help analyze similarity in terms of shared online presence between websites (Ortega & Aguillo, 2008; Vaughan & You, 2006, 2008, 2010). To represent and examine co-link data, numerous studies have been conducted with multi-dimensional scaling (MDS)--studies using MDS to analyze business (Vaughan & You, 2006), university (Thelwall & Wilkinson, 2004), government (Holmberg, 2009), and political domains (Kim, Barnett, & Park, 2010; Romero-Frías & Vaughan, 2010). All these studies found that using MDS to analyze co-links generated worthwhile insights into the data.

### **3.2.4. Natural Language Processing (NLP)**

NLP is a vast area of research that focuses on using computational methods to understand and produce human language content (Hirschberg & Manning, 2015). NLP encompasses a wide range of research topics, two of which are text-based emotion detection and word frequency (B. Liu, 2015).

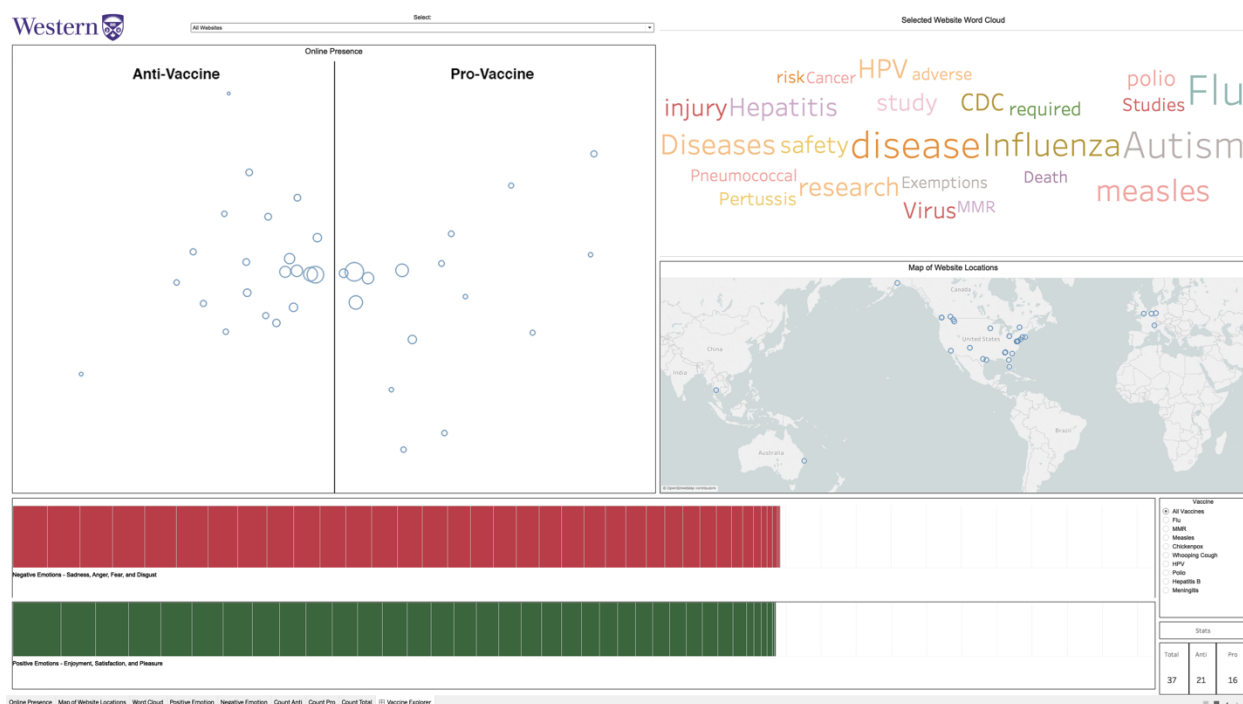
Text-based emotion detection has been examined previously in NLP research (B. Liu, 2015; Ptaszynski, Masui, Rzepka, & Araki, 2014; Rubin, Stanton, & Liddy, 2004; Tokuhisa, Inui, & Matsumoto, 2008). One resource, in particular, that has been developed that makes it possible for researchers to automatically conduct this type of analysis is IBM’s Natural Language Understanding (NLU) API (Vergara, El-Khouly, El Tantawi, Marla, & Lak, 2017). The NLU API (formerly referred to as the AlchemyAPI) has been widely used by many researchers to study topics, sentiments, and emotions in text (Meehan, Lunney, Curran, & McCaughey, 2013; Palomino, Taylor, Göker, Isaacs, & Warber, 2016; Rizzo & Troncy, 2011; Saif, He, & Alani, 2012). The NLU API allows researchers to either input text directly or pull text from URLs of webpages and return a number of different NLP analyses, one of which is emotion analysis. Furthermore, the NLU API can not only detect emotion on the entirety of a text/webpage, but can also return emotion scores for specified target words/phrases (Vergara et al., 2017).

The study of word frequency in text has been examined and used in NLP research (Healey & Ramaswamy, 2011; Katsuki, Mackey, & Cuomo, 2015; McAuley, Leskovec, & Jurafsky, 2012). One of the main concerns for word frequency analysis is how to manage meaningless or unimportant words. In English, like any language, there are many words that are repeated frequently that are not necessarily the key point of interest to a reader. Some of the more obvious examples of these words are “the”, “and”, and “of.” Other types of undesirable words can exist depending on the domain of interest (e.g., dates or numbers). To deal with this issue, the technique of filtering for a list of stop words has been used, and preliminary lists of these words have been created that allow researchers to automatically exclude words that are not of interest (Bird, Klein, & Loper, 2009). To display word frequency data, word clouds have been used successfully (Healey & Ramaswamy, 2011; Katsuki et al., 2015; McAuley et al., 2012). Word clouds display identified words in varying sizes, with larger words being the more frequent. Word clouds are useful because they allow users to quickly see the most prevalent words of a text document and enable them to make quick assessments about what the overall text of a document/website may be discussing.

### 3.3. System Design

The design of VINCENT, displayed in Figure 2, consists of three primary components: the analytics engine, data visualizations, and human-data interactions. In this section, we will discuss these components of the system and explain how the data was collected and managed.

VINCENT was developed in Tableau, version 10.5.

**Figure 2***VINCENT: A Visual Analytic System*

### 3.3.1. Analytics Engine

The analytics engine of VINCENT utilizes webometrics and NLP as its data analysis methods. In this section, we will discuss how, using these methods, data were collected, transformed, and processed. For webometrics, this included leveraging inlink data and geographic location data. For NLP, this included leveraging word frequencies and emotion detection analysis.

The list of 37 vaccine websites (Appendix A) in VINCENT was created based on a list produced for a study on co-link analysis of vaccine websites which included a total of 62 websites (Vaughan & Ninkov, 2018). Websites from that study could be included in VINCENT if they had a central focus on the vaccine debate and a minimum of 200 inlinking domains. The reduction from the original list was primarily due to the elimination of websites that were more minor, websites that had increased their scope beyond just vaccination, and websites that had gone obsolete or merged with another website to form a new website. This list should not be viewed as comprehensive of all vaccine websites, but rather as a sample of some of the more major English-based ones from both sides of the polarized debate.

### 3.3.1.1. Webometrics

Inlink data was collected from each website using MOZ's Link Explorer tool (<https://www.moz.com/link-explorer>). Diverging from some of the previous webometrics research using inlink data, which mostly investigated inlinks coming from pages (Ortega & Aguillo, 2008; Vaughan & You, 2006, 2008, 2010) and sites (Vaughan & Ninkov, 2018), VINCENT uses inlink data about the inlinking domains. Changes in September 2018 to the data provided by MOZ required us to adapt and examine the feasibility of using domain-level inlink data. After comparing domain-level inlink data to data collected for a previous study (Vaughan & Ninkov, 2018), we determined that the domain-level inlink data was a suitable replacement and would be used in the analytics engine of VINCENT.

The shared online presence between the set of websites (Appendix A) was analyzed using MDS. Following similar data analysis techniques to that of previous MDS research (Vaughan & Ninkov, 2018), the inlink data collected on each website was used to create a similarity matrix, which is based on the number of co-links each website shared with one another. Using a computer program originally developed for a previous study (Vaughan & Ninkov, 2018), this co-link data was generated from the collected raw data. Using the output co-link matrix, the data was input into SPSS version 25 and an MDS analysis was conducted. The results of this analysis provided a scatter plot in which each data point was plotted according to the number of co-links they shared, or in other words their shared online presence. Websites that shared more inlinks (and therefore more online presence) were more similar and plotted closer together, while those with fewer inlinks were plotted further away from each other. The goodness of fit between the output scatter plot and the co-link matrix had a stress value of less than 0.05, which suggests a good fit between the two.

Data was also collected regarding the geographic location of the websites. This data was collected through two primary means. The first way of collecting location data was through the sites themselves. Many of the websites had identifying information about the managing owner or organization. The data usually came from an "about us" or "contact us" page and required manual labor to find. For those that did not indicate on their website a location, ICANN WHOIS



registration data was collected. For each of the various collected locations, latitude and longitude coordinates were generated to plot each website on the map of website locations.

### 3.3.1.2. NLP

Word frequency data was collected using the following process. First, each website was analyzed and crawled using InSite5, a software package developed by InSpyder (<https://www.inspyder.com/products/InSite>). With this software, we were able to obtain a CSV export file containing a list of all the words contained on each website, along with the frequency of those word occurrences. After collecting all the raw data about each website, the word frequency lists were filtered to meet the requirements of our analysis. In other words, we wanted only unique words related to the vaccine debate to be displayed. In this effort, we manually created a stop words list to remove irrelevant or common words. The list was built, first, using the Natural Language Toolkit list of stop words for English (Bird et al., 2009). This list of stop words contains some of the most common English words (e.g., “I”, “you”, “too”). From this starting point, the list was expanded to include words that needed to be removed including, but not limited to, letters (e.g., “A”, “B”), dates (e.g., “January”, “Wednesday”), self-reference names (e.g., “NVIC”, “Voices for Vaccines”), people’s names (e.g., “Tom”, “Katie”), Internet words (e.g., “blog”, “post”), and common vaccine debate words (e.g., “vaccines”, “vaccination”). In total, the stop words list, used to refine the word frequency data, consists of 1231 words.

After finalizing each of the website’s individual word frequency list, combined word frequency lists were created for 3 sets of websites: all websites, anti-vaccine websites, and pro-vaccine websites. For each word, the sum of the word frequency was normalized by sum of the total number of words in that set. This generated a proportional count of each word’s presence on the website for each website’s top 25 words. This was a more accurate reflection of the presence of the word on the site rather than simply counting the word frequency totals as some sites had more total words than others. With these proportional word frequencies generated, a list of top 25 words for the 3 sets of websites was also created: all websites, anti-vaccine websites, and pro-vaccine websites.

Text-based emotion detection in the website was conducted with the use of IBM's NLU tool. This tool provides NLP automation through the use of their API and, specifically, can do targeted phrase emotion detection. A user can input text or a URL of a webpage of interest and specify target phrases of interest. The NLU API will return scores for the level of emotion detected for those phrases. Five different emotions (joy, fear, anger, sadness, and disgust) are provided for analysis, which is an overrepresentation of negative emotions (Grimes, 2016). For this system, we did not want to bias our data by over-representing negative emotions. Consequently, the data was cleaned up by merging the 4 negative emotions into one and the labels were changed to reflect a binary of positive emotion (joy) and negative emotion (fear, anger, sadness, and disgust). The vaccines of interest that were examined included: flu, MMR, measles, chicken pox, whooping cough, HPV, polio, hepatitis B, and meningitis. The text was processed using the NLU API's targeted emotion analysis tool. For each of the vaccines, we manually sampled 2 webpages that contained meaningful discussion about the specified vaccine. Several alternate ways of referencing the vaccines were all targeted. For example, with the MMR, targeted phrases included "MMR", "MMR Vaccines", and "MMR Vaccination", among others. The data from each of these different phrases for a vaccine were then merged to reflect the total emotion detected about the specified vaccine.

### 3.3.2. Data Visualizations

VINCENT is comprised of four main visualization components: an online presence map, a word cloud, a map of website locations, and an emotion bar chart. Each of these visualizations represents an important aspect of the websites' information and involves some type of webometrics or NLP data analytics. In this section, each of these visualizations will be discussed, looking at the decisions that were made to represent the data.

#### 3.3.2.1. Online Presence Map

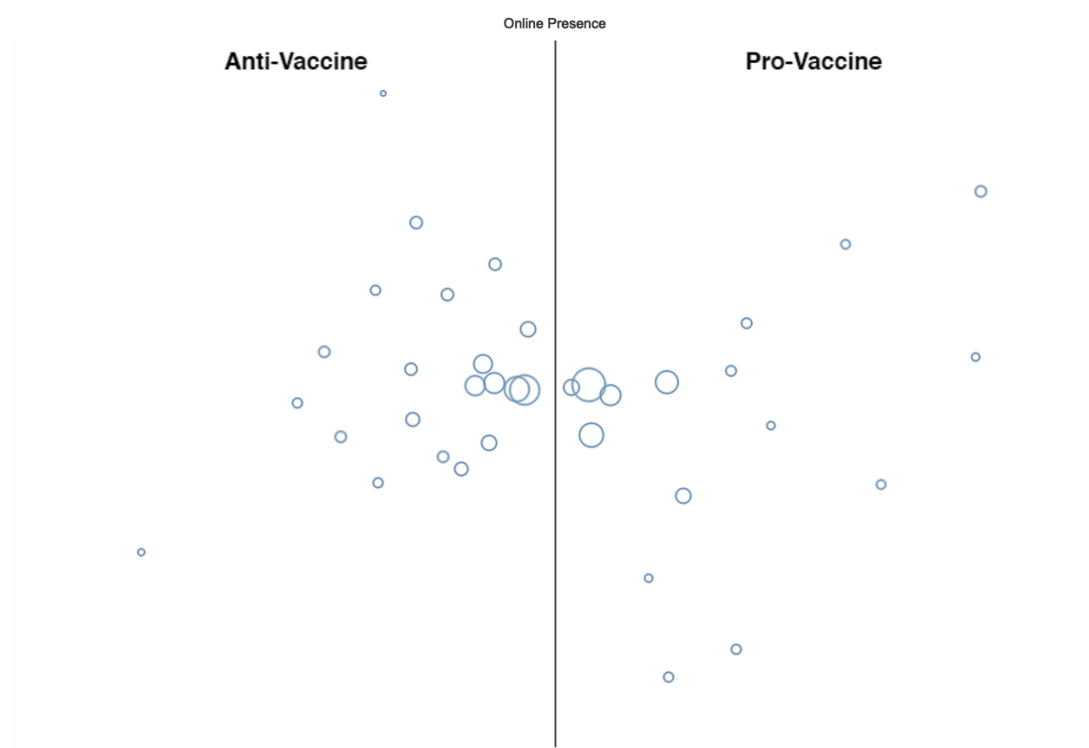
The online presence map, displayed in Figure 3, is a representation of the hyperlink data analyzed from each website. The generated MDS scatter plot map of the websites displays each website in proximity to each other based on their shared online presence. Websites that are plotted closer together share more online presence, while those plotted further away share fewer. Based on this map, polarity between the anti- and pro-vaccine websites was evident, similar to

findings in previous related research (Vaughan & Ninkov, 2018). All anti-vaccine websites ended up on the left side of the map, while all pro-vaccine websites are located on the right side with a space in the middle dividing the two. To display the existence of this polarity, a line dividing the two groups of websites was added to the map with labels for the anti-vaccine and pro-vaccine sides.

Online presence for each website was encoded as a circle representing each of the websites. In this representation, each of the circles was sized based on their total number of inlinking domains. The larger a circle, the more inlinks and, therefore, the larger online presence it has. For reference, the site with the most inlinking domains (9,986) is immunize.org, while the site with the fewest inlinking domains (248) is Vaccine Injury Help Center.

### Figure 3

*MDS Similarity Map*



### 3.3.2.2. Word Cloud

The word cloud, displayed in Figure 4, is a representation of the 25 most common unique words that are related to the vaccine debate from each website or group of websites. Words are sized based on the frequency with which they appeared on the website or group of websites. The bigger a word is on the word cloud, the more frequently it is used on the website, while the smaller a word is, the less frequently it is used. Each word was colored differently to assist with differentiating words from each other.

#### Figure 4

*Word Cloud for all Websites*



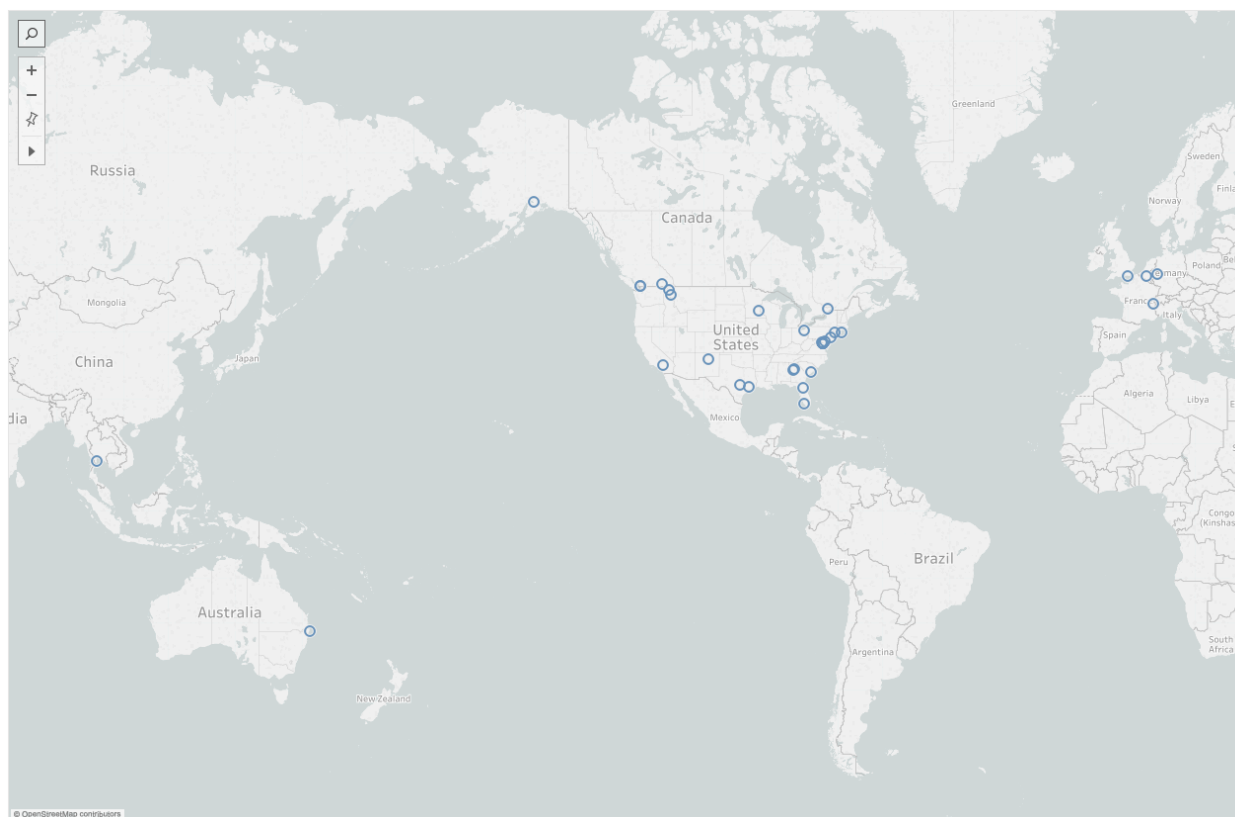
### 3.3.2.3. Map of Website Locations

The map of website locations, displayed in Figure 5, shows a representation of the locations of each website on a world map. Website location is an important piece of data as it allows users of the system to explore the geographic diversity of the websites and identify where clusters of websites may exist. Similar to the online presence map, the website locations use circles to

encode each website. Different from the online presence map, the circles were all sized equally to help the user see the location of each website, and to avoid confusion with excessive overlapping and occlusion of the circles.

**Figure 5**

*Map of Website Locations*



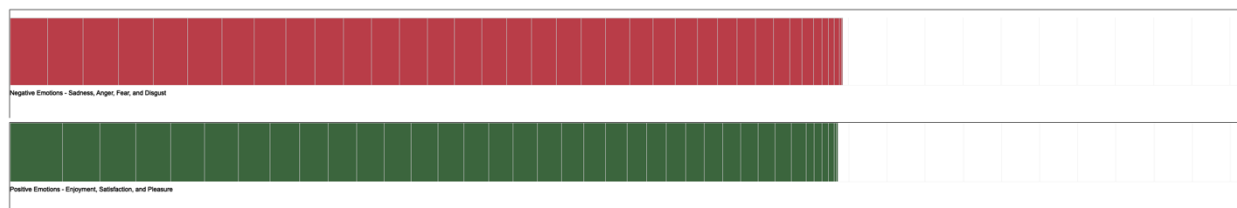
#### 3.3.2.4. Emotion Bar Chart

The emotion bar chart, displayed in Figure 6, represents positive and negative emotions for a selection of each website's text about a set of vaccines. The two bar charts represent the negative (red) and positive (green) emotions detected by the API. Each bar is composed of individual rectangles that refer to individual websites in the set studied. The width of each of these individual rectangles represents the degree of detected emotion on that specific website. The wider the rectangle, the more that emotion is detected. The entire bar is made up of all the smaller rectangles (websites). This bar then represents the overall detected emotion in the text of

the complete website set. The negative and positive bar charts will change in response to the data set that is selected. This will be discussed in more detail in Section 3.3.3.2.

## Figure 6

### *Emotion Bar Chart*



### 3.3.3. Human-Data Interactions

To support users to gain insight into the data and explore the online vaccine debate, many interactions are built into VINCENT. These interactions take place on a global level as well as in the sub-systems of VINCENT. In this section we will explore these interactions and discuss how they will assist users to explore the data.

#### *3.3.3.1. Global System Interactions*

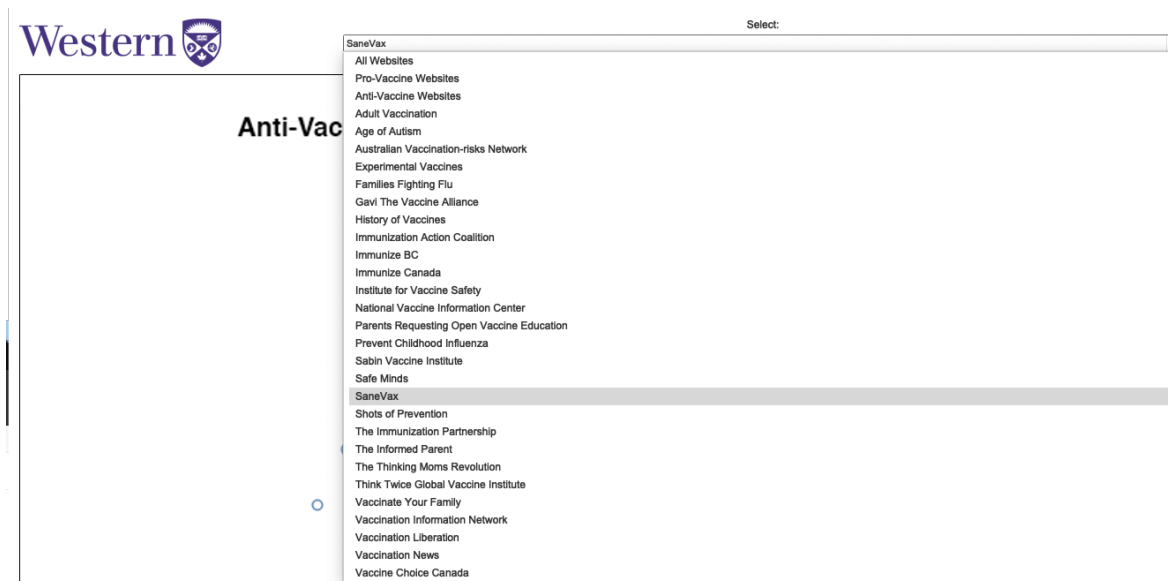
There are several interactions that users can perform on VINCENT that occur at the global system level. These interactions not only affect displayed data at individual, sub-system levels of VINCENT, but also change displayed data at the level of the whole system. Global system interactions in VINCENT include website selection and filtering of websites.

The website selection interaction allows users to focus on a single website. Using this interaction (see Figure 7), users can highlight a single website's data throughout the system in order to determine quickly the website's position on vaccination, online presence, location in the world, and emotion about specific vaccines. Consider the following use case. A user is interested in learning more about the website "SaneVax." They would select this website (Figure 7) from the existing options. VINCENT would then highlight the data points associated with this website, as displayed in Figure 8. For this selected website, the user can immediately find that the website's position is anti-vaccine, that it has strong online presence, that it is located in North Western part

of North America, that it has more negative emotions regarding vaccines than positive, and that it discusses many issues related to HPV (i.e., Cervarix, Gardasil, Cancer, Silgard, HPV).

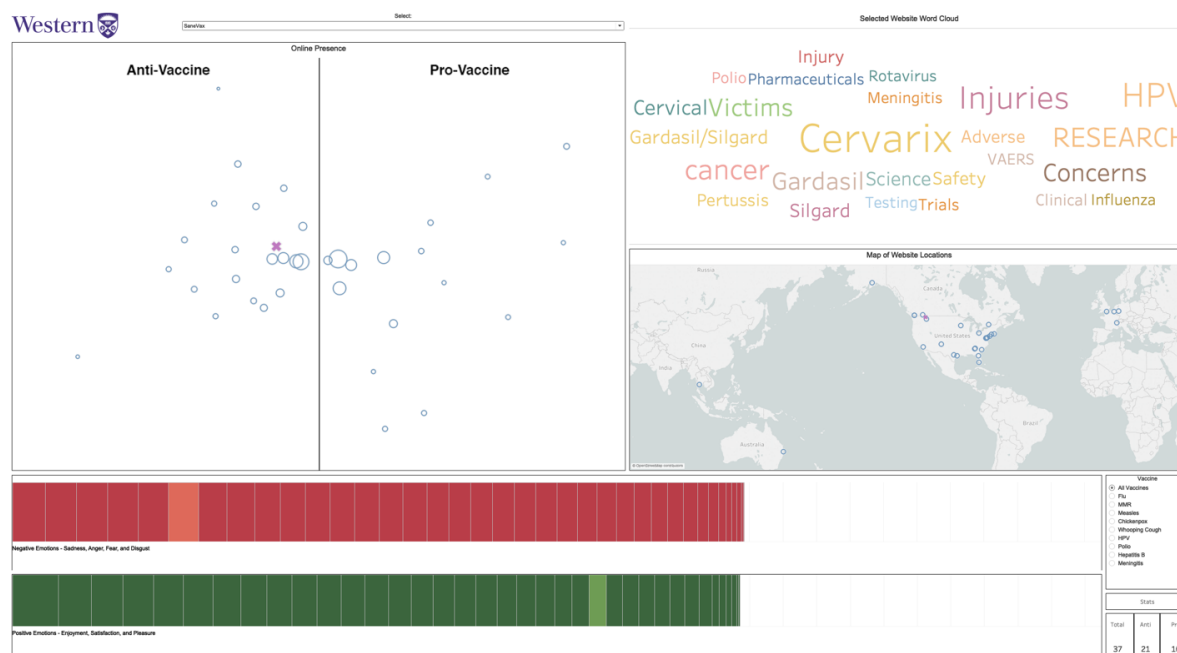
**Figure 7**

*Website Selection Interaction*



**Figure 8**

*VINCENT after Website Selection Interaction*



In addition to the website selection interaction, users have the ability to filter the data to focus on a selected group of websites. Users can highlight and select websites using any of the 3 visualizations, thereby filtering and isolating the data points of a subset of websites. This can be done using the online presence map, map of website locations, or emotion bar chart. Consider a sample use case. A user is curious to learn more about the websites located in North Eastern part of North America. The user goes to the map of website locations and picks websites located in that geographic region. In reaction, the data points on the online presence map and the data of the emotion bar charts are filtered to show only these data points, as displayed in Figure 9. Simultaneously, the stat tracker on the bottom right changes to give the user a numeric count of how many websites they are utilizing now, and how many of each vaccine position is included. The user will quickly see that they have selected 15 websites (10 pro-vaccine and 5 anti-vaccine websites), that the websites are wide ranging in shared online presence, and that they have approximately equal degree of positive and negative emotions associated with the vaccines.

**Figure 9**

*Global Filtered Selection (North Eastern North America)*





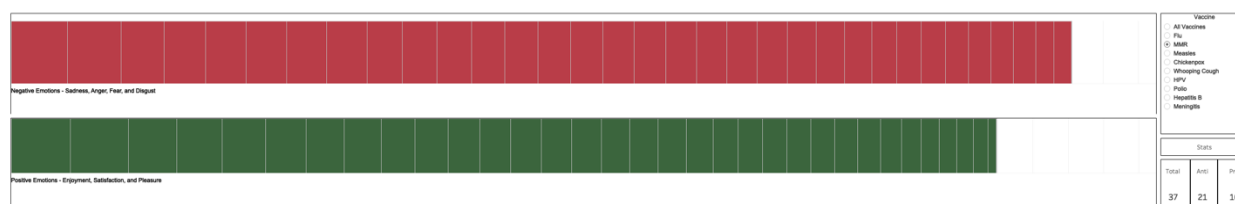
### 3.3.3.2. Sub-System Interactions

There are a number of interactions that can be performed at the sub-systems level of VINCENT. These interactions are focused on isolated elements of the system. They include such interactions as filtering the emotion bar chart to display selected vaccines, hovering display elements to expand an information box, and navigating the map of website locations.

The vaccine selection interaction allows users to filter the displayed data on the emotion bar chart. Upon opening VINCENT, the emotion bar chart displays the overall vaccine emotion data. When a user selects a specific vaccine, the bar chart changes to display only the emotion data that is collected about that specific vaccine. Consider a sample use case. A user is curious about the emotions of the entire set of websites regarding the MMR vaccine. The user would select this vaccine (see right-hand panel in Figure 10), and the bar charts change to display the data. The user can immediately see that there is a greater level of negative emotion on the set of websites than positive emotion regarding the MMR vaccine.

**Figure 10**

*Filtered Vaccine Selection (MMR)*

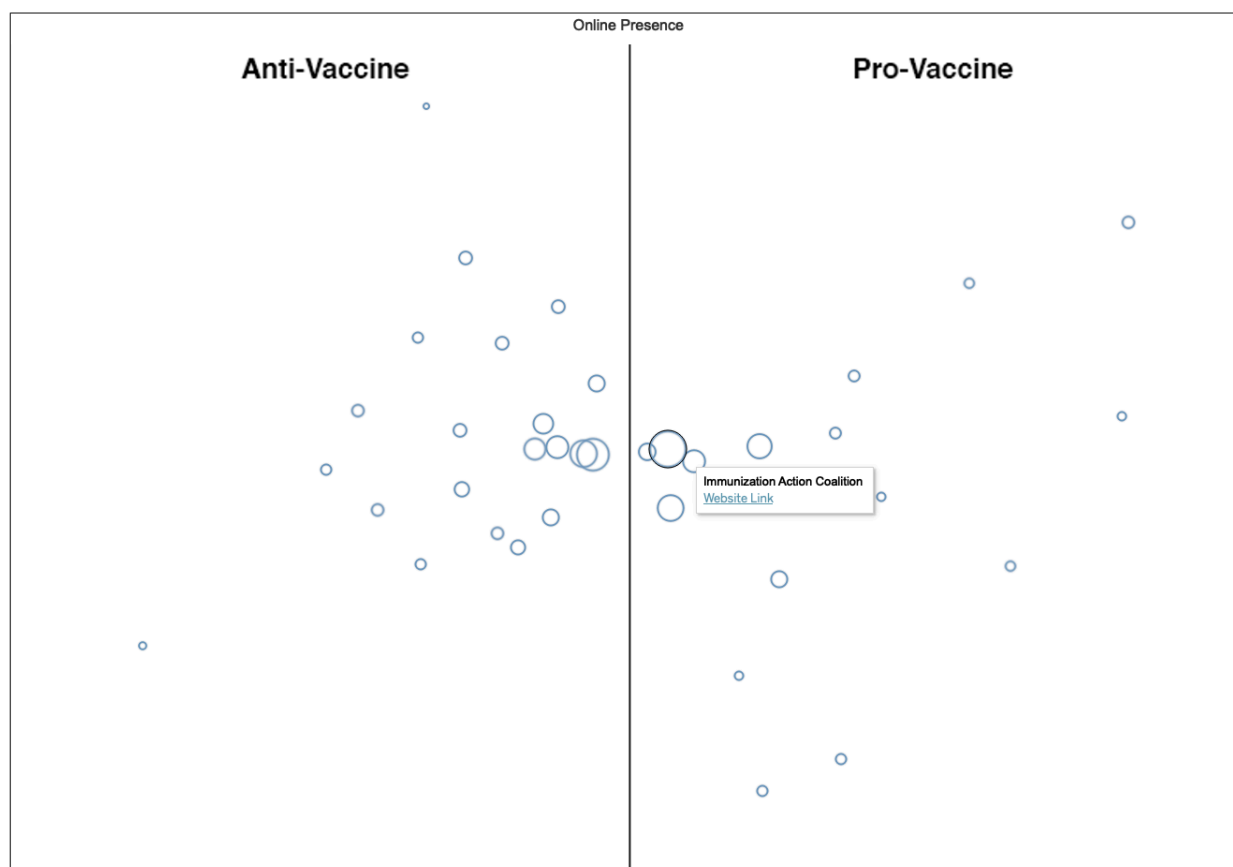


Users also have the option to hover over the online presence map, map of website locations, or emotion bar charts to expand an information box (this is referred to in Tableau as a tooltip) about each specific data point. When a user hovers off the data point, the information box disappears. Again, a sample use case is illustrative of this. A user is interested in identifying which of the pro-vaccine websites have the greatest online presence. To do this, the user would examine the online presence map, determine which token on the pro-vaccine side of the map is the largest, and hover the mouse icon over the token to reveal the information (see Figure 11). In this case, it would be “Immunization Action Coalition.” Similarly, if the user were interested in knowing

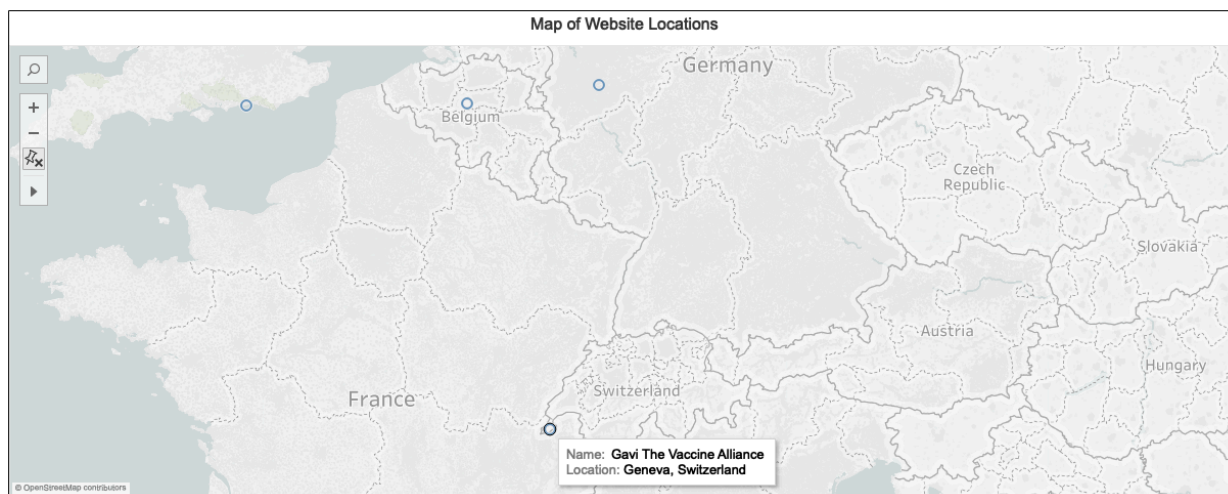
more about a website at a specific location or emotion score, they would hover over those data points to reveal that information.

### Figure 11

*Hover to Expand Information Box*



Finally, on the map of website locations, users have the ability to navigate through the set of websites. On the map of website locations, users can zoom in and out of the map to focus on specific areas. As well, users can click on and drag over the map to move the area of focus. Consider a sample use case. A user is interested in looking at websites in Europe to get a better sense of where exactly they are located. By zooming in on the map and going to Europe (as seen Figure 12), they can clearly identify four websites located there in England, Germany, Belgium, and Switzerland.

**Figure 12***Navigate Map of Website Locations*

### 3.4. Summary and Conclusions

In this paper, we have reported the development of VINCENT, a VAS to help with the investigation of data from websites involved in the online debate on vaccination. VINCENT was created using Tableau, version 10.5. VINCENT incorporates three main sub-systems, each comprised of other sub-systems. An analytics engine made up of webometric (co-link analysis) and NLP (text-based emotion detection) data analysis components; visualization, made up of several different data visualizations; and interaction, made up of a set of different human-data interactions. The development of VINCENT demonstrates that it is feasible to integrate webometrics, natural language processing of website text, data visualization, and human-data interaction into a VAS. VINCENT is novel in its incorporation and integration of the data analysis techniques used (i.e. co-link analysis and text-based emotion analysis) with data visualization and human-data interaction, which had never been previously attempted. VINCENT supports user exploration of data derived from a set of 37 vaccine websites and enables the user to investigate and develop an overall perspective on the vaccine debate. By looking at data from individual websites and groups of websites, a user can identify the breakdown of pro- and anti-vaccine websites, the emotions contained within these websites about specific vaccines, the locations of these websites, and the frequency of vaccine words that appear in these websites. Furthermore, by integrating the data from these different websites,

users can associate the various types of data and uncover patterns that would be otherwise difficult to identify.

Several considerations should go into creating VASes such as VINCENT. First, deciding which tool to use to create the VAS is important. There are advantages and disadvantages to using more programming intensive solutions (such as D3.js) versus more rigid, yet easier to use, toolkit-based solutions (such as Tableau). As well, identifying the appropriate data sources is a challenge that is unique to each project. Online data sources are constantly changing; therefore, it is important for researchers to keep abreast of the current available data. Depending on the resources available to the developer, alternate methods and sources for acquiring proprietary data could improve the value of the system. Next, determining which visualizations are most appropriate for each type of studied dataset is important. For example, the emotion bar charts, presented here, went through several iterations. At first, tree maps were tested but were found to be inadequate at representing certain aspects of the data. Researchers who develop similar VASes need to consider all facets of their data and desired interactions and test various iterations of their system. Finally, incorporating meaningful interactions into the VAS is important. It is necessary to analyze the tasks that users would need to perform, and then determine what combinations of interactions would facilitate the performance of these tasks. In the case of VINCENT, such tasks included comparing websites, identifying groups of websites, and identifying trends in the entire set of websites.

VINCENT was developed to help users make sense of the data from vaccine websites and, ultimately, the online vaccine debate. However, there are many other areas, both within and outside of public health, for which a system such as this could also prove useful. In public health, a similar VAS would be useful for surveillance of other online health debates, such as debates on the efficacy of alternate health claims or debates regarding different medications and drugs. Outside of public health a system similar to VINCENT could prove beneficial in the areas of business, academia, or politics.

One example of such an area that would be well served by a similar VAS is the online discussion about cannabis use. There are diverging positions regarding the risks and benefits of cannabis, and a system similar to VINCENT could enable users to further investigate the debate and make

sense of the data from existing websites. With such a system, users would be able to quickly identify the positions of different websites (i.e., pro- or anti-cannabis, medical or recreational focus on cannabis, and so on), obtain a geographic breakdown of website locations, determine the focus of each website, and identify the detected emotions about various concepts related to cannabis (i.e., “essential oils” or “epilepsy”). Performing tasks such as these could help researchers acquire valuable insight into the online debate on cannabis and determine what (if any) actions could be taken (or policies adopted) to improve public health in this area.

### 3.4.1. Limitations

There were two key limitations to the development of VINCENT. The first set of limitations was related to the data and analysis tools that we used. Social media data could have generated very rich and revealing data for investigation, but these types of data are proprietary and not freely accessible to conduct research of this scale. Moz Link Explorer provided only enough data on inlinks for an adequate co-link analysis at the domain inlink level; getting data for the page- or site-level analysis was not feasible due to the associated cost. As the trend in the area of webometrics is towards collected data becoming increasingly proprietary, researchers need to consider alternative ways of making do with the limited data availability. Additionally, resources like the NLU API are limited in their ability to analyze the websites emotions. Tools like NLU API are essentially only in the infancy of their development. In the future, tools for emotion detection and NLP will certainly improve and be able to achieve a broader range of analysis and better results than are currently possible.

The second set of limitations was related to the interaction capabilities afforded by Tableau as a toolkit. For example, it was not possible in Tableau to allow the filtering interaction to also filter the word cloud selection. Ideally, a user would want to be able to see word clouds of the top 25 words of any subset of websites selected in the other visualizations. However, given the manner by which Tableau allows for the structure of data, and the data management solutions it works with, this was not possible to achieve. The work-around we used for this was to create the website selection interaction that allowed individuals to filter for a specific website throughout VINCENT.

### 3.4.2. Future research

In a follow up paper, we plan to conduct user testing of VINCENT to evaluate whether there is observable benefit to using VINCENT, and, if so, to what extent and in what ways. The findings of this research will lead to the development of best practices for creating similar VASes. They will also help with the identification of potential benefits of VINCENT-like systems that can support exploration of similar public health issues.

### 3.6. References

- Abbott, B. (2019, January 23). Washington State Becomes Latest Hot Spot in Measles Outbreak. *The Wall Street Journal*. Retrieved from <https://www.wsj.com/articles/washington-state-becomes-latest-hot-spot-in-measles-outbreak-11548281172>
- Beigi, G., Hu, X., Maciejewski, R., & Liu, H. (2016). An overview of sentiment analysis in social media and its applications in disaster relief. In *Sentiment analysis and ontology engineering* (pp. 313–340). Springer.
- Bird, S., Klein, E., & Loper, E. (2009). *Natural language processing with Python*. Beijing: O'Reilly.
- Börner, K. (2015). *Atlas of Knowledge: Anyone Can Map*. The MIT Press.
- Bostock, M., Ogievetsky, V., & Heer, J. (2011). D3data-driven documents. *IEEE Transactions on Visualization and Computer Graphics*, *17*(12), 2301–2309. <https://doi.org/10.1109/TVCG.2011.185>
- Brunson, E. K., & Sobo, E. J. (2017). Framing Childhood Vaccination in the United States: Getting Past Polarization in the Public Discourse. *Human Organization*, *76*(1), 38–47.
- Cao, N., Gotz, D., Sun, J., & Qu, H. (2011). DICON: Interactive Visual Analysis of Multidimensional Clusters. *IEEE Transactions on Visualization and Computer Graphics*, *17*(12), 2581–2590. <https://doi.org/10.1109/TVCG.2011.188>
- Cho, I., Wesslen, R., Volkova, S., Ribarsky, W., & Dou, W. (2017). CrystalBall: A Visual Analytic System for Future Event Discovery and Analysis from Social Media Data. In *2017 IEEE Conference on Visual Analytics Science and Technology (VAST)* (pp. 25–35). <https://doi.org/10.1109/VAST.2017.8585658>
- Durbach, N. (2000). ‘They might as well brand us’: working-class resistance to compulsory vaccination in Victorian England. *Social History of Medicine*, *13*(1), 45–63.
- Fox, S., & Rainie, L. (2000). The online health care revolution. Pew Internet & American life project. Retrieved from <https://www.pewinternet.org/2000/11/26/the-online-health-care-revolution/>.

- Grimes, S. (2016). Sentiment, emotion, attitude, and personality, via Natural Language Processing. Retrieved January 20, 2019, from <https://www.ibm.com/blogs/watson/2016/07/sentiment-emotion-attitude-personality-via-natural-language-processing/>
- Han, J., Kamber, M., & Pei, J. (2011). *Data mining: Concepts and techniques*. Amsterdam: Elsevier/Morgan Kaufmann.
- Healey, C., & Ramaswamy, S. (2011). Visualizing twitter sentiment. *Sentiment Viz*. Retrieved from [https://www.csc2.ncsu.edu/faculty/healey/tweet\\_viz/tweet\\_app/](https://www.csc2.ncsu.edu/faculty/healey/tweet_viz/tweet_app/)
- Hirschberg, J., & Manning, C. D. (2015). Advances in natural language processing. *Sciencenat*, 349(6245), 261–266. <https://doi.org/10.1126/science.aaa8685>
- Holmberg, K. (2009). *Webometric network analysis: Mapping cooperation and geopolitical connections between local government administration on the web*. Åbo Akademis förlag-Åbo Akademi University Press.
- Holmberg, K., & Thelwall, M. (2009). Local government web sites in Finland: A geographic and webometric analysis. *Scientometrics*, 79(1), 157–169. <https://doi.org/10.1007/s11192-009-0410-6>
- Hund, M., Böhm, D., Sturm, W., Sedlmair, M., Schreck, T., Ullrich, T., ... Holzinger, A. (2016). Visual analytics for concept exploration in subspaces of patient groups. *Brain Informatics*, 3(4), 233–247. <https://doi.org/10.1007/s40708-016-0043-5>
- Kata, A. (2010). A postmodern Pandora's box: Anti-vaccination misinformation on the Internet. *Vaccine*, 28(7), 1709–1716. <https://doi.org/10.1016/j.vaccine.2009.12.022>
- Kata, A. (2012). Anti-vaccine activists, Web 2.0, and the postmodern paradigm - An overview of tactics and tropes used online by the anti-vaccination movement. *Vaccine*, 30(25), 3778–3789. <https://doi.org/10.1016/j.vaccine.2011.11.112>
- Katsuki, T., Mackey, T. K., & Cuomo, R. (2015). Establishing a link between prescription drug abuse and illicit online pharmacies: analysis of Twitter data. *Journal of Medical Internet Research*, 17(12), e280.



- Keim, D., Andrienko, G., Fekete, J. D., Görg, C., Kohlhammer, J., & Melançon, G. (2008). Visual analytics: Definition, process, and challenges. In *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)* (Vol. 4950 LNCS, pp. 154–175). [https://doi.org/10.1007/978-3-540-70956-5\\_7](https://doi.org/10.1007/978-3-540-70956-5_7)
- Kim, J. H., Barnett, G. A., & Park, H. W. (2010). A hyperlink and issue network analysis of the United States Senate: A rediscovery of the web as a relational and topical medium. *Journal of the Association for Information Science and Technology*, *61*(8), 1598–1611.
- Lewandowsky, S., & Oberauer, K. (2016). Motivated rejection of science. *Current Directions in Psychological Science*, *25*(4), 217–222.
- Liu, B. (2015). *Sentiment analysis: Mining opinions, sentiments, and emotions*. Cambridge University Press.
- Liu, Z., Nersessian, N., & Stasko, J. (2008). Distributed cognition as a theoretical framework for information visualization. *IEEE Transactions on Visualization and Computer Graphics*, *14*(6).
- McAuley, J., Leskovec, J., & Jurafsky, D. (2012). Learning attitudes and attributes from multi-aspect reviews. In *Data Mining (ICDM), 2012 IEEE 12th International Conference on* (pp. 1020–1025). IEEE.
- Meehan, K., Lunney, T., Curran, K., & McCaughey, A. (2013). Context-aware intelligent recommendation system for tourism. In *2013 IEEE International Conference on Pervasive Computing and Communications Workshops, PerCom Workshops 2013* (pp. 328–331). <https://doi.org/10.1109/PerComW.2013.6529508>
- Mitra, T., Counts, S., & Pennebaker, J. W. (2016). Understanding Anti-Vaccination Attitudes in Social Media. In *ICWSM* (pp. 269–278).
- Nair, L., Shetty, S., & Shetty, S. (2016). Interactive visual analytics on Big Data: Tableau vs D3.js. *Journal of E-Learning and Knowledge Society*, *12*(4), 139–150.
- Ninkov, A., & Vaughan, L. (2017). A webometric analysis of the online vaccination debate. *Journal of the Association for Information Science and Technology*, *68*(5), 1285–1294. <https://doi.org/10.1002/asi.23758>

- O’Carroll, P. (2003). Introduction to Public Health Informatics. In P. O’Carroll, W. A. Yasnoff, M. E. Ward, L. H. Ripp, & E. L. Martin (Eds.), *Public Health Informatics and Information Systems* (pp. 3–15). New York, NY, USA: Springer-Verlag.  
<https://doi.org/10.2105/AJPH.34.12.1287-a>
- Ola, O., & Sedig, K. (2014). The challenge of big data in public health: an opportunity for visual analytics. *Online Journal of Public Health Informatics*, 5(3), e223, 1–21.  
<https://doi.org/10.5210/ojphi.v5i3.4933>
- Oliviero, H. (2018, September 4). Whooping cough is making a comeback. Here’s why. *The Toronto Star*. Retrieved from <https://www.thestar.com/life/2018/09/04/whooping-cough-is-making-a-comeback-heres-why.html>
- Ortega, J. L., & Aguillo, I. F. (2008). Visualization of the Nordic academic web: Link analysis using social network tools. *Information Processing & Management*, 44(4), 1624–1633.
- Otterman, S. (2019, January 17). New York Confronts Its Worst Measles Outbreak in Decades. *New York Times*. Retrieved from <https://www.nytimes.com/2019/01/17/nyregion/measles-outbreak-jews-nyc.html>
- Palomino, M., Taylor, T., Göker, A., Isaacs, J., & Warber, S. (2016). The Online Dissemination of Nature–Health Concepts: Lessons from Sentiment Analysis of Social Media Relating to “Nature-Deficit Disorder.” *International Journal of Environmental Research and Public Health*, 13(1), 142.
- Pathak, N., Henry, M. J., & Volkova, S. (2017). Understanding Social Media’s Take on Climate Change through Large-Scale Analysis of Targeted Opinions and Emotions. In *2017 AAAI Spring Symposium Series*.
- Ptaszynski, M., Masui, F., Rzepka, R., & Araki, K. (2014). Emotive or Non-emotive: That is The Question. *ACL 2014*, 59.
- Rizzo, G., & Troncy, R. (2011). Nerd: evaluating named entity recognition tools in the web of data. In *10th International Semantic Web Conference (ISWC’11), Demo Session, Bonn, Germany* (pp. 1–4).

- Romero-Frías, E., & Vaughan, L. (2010). European political trends viewed through patterns of Web linking. *Journal of the Association for Information Science and Technology*, 61(10), 2109–2121.
- Rubin, V. L., Stanton, J. M., & Liddy, E. D. (2004). Discerning emotions in texts. In *The AAAI Symposium on Exploring Attitude and Affect in Text (AAAI-EAAT)*.
- Saif, H., He, Y., & Alani, H. (2012). Semantic sentiment analysis of twitter. *The Semantic Web—ISWC 2012*, 508–524.
- Salomon, G. (1993). No distribution without individuals' cognition: A dynamic interactional view. *Distributed Cognitions: Psychological and Educational Considerations*, 111–138.
- Sedig, K., & Parsons, P. (2013). Interaction design for complex cognitive activities with visual representations: A pattern-based approach. *AIS Transactions on Human-Computer Interaction*, 5(2), 84–113.
- Sedig, K., & Parsons, P. (2016). *Design of Visualizations for Human-Information Interaction: A Pattern-Based Framework. Synthesis Lectures on Visualization* (Vol. 4). <https://doi.org/10.2200/S00685ED1V01Y201512VIS005>
- Sedig, K., Parsons, P., & Babanski, A. (2012). Towards a Characterization of Interactivity in Visual Analytics. *Journal of Multimedia Processing and Technologies, Special Issue on Theory and Application of Visual Analytics*, 3(1), 12–28. <https://doi.org/10.1145/0000000.0000000>
- Shneiderman, B., Plaisant, C., & Hesse, B. W. (2013). Improving healthcare with interactive visualization. *Computer*, 46(5), 58–66.
- Skyttner, L. (2005). *General systems theory: problems, perspectives, practice*. World scientific. ISBN 978-981-256-389-7.
- Stuart, D. (2014). *Web metrics for library and information professionals*. London. Facet.
- Thelwall, M. (2001). Extracting macroscopic information from web links. *Journal of the American Society for Information Science and Technology*, 52(13), 1157–1168.
- Thelwall, M. (2004). *Link Analysis: An Information Science Approach*. Emerald Group Publishing Limited. Retrieved from <http://linkanalysis.wlv.ac.uk/index.html>

- Theilwall, M. (2008). Bibliometrics to webometrics. *Journal of Information Science*, 34(4), 605–621. <https://doi.org/10.1177/0165551507087238>
- Theilwall, M. (2009). Introduction to webometrics: Quantitative web research for the social sciences. *Synthesis Lectures on Information Concepts, Retrieval, and Services*, 1(1), 1–116.
- Theilwall, M., Vaughan, L., & Björneborn, L. (2005). Webometrics. *ARIST*, 39(1), 81–135.
- Theilwall, M., & Wilkinson, D. (2004). Finding similar academic Web sites with links, bibliometric couplings and colinks. *Information Processing & Management*, 40(3), 515–526.
- Theilwall, M., & Zuccala, A. (2008). A university-centred European Union link analysis. *Scientometrics*, 75(3), 407–420.
- Tokuhsa, R., Inui, K., & Matsumoto, Y. (2008). Emotion classification using massive examples extracted from the web. In *Proceedings of the 22nd International Conference on Computational Linguistics-Volume 1* (pp. 881–888). Association for Computational Linguistics.
- Vaughan, L., & Ninkov, A. (2018). A new approach to web co-link analysis. *Journal of the Association for Information Science and Technology*, 69(6), 820–831.
- Vaughan, L., & Wu, G. (2004). Links to commercial websites as a source of business information. *Scientometrics*, 60(3), 487–496.
- Vaughan, L., & You, J. (2006). Comparing business competition positions based on Web co-link data: The global market vs. the Chinese market. In *Scientometrics* (Vol. 68, pp. 611–628). <https://doi.org/10.1007/s11192-006-0133-x>
- Vaughan, L., & You, J. (2008). Content assisted web co-link analysis for competitive intelligence. *Scientometrics*, 77(3), 433–444. <https://doi.org/10.1007/s11192-007-1999-y>
- Vaughan, L., & You, J. (2010). Word co-occurrences on Webpages as a measure of the relatedness of organizations: A new Webometrics concept. *Journal of Informetrics*, 4(4), 483–491.

Vergara, S., El-Khouly, M., El Tantawi, M., Marla, S., & Lak, S. (2017). Building Cognitive Applications with IBM Watson Services: Volume 7 Natural Language Understanding. In *Tech. rep.* (p. 98). IBM Corporation.

Who.int. (2019). Ten health issues WHO will tackle this year. Retrieved February 12, 2019, from <https://www.who.int/emergencies/ten-threats-to-global-health-in-2019>

## Chapter 4 - The Online Vaccine Debate: Study of A Visual Analytics System<sup>4</sup>

Anton Ninkov  
Western University  
Faculty of Information and Media Studies

Dr. Kamran Sedig  
Western University  
Faculty of Information and Media Studies & Department of Computer Science

---

<sup>4</sup> A version of this chapter has been published in:

Ninkov, A., & Sedig, K. (2020). The Online Vaccine Debate: Study of A Visual Analytics System. *Informatics; Multidisciplinary Digital Publishing Institute*. 7(3).

## Abstract

Online debates, specifically the ones about public health issues (e.g., vaccines, medications, and nutrition), occur frequently and intensely, and are having an impact on our world. Many public health topics are debated online, one of which is the efficacy and morality of vaccines. When people examine such online debates, they encounter numerous and conflicting sources of information. This information forms the basis upon which people take a position on such debates. This has profound implications for public health. It necessitates a need for public health stakeholders to be able to examine online debates quickly and effectively. They should be able to easily perform sense-making tasks on the vast amount of online information, such as sentiments, online presence, focus, or geographic locations. In this paper, we report the results of a user study of a visual analytic system (VAS), and whether and how this VAS can help with such sense-making tasks. Specifically, we report a usability evaluation of VINCENT (VIsual aNalytiCs systEm for investigating the online vacciNe debaTe), a VAS previously described. To help the reader, we briefly discuss VINCENT's design in this paper as well. VINCENT integrates webometrics, natural language processing, data visualization, and human-data interaction. In the reported study, we gave users tasks requiring them to make sense of the online vaccine debate. Thirty-four participants were asked to perform these tasks by investigating data from 37 vaccine-focused websites. Half the participants were given access to the system, while the other half were not. Selected study participants from both groups were subsequently asked to be interviewed by the study administrator. Examples of questions and issues discussed with interviewees were: how they went about completing specific tasks, what they meant by some of the feedback they provided, and how they would have performed on the tasks if they had been placed in the other group. Overall, we found that VINCENT was a highly valuable resource for users, helping them make sense of the online vaccine debate much more effectively and faster than those without the system (e.g., users were able to compare websites similarities, identify emotional tone of websites, and locate websites with a specific focus). In this paper, we also identify a few issues that should be taken into consideration when developing VASes for online public health debates.

## 4.1. Introduction

Online debates, specifically the ones about public health issues (e.g., vaccines, medications, and nutrition), occur frequently and intensely, and are having an impact on our world (Kickbusch, 2009; Morphet, Herron, & Gartner, 2019; Velardo, 2015). Many public health topics are debated online, one of which is the efficacy and morality of vaccines (Kata, 2010, 2012). When people examine such online debates, they encounter numerous and conflicting sources of information (Ninkov & Vaughan, 2017). This information forms the basis upon which people take a position on such debates. This has profound implications for public health. It necessitates a need for public health stakeholders to be able to examine online debates quickly and effectively. They should be able to easily perform sense-making tasks on the vast amount of online information, such as sentiments, online presence, focus, or geographic locations. Sense-making is an activity in which a user gradually develops a mental model of an information space (e.g., an online debate) about which they have insufficient knowledge (Klein, Moon, & Hoffman, 2006; Sedig & Parsons, 2013). A sense-making activity is usually comprised of a set of tasks, some of which include: scanning the information space, selecting relevance of items, and examining them in greater detail (Pirulli & Card, 2005). Visual analytic systems (VASes) can help with these sense-making tasks.

VASes are powerful tools that make it possible for users to quickly make sense of complex information presented to them. Made up of three components (data analytics, data visualizations, and human-data interaction), these systems can help users see data in ways that have never before been convenient or, in some cases, possible. For example, VASes enable users to interact with and examine data using methods such as box plots, word clouds, multi-dimensional maps, and stem-leaf plots in a variety of applications such as disaster management (Ragini, Anand, & Bhaskar, 2018), social media reviews (Chang, Ku, & Chen, 2017), or salesforce analysis (Varshney, Rasmussen, Mojsilović, Singh, & DiMicco, 2012).

VINCENT (VISual aNalytiCs systEm for investigating the online vacciNe debaTe) is a VAS designed to help users investigate the online vaccine debate quickly and effectively (Ninkov & Sedig, 2019). It was developed by integrating data analytics (webometrics and natural language processing), data visualizations (scatterplot, bar chart, word cloud, geographic map), and human-



data interaction (filtering, drilling). The system allows users to quickly see and assess websites' online presence, geographic location, focus, and emotion in the text. It is unique in that no other systems have been developed that provides this capability for online public health debates.

In this paper, we report the results of a user study of VINCENT. We gave 34 users tasks that required them to make sense of the online vaccine debate. The users were asked to complete these tasks based on data from 37 vaccine-focused websites (Appendix A). The research questions this study examines is as follows:

- Does VINCENT help users in making sense of the online vaccine debate? Or in other words, do people who use the system:
  - Outperform people without such a system?
  - Find it easier performing analytical and linguistic tasks (e.g., comparing websites similarities, identifying emotional tone of websites, and locating websites with a specific focus) than people without such a system?
  - Have more confidence in their performance—i.e., belief that they found the correct answer—on the tasks than people without such a system?

Additionally, based on user feedback from this study, we will identify a few issues that should be taken into consideration when developing VASes for online public health debates. The remainder of this paper is organized as follows. Section 4.2. discusses the background of online public health debates and VASes. Section 4.3. describes the methodology of this study. Section 4.4. is a summary of the performance results from this study and the response to VINCENT. Finally, Section 4.5. is a discussion of the conclusions, limitations, and future research.

## 4.2. Background

This section discusses the background concepts and terminologies used in this paper. Firstly, a description of online public health debates, more specifically the online vaccine debate, will be provided. This is followed with a discussion about VASes, including why they are important, what their components are, and the means by which they can help users perform tasks.

#### 4.2.1. Online Public Health Debates

Topics related to public health are often discussed and debated online, specifically on the general web and social media. Examples of such topics include vaccination (Kata, 2010, 2012), nutrition (Kovacs, Gillison, & Barnett, 2018; Velardo, 2015), recreational drug use (Bilgrei, 2016), and complementary/alternate medicine use (Zhang et al., 2017). While the implications vary in of these online debates, the underlying methods and mechanisms for transmitting and sharing information have similarities and are inherently interconnected. As previously stated, “The connective power of the Internet brings together those previously considered on the fringe. Members of marginalized groups ... can easily and uncritically interact with like-minded individuals online ... [these] groups have harnessed postmodern ideologies and by combining them with Web 2.0 and social media, are able to effectively spread their messages” (Kata, 2012). As our society has become more connected through the increased use of and access to the Internet, our online public discourse has developed varied ideas about a wide range of health practices, both based in evidence and not.

The online vaccine debate is an example of one of these debates and is an important issue ripe for investigation. As a result of the recent increase in the outbreaks of diseases, such as measles and whooping cough, the anti-vaccination movement is considered by some experts to be an emerging public health problem (Dubé, Vivion, & MacDonald, 2015; Mavragani & Ochoa, 2018). For example, the World Health Organization listed the rise of the anti-vaccination campaign as a top ten health emergency of 2019 (Who.int, 2019). There are many reasons for the persistence of anti-vaccine views, despite the medical community’s unified support of immunization. Increasingly polarized political views and an erosion of trust in scientific findings have produced an environment in which the rejection of scientific conclusions has become more prevalent and accepted among segments of the population (Lewandowsky & Oberauer, 2016). As well, the rise in accessibility to, and widespread use of, the Internet has played a role in amplifying the voice of the anti-vaccination movement (Kata, 2010, 2012). Additionally, as communication technologies have evolved, the public’s attention cycles have become more rapid and driven by the increased information flows (Lorenz-Spreen, Mønsted, Hövel, & Lehmann, 2019). All these factors impact the online vaccine debate as they together lead to people to a less in-depth understanding of these important issues.

The extreme polarization of the vaccine debate is generating a clear divide between anti-vaccine and pro-vaccine groups, as has been revealed through both qualitative classification of inlinks (Ninkov & Vaughan, 2017) and quantitative co-link analysis (Vaughan & Ninkov, 2018). This divide is having harmful effects on the health of the general population. As has been stated, “Providers and policymakers must begin to recognize the jagged, context-dependent, equifinal nature of how parents sort through vaccination-related information or account for their vaccination decisions in order to reverse declining vaccination rates” (Brunson & Sobo, 2017). Some of the specific themes that have galvanized in this polarized debate include those related to autism and vaccines, government conspiracies, and technological developments (Mitra, Counts, & Pennebaker, 2016).

#### **4.2.2. Visual Analytics Systems (VASes)**

People are often victims of information overload in today’s big data environment. It is easy to get lost in and overwhelmed by the voluminous quantity of data. As a result, people struggle to decipher meaning from this sea of data (Keim et al., 2008). VASes that combine human insight with powerful data analytics, data visualizations, and human-data interaction, can alleviate some of these difficulties. Such systems enable potential stakeholders to make sense of data in new ways. By analogy, “Just like the microscope, invented many centuries ago, allowed people to view and measure matter like never before, (visual) analytics is the modern equivalent to the microscope” (Börner, 2015). VASes allow users to see into the data in ways that have never before been possible.

VASes can help users with a variety of cognitive tasks (Sedig & Parsons, 2016). In particular, these systems can be valuable when performing sense-making tasks (Fekete, Jankun-Kelly, Tory, & Xu, 2019; Hohman, Kahng, Pienta, & Chau, 2018). The primary challenges that go along with sense-making tasks are that the relevant information needed is not always easily accessible, stored in the proper format, or located in the proper locations (Marshall & Bly, 2005). In spite of these challenges, people still need to have the ability to rapidly compare and contrast information (Keel, 2007) for which VASes can be particularly useful (Nguyen et al., 2015). While there has been previous studies as to the utility of these systems in healthcare and public health settings,

such research has been fairly limited up to this point, and further investigation is warranted (Caban & Gotz, 2015; Chen & Ebert, 2019; Rind, Wagner, & Aigner, 2019).

VASes are composed of three integrated components: an analytics engine, data visualizations, and human-data interactions (Sedig & Parsons, 2016; Sedig, Parsons, & Babanski, 2012). The analytics engine pre-processes, stores, transforms, and analyzes the data of interest (Han, Pei, & Kamber, 2011). Examples of data analytics techniques that can be integrated into the analytics engine include webometrics and natural language processing (NLP). Data visualizations in VASes involve the visual representations of the information derived from the analytics engine. Visualizations extend the capabilities of individuals to complete tasks by allowing them to analyze data in ways that would be difficult or impossible to do otherwise (Sedig et al., 2012; Shneiderman, Plaisant, & Hesse, 2013). For instance, a scatterplot can be utilized to visually represent coordinates of entities, which helps the user determine, rapidly, the proximity between data points of interest. Human–data interaction is integrated into VASes to allow the user to control the data they access and the means by which the data is processed. Interaction in VASes supports users through distributing the workload between the user and the system during their exploration and analysis of the data (Liu, Nersessian, & Stasko, 2008; Salomon, 1993; Sedig & Parsons, 2016). Specific examples of the numerous human-data interactions that can be incorporated into VASes include filtering, annotating and drilling of data (Sedig & Parsons, 2013), with each interaction supporting different epistemic actions on information by the user.

### 4.3. Methodology

Our study took place between 25 March and 11 April 2019 at a university in Canada. The study was designed to have two sections: the sense-making session and the interview session. The sense-making session consisted of four parts: demographics questionnaire, familiarization period, goal-directed tasks, and post-tasks questionnaire. Select participants from the sense-making session were asked to take part in the interview session, which occurred 2–7 days after the sense-making session and lasted 30 min.

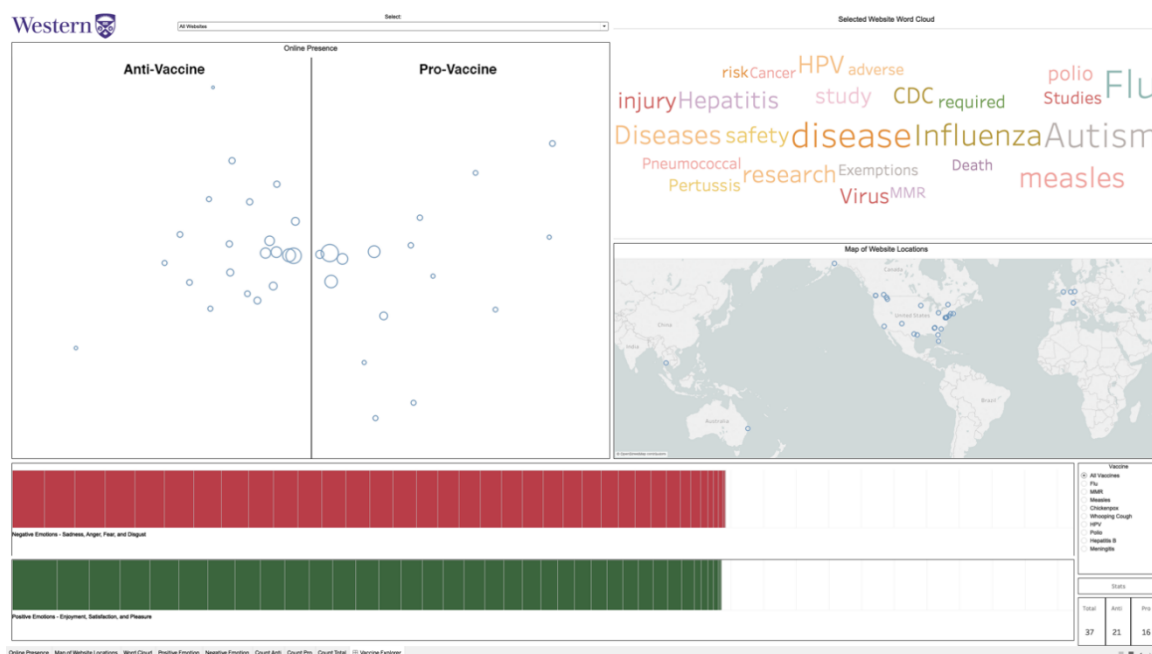
Recruitment took place at the university. In order to be selected, potential participants had to be: at least 18 years of age, currently enrolled as a student at the university, and able to operate a mouse/trackpad and keyboard without any assistance. In total, there were 34 participants for the

sense-making session and 12 were selected for the interview session. Once participants gave formal consent for participating in the study, they were randomly assigned to either the treatment group or the control group. The treatment group was provided VINCENT to complete the tasks. VINCENT incorporated data from a set of 37 vaccine focused websites (discussed in detail in Section 4.3.1.). The very same set of 37 websites were provided to the control group, who only had the use of a web browser and the websites bookmarked in the browser. The list of 37 vaccine websites (Appendix A) was created based on a list produced for a study on the feasibility of co-link analysis for vaccine websites which included a total of 62 websites (Vaughan & Ninkov, 2018). Websites from that previous study could be included in VINCENT if they had a central focus on the vaccine debate and a minimum of 200 inlinking domains. The reduction from the original list was primarily due to the elimination of websites that were more minor, websites that had increased their scope beyond just vaccination, and websites that had gone obsolete or merged with another website to form a new website.

#### **4.3.1. VINCENT: Treatment Instrument**

VINCENT and its components have been discussed in explicit detail in a previous paper (Ninkov & Sedig, 2019). We provide a brief overview of the system to help with understanding this study.

VINCENT (Figure 13) is a VAS that allows users to examine and make sense of data from a set of 37 vaccine-focused websites (listed in Appendix A). These websites range in their positions on vaccines, topics of focus about vaccines, geographic location, and sentiments towards the efficacy and morality of vaccines, both specific ones and vaccines in general. While numerous VASes have been developed and studied previously, VINCENT is novel in that it integrates webometrics (i.e., co-link analysis), NLP (i.e., text-based emotion analysis), data visualization, and human-data interaction (Ninkov & Sedig, 2019). The system is made up of four components: the online presence map, the word cloud, the map of website locations, and the emotion bar charts.

**Figure 13***Screenshot of VINCENT<sup>5</sup>*

Webometrics is the “quantitative study of web-related phenomena” (Thelwall, Vaughan, & Björneborn, 2005). To this end, we have used 2 types of webometrics data: inlinks and geographic locations. Inlinks are hyperlinks directed from an external source to the source of interest (e.g., website A can have an inlink from another website, website B) (Björneborn & Ingwersen, 2004). The inlink data was collected using the MOZ Link Explorer tool (<https://www.moz.com/link-explorer>). Inlink data was used to demonstrate the online presence of the websites. Inlinks were analyzed in two ways: total inlink counts for each website (individual online presence), and a co-link analysis of the shared inlinks between websites (shared online presence). The co-link analysis was conducted using a similar methodology to, and a computer program developed for, a previous study (Vaughan & Ninkov, 2018). Geographic location data was collected using two methods. First, it was done by examining the sites themselves. Many of the websites had personal information that usually came from an “about us” or “contact us” page.

<sup>5</sup> Top left: online presence map. Top right: word cloud. Middle right: map of website locations. Bottom: emotion bar chart

For those that did not indicate on their website a location, WHOIS registration data was collected. For each of the various collected locations, latitude and longitude coordinates were generated to plot each website on the map of website locations.

NLP is a vast area of research that focuses on using computational methods to understand human language content (Hirschberg & Manning, 2015). To this end, we used two types of NLP techniques for website text analysis: term frequency and text-based emotion detection. The word frequency data was collected using InSpyder's InSite 5 (<https://www.inspyder.com/products/InSite>). With this software, we obtained a CSV export file containing a list of all the words on each website, along with the frequency of those word occurrences. We then created a stop-words list that we used to remove erroneous words and only kept un-common words related to the vaccine debate. Text-based emotion analysis was completed using IBM's Natural Language Understanding (NLU) Application Programming Interface (API). With this tool, a user can input text or a URL of a webpage of interest and specify target phrases. The NLU API returns scores for the level of emotion detected for those phrases. The presence of five different emotions (joy, fear, anger, sadness, and disgust) can be analyzed by the tool, which is an overrepresentation of negative emotions (Grimes, 2016). For VINCENT, we did not want to bias our data by over-representing negative emotions. Consequently, the data was cleaned by merging the scores of the 4 negative emotions into one and the labels were changed to reflect a binary of positive emotions (joy) and negative emotions (fear, anger, sadness, and disgust). The vaccines that were examined included: flu, MMR, measles, chickenpox, whooping cough, HPV, polio, hepatitis B, and meningitis.

The interactive data visualizations were developed using the Tableau software (<https://www.tableau.com/>). The online presence map (top left of Figure 13) is a representation of the inlink data analyzed from each website. Online presence is the online attention that a website receives and inlinks can help with its measurement. The scatterplot of the websites was generated using Multi-Dimensional Scaling (MDS) of a co-link analysis of the inlinks. This scatterplot (i.e., online presence map) displays each website in proximity to one another based on their shared online presence. Websites that are plotted closer together share more online presence, while those plotted farther away share less. As well, each website's individual online

presence was encoded using the size of the circle on the map. The larger a circle, the more inlinks and, therefore, the larger online presence it has.

The map of website locations (middle right of Figure 13) displays a representation of the locations of each website on a world map. Similar to the online presence map, the map of website locations uses circles to encode each website. Differing from the online presence map, the circles were all sized equally to help the user see the location of each website, and to avoid confusion with excessive overlapping and occlusion of the circles. For both maps (website locations and online presence), users can select a single website or multiple websites, and it brushes the data throughout the system.

The word cloud (top right of Figure 13) is a representation of the 25 most common, yet unique, words that are related to the vaccine debate from each website or group of websites. Words are sized based on the frequency of their appearance on the website or group of websites. The user can control which word cloud is displayed by using the website selector (top middle of Figure 13).

Finally, the emotion bar chart (bottom of Figure 13) represents the positive and negative emotions found in websites' text towards specific vaccines and vaccines in general. The two bar charts represent the negative (red) and positive (green) emotions detected by the NLU API. Each bar is composed of several rectangles that individually refer to specific websites. The width of each of these individual rectangles represents the degree of detected emotion on that specific website. The wider the rectangle, the more emotional the text is when discussing the selected vaccine. The bar charts change in response to the data that is chosen on the vaccine selector (bottom right of Figure 13).

#### **4.3.2. Sense-Making Session**

People who responded to recruitment were asked to meet for the sense-making session. This session took 45–60 min to complete.



#### 4.3.2.1. Demographics Questionnaire

The first component of the sense-making session was a demographics questionnaire. Each participant responded to questions that asked about their age, education, personal position on vaccines, and familiarity with a variety of topics related to the study (i.e., visual interfaces, the vaccine debate, and vaccine science).

In total, there were 34 participants in the study (17 treatment and 17 control). The mean age of the participants was 24.7 years (standard deviation of 3.9) for the treatment group and 23.1 years (standard deviation of 5.0) for the control group. All participants were university students ranging in their studies from undergraduate to PhD and coming from a variety of different backgrounds and disciplines.

Participants were asked to rank their familiarity with some related topics to the study, ranging from “not familiar” to “very familiar.” With regard to familiarity with visual interfaces, the vaccine debate, and vaccine science, both groups, on average, responded that they were somewhat familiar. Participants were also asked to rate their vaccine stance, between “strongly anti-vaccine” and “strongly pro-vaccine.” The control group had three participants respond “neither pro- or anti-vaccine” while the treatment group had one. The rest of the participants were all either somewhat or strongly pro-vaccine. No participant in this study was anti-vaccine.

#### 4.3.2.2. Familiarization Period

After participants had finished the demographics questionnaire, they were given a 10 min familiarization period before starting the tasks. The familiarization period varied depending on the group to which a participant was assigned. Both groups had the same computer set-up, but different software were available to them.

The control group had ten minutes to familiarize themselves with the websites in the study as well as the layout/functionality of the computer. A web browser window was open and the participants were shown how to access the bookmarks tab consisting of the 37 webpages. They could look through as many webpages as they wanted, were able to open new tabs or windows, or click on hyperlinks to get further information from the webpages.

For the treatment group, the familiarization period was divided into two halves. First, the participants were asked to watch a five-minute video introducing them to VINCENT. This video explained what the system's visualizations represented as well as discussed and demonstrated the various interactions that were built into the system. It did not provide any information about the tasks they were going to perform, ensuring there was no initial advantage for the treatment group over the control group. After watching the introductory video, participants were given another five minutes to use the system freely and familiarize themselves with its functionality.

#### 4.3.2.3. *Tasks*

Following the familiarization period, participants were asked to complete ten tasks. They were given 30 min to complete the tasks (Appendix B), all of which required them to investigate the online vaccine debate as it was presented by the set of 37 websites. These tasks required participants to make sense of various elements of the set of websites, including online presence, shared online presence, geographic location, focus, emotion towards specific vaccines or vaccines in general, and/or a mixture of these. After completing each task, participants were asked to assess how easy the task was, ranging from 1 (very difficult) to 7 (very easy).

All participants were given two pieces of supplementary printed materials: a list of the 37 websites in the study (Appendix A) and a list of definitions for some of the terms that came up in the tasks such as similarity, focus, or online presence (Appendix C). Participants were informed of how much time remained at 15 min and 25 min. Once the 30 min were up, they were allowed to finish responding to a question if it had been started.

#### 4.3.2.4. *Post-Task Questionnaire*

Once participants had finished the tasks, or the thirty minutes were up, the final section of the sense-making session was a post-task questionnaire (Appendices D and E). This questionnaire varied depending on the group in which the participant had been placed. For the control group, participants were asked to assess their confidence in the responses they had given and the easiness of completing all the given tasks (Appendix D). They were also given an opportunity to openly comment on their experiences in the sense-making session. The treatment group was asked the same questions as the control group, but also, was asked about their ability to

understand each of the different data visualizations in VINCENT as well as to connect and control the information from the various the visualizations (Appendix E).

### 4.3.3. Interview Session

The interview sessions were held after the sense-making sessions had been completed. VINCENT was made available for reference during these interviews. The responses from the interviews were analyzed after the study was completed, and the responses were used to help triangulate the results from the sense-making session.

Of the participants that agreed to the interview session during the post-task questionnaire, twelve were selected. These twelve participants were selected in an attempt to reflect the wide range of results that were observed. To this end, three participants were selected from each of the following sub-groups: treatment group users who performed well (participants 6, 9, and 26), treatment group users who performed poorly (participants 5, 20, 32), control group users who performed well (participants 3, 12, 23), and control group users who performed poorly (participants 2, 4, 8).

Participants were asked several questions about their opinions and experiences (Appendix F). They were asked, in a general sense, how they went about completing the tasks, to explain in more detail specific responses they had given, and to compare their experience to what their experience would have been had they been part of the other group. During the interview, the control group was shown the video introducing the system so they could have a basis for some of the comparative questions. The treatment group was not shown the video again, but the system was available to reference during the interview.

## 4.4. Results

This section reports the results of the experimental study. The results will consist of the following three sub-sections: (1) performance results, which includes statistical analysis of the results, comparing how the treatment group performed on the tasks as compared to the control group; (2) response to VINCENT, which includes the participants' feelings towards completing the tasks; (3) usability of VINCENT, which includes the participants response to statements

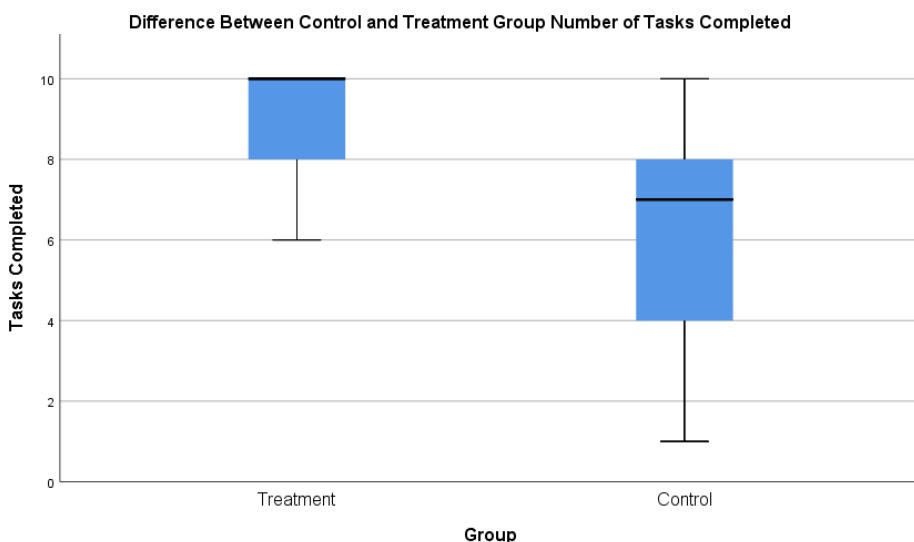
about using the various components of the system. All participants' comments reported in this section are verbatim.

To statistically evaluate the quantitative results, we have used two tests: Mann–Whitney U and Chi-square. Mann–Whitney U tests are used to compare differences between ordinal/continuous variables of two independent groups that have non-normally distributed data. For this study, we used Mann–Whitney U tests to examine if there were significant differences between the two groups with regard to the number of completed tasks, the performance scores on the tasks, and the perceived easiness of and confidence in performing the tasks. Chi-square tests, on the other hand, are used to compare the distribution of nominal variables for independent groups. For this study, we used chi-square to determine if there was a relationship between group and binary task results. It is important to note that not every participant able to complete all 10 tasks. Therefore, we have reported in the results tables the sample size, reflecting how many of the 17 participants in each group completed the task.

#### **4.4.1. Performance Results**

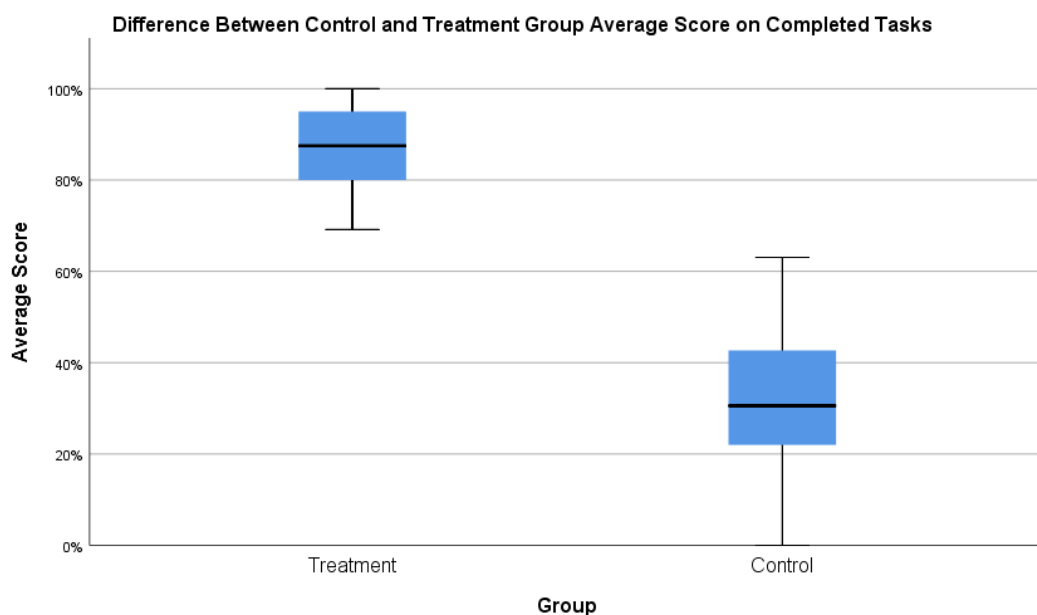
The control group was included to determine by comparison whether VINCENT influenced participants' ability to investigate the online vaccine debate. With regard to completing the tasks, at a descriptive level, the treatment group was able to complete far more than the control group (see Figure 14). The median number of completed tasks for the treatment group was 10/10, while for the control group it was 7/10.

**Figure 14**  
*Boxplot of Tasks Completed*



With regard to the average score (for only the completed tasks), at a descriptive level, the treatment group greatly outperformed the control group (see Figure 15). Every participant in the treatment group outperformed every participant in the control group.

**Figure 15**  
*Boxplot of Average Score*



When comparing the results of the two groups using the Mann–Whitney U test, a significant difference was observed. Table 1 presents the overall statistical analysis of the two groups. Overall, the treatment group was able to complete significantly more tasks and, on the tasks they did complete, were significantly more effective than the control group.

**Table 1**

*Overall Achievement Results*

<b>Group</b>	<b>Median Number of Completed Tasks</b>	<b>Median Percent Score on Tasks Completed</b>
<b>Treatment</b>	10.0	87.5%
<b>Treatment Sample Size</b>	17	17
<b>Control</b>	7.0	30.6%
<b>Control Sample Size</b>	17	17
<b>MWU</b>	72.00	0.00
<b>Significance</b>	$p = 0.010$	$p < 0.000$

In Sections 4.4.1.1. and 4.4.1.2., we discuss the results from the webometrics-based tasks, which helped users assess the websites' online presence (Tables 2 and 3) and geographic locations (Table 4). Tasks that utilized primarily online presence included Tasks 1, 2, and 4, while tasks that utilized primarily geographic locations included Tasks 3 and 10. Then, in Sections 4.4.1.3. and 4.4.1.4., we discuss the results from the NLP-based tasks, which helped users assess websites' focus (Table 5) and emotion in the website text when discussing specific vaccines or vaccines in general (Table 6). Tasks that utilized primarily focus included Tasks 5 and 6 while tasks that utilized primarily emotion included Tasks 7–9.

*4.4.1.1. Online Presence*

In Task 1, we asked participants to identify whether the set had more pro- or anti-vaccine websites, and then specifically how many websites of each vaccine position there were. To complete this task with VINCENT, participants could either highlight all the websites on the online presence map and check for the number of websites that were pro- or anti-vaccine, or otherwise manually count the number of circles on each side of the online presence map. The treatment group was significantly more effective at the task. For Task 1.1, all of the treatment group was able to correctly identify that there were more anti-vaccine websites, while only six

participants in the control group managed to get it correct. The results were similar for Task 1.2, which asked participants to identify the specific number of websites in each vaccine position. In total, 14 of 17 treatment participants were able to correctly identify the exact number of websites for each position, while three of 17 control participants could.

In Task 2, we asked participants to identify both the anti- and pro-vaccine website with the most online presence. To complete this task with VINCENT, users needed to look at the online presence map, identify the biggest circle on the pro- and anti-vaccine side, hover over it and record the website's name. The treatment group was significantly more effective than the control group. Every treatment group participant got both Task 2 sub-tasks correct, while for the control group it was almost the opposite; all but one participant answered all components of the task incorrectly.

In Task 4, participants were asked to give a similarity rating for three pairs of websites. To complete this task with VINCENT participants needed to use multiple different visualizations (online presence map and word cloud). For each pair, they had to see what the two website vaccine positions were, check how far apart they were on the online presence map, and compare their word clouds. The treatment group was significantly more effective at this task than the control group.

**Table 2**

*Online Presence Tasks*

<b>Group</b>	<b>Task 2 Median</b>	<b>Task 4 Median</b>
<b>Treatment Correct</b>	2.0	3.0
<b>Treatment Sample Size</b>	17	17
<b>Control Correct</b>	0.0	1.0
<b>Control Sample Size</b>	17	15
<b>MWU</b>	0.00	47.00
<b>Significance</b>	$p < 0.000$	$p = 0.001$

**Table 3***Online Presence Sub-Tasks*

<b>Group</b>	<b>Task 1.1</b>	<b>Task 1.2</b>	<b>Task 2.1</b>	<b>Task 2.2</b>	<b>Task 4.1</b>	<b>Task 4.2</b>	<b>Task 4.3</b>
<b>Treatment Correct</b>	17	14	17	17	15	14	14
<b>Treatment Sample Size</b>	17	17	17	17	17	17	17
<b>Control Correct</b>	6	3	0	1	7	6	5
<b>Control Sample Size</b>	17	17	17	17	15	15	15
<b>Chi-Square</b>	16.3	14.2	34.0	30.22	6.4	6.1	7.9
<b>Degrees of Freedom</b>	1	1	2	2	1	1	1
<b>Significance</b>	$p < 0.000$	$p < 0.000$	$p < 0.000$	$p < 0.000$	$p = 0.011$	$p = 0.014$	$p = 0.005$

*4.4.1.2. Geographic Locations*

In Task 3, participants searched for the websites that were located outside of North America and needed to identify their country of origin and vaccine position. With VINCENT, users needed to go to the map of website locations, highlight all the websites that were not located in North America, and then record their name, vaccine position, and country of origin. The treatment group was significantly more effective at the task finding the correct websites, locations, and vaccine position with a median accuracy of 100% while the control group had a median accuracy of 50%.

In Task 10, participants had to identify which of the four specified locations presented to them had the highest concentration of each vaccine position. With VINCENT, users needed to highlight the websites in each of the four areas on the map, keep track of how many pro- or anti-vaccine websites there were in each area, and then select the area had the highest concentration of each vaccine position. The treatment group was more effective than the control group, although both groups still fared poorly on the task with the treatment group's accuracy at just over 50% while the control group's accuracy was just under 25%.



**Table 4***Websites' Locations Tasks*

<b>Group</b>	<b>Task 3 Median</b>
<b>Treatment Correct</b>	100%
<b>Treatment Sample Size</b>	17
<b>Control Correct</b>	50%
<b>Control Sample Size</b>	15
<b>MWU</b>	8.00
<b>Significance</b>	$p < 0.000$

*4.4.1.3. Focus*

In Task 5, participants were given four words that were under focus on the vaccine websites. They had to identify if there was a stronger focus on the specified word amongst the pro-vaccine websites or the anti-vaccine websites. With VINCENT, users had to use the website selector to select the anti- and then pro-vaccine word clouds. For each word cloud, they needed to check to see if the given word was there, and if so, record it. The treatment group was significantly more effective at the task than the control group.

In Task 6, participants were given three websites and asked to evaluate the strength of focus on “autism” as strong, weak, or none. With VINCENT, participants needed to use the website selector to select each of the three websites and then scan the word cloud for “autism.” Based on the size of the word (or if it appeared in the word cloud), the user needed to give it a focus rating. The treatment group was more effective than the control as a whole, but only a marginal difference was observed between the results of the two groups. It is worth noting that the treatment group was much more effective at Task 6.3, which presented the participants with a website that did not have any focus on autism.

**Table 5***Word Frequency Tasks*

<b>Group</b>	<b>Task 5 Median</b>	<b>Task 6 Median</b>
<b>Treatment Correct</b>	4.0	2.0
<b>Treatment Sample Size</b>	17	17
<b>Control Correct</b>	2.0	2.0
<b>Control Sample Size</b>	12	12
<b>MWU</b>	32.00	62.00
<b>Significance</b>	$p = 0.001$	$p = 0.057$

*4.4.1.4. Website Text Emotion*

In Task 7, participants had to look at the four websites and determine which had stronger negative than positive emotions associated with the HPV vaccine. For this task, the treatment group needed to use the website selector to choose the website, use the vaccine selector and choose “HPV vaccine”, and then compare the highlighted bars of the emotion bar chart with each other. The treatment group was slightly more effective than the control group on this task, but there was no significant difference observed. The sub-tasks reflected this result, except for Task 7.4 where the treatment group was significantly more effective than the control group. This is an interesting example (vaccines.gov) because it is a pro-vaccine government website which did advocate for and promote the HPV vaccine. However, it also discussed things like the side effects of the vaccine and the people that should not get the vaccine, which is the reason for the higher negative emotion score.

In Task 8, participants had to find the specific website with the strongest positive emotion towards the polio vaccine. With VINCENT, users needed to use the vaccine selector to choose the polio vaccine, and then go to the positive emotion bar on the emotion bar chart, hover over the largest rectangle, and record the name of the website. The treatment group was more effective than the control group on this task. Every treatment group participant responded correctly to this task, while every control group participant responded incorrectly.

In Task 9, participants had to determine which vaccine had the strongest negative emotions associated with it from the anti-vaccine websites. With VINCENT, users needed to use the vaccine selector to choose the identified vaccines and compare the sizes of the negative emotion bars on the emotion bar chart. The treatment group was more effective than the control group on this task.

**Table 6**

*Text-Based Emotion Analysis Tasks*

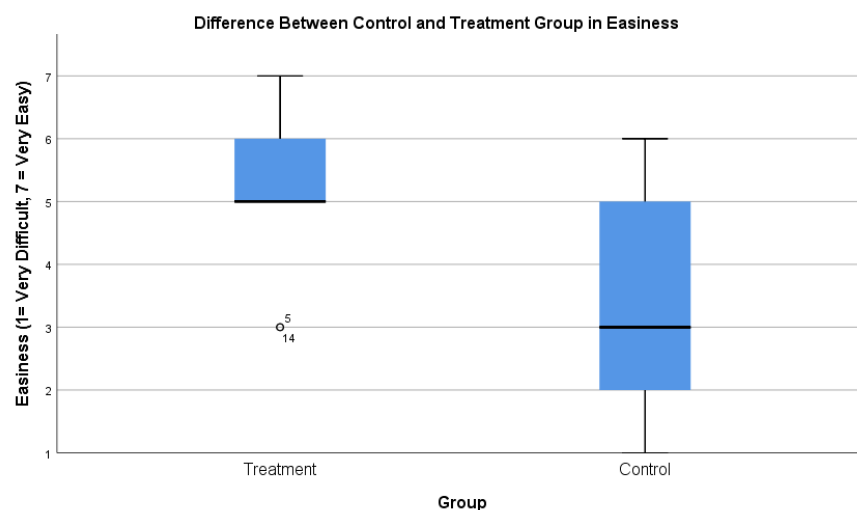
Group	Task 7 Median
Treatment Correct	3.0
Treatment Sample Size	15
Control Correct	2.0
Control Sample Size	11
MWU	62.00
Significance	$p = 0.270$

#### 4.4.2. Response to VINCENT

Overall, the treatment group responded much more positively to the tasks than the control group. At a descriptive level, the treatment group found the tasks much easier to complete (see Figure 16). The median response for how easy they found the tasks was, for the treatment group, easy, while for the control group it was somewhat difficult.

**Figure 16**

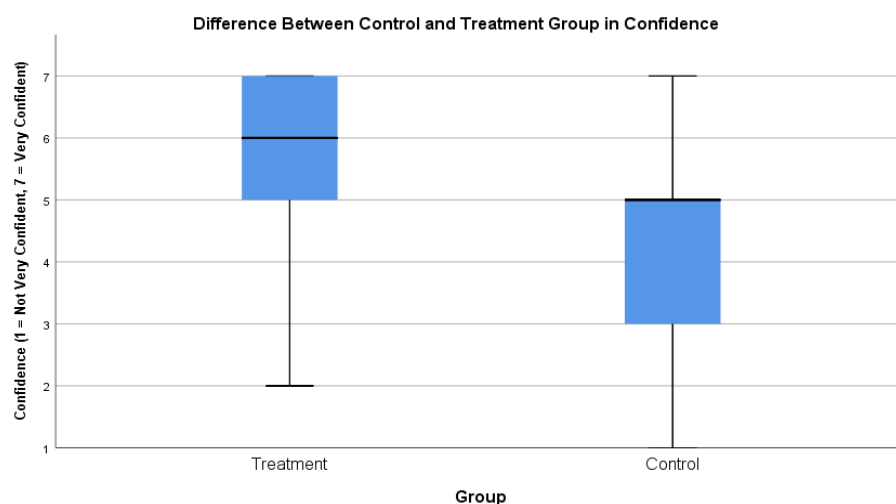
*Box Plot of Easiness*



As well, at a descriptive level, the treatment group was much more confident in their responses (see Figure 17). The majority of responses from the treatment group ranged from extremely confident to somewhat confident, while for the control group, the majority of responses were somewhat confident to somewhat not confident.

**Figure 17**

*Box Plot of Confidence*



Integrating the webometrics and NLP components made sense to the treatment group, as they responded that they found it easy to connect the information across multiple visualizations. This further indicates that the system was usable both overall and not just with regard to the various components of the system, as will be described subsequently.

The treatment group found most of the tasks to be straightforward with the use of VINCENT. They had to identify what the task was asking them to find, search the VAS for the corresponding information, analyze the information (if required), and develop/identify the appropriate response. It was challenging when the system did not match their mental model of how the interaction should work or if the task required them to go beyond simply finding the information in the system and required them to evaluate the information presented to them in further detail.

Comparing the two groups' responses to the tasks, a significant difference was observed (Table 6). The treatment group found the tasks to be much easier to complete and were much more confident in their responses than the control group.

**Table 7**

*Overall Easiness of and Confidence in Response to Tasks*

<b>Group</b>	<b>Median Easiness Completing All Tasks</b>	<b>Median Confidence in Responses to all Tasks</b>
<b>Treatment</b>	5.0	6.0
<b>Treatment Sample Size</b>	17	17
<b>Control</b>	3.0	5.0
<b>Control Sample Size</b>	17	17
<b>MWU</b>	38.50	63.00
<b>Significance</b>	$p < 0.000$	$p = 0.004$

Two main observations were identified from the interview sessions that helped explain, in a general sense, these responses. First, the amount of time they had to complete the tasks was an important factor. VINCENT helped the participants deal with the vast amount of information quickly. Using the system, the treatment group could easily and rapidly find the information they needed to complete the tasks while the control group, on the contrary, found it very time consuming to go through the websites and get the information they needed.

*Participant 8: "I didn't have enough time ... if it was 10 websites, obviously, I would have done way better than 37 websites. I would have been able to look over them all in depth. However, 37 websites were a lot and I had to look through every single one"*

*Participant 12: "I just found it pretty challenging overall just because the amount of websites I was given and the fact that I had to analyze them beyond just the first page or the domain"*

*Participant 2: "The tasks that required me to go through all the websites and determine are they pro-vaccine or anti-vaccine (among other things) ... it was super hard."*

Second, VINCENT made it easier (and in some cases, possible) for participants to analyze and evaluate the information required to complete the tasks. The system offloaded much of the

analysis and allowed the treatment group to visualize the data from the websites in ways that helped them to easily see patterns and make judgments about the data. One treatment group participant highlighted this sentiment.

*Participant 9: “I was confident (in my responses) because ... I would read the question and straight after I looked at the visualization ... (and would) find what I was looking for. The terminology of the question was right there in the visualization. So, it was not like I had to do any further research. (The task) asked something, I clicked it, I hovered, I did some maneuvering around and navigating around the visualization and it was right there. Nothing was hidden, everything was just there.”*

We will discuss the response to the various components of VINCENT in the following subsections. In Section 4.4.2.1. and Section 4.4.2.2., we discuss the user response to the system with regard to the webometrics-based tasks, which helped users assess websites’ online presence (Table 8) and geographic location (Table 9). These included Tasks 1, 2, 3, 4, and 10. In Section 4.4.2.3. and Section 4.4.2.4., we discuss the user response to the system with regard to the NLP-based tasks, which helped users assess websites’ focus (Table 10) and website text emotion when discussing specific vaccines or vaccines in general (Table 11). These included Tasks 5–9.

#### *4.4.2.1. Online Presence*

The treatment group found Task 1 to be significantly easier to complete than the control group. The treatment group quickly understood how to read the online presence map to get the information they needed. The control group did poorly on this task, with the majority responding incorrectly that the set of websites had more pro-vaccine than anti-vaccine websites. To complete this task, the control group had to find ways to investigate the set of websites quickly and effectively. The control group cited several reasons they struggled to do this, including that there were too many websites they had to assess, and they could not find appropriate identifying factors or indicators of what makes a website pro-vaccine or anti-vaccine.

*Participant 8: “I would determine if they were pro- or not pro- by their layout or their about section ... It was kind of hard to keep track of every website within the time frame.”*

*Participant 12: “The first task, I found pretty challenging because there were over 30 websites and the titles themselves, some of them didn’t really give away whether they were pro- or anti-vaccine. So, I had to basically click through all of them and then, even on the cover page, I was sometimes not even sure. Then... so I would have to explore the website and that took a really long time.”*

*Participant 12: “Sometimes it was clear from the outside what the bias was, for instance the title of the webpage often communicated what the stance was, but that can be misleading. The quality of the webpage, a lot of the anti-vaccine websites looked like they were hastily put together whereas the pro-vaccine websites were usually government organizations and often times that was a hint, but ultimately it is the words that count.”*

After seeing VINCENT, the control group discussed how the system would have helped them complete the task and saw how their original perceptions were inaccurate. As well, the treatment group discussed how they felt they would have fared on the task without the system. Some factors they observed that made the task more difficult for the control group included: pre-existing biases, difficulty quickly judging websites, and the juxtaposition of the websites affecting their determination of the website’s stance (i.e., a very anti-vaccine website next to a somewhat anti-vaccine website made the latter appear less anti-vaccine).

*Participant 8: “So, I thought (Australian Vaccination-Risks Network) was pro-vaccination but really it is anti-vaccination, I didn’t get the whole vibe or the whole message of it being anti-vaccination ... I guess I just didn’t go as in depth as other websites ... This system would have helped since it not only marked it as anti-vaccination but I could see (for example) what negative emotions it had”*

*Participant 20: “I don’t think I’d make as objective a decision (without the system) on which websites are anti- or pro- as with (the) system, but also going back to back one website may seem more pro- or anti-vaccine because it was just after another type of website. If I had looked at a really anti-vaccine website and then looked at another anti-*

*vaccine website but it was more mild, I may have personally put it in pro-vaccine category because of my own personal experience.*

*Participant 26: “(Without the tool I would look at) how I trusted the name of the website ... I know what government websites would be called, I know I can trust them in general, and I would put it as a pro-vaccine website ... whereas something like Vaxxter, I’d be instantly questionable and think it’s anti-vaccine ... it’s not a real word, its playing on catchphrasiness and that is a common thing with dubious websites, but then there is one, at the same time, called GAVI vaccine alliance ... which if you were to ask me right of the bat if its pro- or anti, I’d say its anti- ... but ... I actually found out its a pro- one, so (my assumptions) don’t (always) work.”*

The treatment group found Task 2 to be significantly easier than the control group as well. The treatment group was quickly able to understand that the size of the circles reflected online presence of each website and were quickly and easily able to identify the two websites of interest. The control group struggled with the task and found it more difficult to complete. The difficulty was due to the control group not identifying successful ways to judge online presence quickly. They tended to depend on superficial aspects of the website in an attempt to make these determinations, including the look of the website, the content, or the amount of built in interaction.

*Participant 2: “Online presence for me was the quality of the content and representation ... If I’m going through a website that has nothing in it and is just 1 page, that is for me, I don’t think that is going to have much attention or online presence than a website that has a blog and different authors write in it and it is, for example, interactive, you can go and comment and different posts etc. ...”*

*Participant 8: “If I were a mother, I would choose websites that were the most family related. So, in that sense, I would choose those ones as the ones that had the strongest online presence”*

The treatment group found Task 4 to be slightly easier than the control. However, no significant difference was observed between the two groups’ results. This was a task that required going



beyond just finding information on VINCENT, requiring participants to compare and analyze the information. Some participants in the treatment group highlighted the reasons they felt the task was challenging even with VINCENT, ranging from uncertainty on how to read the online presence map as well as being distracted by the amount of information they needed to assess and report.

*Participant 6: “The horizontal axis on the (online presence map) I roughly interpret as the farther left it is the more anti-vaccine it is, the farther right it is the more pro-vaccine it is ... But I don’t know what the vertical axis is telling but it seems like it is a really useful amount of real estate. If it has the opportunity to tell me something, that’d be fantastic ... Because I was putting myself in a mental state of ‘what does the vertical axis mean’, that took a lot of time for me to figure out what I thought was going on”*

*Participant 9: “I wish I could keep both Xs on the map to help with compare and contrast”*

*Participant 5: “I didn’t feel super confident about this because I think I went off in too much detail talking about all the differences and maybe the negative and positive emotions thing tripped me up...”*

*Participant 20: “Because there is so much information here, I wrote -as you can see here- a lot. And I feel like for me it was more difficult because I wanted to write more and there wasn’t enough time to do so.”*

The control group found Task 4 easier to complete than they had found Tasks 1 or 2. An important reason they found this comparatively easier to the previous tasks was that instead of looking through all of the websites, this task only required them to focus on and compare two websites at a time. In the eyes of the control group, this was much more manageable and gave them an opportunity to look more closely at the information they had to assess.

*Participant 3: “To me it was an easy task to complete in terms of the other tasks because I was only comparing 2 websites at a time”*

*Participant 12: “With the specific websites, I was prompted to look at one aspect of them and since it was only about 2 at a time it was easier to remember what I looked at on the first website and then compare that to the next one versus (looking at all the website). I had already forgotten what I looked at 2 websites ago and since I wasn’t looking for anything specific it was sort of overall the feel of the website, that was a lot harder”*

**Table 8**

*Online Presence Tasks Easiness*

<b>Group</b>	<b>Task 1</b>	<b>Task 2</b>	<b>Task 4</b>
<b>Treatment</b>	7.0	7.0	6.0
<b>Treatment Sample Size</b>	17	17	17
<b>Control</b>	4.0	5.0	2.0
<b>Control Sample Size</b>	17	17	15
<b>MWU</b>	40.00	14.00	87.50
<b>Significance</b>	$p < 0.000$	$p < 0.000$	$p = 0.196$

#### 4.4.2.2. Geographic Locations

The treatment group found Task 3 significantly easier to complete than the control group. With VINCENT, participants understood how the map worked and how to locate the websites to uncover the information. Participants reinforced this finding in the interviews, explaining how the map made sense to them and that it was easy for them to see and evaluate the information.

*Participant 26: “There are visual spaces in the app that are definitely more approachable ... For example, the (geographic) map, most people have a mental model of how a map works, so they see the map and they see locations dotted on a map and they can easily approach this and get instant context”*

*Participant 3: “I would have (felt it was) the easiest and had the highest confidence (in my response with the system) because even just seeing here, you can see there are 6 countries and I would immediately be able to get the information I needed.”*

The control group struggled to find the information they needed. A common strategy was to look to see if there was any indication about the geographic location from the name of the website or the top-level domain (e.g., .uk, .au). One participant (participant 2) had a computer science background and mentioned how they used these skills to help do this task, specifically using WHOIS to help locate the websites. But even this method was only somewhat effective, as they were limited in time and could only search the websites that they suspected of being located outside of North America.

*Participant 2: “One task ... (asked me to determine) which country is this website coming from. So, in order to do that, I looked up some of the websites from WHOIS. Some of the websites didn’t include their address or postal code in their about page or any other page, so I had to look up online to see from which country is this website coming from.”*

*Participant 8: “I looked at the URL. I think .com is North America, so there are some that are .eu or .uk so I thought those would be Europe our United Kingdom. So that is what I put as the answer ... This question would have been easier (with the system) with the system because the map shows where the website is located.”*

*Participant 3: “For example, looking at the locations, the best I could do was try and look at the ending of the URL, the domain, and then go to that website and see if I could find out anything about where it was from.”*

The treatment group also found Task 10 to be easier to complete than the control group. In general, the treatment group seemed to understand what they were looking for and how to interact with the system to get that information. One aspect of this task that both the treatment and control groups identified having difficulty understanding was some of the geographic terminology used. “Midwestern USA” specifically seemed to cause confusion amongst participants. Some participants expressed their confusion about what this area meant. Participants said they would have benefited from having geographic regional labels added to the map to help them keep track of and identify the various regions.

*Participant 20: “I would have kind of recognized western North America is here, eastern North America is here, Europe, but Midwestern USA, I don’t know what that means ... If you had those questions and the labels were on the map (it would have helped)”*

*Participant 9: “I felt like the majority of the websites I was looking at fell in both (Midwest and Western USA) so I couldn’t specify which one”*

One treatment group participant highlighted why VINCENT was useful for this task or any other task that required examining groups of websites by their geographic locations. With the system, the user can quickly put websites together based on geography and see if any relationships exist between this and vaccine position, focus, or emotions regarding vaccines.

*Participant 6: “I think of geographical terms ... so I want to know what these clusters of websites have in common because of their geographical proximities and simply by highlighting them there happens to be in the pacific northwest a strong anti-vaccine tendency at least from the sample we have available, which is interesting to me”*

**Table 9**

*Geographic Location Tasks Easiness*

<b>Group</b>	<b>Task 3</b>
<b>Treatment</b>	6.0
<b>Treatment Sample Size</b>	17
<b>Control</b>	2.0
<b>Control Sample Size</b>	15
<b>MWU</b>	36.50
<b>Significance</b>	$p < 0.000$

#### 4.4.2.3. Focus

The treatment group found Task 5 to be slightly easier to complete than the control group, but the difference was not significant between the two groups. For the control group, this task required participants to have completed a general overview of both groups of websites. However, some participants in the control group expressed that they relied primarily on their previous knowledge of the vaccine debate to connect which words they thought were more likely to be pro- or anti-vaccine focused. For example, in Task 5.1 (the only one that did not have a significant difference in results between the two groups), participants had to determine whether

the word “cancer” had a stronger focus in the anti- or pro-vaccine group. One control group participant explained how they applied their knowledge of the vaccine debate to the complete Task 5.

*Participant 12: “This (task required) skimming through some of the main websites, but also a huge part of it was also my previous knowledge on the pro- & anti-vaccine-debate and making assumptions whether these websites would have more focus on these various issues ... I would say this question was a lot of assuming, because I know that mumps, that is something you can prevent with vaccines, so I would assume pro-vaccine, and same with virus. But cancer ... (is something) that people that don’t believe in vaccines (would say) causes, so I would just assume that they would talk about those on anti-vaccine websites and talk about the dangers here”*

The treatment group found Task 6 to be significantly easier to complete than the control group. For the treatment group, the task required them to make judgments on the words in the word clouds. If the word showed up and was large, it was an indication that it had a strong focus. If the word showed up but was small, then it was an indication that it had a weak focus. If the word did not show up at all, it was none. The treatment group expressed that identifying the strong focus or no focus words was easier than identifying weak focus words, as some had difficulty analyzing and evaluating the smaller words in the word cloud. This was reflected in the results for Sub-Task 5.1, which required them to evaluate a smaller word. It was the task with which the treatment group fared the worst.

*Participant 5: “I was able to easily differentiate between the ones that were like used the most but I think I found it somewhat difficult ... that I was not able to tell the difference between the smaller words”*

The control group had to make the assessments by reading through the websites. Again, this was challenging given the amount of content and time that they had. Participants would use strategies such as relying on their previous knowledge, going to about pages or using search features to find the words on the website.

*Participant 3: “One of the hardest questions was about the focus of the website. I found that hard because I could only look at so many words for each website, and typically I would just look at the homepage of those websites and see if any of those words popped out at me.”*

*Participant 12: “I went on each website and I used the search tool and I just searched the word “Autism” and saw how many hits came up ... and just going through the homepage, usually a lot of the anti- ones will have autism on their first page because that is like the main problem people have with vaccines ... but the main strategy for me was using the search button on the websites”*

*Participant 3: “Using this tool would have been much easier because I can see the world autism right there and I would be able to ... choose the websites and see if it comes up. So, it comes up and I was very wrong (about my previous answer), which shows it wasn’t as easy as I thought it was”*

*Participant 23: “Typically the sites would have something on the front page that would give relevant information out—stories you could click on and read more about. If in those it mentioned something like mumps or Gardasil vaccine, then that would be an indication of a strong focus ... when looking at a page I would find what they focus on by looking at what they immediately present.”*

**Table 10**

*Focus Tasks Easiness*

<b>Group</b>	<b>Task 5</b>	<b>Task 6</b>
<b>Treatment</b>	6.0	7.0
<b>Treatment Sample Size</b>	17	17
<b>Control</b>	5.5	5.5
<b>Control Sample Size</b>	12	12
<b>MWU</b>	69.0	29.5
<b>Significance</b>	$p = 0.130$	$p = 0.001$

#### 4.4.2.4. Website Text Emotion

The treatment group found Tasks 7, 8 and 9 to be significantly easier to complete than the control group. They found that it was easy to read the emotion bar charts and make comparisons

between the positive and negative emotions, identify a specific website's emotion towards a vaccine, or assess the emotion of a sub-group of websites. One control group participant highlighted the reason VINCENT made these tasks dealing with emotion easier, especially considering what the task would have been like without the VAS.

*Participant 9: "I don't think it would be easy (to evaluate emotion) at all without the system because nobody writes on a website "I feel strongly negatively about XYZ" it's never there. You really need to read through and, again, what you determine might not be at all what they are trying to say, so I wouldn't be able to confidently answer this. And the visualization is perfect. It tells you right there"*

While the treatment group found the tasks easy to complete, especially compared to the control groups' experience, there was some confusion using VINCENT's emotion bar charts as noted in the excerpts below. These included difficulty with differentiating each section of the bar chart from one another or a need to highlight corresponding negative and positive emotions simultaneously. These difficulties shed light on why the performance scores of the treatment group on Task 7 were only slightly more effective than the control group.

*Participant 5: "I was having difficulties differentiating the smaller ones to the right. Maybe it's the color?"*

*Participant 32 "If I hover on this one, it would be nice if the relative one hovered as well so I can compare very easily"*

**Table 11**

*Website Text Emotion Tasks Easiness*

<b>Group</b>	<b>Task 7</b>	<b>Task 8</b>	<b>Task 9</b>
<b>Treatment</b>	6.0	7.0	6.0
<b>Treatment Sample Size</b>	15	13	12
<b>Control</b>	5.0	3.0	2.0
<b>Control Sample Size</b>	11	6	7
<b>MWU</b>	42.50	16.50	15.50
<b>Significance</b>	$p = 0.031$	$p = 0.007$	$p = 0.011$

### 4.4.3. Usability of VINCENT

In the post-task questionnaire, the 17 treatment group participants were asked to respond to statements directly about the usability of VINCENT, both with regard to specific components and the overall system. The results of the responses to each statement are displayed in Table 11 at the end of this section. In this section, we will discuss the results and some observations for each of the responses to VINCENT's usability.

With regard to the general use of the system, the participants found the VINCENT to be easily usable. As reflected in the responses to Statements 9, 10, and 11, the majority of participants agreed, in some capacity, with this sentiment. As well, when asked about each of the specific components of the system, the majority also felt that the system was usable (each discussed in more detail later in this section). Participants who used VINCENT felt that it was easy to connect the information across the various visualizations together. As well, these participants also said that they found it easy to control the visualizations to see what they wanted to know. Each section of the system worked in cohesion with one another to allow them to investigate each component separately and then put the information together when needed, as is apparent from feedback from Participant 9 and 20.

*Participant 9: "I think it is a very user-friendly tool and one that people will easily be able to extract from the information they need."*

*Participant 20: "Moving around in the map, seeing emotions associated with specific vaccines, having the system link most of what I was trying to do together in all windows (I thought worked well)."*

Additionally, these participants found it helpful to have the various types of information from the websites integrated together into one system that they controlled. This means that these participants found the interactions in VINCENT easy to use, such as filtering groups of websites or individual websites. This is highlighted in feedback from Participants 26 and 27.



*Participant 26: “Associating the various online presences with positions with the map feature was interesting. It allowed me to rapidly understand where the website was located, and if that had anything to do with the content of their position.”*

*Participant 27: “The integration of various types of information into one system I can control.”*

With regard to Statement 3 (ability to clearly understand and evaluate the various website’s online presence with the online presence map), all but one of the participants who used VINCENT agreed. This highlights that encoding the individual online presence of websites into the size of its representative circle on the map worked well. These participants understood this representation and could use the information to complete the tasks. Participant 29 highlights this sentiment in their feedback.

*Participant 29: “The graph with the online presence ... worked quite well and were generally clear and easy to use.”*

Also dealing with the online presence map, the majority of the treatment group said that they also agreed with Statement 4 (ability to clearly understand and evaluate the shared presence of multiple websites on the online presence map). All but three participants responded that they agreed at some level. While the participants seemed to find that display of shared presence as the proximity of the data points usable, it is important to note that some of the feedback, such as from Participants 31 and 32 below, suggested that the functionality for comparing specific websites to one another could have been improved.

*Participant 9: “I think it may be useful to incorporate a check box or something that will help pick 2 or more websites to investigate together - I found myself hovering over the circles to find the websites and then trying to draw a square around them to compare and contrast.”*

*Participant 31: “Comparing the presence and position of 2 selected websites (was confusing).”*

For Statement 5 (ability to clearly understand and evaluate the focus of individual websites and groups of websites using the word cloud), the treatment group participants agreed with the sentiment that they found the word cloud to be usable. It was clear that the drop-down selector was how they controlled what data would show up in this area of the system. However, participants did highlight challenges in the usability because the drop-down selector was disconnected from the filtering of the map, as mentioned by Participant 20.

*Participant 20: "I was confused as to why the word cloud did not reflect my selected website when I clicked through the online presence map or the Location Map. Instead, to see the word cloud for a specific website I had to select that website manually through a drop down menu."*

With regard to Statement 6 (ability to clearly understand and evaluate the geographic dispersion of individual websites and groups of websites using the map of website locations), all but one of the treatment group participants found the map of website locations easily usable. Participants understood how to select websites on the map and zoom in and out to see specific locations. There may have been a bit of a learning curve to developing this understanding of how the zoom feature and search function worked, as was noted by Participant 1.

*Participant 1: "At one point I found the Map of Website Locations a bit tricky to control the zoom and search function, but I got the hang of it after playing around for a minute."*

By majority, the treatment group participants agreed with Statement 7 (ability to understand and evaluate the distribution of emotion that several websites collectively had towards various vaccines) and found that they could use the emotion bar charts to see the collective sentiment towards various vaccines. Compared to the other components of the system, however, these participants had the least strong agreement with this statement. This may be a result of some confusion with the functionality of the emotion bar chart, noted by several participants in the comments. For example, tasks that involved selecting multiple websites on other parts of VINCENT and then comparing the emotions detected were among the most challenging from a usability perspective. One participant noticed that there was a lag in the loading of the data, which may have been a hinderance to its usability for other users as well.

*Participant 15: “I found the positive and negative emotions slightly confusing as sometimes when a site has a very low emotion I could not always tell that it was there when comparing it with another website.”*

*Participant 27: “(The) computer lags when clicking vaccine types by disease”*

Finally, continuing with the emotion bar charts, all but one of the treatment group participants also agreed with statement 8 (ability to understand and evaluate the distribution of emotion that a single website had towards various vaccines). These participants found it straightforward to use the emotion bar charts to examine individual website’s emotions. These participants found it clear that they could select the vaccine they wanted to look at using the vaccine selector, and then click on specific websites or select them from the drop-down menu to examine individual websites. Several participants mentioned the positive experience using this feature in the comments, such as Participant 26.

*Participant 26: “It was also powerful to quickly filter each vaccine's emotions on a website by website basis.”*

**Table 11***Usability of VINCENT*

Statement	Strongly Disagree	Disagree	Somewhat Disagree	Neither Agree or Disagree	Somewhat Agree	Agree	Strongly Agree	Median Response
3. I was able to clearly understand and evaluate the various individual website's online presence using the online presence map.	0 (0%)	1 (5.8%)	0 (0%)	0 (11.8%)	2 (11.8%)	4 (23.6%)	10 (58.8%)	Strongly Agree
4. I was able to clearly understand and evaluate the shared online presence of multiple websites using the online presence map.	0 (0%)	1 (5.8%)	1 (5.8%)	1 (5.8%)	0 (0%)	5 (29.4%)	9 (52.9%)	Strongly Agree
5. I was able to clearly understand and evaluate the focus of individual websites and groups of websites using the word cloud.	0 (0%)	1 (5.8%)	0 (0%)	1 (5.8%)	1 (5.8%)	6 (35.2%)	8 (47.1%)	Agree
6. I was able to clearly understand and evaluate the geographic dispersion of the various websites on the map of website locations.	0 (0%)	1 (5.8%)	0 (0%)	0 (0%)	1 (5.8%)	4 (23.6%)	11 (64.7%)	Strongly Agree
7. I was able to clearly understand and evaluate the distribution of emotion that several websites collectively had towards various vaccines.	0 (0%)	0 (0%)	1 (5.8%)	1 (5.8%)	5 (29.4%)	2 (11.8%)	8 (47.1%)	Agree
8. I was able to clearly understand and evaluate the distribution of emotion that a single website had towards various vaccines.	0 (0%)	0 (0%)	0 (0%)	1 (5.8%)	1 (5.8%)	5 (29.4%)	10 (58.8%)	Strongly Agree
9. I found it easy to connect the information across the various visualizations together.	0 (0%)	0 (0%)	0 (0%)	0 (0%)	5 (29.4%)	5 (29.4%)	7 (41.2%)	Agree
10. I found it easy to control the visualizations to see what I wanted to know.	1 (5.8%)	0 (0%)	0 (0%)	1 (5.8%)	5 (29.4%)	5 (29.4%)	5 (29.4%)	Agree
11. I found it helpful to have the various types of information integrated into one system that I controlled.	0 (0%)	0 (0%)	0 (0%)	1 (5.8%)	2 (11.8%)	5 (29.4%)	9 (52.9%)	Strongly Agree

## 4.5. Discussion and Conclusions

This section discusses the conclusions of this study. We will first discuss the overall conclusions in Section 4.5.1. We will follow this with a discussion about conclusions specifically with regard to the various components of the system in Sections 4.5.2. (webometrics) and 4.5.3. (NLP). Next, in Section 4.5.4., we will go over the considerations that for developing future VASes for online public health debates. Finally, in Sections 4.5.5. and 4.5.6., we will discuss the limitations of the study and future research accordingly.

### 4.5.1. Overall

Overall, the study found that VINCENT was a valuable resource for users when making sense of the online vaccine debate. Participants who used the system were more effective at the prescribed tasks, found the tasks easier, and were more confident in their responses than those who did not use the system. By integrating webometrics, NLP, data visualization, and human-data interaction, the system enabled users to make sense of the presented data quickly and effectively. Some participants in the treatment group highlighted the general importance of the system in their ability to complete the tasks.

*Participant 32: “The tasks were very easy to answer with the visualization ... it was very easy to find the required information”*

*Participant 20: “I thought it was really neat to see the different results on screen at the same time. The location, the emotional affects, where the websites sat pro- or anti-vaccine, and the word bank, it gave you so much information at the same time and it kind of allows you to more easily draw a conclusion about a website being able to see it all at once rather than having to research this, sit down and figure it out—(I) wouldn't make the same conclusions or draw it all together at once on my own.”*

The amount of information, its complexity, and the amount of time participants had to assess it were the important factors identified to explain why the treatment group was more effective at the tasks than the control group. VINCENT made it easy to quickly make sense of the data in the

vaccine debate. One participant highlighted this, explaining how they found the system useful and fun.

*Participant 9: “There are a lot of sources out there, and you really can’t go through all of them. The way you have them put all together in one visualization where I can see everything, it just looks very helpful. It’s not something I would do research on my own, but it was fun to navigate with it.”*

The treatment group understood how to interact with the system to get the information they needed. They also saw the potential for VINCENT to help them explore this information space further and more accurately than they could on their own. One control group participant discussed this idea after seeing the system for the first time.

*Participant 12: “I think if I had that system, I definitely would have been more confident in my answers. I definitely would have been able to complete them faster and I wouldn’t be so uncertain about almost everything I put on the questions and probably felt more confident in them. Overall, I think the system definitely lays it out for you in a really simple manner so that all that information is accessible to you within the click of a button vs. having to do it manually and go through all the websites and make your own judgement”*

#### **4.5.2. Webometrics**

The webometrics components of VINCENT (maps of geographic location and online presence) were extremely valuable resources for participants to complete the online presence and geographic location tasks. The treatment group was much more effective at identifying the website vaccine position than the control group. The latter struggled to assess quickly the information on the websites and to make surface level judgments about those websites. Further, the control group would frequently misjudge an anti-vaccine website as pro-vaccine, a finding worthy of further investigation.

The implications for this finding alone indicate that unaided by the system, people can struggle to make sense of the overall message that vaccine websites are presenting to them. For example, a website may advocate for parental vaccine choice; we found that participants could

misinterpret this as an indication that the website is pro-vaccine. The system was necessary for completing such a task because it showed, supported by an analysis of inlinks, how much shared online presence the websites had with one another, and therefore provided further insights into whether websites were actually pro- or anti-vaccine. The users could then corroborate the findings by looking at the websites' emotion or focus data.

It was interesting to note that some participants who were interviewed, both from the treatment group and control group, felt confident that they were or would be able to accurately perform Task 1 with or without the system, as the excerpts below demonstrate.

*Participant 3: "I think what I meant by it was easy was that it was easy once I was looking at the website to determine if it was pro- or anti-vaccine. But it wasn't completely easy because there were so many to go through"*

*Participant 5: "For websites, it would have been a lot longer a process. I probably would have used the paper you gave me and just counted them pro- and anti-, but it would have been more lengthy. I still think I would have felt just as confident because I'm literate and I can see what is going on, but it would have been much longer of a process but just as confident"*

*Participant 9: "Without the system, it would have been easy but time consuming ... I'd also be confident—100% with the system and 90% without"*

*Participant 8: "It wasn't a hard task to do, just very time consuming and required putting some thought into it."*

In other words, some participants were not aware of how poorly they did or would do on Task 1 without VINCENT, further demonstrating how beneficial the system was in helping the user make sense of the debate. There is often ambiguity and a lack of clarity about what the information from these websites is trying to convey, which takes more effort to uncover.

The information displayed by the online presence map was very clear to users, as demonstrated by their ability to perform these tasks accurately. One issue that caused confusion, however, was the meaning and interpretation of the axes. The coordinates of the websites on the scatterplot

were generated using MDS, which plots the data points with regard to their similarity to one another. With MDS, the proximities of the data points are more important than their location independently on the axes. These axes are not always well defined, and it is up to the reader of the map to discern their meaning from the scatterplot. In this case, the researchers could not infer exactly what the y-axis was reflective of and therefore did not indicate it on the map. In order to limit confusion, it would be important to further explain this in the instructional video.

The map of website locations was clear to the users, as they were able to navigate it accurately and easily. However, many users were not familiar with some of the specific geographic terms of the task. Further geographic information should be provided in the system, and it should not be assumed that these geographic areas are common knowledge. Color coding the areas so the user can see them clearly or adding the areas to the information box are possible solutions.

#### 4.5.3. NLP

The NLP components of VINCENT (the emotion bar chart and word cloud) were also extremely valuable for users. The treatment group was able to properly interact with the emotions bar chart and get the information they needed. Some users mentioned, however, that they found it somewhat confusing to use and were unsure if they properly understood the generated information. Comparing a single website's positive and negative emotion regarding a vaccine, for example, was impeded because the bars did not line up on top of each other. Some participants in the treatment group highlighted why they struggled to use the emotion bar chart, citing reasons like: difficulty activating the emotion bar chart correctly or struggling to compare the positive and negative emotions of a single website.

*Participant 6: "The (emotion analysis) on the bottom. I didn't know quite how to read that. When I select an individual website from the main interface here, I didn't quite understand how the changes work."*

*Participant 26: "(The emotion panel) is very powerful as long as it is activated correctly, and that was one of my issues—I didn't know how to properly activate it. But it is powerful itself"*



*Participant 6: “I think I may have looked at (the task) too quickly. Since the bars aren’t lining up (it was difficult to assess) ....”*

The activation of the word cloud was straightforward for users. Some mentioned, however, that assessing the size of the smaller words in the word cloud was challenging. As the words got smaller, it was more difficult to differentiate the words from one another or determine if the size of one word was bigger or smaller than another. The word cloud could either be expanded to take up more space on the display so that the sizes of the words are easier to differentiate, or another method of visualizing word frequency could be considered.

#### **4.5.4. Considerations for Developing VASes for Online Public Health Debates**

The study suggests that several considerations should go into developing future VASes for online public health debates. These considerations include: find ways to reduce the users’ effort in evaluating the data, plan around users’ mental models to determine how the VASes interactions should work, and include more supporting information for users to accurately assess the data in the VAS.

When developing these systems, careful consideration must be given to the difficulty users have evaluating information. It is important to reduce the analysis effort required to use the system by adopting well-thought-out design and additional data analytics methods. We found that for tasks requiring participants to analyze data in depth and go beyond simply locating specific pieces of information in the system, the treatment group’s performance was poorer than on tasks that required the users to conduct less in-depth analysis. For example, in Task 7 (one in which treatment group did not significantly outperform the control group), participants were asked to evaluate if various websites had stronger positive or negative emotions with regard to the HPV vaccine. While the treatment group was generally able to locate on the emotion bar charts the appropriate information needed for this task, and subsequently connect the pieces of information together, making the proper judgment with regard to exactly how much the bars differed was challenging. Implementing ways for users to automatically compare this data would limit the potential to misunderstand the charts and information presented.

Development of these systems requires careful consideration of previous mental models of users. Adopting strategies that meet these models to avoid confusion is important. While many of VINCENT's design features did this well, some aspects of the system did not match these mental models. For example, users could interact with the system to select websites directly on the visualizations or select specific websites via a drop-down website selector tool. Participants were, at times, confused by these interactions because they split the selection process into two separate interactions (this was a limitation of the tool used to develop VINCENT). Integrating all of the selection options together would have made the system more usable for participants. Furthermore, participants mentioned that they expected that using the website selector tool, they could select more than one website at a time (specifically for the comparison tasks). Creating systems that function more intuitively for users would reduce the time and effort it takes users to learn how to use the VAS and do the tasks required of them.

It is important that these systems convey information with enough supporting information for users to properly perform sense-making tasks. For example, in Task 10, participants were asked to evaluate the concentration of pro- and anti-vaccine websites according to several specified geographic regions. Some of these regions (specifically "Midwestern USA") were not clear to the users. Supporting information must be included in these systems so that the users are not limited in their ability to perform tasks.

#### **4.5.5. Limitations**

There were several limitations to this study. First, the amount of time which we could reasonably ask of participants was limited. Participants from both groups mentioned in the interviews and the post-task questionnaires that time was a factor in their ability to complete the tasks. Had participants had additional time, the results could have been affected. However, it was clear that when under a time constraint, VINCENT enabled participants to complete more tasks.

As well, participants in this study were required to be university students in Canada, which limited the diversity of the population of the study. Testing the system on and against users with more knowledge and experience with the public health issue of interest may have yielded different results than what we found in our study.

The tools used to design the visualizations and interactions of VINCENT limited the functionality and, subsequently, the effectiveness of the system. The separated selector system (highlighting a website or selecting it from the dropdown menu), filtering of the emotion bar charts, and the inability to select more than one website at a time from the dropdown list were all mentioned as setbacks for the treatment group during this research. Developing a system from scratch or using another visualization development tool (like D3.js) may have alleviated some of these design limitations. There are currently no visual analytics systems that examine online public health debates. As a result, it is difficult to compare the tool developed here to other existing research.

Finally, the tools and data used for the data analytics of VINCENT also carried limitations. For example, the online presence map relied on domain-level inlinks from MOZ's Link Explorer. Further analysis using different levels of inlinks (site- or page-level) could improve the data analytics. For geographic location data, we would rely on WHOIS registration data if we could not find location information on the websites. While this was a useful tool for identifying locations, it can also be misleading if the website is registered in one location but is hosted or aimed at an audience in another location. Furthermore, the NLU API was used as the method of analyzing website text. This out of box text-based emotion analysis was useful for this study, but more reliable results could potentially be achieved using a customized NLP tool that had been trained on the text of the domain of interest. For example, BioBERT is an NLP tool that has been trained on large-scale biomedical corpora and could be useful for these types of public health related tasks (Lee et al., 2019).

#### **4.5.6. Future Research**

This study found that VINCENT was a valuable resource for users investigating the online vaccine debate, a noted public health issue of our time. Further research is needed to examine how systems like this can be applied in other areas of debate, both within public health and in other domains. Those with a vested interest in making sense of public health topics, for example cannabis or alternative health practices, as well as topics from other domains (e.g., academia, business, or politics) could benefit from the development of similar systems.

Furthermore, future research should look at using alternate methods of data analytics, data visualizations and human-data interactions to those utilized in this study. Social media could also be an important medium for further analysis of online public health debates. As well, social network analyses for examining and visualizing shared online presence in place of MDS, used here, could result in more effective user performance on the sense-making tasks. By examining alternate methods for developing VASes for online public health debates, future systems can be developed with a clearer understanding of which methods are best for users who need to make sense of online public health debates.

## 4.6. References

- Bilgri, O. R. (2016). From “herbal highs” to the “heroin of cannabis”: Exploring the evolving discourse on synthetic cannabinoid use in a Norwegian Internet drug forum. *International Journal of Drug Policy*, 29, 1–8.
- Björneborn, L., & Ingwersen, P. (2004). Toward a basic framework for webometrics. *Journal of the American Society for Information Science and Technology*, 55(14), 1216–1227. <https://doi.org/10.1002/asi.20077>
- Börner, K. (2015). *Atlas of Knowledge: Anyone Can Map*. The MIT Press.
- Brunson, E. K., & Sobo, E. J. (2017). Framing Childhood Vaccination in the United States: Getting Past Polarization in the Public Discourse. *Human Organization*, 76(1), 38–47.
- Caban, J. J., & Gotz, D. (2015). Visual analytics in healthcare – opportunities and research challenges. *Journal of the American Medical Informatics Association*, 22(2), 260–262. <https://doi.org/10.1093/jamia/ocv006>
- Chang, Y. C., Ku, C. H., & Chen, C. H. (2019). Social media analytics: Extracting and visualizing Hilton hotel ratings and reviews from TripAdvisor. *International Journal of Information Management*, 48, 263-279.
- Chen, M., & Ebert, D. S. (2019). An ontological framework for supporting the design and evaluation of visual analytics systems. In *Computer Graphics Forum* (Vol. 38, pp. 131–144). Wiley Online Library.
- Dubé, E., Vivion, M., & MacDonald, N. E. (2015). Vaccine hesitancy, vaccine refusal and the anti-vaccine movement: influence, impact and implications. *Expert Review of Vaccines*, 14(1), 99–117.
- Fekete, J.-D., Jankun-Kelly, T. J., Tory, M., & Xu, K. (2019). Provenance and Logging for Sense Making (Dagstuhl Seminar 18462). Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik.
- Grimes, S. (2016). Sentiment, emotion, attitude, and personality, via Natural Language Processing. Retrieved January 20, 2019, from <https://www.ibm.com/blogs/watson/2016/07/sentiment-emotion-attitude-personality-via-natural-language-processing/>
- Han, J., Pei, J., & Kamber, M. (2011). *Data mining: concepts and techniques*. Elsevier.
- Hirschberg, J., & Manning, C. D. (2015). Advances in natural language processing. *Science*, 349(6245), 261–266. <https://doi.org/10.1126/science.aaa8685>

- Hohman, F. M., Kahng, M., Pienta, R., & Chau, D. H. (2018). Visual analytics in deep learning: An interrogative survey for the next frontiers. *IEEE Transactions on Visualization and Computer Graphics*.
- Kata, A. (2010). A postmodern Pandora's box: Anti-vaccination misinformation on the Internet. *Vaccine*, 28(7), 1709–1716. <https://doi.org/10.1016/j.vaccine.2009.12.022>
- Kata, A. (2012). Anti-vaccine activists, Web 2.0, and the postmodern paradigm - An overview of tactics and tropes used online by the anti-vaccination movement. *Vaccine*, 30(25), 3778–3789. <https://doi.org/10.1016/j.vaccine.2011.11.112>
- Keel, P. E. (2007). EWall: A visual analytics environment for collaborative sense-making. *Information Visualization*, 6(1), 48–63.
- Keim, D., Andrienko, G., Fekete, J. D., Görg, C., Kohlhammer, J., & Melançon, G. (2008). Visual analytics: Definition, process, and challenges. In *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)* (Vol. 4950 LNCS, pp. 154–175). [https://doi.org/10.1007/978-3-540-70956-5\\_7](https://doi.org/10.1007/978-3-540-70956-5_7)
- Kickbusch, I. (2009). Health literacy: engaging in a political debate. *International Journal of Public Health*, 54(3), 131–132.
- Klein, G., Moon, B., & Hoffman, R. R. (2006). Making sense of sensemaking 1: Alternative perspectives. *IEEE Intelligent Systems*, 21(4), 70–73.
- Kovacs, B. E., Gillison, F. B., & Barnett, J. C. (2018). Is children's weight a public health or a private family issue? A qualitative analysis of online discussion about National Child Measurement Programme feedback in England. *BMC Public Health*, 18(1), 1295.
- Lee, J., Yoon, W., Kim, S., Kim, D., Kim, S., So, C. H., & Kang, J. (2019). Biobert: pre-trained biomedical language representation model for biomedical text mining. *ArXiv Preprint ArXiv:1901.08746*.
- Lewandowsky, S., & Oberauer, K. (2016). Motivated rejection of science. *Current Directions in Psychological Science*, 25(4), 217–222.
- Liu, Z., Nersessian, N., & Stasko, J. (2008). Distributed cognition as a theoretical framework for information visualization. *IEEE Transactions on Visualization and Computer Graphics*, 14(6).

- Lorenz-Spreen, P., Mønsted, B. M., Hövel, P., & Lehmann, S. (2019). Accelerating dynamics of collective attention. *Nature Communications*, *10*(1), 1759.
- Marshall, C. C., & Bly, S. (2005). Saving and using encountered information: implications for electronic periodicals. In *Proceedings of the Sigchi conference on human factors in computing systems* (pp. 111–120). ACM.
- Mavragani, A., & Ochoa, G. (2018). The internet and the anti-vaccine movement: tracking the 2017 EU measles outbreak. *Big Data and Cognitive Computing*, *2*(1), 2.
- Mitra, T., Counts, S., & Pennebaker, J. W. (2016). Understanding Anti-Vaccination Attitudes in Social Media. In *ICWSM* (pp. 269–278).
- Morphett, K., Herron, L., & Gartner, C. (2019). Protectors or puritans? Responses to media articles about the health effects of e-cigarettes. *Addiction Research & Theory*, 1–8.
- Nguyen, P. H., Xu, K., Wheat, A., Wong, B. L. W., Attfield, S., & Fields, B. (2015). Sensepath: Understanding the sensemaking process through analytic provenance. *IEEE Transactions on Visualization and Computer Graphics*, *22*(1), 41–50.
- Ninkov, A., & Sedig, K. (2019). VINCENT: A visual analytics system for investigating the online vaccine debate. *Online Journal of Public Health Informatics*, *11*(2).
- Ninkov, A., & Vaughan, L. (2017). A webometric analysis of the online vaccination debate. *Journal of the Association for Information Science and Technology*, *68*(5), 1285–1294. <https://doi.org/10.1002/asi.23758>
- Pirolli, P., & Card, S. (2005). The sensemaking process and leverage points for analyst technology as identified through cognitive task analysis. In *Proceedings of international conference on intelligence analysis* (Vol. 5, pp. 2–4). McLean, VA, USA.
- Ragini, J. R., Anand, P. M. R., & Bhaskar, V. (2018). Big data analytics for disaster response and recovery through sentiment analysis. *International Journal of Information Management*, *42*, 13–24.
- Rind, A., Wagner, M., & Aigner, W. (2019). Towards a Structural Framework for Explicit Domain Knowledge in Visual Analytics. *ArXiv Preprint ArXiv:1908.07752*.
- Salomon, G. (1993). No distribution without individuals' cognition: A dynamic interactional view. *Distributed Cognitions: Psychological and Educational Considerations*, 111–138.

- Sedig, K., & Parsons, P. (2013). Interaction design for complex cognitive activities with visual representations: A pattern-based approach. *AIS Transactions on Human-Computer Interaction*, 5(2), 84–113.
- Sedig, K., & Parsons, P. (2016). *Design of Visualizations for Human-Information Interaction: A Pattern-Based Framework. Synthesis Lectures on Visualization* (Vol. 4).  
<https://doi.org/10.2200/S00685ED1V01Y201512VIS005>
- Sedig, K., Parsons, P., & Babanski, A. (2012). Towards a Characterization of Interactivity in Visual Analytics. *Journal of Multimedia Processing and Technologies, Special Issue on Theory and Application of Visual Analytics*, 3(1), 12–28.  
<https://doi.org/10.1145/0000000.0000000>
- Shneiderman, B., Plaisant, C., & Hesse, B. W. (2013). Improving healthcare with interactive visualization. *Computer*, 46(5), 58–66.
- Thelwall, M., Vaughan, L., & Björneborn, L. (2005). Webometrics. *ARIST*, 39(1), 81–135.
- Varshney, K. R., Rasmussen, J. C., Mojsilović, A., Singh, M., & DiMicco, J. M. (2012). Interactive visual salesforce analytics.
- Vaughan, L., & Ninkov, A. (2018). A new approach to web co-link analysis. *Journal of the Association for Information Science and Technology*, 69(6), 820–831.
- Velardo, S. (2015). The nuances of health literacy, nutrition literacy, and food literacy. *Journal of Nutrition Education and Behavior*, 47(4), 385–389.
- Who.int. (2019). Ten health issues WHO will tackle this year. Retrieved February 12, 2019, from <https://www.who.int/emergencies/ten-threats-to-global-health-in-2019>
- Zhang, S., Qiu, L., Chen, F., Zhang, W., Yu, Y., & Elhadad, N. (2017). We make choices we think are going to save us: Debate and stance identification for online breast cancer CAM discussions. In *Proceedings of the 26th International Conference on World Wide Web Companion* (pp. 1073–1081).



## Chapter 5 - Online public health debates: A framework-based approach to visual analytics<sup>6</sup>

Anton Ninkov  
Western University  
Faculty of Information and Media Studies

Dr. Kamran Sedig  
Western University  
Faculty of Information and Media Studies & Department of Computer Science

---

<sup>6</sup> A version of this chapter was submitted for publication:

Ninkov, A., & Sedig, K. (in review, 2020) Online public health debates: A framework-based approach to visual analytics. *Online J. Public Health Inform.*

## Abstract

Nowadays, many people are deeply concerned about their physical well-being and invest time and effort into investigating health-related topics. In response to this, many online websites and social media profiles have been created, resulting in a plethora of information on such topics. Much of this information is oftentimes conflicting; this results in online camps that have different positions and arguments on a topic. In this paper, we refer to the collection of all such positionings and entrenched camps on a topic as an online public health debate. The information people encounter regarding such debates can ultimately influence how they make decisions, what they believe, and how they act. Therefore, there is a need for public health stakeholders (i.e., people with a vested interest in public health issues) to be able to make sense of online debates quickly and accurately.

In this paper, we present a framework-based approach for investigating online public health debates called ODIN (Online Debate entIty aNalyzer). We first present the concept of online debate entities (ODEs), which is a generalization for those that participate in online debates (e.g., websites and Twitter profiles). We then present ODIN, in which we identify, define, and justify seven ODE attributes that we consider important for making sense of these debates. Next, we provide an overview of four online public health debates (vaccines, statins, cannabis, and dieting plans) using ODIN. We then present four framework-based visual analytics systems (VASes) designed to help stakeholders to quickly make sense of these online public health debates.

## 5.1. Introduction

Nowadays, many people are deeply concerned about their physical well-being and invest time and effort into investigating health-related topics. In response to this, many online websites and social media profiles have been created, resulting in a plethora of information on such topics (Kitchens, Harle, & Li, 2014; Swar, Hameed, & Reychav, 2017). Although some of this information is very useful, making sense of this information

with its large quantity and varying quality is difficult and can confuse people (Seymour, Getman, Saraf, Zhang, & Kalenderian, 2015). One of the reasons for this difficulty is that people are often confronted with various sources of conflicting information (Truumees et al., 2020; Yoon, Sohn, Choi, & Jung, 2017). Oftentimes, such conflicting information creates online camps that have different positions and arguments. In this paper, we refer to the collection of all such positionings and entrenched camps on a topic as an online public health debate.

The information that people encounter when examining online public health debates can influence how they come to make decisions, what they believe, and how they act (Fox & Duggan, 2013; Miller & Bell, 2012). The ability to quickly examine and make sense of online health debates can help the general public assess the balance of such debates. Beyond the general populace, other stakeholders, such as public health practitioners and policy makers, may also need to make sense of such online debates quickly and accurately.

Currently, however, making sense of such online debates is not straightforward. This is the result of several factors. Here, we highlight four of these factors: 1) number of sources—there are many sources of information about each health issue; 2) distribution of sources—information sources are distributed across the Internet on different websites, including a multiplicity of social media platforms; 3) veracity of information—it is difficult to examine the veracity of information originating from different sources; and, 4) positioning of sources—it is not easy to immediately understand the sentiments expressed by different sources and how they position themselves.

Computational tools can help alleviate some of the difficulties encountered above (Jonassen, 1995; Z. Liu, Nersessian, & Stasko, 2008; Sedig, Klawe, & Westrom, 2001). A subset of computational tools that can help stakeholders make sense of complex information, such as the information encountered in public health debates, is visual analytics systems (VASes) (Sedig & Parsons, 2016). By integrating data analytics, data visualizations, and human-data interaction, VASes can facilitate sense-making activities

(Pirolli & Card, 2005; Sedig & Parsons, 2016). It is important for such systems to help stakeholders not only perform the necessary sense-making activities, but also be human-centered, fitting the needs of the users. Without a proper understanding of the structure and elements of online public health debates, however, designing and building of such systems would be haphazard rather than systematic in approach. To create human-centered tools in healthcare, we must have appropriate frameworks (Sedig, Naimi, & Haggerty, 2017).

To help design VASes that facilitate making sense of online public health debates, in this paper, we present and discuss a framework that we have developed called ODIN (Online Debate entlTy aNalyzer). This framework is for generalizing online public health debates and is based on a construct, which we call Online Debate Entities (ODEs). ODEs are the various sources of information--that is, organizations and people with an online presence debating public health topics (e.g., websites, Twitter profiles, Facebook users, and/or Reddit users). The ODIN framework helps with the analysis of various attributes of ODEs which are needed to permit stakeholders to quickly make sense of any online debate. In this framework, we identify and define seven attributes: presence, shared presence, geographic location, registrant, age, focus, and sentiments. Using four examples of online public health debates (vaccines, cannabis, statins, and dieting plans), we demonstrate how ODIN can be used not only for systematizing the analysis of ODEs, but also for helping design framework-based VASes that facilitate stakeholders' investigations of other online public health debates.

The remainder of this paper is organized as follows. Section 5.2., Background, discusses information spaces, sense-making, VASes, and online public health debates. Section 5.3. presents ODIN and discusses how it helps with analyzing online debates. Section 5.4. describes the four online debate cases using ODIN. Section 5.5. demonstrates how ODIN-based VASes can be developed and how they facilitate making sense of online debate. Section 5.6. presents some conclusions and discusses the limitations of this approach and describes future work.

## 5.2. Background

This section discusses the background concepts and terminologies used in this paper. We begin with information spaces and sense-making. We follow this with a discussion on VASes, why they are important, what they are made of, and examples of how they can help users to complete sense-making activities. Finally, we go over online public health debates and examine 4 specific cases: vaccines, cannabis, statins, and dieting.

### 5.2.1. Information Spaces & Sense-Making

Information spaces are bodies of information that are thought of as having spatial characteristics (Sedig & Parsons, 2016). Compared to the related concept of “data”, which refers to information that has already been discerned and recorded, an information space is a useful idea for visual analytics research because it allows the freedom to conceptualize unstructured or abstract information. For online debates, these concepts are important because they are unstructured information where data is not necessarily readily available (Moreno, Ozogul, & Reisslein, 2011; Snyder, 2014; J. C. Thomas, Diamant, Martino, & Bellamy, 2012). Information spaces are made up of information items (e.g., entities, properties and relationships) that exist at various levels of granularity. Making sense of information spaces is an example of a cognitive activity.

Cognitive activities are part of everyday, modern life. Cognitive activities that are information-intensive and involve intense human cognition can be further described as complex cognitive activities (Ericsson & Hastie, 1994; Funke, 2010). Complex cognitive activities have two distinct characteristics: 1) they require the use of complex psychological processes and 2) they exist in the presence of complex conditions (Knauff & Wolf, 2010). An example of a complex cognitive activity is making sense of online public health debates.

Sense-making is an activity in which people gradually develop mental models of an information space about which they have insufficient knowledge (Klein, Moon, &

Hoffman, 2006; Sedig & Parsons, 2013). Sense-making often includes a set of tasks, some of which can include: scanning the information space, selecting relevance of items, and examining items in more detail (Pirolli & Card, 2005). Sense-making can require people to take complex information and uncover meaning from it that otherwise could go unnoticed. Sense-making activities are often problems that are ill-structured and open ended (Buchel & Sedig, 2016). This includes, for example, investigating online debates (Ninkov & Sedig, 2019). Sense-making involves users establishing goals, discovering an information space's structure, and determining what questions to ask as well as how the answers to those questions should be organized (Russell, Stefik, Pirolli, & Card, 1993). A challenge when completing sense-making activities is that relevant information for completing required tasks is not always easy to access, stored in the proper format, or located in the correct locations (Marshall & Bly, 2005). Visual analytics systems have been developed to support users in sense-making activities in a variety of information spaces (Fekete, Jankun-Kelly, Tory, & Xu, 2019; Hohman, Kahng, Pienta, & Chau, 2018).

### 5.2.2. VASes

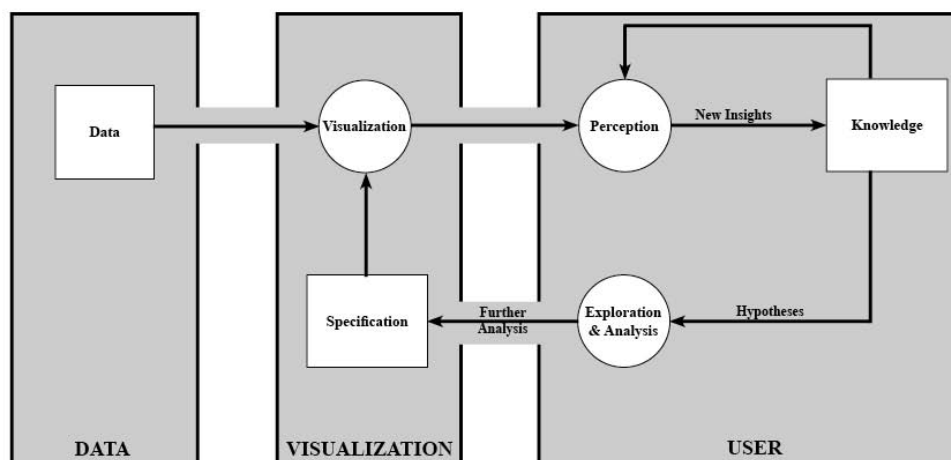
In today's environment of big data, people are often the victims of information overload. They can get lost in and overwhelmed by the voluminous data presented to them, and consequently struggle to decipher meaning (Keim et al., 2008). With VASes that combine human insight with powerful data analytics, data visualizations and human-data interaction, users can be alleviated of some of these problems they face. VASes can enable stakeholders to make sense of data in ways that have never before been convenient or possible. "Just like the microscope, invented many centuries ago, allowed people to view and measure matter like never before, (visual) analytics is the modern equivalent to the microscope" (Börner, 2015). VASes combine data analytics with interactive visualizations to synthesize, analyze, and facilitate high-level cognitive activities, like sense-making, that involve investigating data (Keim, Kohlhammer, Ellis, & Mansmann, 2010; J. J. Thomas & Cook, 2005).

VASes are composed of three integrated components: 1) an analytics engine, 2) data visualizations, and 3) human-data interactions (Sedig & Parsons, 2016; Sedig, Parsons, & Babanski, 2012). The analytics engine pre-processes and stores data (e.g., data cleaning & fusion), transforms it (e.g., normalization), and analyzes it (e.g., multi-dimensional scaling, emotion analysis) (Han, Pei, & Kamber, 2011). Data visualizations in a VAS can be visual representations of the information derived from the analytics engine.

Visualizations extend the capabilities of individuals to complete tasks by allowing them to analyze data in ways that would be difficult or impossible to do otherwise (Sedig et al., 2012; Shneiderman, Plaisant, & Hesse, 2013). Human-data interaction is used in VASes to allow the user to control the data they see and the way the data is processed.

Interaction in VASes supports users through distributing the workload between the user and the system during their exploration and analysis of the data (Z. Liu et al., 2008; Salomon, 1993; Sedig & Parsons, 2016).

Sense-making requires people to rapidly compare and contrast information items unhindered by information formats or locations (Keel, 2007), for which VASes can be particularly useful. The sense-making loop for VASes can help us understand the process that users go through using these systems (Figure 18) (Keim et al., 2008; Van Wijk, 2005). Data is analyzed and fed into the VAS. To make sense of the data, the user first perceives it through the visualization, and then gathers new insights from it by conducting interactions on the VAS. Based on the new knowledge they acquire, the user can then further analyze and explore the information in new ways by specifying the VAS to what they need to see, which vary depending on the tasks at hand. The loop repeats as the user completes sense-making tasks and generates further insights from the data.

**Figure 18***Sense-Making Loop*<sup>7</sup>

### 5.2.3. Online Public Health Debates

The term “online debate” is used in this paper to refer to a concept or topic that is widely discussed on the Internet and has two or more differing points of view. While some online debates exist on a micro-level (i.e., instances of a topic being debated back and forth between two ODEs on a web forum) (Herring, Job-Sluder, Scheckler, & Barab, 2002; Nicholson & Leask, 2012; Oraby et al., 2017), in this paper we are focused on online debates that manifest themselves on a macro-level (i.e., prominent ODEs promoting information and viewpoints on various positions of a debate) (Getman et al., 2018; Ninkov & Sedig, 2019; Ninkov & Vaughan, 2017).

Some well-documented examples of topics that are debated online include: gun control (Sridhar, Foulds, Huang, Getoor, & Walker, 2015; Walker, Tree, Anand, Abbott, & King, 2012), vaccines (Kata, 2010, 2012; Ninkov & Sedig, 2019; Ninkov & Vaughan, 2017), abortion (Bloch, 2007; R. L. Hill, 2017), and climate change (Collins & Nerlich, 2015;

<sup>7</sup> The Sense-Making Loop has been based on (Keim et al., 2008; Van Wijk, 2005)



Howarth & Sharman, 2015). Online debates are important for researchers to investigate because they have been shown to influence the way people perceive an issue and have an impact on their real-world decisions (Kickbusch, 2009; Morphett, Herron, & Gartner, 2019; Velardo, 2015). People often rely on the Internet to gather the information they need to form their opinions; this practice extends to topics regarding what is best for their health.

Every online debate is an information space made up of information items, including ODEs and their attributes. The ODEs that participate in a debate have contrasting views and opinions to one another. ODEs include websites or social media profiles (e.g., Twitter) of people, organizations, companies, and government agencies. Each of these ODEs have important characteristics that help position them within the information space. Visually representing this information is important for stakeholders to be able to make sense of it.

Next, we examine the background and applications of four example online public health debates. These debate topics include: vaccines (Getman et al., 2018; Nicholson & Leask, 2012; Seeman & Rizo, 2009), cannabis (Bilgrei, 2016; Hasan & Ng, 2014), statins (Huesch, 2017; Navar, 2019), and dieting plans (Jauho, 2016; Mazzi, 2018).

#### *5.2.3.1. Vaccines*

There is little debate within the medical community on the efficacy of vaccines. Some public health experts claim that there have been up to 103 million prevented contagious diseases since 1924 as a result of vaccines (Van Panhuis et al., 2013). Despite the documented successes of the practice, there are still some in the medical community who have concerns about the risks or side effects of vaccines, and question the practice of mandatory vaccination (Mawson, Ray, Bhuiyan, & Jacob, 2017). The most infamous example of this is in a 1998 article published in *The Lancet* (now redacted) by Andrew Wakefield that reported a connection between the measles, mumps, and rubella vaccine and developmental disorders in children (Wakefield et al., 1998). The findings of this

study still influence the discussion around vaccination today – Wakefield is a prominent figure throughout the anti-vaccine community. The debate about vaccines is lively within the general public and has become even more contentious over time. It is well documented, with a recent rise in the number of unvaccinated children and the subsequent re-emergence of previously near eradicated diseases, such as measles (Cherney & McKay, 2019; Kwai, 2019). Because of this trend, the anti-vaccination movement has been considered an emerging public health problem by some experts (Dubé, Vivion, & MacDonald, 2015; Mavragani & Ochoa, 2018). The World Health Organization even listed the rise of the anti-vaccination campaign as a top ten health emergency in 2019 (Who.int, 2019).

#### 5.2.3.2. Cannabis

Cannabis is the subject of debate both within the medical community and the general public. The cannabis debate occurs with regard to its medical applications, the movement towards legalization for recreational purposes, and the risks associated with its use both medically and recreationally. Cannabis has been examined by researchers for its potential medical uses in helping patients with conditions such as chronic pain (K. P. Hill, 2015), epilepsy (Maa & Figi, 2014), and Parkinson's disease (Lotan, Treves, Roditi, & Djaldetti, 2014). The drug has shown some promise to help with these conditions, and potentially has uses for some elusive illnesses like fibromyalgia (Fiz, Durán, Capellà, Carbonell, & Farré, 2011) and insomnia (Russo, Guy, & Robson, 2007). While these beneficial applications of cannabis have been studied, there is also a great deal of research detailing the harm that cannabis can pose to users. Some reasons for concern include: addiction (Maldonado, Berrendero, Ozaita, & Robledo, 2011; Wenger, Moldrich, & Furst, 2003), infertility (Rajanahally et al., 2019), and adverse effects on cognitive motor functions (Volkow, Baler, Compton, & Weiss, 2014). As well, there are concerns that cannabis can act as a gateway into other more dangerous drugs like heroin or cocaine (Fergusson, Boden, & Horwood, 2006; Secades-Villa, Garcia-Rodríguez, Jin, Wang, & Blanco, 2015). There is disunity in the medical community about the correct position on cannabis

and whether the risks of the drug outweigh its potential benefits. Current literature about the potential applications or risks of the drug discusses the limited research there has been thus far as a result of the legal constraints, and calls for further research on the topic urgently (National Academies of Sciences, Engineering, and Medicine, 2017).

#### 5.2.3.3. *Statins*

The debate about cholesterol management and the use of statins is a public health issue within the medical community (Godlee, 2016; Redberg & Katz, 2016). Statins are debated with regard to two important concepts: 1) who should receive the drug and 2) how common are its side effects (J. A. Hill, 2019). There are some medical professionals that believe statins should only be used to help manage cholesterol for secondary prevention, meaning that after an event such as a heart attack, statins are administered to a patient to help control or prevent a future event (Abramson, Rosenberg, Jewell, & Wright, 2013). Others believe that statins can play a wider role in managing cholesterol via primary prevention, meaning that patients with high cholesterol, who are at risk for an event, take statins regularly to stop the event from happening. The reasons for these two diverging positions are varied. There is a range of concern about the effectiveness of statins in preventing events compared to the potential risks involved with the side effects that can be caused by taking the medication. A 2019 study found that 10% of patients who were eligible for statins declined the medication citing side effects of the drug as their primary concern (Bradley et al., 2019). Some of the most commonly reported side effects of statins are muscle aches & pains, muscle breakdown (rhabdomyolysis), diabetes, and liver failure (Newman et al., 2019). Patients' concern regarding these side effects is believed to also influence how they might react to the drug itself. The term "nocebo" was recently used by the American Heart Association to explain a growing phenomenon of patients experiencing side effects as a result of their expectations of a side effect from the drug and not necessarily from the drug itself (Newman et al., 2019).

#### 5.2.3.4. *Dieting Plans*

What and how much people eat is a critical part of their health, and dieting plans strategies for weight loss and health management are an important public health debate. There is a rich history informing the debate on dieting. Origins of the practice date back over 2000 years (Foxcroft, 2012). Specific trends and dieting plans themselves have evolved over time. The practice of dieting is expansive, and, depending on the diet and context, they can be harmful, ineffective, or beneficial (Lowe & Timko, 2004). There are many dieting plans promoted as the right choice for a healthy diet. For example, a vegan diet encourages people to abstain from eating any animal products (e.g., milk, eggs, meat). Not only is this encouraged for personal health and weight loss (Barnard et al., 2006; Huang, Huang, Hu, & Chavarro, 2016), but it has also been promoted as a way to combat climate change (Cleveland & Gee, 2017; Springmann, Godfray, Rayner, & Scarborough, 2016). Another example is the Keto Diet, a popular diet that encourages people to eat fewer carbohydrates (e.g., bread, rice, pasta) and more fats (e.g., eggs, avocado, cheese). The diet has been promoted as a way to help people manage weight loss (Abbasi, 2018; Bueno, de Melo, de Oliveira, & da Rocha Ataide, 2013; Paoli, 2014) as well as epilepsy in some situations (Masino & Rho, 2019; Neal et al., 2008). Examples of other popular, competing dieting plans include: South Beach Diet (Ariagno, 2018), Atkins Diet (Matarese & Harvin, 2018), and Paleo Diet (Andromalos, 2018).

### 5.3. ODIN

In this section, we present a framework to analyze ODEs called ODIN (Online Debate Entity Analyzer). Online debates occur on the general web and on social media platforms (e.g., Twitter) as a result of contrasting or conflicting information about a topic. They manifest themselves in the lack of agreement between ODEs about what views on a topic or issue are correct or should be accepted. In a general sense, online debates are polarized around two or more opposing positions (pro- or anti- “phenomenon”). However,

contained within each position are nuanced sub-positions, which share similarities with one another but may also be at odds with one another based on a variety of other factors.

Every online debate is an information space. Each ODE and its attributes are information items that help position it within the space. In this paper, we identify, define, and justify seven attributes that we consider to be important for making sense of online debates.

These attributes include presence, shared presence, geographic location, registrant, age, focus, and sentiments. Next, we demonstrate, by citing literature, how online debates can be analyzed and measured. We then provide an example of how each attribute supports making sense of online debates. At the end of the section we provide Table 12, which summarizes the attributes and can be used as a reference for ODIN. In Section 5.4., we use ODIN and show how it can be applied in analyzing 4 online public health debates (vaccines, cannabis, statins, & dieting plans).

All ODEs have a presence, that is, the attention received by the ODE. The more presence an ODE has, the more popularity and/or authority it holds in a debate. Presence can be quantified by various metrics. On the general web, inlinks (Ninkov & Vaughan, 2017; Thelwall, 2004; Thelwall, Sud, & Wilkinson, 2012), website traffic (Baka & Leyni, 2017; Brumshteyn & Vas'kovskii, 2017), and website rankings from online resources such as Alexa (alexa.com), MOZ (moz.com), or Majestic (majestic.com) are all powerful tools for measuring presence. On Twitter, metrics like followers (McCoy, Nelson, & Weigle, 2017; Triemstra, Poeppelman, & Arora, 2018) or follower/following ratios (Anger & Kittl, 2011; Borgmann et al., 2016) have been used as indicators of presence. With these metrics, it is possible to quickly make sense of the popularity and authority of ODEs. If ODEs have a strong presence, the information they are sharing about the online debate has greater reach than ODEs with a weaker presence.

Presence can also be examined to compare the similarity of ODEs. The more shared presence ODEs have, the more likely they are similar to one another in their views and position. Shared presence has been determined using co-link (general web) (Thelwall et al., 2012; Vaughan & Ninkov, 2018) or co-follower (Twitter) (Mosleh, Pennycook,

Arechar, & Rand, 2019; Shi, Mast, Weber, Kellum, & Macy, 2017) analyses. Using these metrics can help stakeholders make sense of an online debate's structure quickly. If multiple ODEs share a lot of presence, they are likely to be similar to each other (Holmberg, 2009; Thelwall & Wilkinson, 2004). Because of this, stakeholders can use shared presence data to quickly see ODEs that are similar to one another and identify potential debate position clusters. On the contrary, stakeholders can also use shared presence to see ODEs that are not similar to one another and identify ODEs that may appear on the surface to have one position but really have different positions at a deeper level.

The registrant is the person or organization that owns an ODE. The registrant is not always easy to determine based on the name and surface level appearance of an ODE; this makes uncovering the background of an ODE's registrant an important task. Content analyses along with a registrant classification system can be used to collect this data (Ninkov & Vaughan, 2017; Rothenfluh & Schulz, 2018). As well, for the general web, services like WHOIS (<https://lookup.icann.org/>) that provide information on the registration of a website have been used to record this information (Cetin, Ganan, Korczynski, & van Eeten, 2017; S. Liu, Foster, Savage, Voelker, & Saul, 2015). Knowing the identity of an ODE registrant makes it easier to make sense of an online debate. Not only does it allow stakeholders to quickly see the obvious biases or conflict of interests, but the relationships between multiple ODEs can also become more transparent with this information. For example, a debate may have ODEs that share a position and have similar sentiments towards a variety of issues, but the registrant is the same person or organization. This would reveal that this position may be overrepresented.

The geographic location of an ODE describes where it is located in the real world. This can be measured by conducting a content analysis of the ODEs (Halavais, 2000; Holmberg & Thelwall, 2009; Ninkov & Sedig, 2019). For the general web, services like WHOIS (<https://lookup.icann.org/>) that provide registration information for websites have been used (Janc, 2016; Ninkov & Sedig, 2019; Tsou et al., 2013), while on Twitter,

location information of tweets are provided by the website itself (Stefanidis et al., 2017; Waseem & Hovy, 2016). Geographic location information is important because when combined with the locations of other ODEs, they make it easy to make sense of the real-world dispersion of an online debate and its positions. For example, specific positions of an online debate may be spread out throughout the world or have regionalization. There could be important reasons for geographic clusters, such as differences in local public health policies or traditions.

The age of an ODE is the length of time that it has been registered. For the general web, the Internet Archive's Wayback Machine (<http://web.archive.org/>) makes it possible to see how long a website has been active. On Twitter, the age of the profile is visible on its homepage. The age of ODEs can be used as a way to assess the evolution of an online debate (Jain & Gupta, 2016; Kefi & Perez, 2018; Kend & Goode, 2018; Zheng et al., 2019). If ODEs of one position are collectively older than ODEs of another, it can be an indication that the debate has evolved over time and the new position is becoming increasingly prominent.

ODEs share content that is necessary for stakeholders to examine to make sense of a debate. The focus of this content can indicate an ODE's position in a debate. Frequently occurring words and phrases related to a debate reveals the focus of an ODE (Fan & Welch, 2016; Ninkov & Sedig, 2019; Ruiz & Barnett, 2015; Waller, Hess, & Demetrious, 2016). As well, natural language processing methods such as topic modelling (Skeppstedt, Kerren, & Stede, 2018; Vilares & He, 2017) or content analysis (Austgulen, 2014; O'Connor, 2017) have been used to identify the focus of ODEs. The focus is important to quickly make sense quickly of the content shared by an ODE or group of ODEs. For example, if various websites appear to have a similar position on a debate, but differ in their focus, they may share an overall position (anti- or pro-) but diverge from one another in their sub-position. As well, if multiple ODEs share a common focus but differ in their sentiments about them, this would be an indication that they are on opposite ends of a debate.

Finally, ODEs share sentiments with regard to various topics. These sentiments can take many forms and give stakeholders important context to help make sense of a debate. Sentiments can be shared through the text, images, audio, and video elements of a website or social media profile. These sentiments can describe the overall tone of an ODE or the specific topics they discuss. ODEs' text can be analyzed using content analysis (Chew & Eysenbach, 2010; Moran, Lucas, Everhart, Morgan, & Prickett, 2016) as well as various natural language processing tools to assess metrics like polarity (i.e., positive and negative sentiment) (Barnaghi, Ghaffari, & Breslin, 2016; Fang & Zhan, 2015; B. Liu, 2015) or emotion (e.g., joy, fear, anger, sadness) (Grimes, 2016; Hirschberg & Manning, 2015; Yu & Ho, 2014). With this information easily accessible, stakeholders can determine the position of an ODE quickly. As well, this information can help stakeholders determine the polarizing issues within a debate. For example, if multiple ODEs mention an issue frequently and some of them share negative sentiments while others share positive ones, it would be an indication of a polarizing issue within the debate.

**Table 12**

*Attributes of ODEs*

Attribute	Definition	Measurement - General Web	Measurement – Twitter
<b>Presence</b>	Attention received by an ODE, which indicates its popularity and authority.	Inlinks Website traffic Website ranking	Followers Follower/Following Ratio
<b>Shared Presence</b>	Presence various ODEs share with one another	Co-Link Analysis	Co-Follower analysis
<b>Registrant</b>	Person and/or organization that registered an ODE	Registration information Content analysis	Profile registration name Content analysis
<b>Geographic Location</b>	Geographic location of an ODE's registration	Registration information Content analysis	Profile registration location Content analysis
<b>Age</b>	Time since ODE was created	Internet Archive	Reported creation date
<b>Focus</b>	Frequently mentioned topics and concepts of an ODE	Word/term frequency Topic Modelling tools Content analysis	Word/term frequency Topic Modelling tools Content analysis
<b>Sentiments</b>	Feelings, attitudes, and emotions of an ODE's text and multimedia content	NLP tools Content analysis	NLP tools Content analysis



## 5.4. Online Public Health Debates – Four Case Studies

In this section, we conduct an overview of four online debates: vaccines, cannabis, statins, and dieting plans. For each online debate, we discuss how it manifests itself online. We do this using ODIN, described in Section 5.3. We show how this framework can be used to make sense of these online debates. It should be noted that in this section, we are not attempting to suggest that this is a thorough investigation into the complex content of these debates, or to conclude what the “right” position in each debate is. Rather, we are attempting to demonstrate that these debates exist online and can be made sense of through an examination of the attributes described in ODIN.

### 5.4.1. Vaccines

The online debate about vaccines consists of pro- and anti-vaccine ODEs. Within each position there are sub-positions that approach the debate differently (Vaughan & Ninkov, 2018). Some anti-vaccine sub-positions include parents against vaccines (e.g., The Informed Parent and ThinkTwice Global Vaccine Institute), autism concern (e.g., Age of Autism), or general anti-vaccine (e.g., National Vaccine Information Center and Vaxxter). Some of the pro-vaccine sub-positions include flu vaccine focus (e.g., Families Fighting the Flu), vaccines for children (e.g., Shot of Prevention and Voices for Vaccines) and general pro-vaccine (e.g., Immunization Action Coalition and GAVI: The Vaccine Alliance).

Vaccine ODEs vary in their presence. Anti-vaccine ODEs like the National Vaccine Information Center (NVIC) and Age of Autism as well as pro-vaccine ODEs like Immunization Action Coalition or GAVI: The Vaccine Alliance have a strong presence. On the other hand, ODEs such as The Informed Parent (anti-vaccine) or Shot of Prevention (pro-vaccine) have a weak presence (Ninkov & Sedig, 2019). Based on presence data, it is possible to quickly identify ODEs that are popular and hold authority in the vaccine debate. One example of this is NVIC, which has both a strong presence in

the debate on the general web (Ninkov & Sedig, 2020) and Twitter (@NVICLoeDown has over 12,000 followers).

As well, shared presence of vaccine ODEs quantifies which ODEs share similar online presence, and which share little to no presence at all. The difference in shared presence is clear in the vaccine debate: ODEs that share either a pro- or anti-vaccine position will also share more presence with one another (Ninkov & Sedig, 2020; Vaughan & Ninkov, 2018). Stakeholders can use this information to help determine which ODEs belong to which positions. Within the anti- and pro-vaccine positions, other sub-positions can be identified using shared presence as well. For example, two ODEs that have a very similar presence include Shot of Prevention and Voices for Vaccines (Ninkov & Sedig, 2019; Vaughan & Ninkov, 2018). Upon further examination, clearly both these websites follow the sub-position of child vaccination.

Vaccine ODEs are owned and operated by various registrants, each with their own biases and motivations. For example, ThinkTwice Global Vaccine Institute's registrant (an anti-vaccine ODE) is New Atlantean Press, a holistic health publisher. New Atlantean Press uses this ODE to promote anti-vaccine information as well as promote the sale of their books. Additionally, some vaccine ODEs are operated under the same registrant. Shot of Prevention and Vaccinate Your Family are both registered under the same registrant and likely share the same motivations and views on the debate. The registrant of vaccine ODEs can also be examined using classification systems that place the registrant into various categories. A 2017 study showed there was a significant difference in the classification of ODE registrants (those that linked to a known vaccine ODE) between anti- and pro-vaccine. The ODEs from the pro-vaccine sample were more likely to have registrants such as medical entities while registrants of ODEs from the anti-vaccine sample were more likely to be individuals (Ninkov & Vaughan, 2017).

The geographic locations of vaccine ODEs reveal the distribution of the debate globally. Overall, there is a concentration of English-language vaccine ODEs in North America and Europe (Ninkov & Sedig, 2019). For example, The Informed Parent is an anti-

vaccine ODE that is located in England and Age of Autism, another anti-vaccine ODE, is located in Virginia. There are also clusters of ODEs that share both their vaccine position and location in various regions of the world. An example of this is in North Western North America, where there is a concentration of anti-vaccine ODEs (Ninkov & Sedig, 2019).

There has been online debate about the efficacy and morality of vaccines since near the beginning of the Internet. A 2017 study (Ninkov & Vaughan, 2017) found that the age of vaccine ODEs sampled were from 18 years to 6 months, with an average age of 8 years. Noticeably, there was no significant difference found between the mean ages of ODEs that were anti- and pro-vaccine. This indicates that both positions of the debate have been represented online for about equally as long.

Vaccine ODEs tend to focus on different topics depending on their debate position. Collectively, anti-vaccine ODEs have a strong focus on side effects (e.g., “autism”, “death”, “injury”) and vaccine choice (e.g., “exemptions”, “choice”, “required”) while pro-vaccine ODEs have a strong focus on vaccine-preventable diseases (e.g., “influenza”, “measles”, “flu”) and vaccine organizations (e.g., “Center for Disease Control”, “National Foundation for Infectious Diseases”). On an individual basis, the focus of an ODE reveals the specific components of the debate that an individual might be interested in. Sabin Vaccine Institute focuses on topics like “development” and “uptake” (Ninkov & Sedig, 2019). Upon further investigation, this is an appropriate indication of the ODE’s mission to create and spread new vaccines throughout the world.

The sentiments and emotions that ODEs share about vaccine topics, both in a general sense and about related concepts (e.g., specific vaccines, side effects), vary. Anti-vaccine ODEs often share negative sentiments about vaccines that use emotions like fear or anger (Ninkov & Sedig, 2019). Pro-vaccine ODEs, on the other hand, often share more positive sentiments about vaccines that use emotions like joy and pleasure when discussing vaccines (Ninkov & Sedig, 2019). Two examples which highlight this difference between the sentiments of the two sides of the debate are seen in Figures 19 and 20. In Figure 19,

news items on the anti-vaccine website Vaxxter are displayed. Phrases like “embroiled in fight with state’s laws”, “ineffective vaccines” or “vaccinated 6-year-old dies” all highlight the negative sentiments shared about vaccines by this ODE. On the other side of the debate, in Figure 20, the homepage of Voices for Vaccines is displayed. Phrases like “pro-vaccine”, “importance of on-time vaccination”, and “we want to protect all children” all highlight the positive sentiments shared about vaccination by this ODE. On a larger scale, the sentiments shared regarding vaccine topics are important for evaluating the content and views of the different sides of the debate.

**Figure 19**

*Vaxxter Website “News” Page<sup>8</sup>*

The screenshot displays the Vaxxter website's news page. It features three main news articles on the left and a sidebar on the right with navigation and utility options.

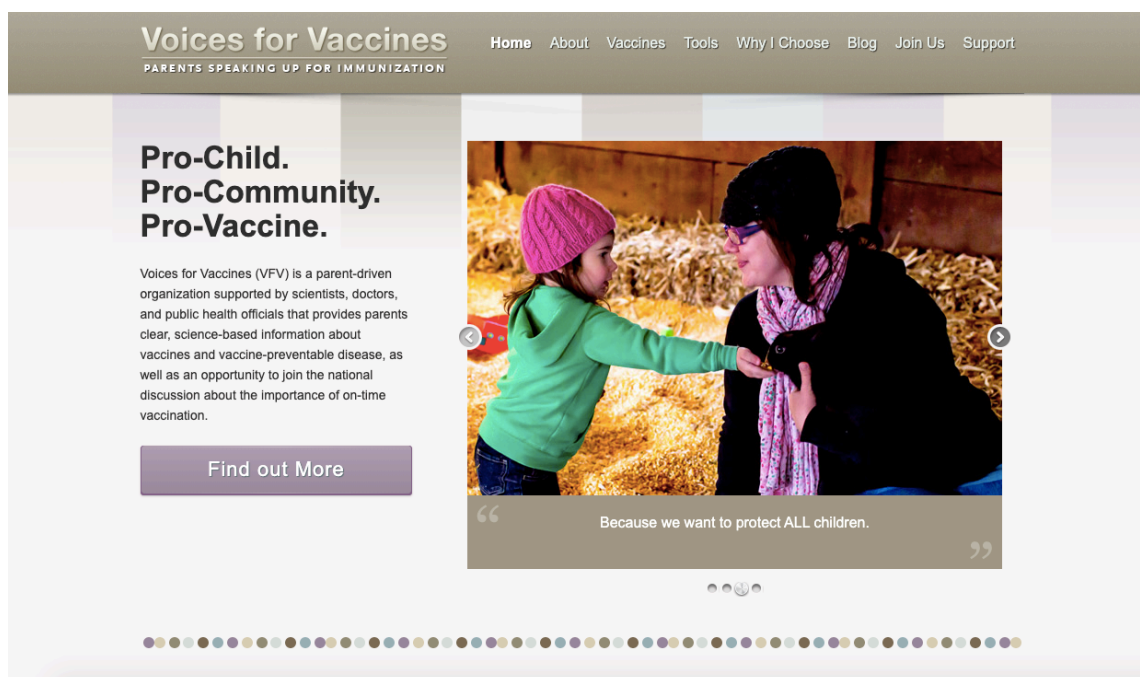
**Article 1:** **SCHOOL POLICIES, VACCINE STORIES, VAX EXEMPTIONS 01/13/2019**  
**Mississippi Vaccine Choice Activist Embroiled In Fight With State’s Laws**  
 Mary Jo Perry of Mississippi Parents for Vaccine Rights (MPVR) wants parents to know that she stands behind their medical decisions, at least in cases of vaccine choice. Her organization, MPVR, continues to push hard...

**Article 2:** **VACCINE STORIES 01/13/2019**  
**23,000 Californians Given ‘Ineffective Vaccines:’ Officials Are Trying To Track Them Down**  
 Ventura County California health officials are furiously attempting to contact the parents of children who received “ineffective vaccines.” Health officials claim the vaccines are now deemed ineffective because the county workers responsible for them stored...

**Article 3:** **FLU SHOTS, SIDE EFFECTS - VAX, VACCINE - OTHER, VACCINE STORIES 12/12/2018**  
**Healthy Vaccinated 6-Year-Old Dies From Flu In Connecticut**  
 A Connecticut couple is living the worst possible parenting nightmare. Christy Pugh and David Splan’s daughter, Emma, passed away after experiencing myocarditis, an inflammation of the heart that inhibits its ability to pump. Her death is...

**Sidebar:**  
 SEARCH ...  
 RECENT POSTS  
 SPECIAL ANNOUNCEMENT REGARDING SPRING 2020 BOOT CAMP  
 Informed Consent – Losing Ground in America  
 Ebola: Ignoring Nutrition and Playing Politics  
 Fighting for Your Rights in Ohio  
 Vaccines 2020: Big Pharma’s Admissions of Fraud  
 CATEGORIES  
 Vaccine Stories  
 TRANSLATE THIS SITE  
 Select Language  
 Powered by Google Translate  
 JOIN OUR LIST  
 By joining Vaxxter’s email list, you will receive science-

<sup>8</sup> <https://vaxxter.com/category/uncategorized/vaccine-tales/>

**Figure 20***Voices for Vaccines' Website Homepage<sup>9</sup>*

### 5.4.2. Cannabis

The online debate about cannabis consists of ODEs with pro- and anti-cannabis views. There appear to be far more ODEs that are pro-cannabis, and within the pro-cannabis group there are a range of sub-positions including those focused on: cannabis legalization (e.g., NORML and Marijuana Policy Project), medical cannabis use (e.g., American Alliance for Medical Cannabis [AAMC], CannabisMD, Society of Cannabis Clinicians), cannabis lifestyle (e.g., High Times), or government cannabis administration (e.g., Ontario Cannabis Store [OCS]). On the other hand, the anti-cannabis position have ODEs that all show concern for the various effects of cannabis use on people and society (e.g.,

<sup>9</sup> <https://www.voicesforvaccines.org/>

Citizens Against Legalizing Marijuana [CALM] and Smart Approaches to Marijuana [SAM]).

Cannabis ODEs vary in their presence. There appear to be more pro-cannabis ODEs in count, many of which seem to have a stronger presence than their anti-cannabis counterparts. Pro-cannabis ODEs like NORML and High Times have approximately 292,000 and 724,000 Twitter followers. Anti-cannabis ODEs, on the other hand like CALM and SAM have far fewer, with approximately 450 and 6,000 Twitter followers. The shared presence of cannabis ODEs varies as well. Many of the anti-cannabis ODEs link to each other's websites and have specific sections dedicated to "other organizations." As well, government cannabis administration ODEs in Canada appear to have a strong shared presence with each other, as each province runs its own cannabis distribution network.

Cannabis ODEs are owned and operated by various registrants, each with their own biases and motivations. For example, there are some ODEs that are operated by non-profit organizations such as NORML, SAM, or AAMC. Some registrants may have biases that influence the content they share on the debate. High Times lifestyle magazine is operated by High Times Holding Company, whose stocks are publicly traded and connected with other cannabis retail companies. The content they share supports a pro-cannabis lifestyle but may also promote views that encourage people to purchase the products that this company sells.

The geographic locations of cannabis ODEs reveal the distribution of the debate globally. Overall, there appears to be a concentration of English-language cannabis ODEs in North America and Europe, with even more in areas where cannabis has been legalized like Canada and areas of the United States. NORML is a cannabis legalization ODE that is located in Washington D.C. while SAM, an anti-cannabis ODE, is located nearby in Virginia. There are also clusters of ODEs that share both their vaccine position and location in various regions of the world. An example of this is in California, where there

appears to be a concentration of medical cannabis ODEs (e.g., Society of Cannabis Clinicians).

The debate around cannabis has existed online for a long time and appears to have become more prevalent over time. Older cannabis ODEs have been operating since the late 1990s (e.g., High Times, NORML), while newer ones have been created as recently as 2018 (e.g., OCS). This debate has evolved over time, with shifts seen in the discussion with regard to the forms of cannabis consumption that are promoted as well as the legal implications of cannabis laws (Månsson, 2014; Meacham, Paul, & Ramo, 2018; Mitchell, Sweitzer, Tunno, Kollins, & McClermon, 2016). As the trend of legalization in various regions of the world continues, new ODEs will likely emerge that discuss regional issues related to the use of cannabis.

Cannabis ODEs focus on various topics related to the debate depending on their position. The anti-cannabis ODEs appear to focus heavily on topics related to side-effects such as “mental illness”, “addiction”, and “lung cancer.” Within the pro-cannabis ODEs, the focus appears to have some variance depending on the sub-position of the ODE. The cannabis legalization ODEs appear to focus on topics such as “legalization”, “law”, “policy”, and “testing.” The medical cannabis use position appears to be focused on topics such as “healing”, “compassion”, as well as a variety of ailments that cannabis can treat (e.g., “sleep”, “pain”, “depression”). The cannabis lifestyle position appears to be focused on topics such as “strands”, “shop”, and “culture.” Finally, government cannabis administrators appear to be focused on topics such as product types offered (e.g., “extracts”, “flower”, “edibles”) and topics related to consumption (e.g., “methods”, “safety”, and “learn”).

The sentiments and emotions shared by cannabis ODEs also vary depending on their position in the debate. The anti-cannabis ODEs appear to evoke fear and sadness when promoting their views. SAM has a page in their website titled “The Victims of Marijuana” where they discuss real life cases of individuals who have died with relation to cannabis use. In Figure 21, an example of a reported victim of cannabis is shown. In



this example, the image used is of the deceased person’s father sitting in front of a memorial with a description in the text of how they died from heroin addiction. They connect the person’s cannabis use to them taking this harder drug. Among the pro-cannabis positions, many emotions and sentiments are shared. One pro-cannabis ODE (cannabis legalization) appears to evoke a more neutral, calm, and relaxed tone. For example, in Figure 22, the homepage for NORML is displayed. Here we see the content shared is mostly of a neutral tone, such as “NORML Responds” and “marijuana laws and penalties.” There are also phrases that appear to share a more pleasant tone, such as “help legalize” and “legalization has been a success.”

## Figure 21

SAM Website “The Victims of Marijuana” Page<sup>10</sup>



### Jeffrey Veatch

In 2008 Jeffrey’s son Justin died after snorting heroin. Although Justin’s story ended with heroin, it started with marijuana.

Justin Veatch was a talented young musician and songwriter. When he was 14, he began experimenting with marijuana. His parents were concerned and took him to a counselor to get help. Jeffrey says the counselor didn’t seem too concerned with Justin’s marijuana use and so they let it slide.

Justin’s marijuana use coincided with anxieties he did not discuss with his parents. Instead, Justin self-medicated by adding other recreational drugs and ultimately prescription opiates. He became dependent upon the opiates and when he could no longer access strong enough pills, he turned to snorting street heroin.

Jeffrey Veatch agrees that not all young people who use marijuana will end up using other drugs. But he knows in his son’s case, marijuana was the trigger that caused him to lose his inhibitions and experiment with other, more dangerous drugs that ultimately lead to his death.

Jeffrey Veatch was a network news writer for more than 40 years and won a National Writers’ Guild Award in 2007. But in September 2008, his life was forever changed when Justin died.

Today, Jeffrey Veatch goes to high schools across the country, telling his son’s story. He says that in every place he visits, he is told there are a handful of young people who have been impacted by marijuana in much the same way Justin was. Justin Veatch’s story must be heard by lawmakers considering expanded use and access to this drug as our nation is in the grips of an addiction epidemic.

To learn more about Justin’s story, visit <https://thejustinveatchfund.org>

<sup>10</sup> <https://learnaboutsam.org/victim-stories/>



Figure 22

NORML Website Homepage<sup>11</sup>

**NORML**  
Working to reform marijuana laws

Home Action Center About Marijuana State Info Legal Issues Library News Releases Blog About NORML Support

**Poll: Respondents in Adult-Use Marijuana States Say Legalization Has Been Successful**  
Read more »

Do your part to **HELP LEGALIZE MARIJUANA!** (how?)

ACT! DONATE  
LAWYERS SHOP  
JOIN NORML

**NORML NEWSLETTER**  
Sign up to receive legislative alerts, news & analysis from NORML.

**NORML Blog, Marijuana Law Reform »**  
Working to reform marijuana laws

**DEA Assisting “to the Maximum Extent Possible” in the Federal Law Enforcement Response to Nationwide Protests — NORML Responds**  
READ MORE »  
by Justin Strekal, NORML Political Director

As first reported by BuzzFeed News, the federal Drug Enforcement Administration (DEA) is expanding its law enforcement powers so that it can better assist “to the maximum extent possible in the federal law enforcement response” to the wave of ongoing, nationwide protests that have followed the killing of George Floyd by a member of the Minneapolis police.

California Marijuana Laws & Information  
Choose a different state: CA

MARIJUANA LAWS & PENALTIES »  
ARRESTS AND CROP DATA »  
CALIFORNIA NORML CHAPTERS »  
MARIJUANA LAWYERS IN CALIFORNIA »

NORML's online network

### 5.4.3. Statins

The online debate about statins consists of ODEs with pro-statin, anti-statin, and neutral views. There appears to be a greater amount of ODEs that are pro-statin than anti-statin. Within the pro-statin group, there also appears to be a variety of sub-positions including those focused on: non-profit medical organizations (e.g., American Heart Association, British Heart Association), general advocacy organizations (e.g., Take Cholesterol to Heart), and pharmaceutical companies (e.g., Kowa Pharmaceuticals). The neutral ODEs are government organizations such as Public Health Agency of Canada or the Center for Disease Control. While they do not warn of the risks of statins, like the anti-statin ODEs

<sup>11</sup> <https://norml.org/>

do, they do not recommend them as the solution for high cholesterol either. These ODEs share the latest information regarding statins for healthcare professionals. The anti-statin ODEs share the common viewpoint that taking statins for primary prevention is inadvisable. No ODEs were found that warned against statins for secondary prevention. Examples of anti-statin ODEs include Dr. Joseph Mercola, Alliance for Natural Health, The Health Examiner, and Dr. Aseem Malhorta.

The presence of statin ODEs ranges in both extremes. Non-profit medical organizations ODEs like the American Heart Association (291,000 Twitter followers) and British Heart Association (323,500 Twitter followers) have very strong presence. The American Heart Association also has many associated ODEs that focus on particular portions of the populations (e.g., Go Red For Women) or regions (e.g., Midwest and Eastern). As well, the anti-statin ODE Dr. Joseph Mercola (290,000 Twitter followers) has a strong online presence. There are also statin ODEs with very weak presence, such as the anti-statin ODE The Health Examiner (22 Twitter followers). With regard to the shared presence, there appears to be connections between ODEs of the same position. For example, on Dr. Joseph Mercola's website, there are several references, interviews, and links to another like-minded individual (Dr. Aseem Malhorta) who also operates an anti-statin ODE. All of the associated ODEs to the American Heart Association are linked on Twitter and often interact with one another. This is an indication that there is likely be a strong shared presence between these ODEs.

Statin ODEs are owned and operated by a variety of people and organizations. There are some ODEs that are operated by individual people, such as the two anti-statin ODEs run by Dr. Joseph Mercola and Dr. Aseem Malhorta. Other ODEs have registrants that are less straightforward in relation to their title that could have implications on their biases. For example, the ODE "Take Cholesterol to Heart" appears to be an independent advocacy organization for statins that is run or operated by Howie Mandel, the spokesperson for the organization. A deeper investigation into this ODE reveals that it is, in fact, owned by Kowa Pharmaceutical, a company that creates statins. This could bias

what this organization recommends to the public regarding appropriate statin use, and it is important to consider.

Statin ODEs are located in a variety of places all over the world. Many countries have their own non-profit heart health organizations that promote the use of statins. Some of these countries include the United States (the American Heart Association), Britain (British Heart Association), Canada (Heart and Stroke Foundation), and Australia (Australian Heart Foundation). As well, in these countries, there are often government run health organizations that share information on statins for the public. For example, in Canada there is the Public Health Agency of Canada while in the United States there is the Centre for Disease Control.

The discussion around statins has been ongoing since the drugs commercial release in the 1980s (Endo, 2010). As a result, the age range of statin ODEs is large, starting at the beginning of the Internet lasting until today. Some of the oldest ODEs date back to the late 1990s (British Heart Foundation in 1999, Dr. Joseph Mercola in 1997) and early 2000s (American Heart Foundation in 2003). On the other hand, some statin ODEs have been created very recently, such as the advocacy group Take Cholesterol to Heart in 2017 or The Health Examiner in January 2020. There appears to be no slowdown in the discussion of the appropriate role of statins in cholesterol management.

Statin ODEs focus on various topics related to the debate, which vary depending on their position. Anti-statin ODEs appear to focus heavily on side effects (e.g., “neurological”, “muscle pain”, “liver damage”) and alternate methods for managing cholesterol (e.g., “dietary approach”, “omega-3 fatty acids”, “garlic”). Neutral ODEs appear to mostly discuss resources that people can use to find out further information on statins (e.g., “healthcare professional”, “drug product database”, “FDA”). Pro-Statin ODEs appear to be focused on the benefits of taking statins (e.g., “reduce risk”, “fight inflammation”) and, associated with this, the risks of having high cholesterol (e.g., “narrowed arteries”, “heart attack”, “stroke”).

The sentiments and emotions shared by statin ODEs vary depending on their position in the debate. At one end of the debate, the anti-statin ODEs appear to evoke fear when it comes to statins. For example, displayed in Figure 23, an article published on Dr. Joseph Mercola's website discusses the risks of taking statins associated with mental health issues. In the article, they discuss risks such as "depression", "anxiety", and "suicide." As well, the article further invokes fear in that it claims that "the scientific community is marked by a significant lack of interest in investigating the effects on personality." On the other end of the debate, the pro-statin ODEs tend to invoke positive emotions like joy when discussing statins and fear when discussing the implications of high cholesterol. For example, in Figure 24 & 25, on the homepage of the ODE Take Cholesterol to Heart, there is a very pleasant picture of comedian Howie Mandel with text about how he likes inspiring individuals to take statins and that it is better than getting a laugh. On another page, discussing the risks of high cholesterol, there is another image of Mandel, this time looking concerned with text about how many people have high cholesterol, and how it can be an invisible problem until the consequences are too late to manage.

## Figure 23

*Dr. Joseph Mercola Website Article on Statins<sup>12</sup>*

**MERCOLA**  
Take Control of Your Health

All Find Answers to Your Health Questions

Sign in Join

**STORY AT-A-GLANCE**

- > Statin drugs are associated with an increased risk of depression and anxiety. Researchers also find they increase acts of aggression and violence, raise the risk of suicide and compromise cognition
- > The number of people taking the drugs increased to 35 million in the U.S. following new recommendations in 2016
- > The scientific community is marked by a significant lack of interest in investigating the effects on personality
- > Other known side effects include musculoskeletal disorders, atherosclerosis, cataracts, neurodegenerative diseases and osteoporosis
- > Your protection may be increased by evaluating your risk of heart disease, understanding your cholesterol numbers and analyzing your iron levels and overall diet

Data from the CDC<sup>1</sup> in 2017 show heart disease causes one death every 37 seconds in America and that it is the leading cause of death in the U.S. It created a financial burden of \$219 billion in 2014 and 2015. Every 40 seconds someone has a heart attack. Those at higher risk are smokers and those who have high blood pressure, high blood cholesterol and/or diabetes.<sup>2</sup>

Top

## Figure 24

*Take Cholesterol to Heart Website Homepage<sup>13</sup>*

Take Cholesterol to Heart

Cholesterol Facts All Statins Are Not the Same Discover Your Statin Status ACTION Survey Resources Press Room Talking With Your Doctor Learn About Howie's Journey

Supporting patients managing high cholesterol during the COVID-19 pandemic  
See more →

**“If I can inspire people to talk to their doctor about their statin, that’s better than getting a laugh.”**

See actor and comedian Howie Mandel share his experience with high cholesterol and being open with his doctor.

<sup>12</sup> <https://articles.mercola.com/sites/articles/archive/2020/01/29/statin-devastating-effects-on-brain.aspx>

<sup>13</sup> <https://www.takecholesteroltoheart.com/>

**Figure 25**

*Take Cholesterol to Heart Website, Cholesterol Facts Page<sup>14</sup>*

**“When I have a cold, I know my chest is congested. With cholesterol, I don’t know if my arteries are congested. And that’s more dangerous.”**

**High cholesterol affects more than 95 million Americans.**

Cholesterol is a waxy, fat-like substance that can be found in every cell of the body. Some cholesterol is produced by your liver, while other cholesterol comes from the food you eat. Cholesterol moves through your bloodstream and is made by the body in order to keep cells healthy. Cholesterol is also needed for the production of hormones, vitamin D, and other substances.

**Why is high cholesterol considered dangerous?**

#### 5.4.4. Dieting Plans

The online debate about dieting plans is made up of many different positions. While there are some anti-dieting plans out there, the bulk of the debate appears to be among pro-dieting plans, specifically about which plans are best. Within the pro-diet group there are a range of sub-positions that promote specific diets such as ketogenic diet (e.g., The Charlie Foundation, Matthew’s Friends, Keto Resources and My Keto Kitchen), vegan diet (e.g., Vegan, The Vegan Society, Vegan Action, The Vegetarian Resource Group, and Veg Source), south beach diet (e.g., South Beach Diet, South Beach Diet 101, and South Beach Diet Club), and Atkins diet (e.g., Atkins and Atkins Nutritionals). On the other side of the debate, the anti-diet position is aligned in their view that the problem isn’t about getting people to eat healthier, but rather that society needs to accept people for whatever their body types and eating habits are (e.g., Association for Size Diversity

<sup>14</sup> <https://www.takecholesteroltoheart.com/cholesterol-heart-disease>

and Health, National Association to Advance Fat Acceptance, No Lose, and Fat!So?). In promoting body acceptance, these ODEs suggest that the risks of dieting plans (physical and mental) is more harmful on society than accepting obesity and different body types.

With regard to presence, it appears that the pro-dieting plan ODEs have a stronger presence than the anti-dieting plan ODEs. Most of the anti-dieting plan ODEs have weak presence when compared to the pro-dieting plan ODEs. For example, the Association for Size Diversity and Health is one of the most present anti-dieting plan ODE, with 3,345 Twitter followers. One of the most present pro-dieting plan ODEs, on the other hand, is Vegan with over 250,000 Twitter followers. Among the pro-dieting plan positions and their ODEs, there is quite a lot of variability in their presence. Continuing with the vegan diet as an example, it appears that there are numerous ODEs with a strong presence. Aside from the one just mentioned, the Vegan Society is another example (230,100 Twitter followers). However, there are also many ODEs that promote this position that have a weaker presence. One example of this is the ODE Vegan Action, with 2,691 Twitter followers. With regard to shared presence, there appears to be connections between ODEs of the same position or between ODEs of similar diet plans. For example, in The Charlie Foundation website, there is a section where they list friends of the organization. In this section, other keto diet ODEs, including Matthew's Friends, are listed. This would indicate a strong shared presence between these organizations. The keto diet and south beach diet are very similar to one another in terms of their dieting plans and therefore the ODEs promoting these positions likely share a strong online presence. On the ODE South Beach Diet, for example, there are numerous references and links to keto diet material highlighting this.

Dieting plan ODEs are owned and operated by a variety of people and organizations. ODEs with a pro-dieting plan position and strong presence, such as Atkins and South Beach Diet, are often owned by organizations that created the diet and sell the related publications, materials, and food. Most dieting plan ODEs that are not directly connected to a business, such as the Association for Size Diversity and Health and National

Association to Advance Fat Acceptance, are operated by non-profit organizations of the same name.

Dieting plan ODEs are located throughout the world, but there appears to be a large concentration of the English language ones in North America. ODEs such as Atkins (Denver, Colorado), South Beach Diet (Washington, Pennsylvania), and Matthews Friends (Surrey, BC) are all located in this region. While there are examples of ODEs located outside of North America, such as The Vegan Society (Birmingham, England), it appears to be less likely for this to be the case.

The age of dieting plan ODEs is wide ranging. There have been ODEs for various dieting plans since the late 1990s (e.g., Atkins, The Vegan Society). However, specific dieting plans have come into existence or become more popular at various periods of time. By examining the age of the ODEs associated, the evolution of the debate about the best dieting plans is clearer. The South Beach Diet, which was created in 2003, highlights this. The official ODE of this dieting plan's creator (South Beach Diet) launched in 2003. Another unassociated but also popular ODE of the same position was then launched in 2004 (South Beach Diet 101) and as the diet became more popular, more ODEs were developed. The diet is still popular to this day, with new ODEs of this position created as recently as 2018 (South Beach Diet Club).

The focus of dieting plan ODEs varies among the different position. Anti-dieting plan ODEs tend to focus on topics such as "fat acceptance" or "healthy at every size." Pro-dieting plan ODEs have different focuses depending on the dieting plan requirements. For example, Atkins diet, ketogenic diet, and south beach diet ODEs all appear to focus heavily on "carbohydrates", "low carb", or "sugar levels." This is because these diets promote reducing the amount of carbohydrates one eats and replacing it with options that will promote weight loss. However, within these positions there appear to be some variability between the ODEs. For example, the ketogenic diet has some ODEs that are focused almost entirely on the topic of "weight loss" or "body shaping (e.g., Keto Resources, My Keto Kitchen). Other ketogenic diet ODEs (e.g., The Charlie Foundation)



promote the keto diet for its applications in helping individuals with epilepsy and focus on different topics like “metabolic disorders” or “neurological conditions.” Vegan diet ODEs tend to have a very different focus from the other dieting plans mentioned. While the health benefits of the dieting plan are often mentioned, there is a strong focus typically on justice for animals and to combat climate change.

The sentiment of dieting plan ODEs varies for each position. For example, the anti-dieting plan ODEs appear to contain a mixture of negative sentiments (regarding dieting plans and the culture surrounding them) and positive sentiment (accepting people as they are). Figure 26 shows the Healthy at Every Size webpage of the Association for Size Diversity and Health. On this page, positive sentiments are shared with focuses on things such as “weight inclusivity” and “eating for well-being.” At the same time, there appear to be negative sentiments shared when referring to society’s status for health of individual such as “judge” and “oppress.” Pro-dieting plan ODEs appear to share much more positive sentiments in general. These ODEs try to give a sense of encouragement and often share personal stories of weight loss or body change that are meant to encourage other individuals to feel positive about the possibility of losing weight or becoming healthy as well. Vegan diet ODEs appear to share different sentiments than the other diets mentioned. For example, Figure 27 shows the lifestyle page of the ODE Vegan, which is very positive. There is a happy picture of an individual and a dog and in the text, they refer to choosing a vegan diet as a “compassionate choice” to protect animals. In another section of the website, however (Figure 28), negative sentiments (e.g., anger and sadness) are shared on animal cruelty. These sentiments are generated from the images of chickens being mass killed by foam and the cruelty of the meat industry is described in the text with words such as “wasting” birds or referring to the foam as “chicken killing foam.”

Figure 26

Association for Size Diversity and Health Website, *Healthy at Every Size* Page<sup>15</sup>

The screenshot shows the website for the Association for Size Diversity and Health (ASDAH). The header includes the ASDAH logo and navigation links: ABOUT, BLOG, MEDIA, and CONTACT. A left sidebar contains a menu with links to Home, HAES® Approach, Blog, Resources, Projects, and Membership, along with buttons for 'JOIN ASDAH' and 'FIND HAES EXPERT'. The main content area features a blue banner with the title 'HAES® Principles' and a 'Printer Friendly' icon. Below the banner, the text reads: 'The Health At Every Size® Approach'. The main text explains that ASDAH affirms a holistic definition of health, not just the absence of illness, and that health exists on a continuum. It also describes the HAES approach as one that addresses broad forces supporting health and honors social connections. At the bottom, a list of four principles is provided:

- 1. Weight Inclusivity:** Accept and respect the inherent diversity of body shapes and sizes and reject the idealizing or pathologizing of specific weights.
- 2. Health Enhancement:** Support health policies that improve and equalize access to information and services, and personal practices that improve human well-being, including attention to individual physical, economic, social, spiritual, emotional, and other needs.
- 3. Respectful Care:** Acknowledge our biases, and work to end weight discrimination, weight stigma, and weight bias. Provide information and services from an understanding that socio-economic status, race, gender, sexual orientation, age, and other identities impact weight stigma, and support environments that address these inequities.
- 4. Eating for Well-being:** Promote flexible, individualized eating based on hunger, satiety, nutritional needs, and pleasure, rather than any externally regulated eating plan focused on

<sup>15</sup> <https://www.sizediversityandhealth.org/>

Figure 27

*Vegan Website, Lifestyle Page<sup>16</sup>*

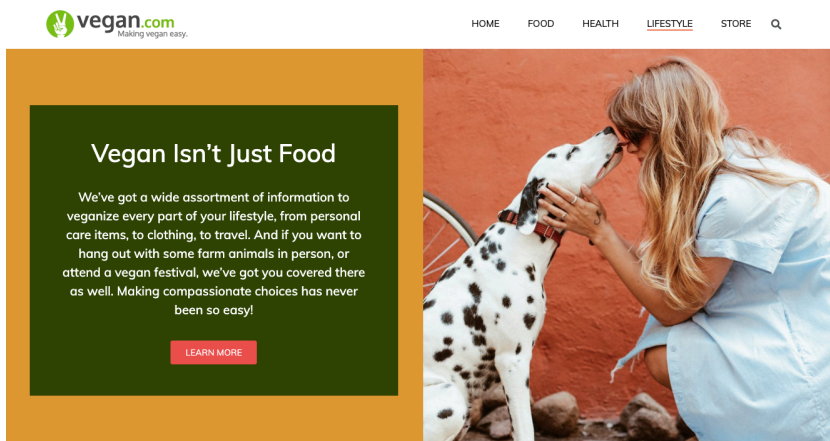
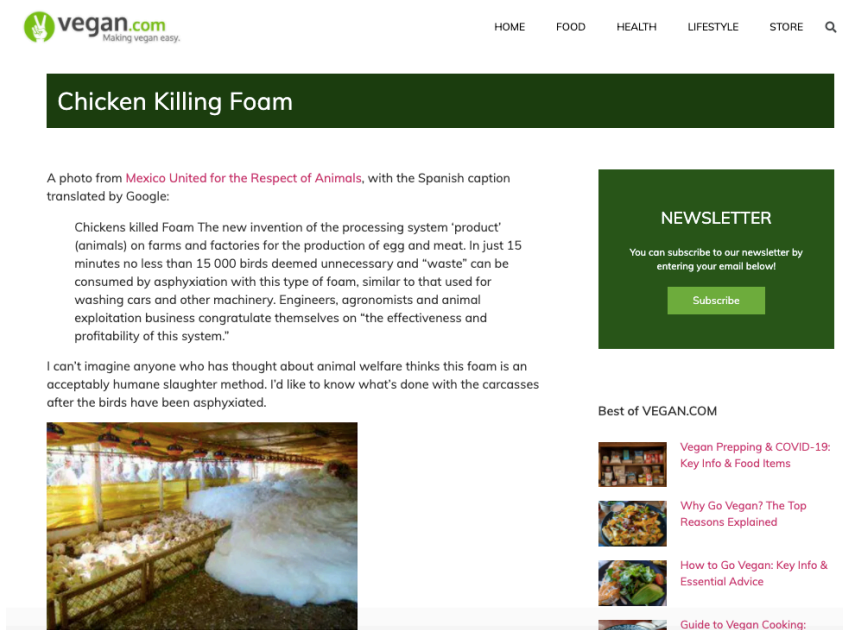


Figure 28

*Vegan Website, Chicken Killing Foam Page<sup>17</sup>*



<sup>16</sup> <https://www.vegan.com/lifestyle/>

<sup>17</sup> <https://www.vegan.com/chicken-killing-foam/>

## 5.5. ODIN-Based Design of Visual Analytics Systems

Making sense of online debates requires stakeholders to complete tasks that involve investigating the attributes outlined in ODIN, presented in Section 5.3. In this section, we discuss how ODIN-based VASes can enable stakeholders to complete these tasks to make sense of ODE attributes. Using four scenarios, one for each of the online public health debates described in Section 5.4., we demonstrate how ODIN can be incorporated into the design of VASes that facilitate making sense of these online public health debates. To demonstrate this, in this section, we present four ODIN-based VASes: an existing tool (VINCENT) and three mock-ups.

The existing ODIN-based VAS, VINCENT, was developed to help stakeholders make sense of the online vaccine debate (Ninkov & Sedig, 2019, 2020). The design and empirical evaluation of VINCENT suggest that such systems can be very useful. A study of VINCENT showed that participants who used the system to complete tasks were quicker, more confident in their abilities, and more accurate in their responses than participants who did not have the system. Therefore, we use the same framework-based approach to design the other VASes.

The other three ODIN-based VASes are mock-ups of potential tools for making sense of the online statin, cannabis and diet debates. It is important to note that these three VASes and scenarios are not based on real data but are rather conceptualizations of what the data could be. We use these conceptualizations to demonstrate how the integration of these attributes in VASes could help stakeholders make sense of the debates.

### 5.5.1. Vaccines

**Figure 29**

*VINCENT, a VAS for Vaccines<sup>18</sup>*

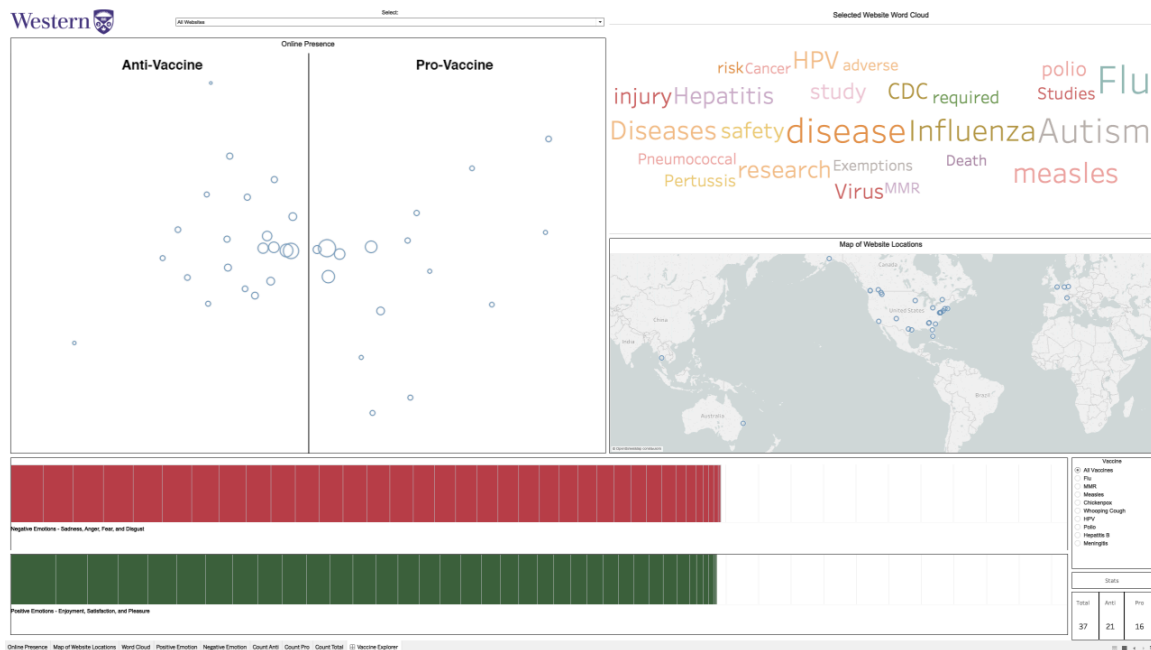


Figure 29 shows VINCENT, a VAS developed to help stakeholders make sense of the online debate about vaccines (Ninkov & Sedig, 2019, 2020). In VINCENT, vaccine ODEs have been identified, and attribute data including presence, shared presence, focus, geographic location, and sentiment have been integrated. VINCENT is composed of four main components: 1) the Presence Map (top left), 2) the Word Cloud (top right), 3) the Map of Website Locations (middle right), and 4) the Emotion Bar Charts (bottom). The Presence Map is a representation of the hyperlink or follower data analyzed from each website. Shared presence is determined by analyzing co-occurrences (links or followers) between ODEs using Multi-Dimensional Scaling (MDS) (Thelwall & Zuccala, 2008; Vaughan & Ninkov, 2018; Vaughan & You, 2008). With this method of analysis, ODEs'

<sup>18</sup> See Chapters 3 & 4

inlinks or followers are analyzed with a co-occurrence analysis and visualized using MDS. This technique plots each ODE on the map as a circle positioned in proximity to other circles based on the similarity of the ODEs' presence (i.e., how much shared presence they have). If the circles are plotted closer together, then the ODEs they represent share more presence, and vice versa. Each ODE's circle on this map has been sized to reflect its individual presence. The bigger the circle, the greater presence it has. The Word Cloud is a representation of the 25 most common unique words that are related to the vaccine debate from each website. Words are sized based on the frequency with which they appeared on the website or group of websites. The bigger a word in the Word Cloud is, the more frequently it is used on the website. The Map of Website Locations shows a representation of the locations of each website on a world map. Similar to the online Presence Map, the Map of Website Locations uses circles to encode each website, but the circles have all been sized equally to help the user see the location of each website clearer. The Emotion Bar Charts represent positive and negative emotions for a selection of each website's text about a set of vaccines, selected on the right side of the bar charts. The Emotion Bar Charts represent the negative (red) and positive (green) emotions detected by IBM's Natural Language Processing API (Grimes, 2016). Each bar is made up of all the smaller rectangles (ODEs). The bar represents the overall detected emotion in the text of the complete set of ODEs. The width of each rectangle within each bar chart represents the degree of detected emotion in that ODE's content. The wider the rectangle, the more the emotion is detected.

**Scenario:** A public health stakeholder suspects that there is a lot of anti-vaccine sentiment coming from North Western North America. They want to investigate this issue by identifying which part of North America has the highest concentration of anti-vaccine ODEs. If the North West does have the highest concentration of anti-vaccine ODEs, the stakeholder wants to know which vaccines, in particular, have the strongest negative emotions associated with them.

The stakeholder uses the VAS by first navigating to the Map of Website Locations and selecting the ODEs in the region they want to investigate. After highlighting the ODEs for each of the various regions of interest (North West, North East, South East, Mid West, South West) and checking the position of the ODEs, they confirm that there are the highest concentration of anti-vaccine ODEs in the north west. The stakeholder then selects each of the listed vaccines and compares the Emotion Bar Charts. They find that the ODEs had the strongest negative emotions associated with the measles, mumps, and rubella vaccine. As the stakeholder continues to investigate the ODEs using the Word Cloud, they quickly find that some of the anti-vaccine ODEs in the north west, like Vaccination News and Vaccine Choice Canada, have a strong focus on the issue of autism. Given the popular belief among anti-vaccine groups that the MMR vaccine is linked to autism (Nicholson & Leask, 2012), the prevalence of this sub-position is notable.

5.5.2. Statins

Figure 30

VAS for Statins

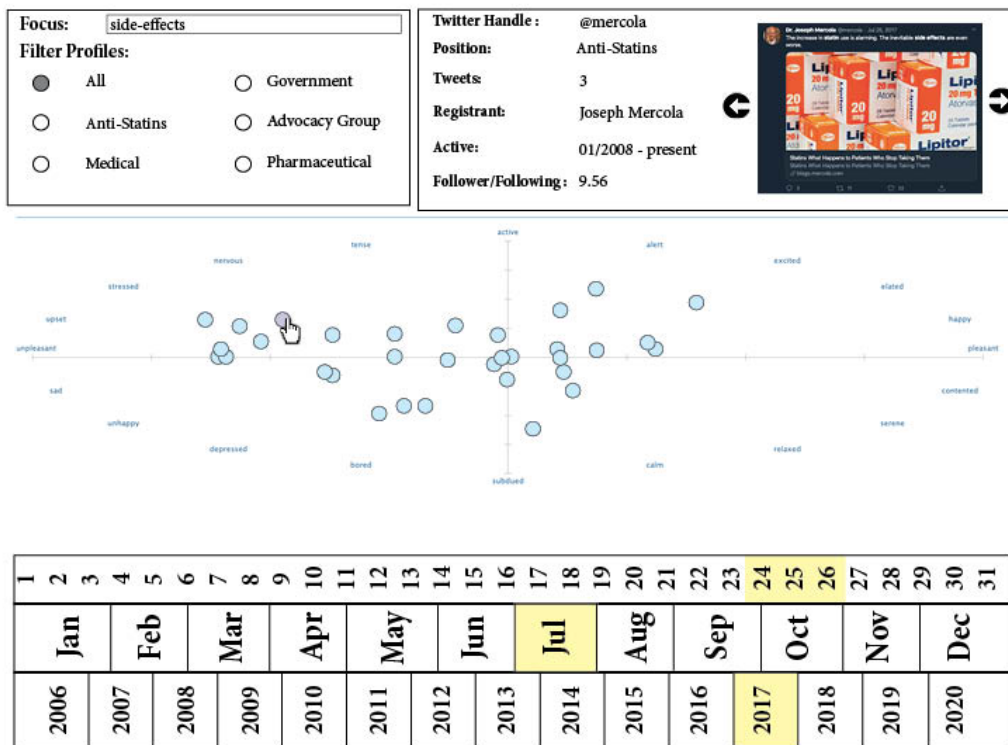


Figure 30 shows a VAS envisioned to help stakeholders make sense of the online debate about statins. In this VAS, statin ODEs from Twitter have been identified and attribute data including age, position, registrant, presence, focus, and sentiments have been integrated. The VAS consists of four components: 1) the Focus Panel (top left), 2) the ODE Information Box (top right), 3) the Sentiment Map (middle), and 4) the timeline (bottom). The Focus Panel is made up of the focus input, where a user can specify words or phrases that they want to know more about. As well, in this panel the stakeholders can filter the ODEs to display a specified position. The ODE information box reveals information on individual ODEs as they are selected by the stakeholders in the Sentiment



Map. This information includes the ODE's Twitter handle, position, number of tweets, registrant, years active, and presence (follower/following ratio). The Sentiment Map displays all the ODEs that match the criteria specified in the Focus Panel within the time frame selected on the timeline. Each ODE is represented on the map with a circle. The placement of the circle is based on a sentiment analysis of tweets containing the specified focus (this method of sentiment analysis has been adopted from (Healey & Ramaswamy, 2011)). The timeline shows the years and months since the beginning of Twitter (2006) to the present. The yellow highlighted sections display the selected time period by the stakeholders.

Stakeholder investigates the debate using this VAS by first specifying the period they are interested in investigating in the timeline. They then choose the focus of the debate they wish to know more about and the position(s) of the ODEs that they want to include in their search. With the specifications made, the ODEs are scanned and those that match the criteria are then plotted onto the Sentiment Map based on the aggregation of the sentiment of their focus-related tweets. Stakeholders identify specific ODEs on the Sentiment Map for which they want to learn more about. By choosing a circle of an ODE on the map, they reveal its content in the ODE information box. When new data points are chosen or the date range is changed, the data in the system updates accordingly and the user can begin examining the new data presented.

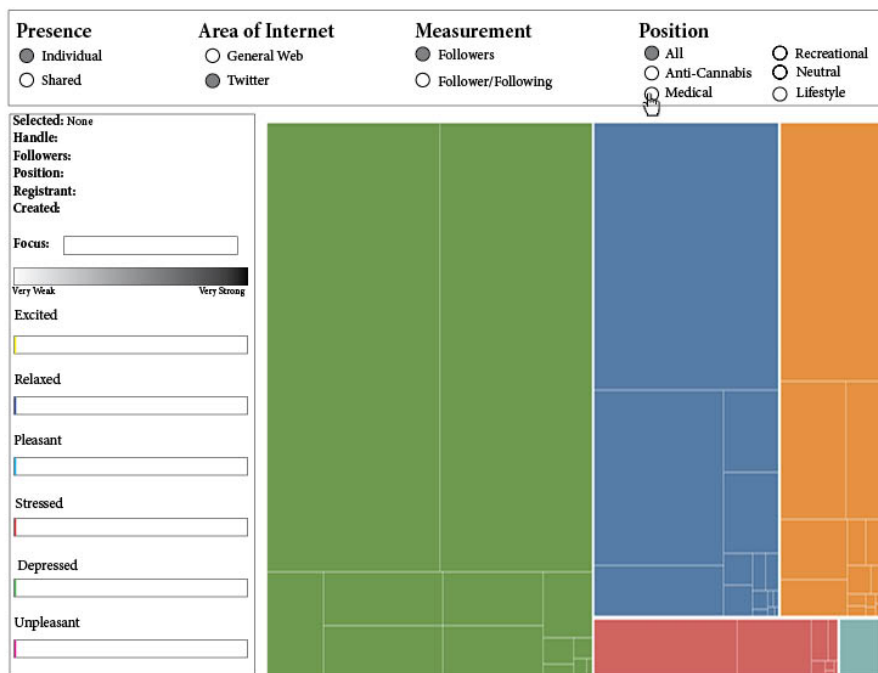
**Scenario:** A public health stakeholder wants to make sense of what happened in the summer of 2017, when an event occurred that the stakeholder hypothesizes may have led to an uptick in hesitancy of patients choosing to take statins based on the fear of side-effects. They use the VAS in Figure 30 to help them make sense of the debate as it occurred on Twitter. They first specify the period of time they want to investigate (July 24-26, 2017) and the focus ("side-effects"). The VAS then processes the data from the ODEs active during that time and outputs a Sentiment Map based on the content they shared related to the focus. The stakeholder then selects the data points of interest on the Sentiment Map and reveals further information about each ODE.

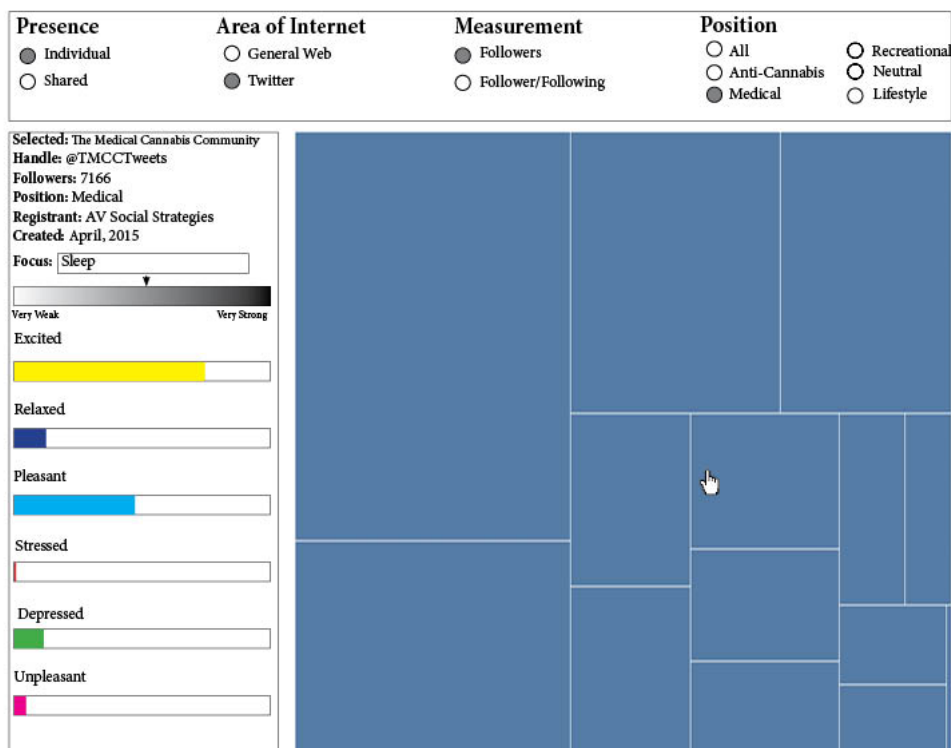
The stakeholder has come across a particular ODE (Dr. Joseph Mercola) that they want to know more about. The tweets analyzed have a nervous and unpleasant sentiment towards the focus of “side-effects.” By selecting the ODE, the stakeholder can see that its position is anti-statin, that it has been active on Twitter since 2008 and that it has a strong presence, as indicated by the follower/following ratio. In the time period they have selected, they published 3 tweets that contained reference to “side-effects.” They can also scroll through these tweets in the information box and evaluate them. By repeating this process through selecting other ODEs on the Sentiment Map, or by specifying another time period they are interested in, the stakeholder can start to make sense of how the debate was shaped on Twitter during this event.

### 5.5.3. Cannabis

**Figure 31**

*VAS for Cannabis*



**Figure 32***VAS for Cannabis with Selections Made*

Figures 31 and 32 shows a VAS envisioned to help stakeholders make sense of the online debate about cannabis. In this VAS, cannabis ODEs have been identified and attribute data including presence, shared presence, focus, sentiment, age and registrant have been integrated. The VAS consists of three components: 1) the Selection Panel (top), 2) the Presence Map (bottom right), and 3) the Focus and Information Panel (bottom left). The Selection Panel allows the user to control how presence is evaluated, including the type (shared or individual), the area (general web or twitter), and the measurement method. As well, users can choose to filter the ODEs in the VAS using this panel. In the Presence Map, based on the presence type selected, either a tree map of individual presence (Figures 31 & 32) or a (MDS) map of shared presence (Figure 33 in Section 5.5.4.) is

populated. The VAS is set to individual presence for the example in this section and the tree map is displayed. The bigger the cell of the tree map, the more presence the ODE it represents has. Finally, the Focus and Information Panel allow the user to drill into specific ODEs by selecting on them from the Presence Map. When an ODE is selected, its name, registrant, age, and position become immediately available. Stakeholder can then investigate the ODE's sentiment regarding a particular focus by choosing a word or phrase. The system then analyzes the related content from the ODE and returns the emotions detected, as well as the strength of the ODE's focus on the word or phrase.

Stakeholders investigate the debate with this VAS by first selecting the level that they want to analyze presence (individual or shared). Next, they choose the area of the Internet they want to analyze (general web or Twitter) and then set the measurement technique they want the system to use. As well, in this panel, they determine whether or not they want to filter the ODEs and if so, how. The resulting tree map shows the presence of the ODEs based on the specifications. The cells of the tree map are organized, and color coordinated based on their position. Stakeholders can investigate the debate by selecting specific cells to reveal ODE information and repeating this process with all the relevant ODEs to their investigation.

**Scenario:** A public health stakeholder has come across claims on social media that suggest using cannabis can help individuals who struggle with sleep. The stakeholder has investigated the literature on the subject and have found some evidence to support this claim. They now need to investigate if this claim is shared widely on Twitter, and specifically among the medical cannabis ODEs. They use the VAS to examine the online debate on Twitter. In Figure 31, the VAS has been set for individual presence on Twitter using follower/following metrics without any filter of position or ODE selected. The stakeholder continues their investigation by filtering the VAS so only ODEs with a medical position are shown (Figure 32). They also investigate by selecting the individual cells to reveal further information about the ODE.

In Figure 32, the ODE that has been selected is one with a moderately strong presence called “The Medical Cannabis Community.” The stakeholder can quickly see in the information panel on the left the ODE’s Twitter handle, presence, position, registrant, and age. They can also input “sleep” into the focus and see the sentiments they have on the topic. In this case, they have found that the ODE has a moderately strong focus on the topic. As well, they can see that overall, the tweets from the ODE with this focus have an excited and pleasant sentiment. By repeating this process with other ODEs of interest, both with the same position or not, the stakeholder can quickly see how the discussion of sleep in the online cannabis debate manifests itself.

### 5.5.4. Dieting Plans

**Figure 33**

*VAS for Dieting Plans*



Figure 33 shows a VAS envisioned to help stakeholders make sense of the online debate about dieting plans. In this VAS, dieting plan ODEs have been identified and the attribute data described in Section 5.5.3. have been integrated in a similar fashion. The VAS has the three main components: 1) Selection Panel (top), 2) the Presence Map (bottom right) and 3) the information and Focus Panel (bottom left). In this example, the VAS has been set to analyze shared rather than individual presence. The shared presence is analyzed using co-link analysis and MDS, the same method described in Section 5.5.1. To make the information easier to interpret, the position of the ODEs has also been integrated into the MDS map by colouring the circles according to position, and then highlighting the areas on the map that contain ODEs with those positions.

Stakeholders investigate the debate using this VAS by first selecting (in the Selection Panel) the area of the Internet, measurement level of co-occurrence analysis (page or domain), and position they want to investigate. For co-occurrence analysis, depending on the level of measurement the stakeholders select (Vaughan & Ninkov, 2018), the ODEs' inlinks are analyzed and plotted accordingly in the Presence Map. As in Figure 33, stakeholders can investigate the ODEs by selecting circles on the map to reveal details in the Focus and Information Panel.

**Scenario:** A public health stakeholder is comparing diet positions to one another. Specifically, they are interested in identifying 1) diet communities that have greater similarity to each other and 2) how those diet communities view one another. They use the VAS to investigate the online debate on the general web. In Figure 33, the VAS has been set for shared presence on the general web using page-level co-link analysis without any filter of position. The stakeholder begins to select individual cells of interest to reveal further information about the ODE in the Focus and Information Panel. They notice that South Beach Diet and Keto Diet are close to one another on the map, which makes sense since both dieting plans promote reducing carbs as an important component of the diet.

In Figure 33, the stakeholder selects the ODE “Keto Resource.” The stakeholder can quickly identify, in the Focus and Information Panel, the ODE’s Twitter handle,

presence, position, registrant, and age. In this case, the stakeholder is interested in the “South Beach.” The stakeholder finds that the selected ODE has a somewhat strong focus on the phrase, indicating the dieting plans are likely associated with one another. As well, they can see that, overall, the related content on the website has an excited and pleasant sentiment. By repeating this process with other ODEs of interest, both with the same position or not, the stakeholder can quickly see how the discussion of sleep in the online diet debate manifests itself.

## 5.6. Discussion

In this paper, we have presented a framework for analyzing ODEs called ODIN. This framework consists of seven ODE attributes, discussed in detail in Section 5.3. For each attribute, we have described the types of data collection and analysis methods that can be used. We then demonstrated how the framework can be applied in the analysis of four online public health debates (vaccines, cannabis, statins, & dieting plans).

We used the analyses to guide the development of four framework-based VASes. These VASes were designed to help stakeholders make sense of online public health debates. Online debates are complex information spaces and making sense of them can be challenging. Without the assistance of VASes, it is difficult, or in some cases impossible, to quickly assess the various attributes of these debates. There are many ways in which ODIN-based VASes can implement attributes from the framework, as we have discussed in this paper. One of these VASes presented was VINCENT, a VAS that has already been developed and shown to be very useful for stakeholders (Ninkov & Sedig, 2020). As well, we described three other mock-ups of ODIN-based VASes to show other potential applications of the framework for analyzing online public health debate.

The framework presented in this paper is meant to be generalizable so that it can apply to other debates, either related to public health or in other areas (e.g., political debates). As long as the discussion of interest consists of multiple opposing positions, has ODEs that

can be identified, and has data about the attributes that can be analyzed, the framework can be implemented.

### **5.6.1. Limitations & Future Research**

While we have included seven ODE attributes in ODIN that we consider to be important for making sense of online debates, this list can likely be expanded with other ODE attributes. Examining other online debates (both about public health topics or other topical areas of debate) and including a large range of social media platforms such as Facebook or Reddit would likely result in the inclusion of other attributes.

We have focused our discussion of online debates with focus on two specific manifestations of the Internet: the general web and Twitter. However, there are many other social media platforms (e.g., Facebook, Instagram, Reddit) in which these debates also occur. Each area of the web has unique qualities and uses different labels and structures for the various attributes described in ODIN. While we have not discussed each of these social media platforms in detail individually, the attributes described are generalizable so that they can be applied or adapted to these other social media websites. The application of such attributes to such alternate social media websites could be a fruitful area of investigation.

It is also important to note that the data that researchers can access is constantly changing. Considerations of what the best data sources are, as well as how data may be manipulated is important. For example, Twitter follower data can be manipulated through users buying followers (Aggarwal, Kumar, Bhargava, & Kumaraguru, 2018). While methods exist (e.g., follower/following ratio) that attempt to expose those that artificially boost their social media statistics, improved approaches for identification of profiles with inflated numbers are needed to gauge a more accurate view of presence.

Finally, there is also a need for the examination of alternate methods of data analytics, data visualizations, and human-data interactions for ODIN-based VASes. For example, in this paper we have discussed and implemented MDS for integrating shared online



presence in VASes. However, other methods of social network analysis could have been utilized. As more VASes are created to meet the challenges of making sense of online public health debates, best practices for their design can be catalogued.

## 5.7. References

- Abbasi, J. (2018). Interest in the ketogenic diet grows for weight loss and type 2 diabetes. *Jama*, 319(3), 215–217.
- Abramson, J. D., Rosenberg, H. G., Jewell, N., & Wright, J. M. (2013). Should people at low risk of cardiovascular disease take a statin? *Bmj*, 347, f6123.
- Aggarwal, A., Kumar, S., Bhargava, K., & Kumaraguru, P. (2018). The follower count fallacy: detecting twitter users with manipulated follower count. In *Proceedings of the 33rd Annual ACM Symposium on Applied Computing* (pp. 1748–1755).
- Andromalos, L. (2018). The Paleo Diet. In *Clinical Guide to Popular Diets* (pp. 71–86). CRC Press.
- Anger, I., & Kittl, C. (2011). Measuring influence on Twitter. In *Proceedings of the 11th international conference on knowledge management and knowledge technologies* (pp. 1–4).
- Ariagno, M. (2018). The South Beach Diet. In *Clinical Guide to Popular Diets* (pp. 87–98). CRC Press.
- Austgulen, M. H. (2014). Environmentally sustainable meat consumption: An analysis of the Norwegian public debate. *Journal of Consumer Policy*, 37(1), 45–66.
- Baka, A. B. A., & Leyni, N. (2017). Webometric study of world class universities websites. *Qualitative and Quantitative Methods in Libraries*, 105–115.
- Barnaghi, P., Ghaffari, P., & Breslin, J. G. (2016). Opinion mining and sentiment polarity on twitter and correlation between events and sentiment. In *2016 IEEE Second International Conference on Big Data Computing Service and Applications (BigDataService)* (pp. 52–57). IEEE.
- Barnard, N. D., Cohen, J., Jenkins, D. J. A., Turner-McGrievy, G., Gloede, L., Jaster, B., ... Talpers, S. (2006). A low-fat vegan diet improves glycemic control and cardiovascular risk factors in a randomized clinical trial in individuals with type 2 diabetes. *Diabetes Care*, 29(8), 1777–1783.

- Bilgri, O. R. (2016). From “herbal highs” to the “heroin of cannabis”: Exploring the evolving discourse on synthetic cannabinoid use in a Norwegian Internet drug forum. *International Journal of Drug Policy*, 29, 1–8.
- Bloch, J. P. (2007). Cyber wars: catholics for a free choice and the online abortion debate. *Review of Religious Research*, 165–186.
- Borgmann, H., Woelm, J.-H., Merseburger, A., Nestler, T., Salem, J., Brandt, M. P., ... Loeb, S. (2016). Qualitative Twitter analysis of participants, tweet strategies, and tweet content at a major urologic conference. *Canadian Urological Association Journal*, 10(1–2), 39.
- Börner, K. (2015). *Atlas of Knowledge: Anyone Can Map*. The MIT Press.
- Bradley, C. K., Wang, T. Y., Li, S., Robinson, J. G., Roger, V. L., Goldberg, A. C., ... Peterson, E. D. (2019). Patient-reported reasons for declining or discontinuing statin therapy: insights from the PALM registry. *Journal of the American Heart Association*, 8(7), e011765.
- Brumshteyn, Y. M., & Vas'kovskii, E. Y. (2017). Analysis of the webometric indicators of the main websites that aggregate multithematic scientific information. *Automatic Documentation and Mathematical Linguistics*, 51(6), 250–265.
- Buchel, O., & Sedig, K. (2016). From data-centered to activity-centered geospatial visualizations. *Geospatial Research: Concepts, Methodologies, Tools, and Applications: Concepts, Methodologies, Tools, and Applications*, 1, 246.
- Bueno, N. B., de Melo, I. S. V., de Oliveira, S. L., & da Rocha Ataide, T. (2013). Very-low-carbohydrate ketogenic diet v. low-fat diet for long-term weight loss: a meta-analysis of randomised controlled trials. *British Journal of Nutrition*, 110(7), 1178–1187.
- Cetin, O., Ganan, C., Korczynski, M., & van Eeten, M. (2017). Make notifications great again: learning how to notify in the age of large-scale vulnerability scanning. In *Workshop on the Economy of Information Security*.
- Cherney, M., & McKay, B. (2019). ‘All You Need Is One Person on a Plane’: Stifling a Lethal Measles Outbreak. *The Wall Street Journal*.

- Chew, C., & Eysenbach, G. (2010). Pandemics in the age of Twitter: content analysis of Tweets during the 2009 H1N1 outbreak. *PloS One*, 5(11).
- Cleveland, D. A., & Gee, Q. (2017). Plant-based diets for mitigating climate change. In *Vegetarian and plant-based diets in health and disease prevention* (pp. 135–156). Elsevier.
- Collins, L., & Nerlich, B. (2015). Examining user comments for deliberative democracy: A corpus-driven analysis of the climate change debate online. *Environmental Communication*, 9(2), 189–207.
- Dubé, E., Vivion, M., & MacDonald, N. E. (2015). Vaccine hesitancy, vaccine refusal and the anti-vaccine movement: influence, impact and implications. *Expert Review of Vaccines*, 14(1), 99–117.
- Endo, A. (2010). A historical perspective on the discovery of statins. *Proceedings of the Japan Academy, Series B*, 86(5), 484–493.
- Ericsson, K. A., & Hastie, R. (1994). Contemporary approaches to the study of thinking and problem solving. In *Thinking and problem solving* (pp. 37–79). Elsevier.
- Fan, S. C., & Welch, J. M. (2016). Content analysis of virtual reference data: Reshaping library website design. *Medical Reference Services Quarterly*, 35(3), 294–304.
- Fang, X., & Zhan, J. (2015). Sentiment analysis using product review data. *Journal of Big Data*, 2(1), 5.
- Fekete, J.-D., Jankun-Kelly, T. J., Tory, M., & Xu, K. (2019). Provenance and Logging for Sense Making (Dagstuhl Seminar 18462). Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik.
- Fergusson, D. M., Boden, J. M., & Horwood, L. J. (2006). Cannabis use and other illicit drug use: testing the cannabis gateway hypothesis. *Addiction*, 101(4), 556–569.
- Fiz, J., Durán, M., Capellà, D., Carbonell, J., & Farré, M. (2011). Cannabis use in patients with fibromyalgia: effect on symptoms relief and health-related quality of life. *PloS One*, 6(4), e18440.
- Fox, S., & Duggan, M. (2013). Health online 2013. *Health*, 2013, 1-55.

- Foxcroft, L. (2013). *Calories and Corsets: A history of dieting over two thousand years*. London: Profile.
- Funke, J. (2010). Complex problem solving: A case for complex cognition? *Cognitive Processing, 11*(2), 133–142.
- Getman, R., Helmi, M., Roberts, H., Yansane, A., Cutler, D., & Seymour, B. (2018). Vaccine hesitancy and online information: The influence of digital networks. *Health Education & Behavior, 45*(4), 599–606.
- Godlee, F. (2016). Statins: we need an independent review. *British Medical Journal Publishing Group*.
- Grimes, S. (2016). Sentiment, emotion, attitude, and personality, via Natural Language Processing. Retrieved January 20, 2019, from <https://www.ibm.com/blogs/watson/2016/07/sentiment-emotion-attitude-personality-via-natural-language-processing/>
- Halavais, A. (2000). National borders on the world wide web. *New Media & Society, 2*(1), 7–28.
- Han, J., Pei, J., & Kamber, M. (2011). *Data mining: concepts and techniques*. Elsevier.
- Hasan, K. S., & Ng, V. (2014). Why are you taking this stance? Identifying and classifying reasons in ideological debates. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)* (pp. 751–762).
- Healey, C., & Ramaswamy, S. (2011). Visualizing twitter sentiment. *Sentiment Viz*. Retrieved from [https://www.csc2.ncsu.edu/faculty/healey/tweet\\_viz/tweet\\_app/](https://www.csc2.ncsu.edu/faculty/healey/tweet_viz/tweet_app/)
- Herring, S., Job-Sluder, K., Scheckler, R., & Barab, S. (2002). Searching for safety online: Managing "trolling" in a feminist forum. *The Information Society, 18*(5), 371–384.
- Hill, J. A. (2019). Why are we still in the middle of a 'statins war'?
- Hill, K. P. (2015). Medical Marijuana for Treatment of Chronic Pain and Other Medical and Psychiatric Problems A Clinical Review. *JAMA, 313*(24), 2474–2483. <https://doi.org/http://dx.doi.org/10.1001/jama.2015.6199>

- Hill, R. L. (2017). The political potential of numbers: data visualisation in the abortion debate. *Women, Gender & Research*, 26(1), 83–96.
- Hirschberg, J., & Manning, C. D. (2015). Advances in natural language processing. *Scienencat*, 349(6245), 261–266. <https://doi.org/10.1126/science.aaa8685>
- Hohman, F. M., Kahng, M., Pienta, R., & Chau, D. H. (2018). Visual analytics in deep learning: An interrogative survey for the next frontiers. *IEEE Transactions on Visualization and Computer Graphics*.
- Holmberg, K. (2009). *Webometric network analysis: Mapping cooperation and geopolitical connections between local government administration on the web*. Åbo Akademis förlag-Åbo Akademi University Press.
- Holmberg, K., & Thelwall, M. (2009). Local government web sites in Finland: A geographic and webometric analysis. *Scientometrics*, 79(1), 157–169. <https://doi.org/10.1007/s11192-009-0410-6>
- Howarth, C. C., & Sharman, A. G. (2015). Labeling opinions in the climate debate: A critical review. *Wiley Interdisciplinary Reviews: Climate Change*, 6(2), 239–254.
- Huang, R.-Y., Huang, C.-C., Hu, F. B., & Chavarro, J. E. (2016). Vegetarian diets and weight reduction: a meta-analysis of randomized controlled trials. *Journal of General Internal Medicine*, 31(1), 109–116.
- Huesch, M. D. (2017). Commercial online social network data and statin side-effect surveillance: a pilot observational study of aggregate mentions on facebook. *Drug Safety*, 40(12), 1199–1204.
- Jain, A. K., & Gupta, B. B. (2016). Comparative analysis of features based machine learning approaches for phishing detection. In *2016 3rd international conference on computing for sustainable global development (INDIACom)* (pp. 2125–2130). IEEE.
- Janc, K. (2016). A global approach to the spatial diversity and dynamics of internet domains. *Geographical Review*, 106(4), 567–587.
- Jauho, M. (2016). The social construction of competence: Conceptions of science and expertise among proponents of the low-carbohydrate high-fat diet in Finland. *Public Understanding of Science*, 25(3), 332–345.

- Jonassen, D. H. (1995). Computers as cognitive tools: Learning with technology, not from technology. *Journal of Computing in Higher Education*, 6(2), 40–73.
- Kata, A. (2010). A postmodern Pandora's box: Anti-vaccination misinformation on the Internet. *Vaccine*, 28(7), 1709–1716. <https://doi.org/10.1016/j.vaccine.2009.12.022>
- Kata, A. (2012). Anti-vaccine activists, Web 2.0, and the postmodern paradigm - An overview of tactics and tropes used online by the anti-vaccination movement. *Vaccine*, 30(25), 3778–3789. <https://doi.org/10.1016/j.vaccine.2011.11.112>
- Keel, P. E. (2007). EWall: A visual analytics environment for collaborative sense-making. *Information Visualization*, 6(1), 48–63.
- Kefi, H., & Perez, C. (2018). Dark Side of Online Social Networks: Technical, Managerial, and Behavioral Perspectives.
- Keim, D., Andrienko, G., Fekete, J. D., Görg, C., Kohlhammer, J., & Melançon, G. (2008). Visual analytics: Definition, process, and challenges. In *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)* (Vol. 4950 LNCS, pp. 154–175). [https://doi.org/10.1007/978-3-540-70956-5\\_7](https://doi.org/10.1007/978-3-540-70956-5_7)
- Keim, D., Kohlhammer, J., Ellis, G., & Mansmann, F. (2010). *Mastering the information age solving problems with visual analytics*. Eurographics Association.
- Kend, M., & Goode, S. (2018). The Effect of Website Age on Reported Cash Flows.
- Kickbusch, I. (2009). Health literacy: engaging in a political debate. *International Journal of Public Health*, 54(3), 131–132.
- Kitchens, B., Harle, C. A., & Li, S. (2014). Quality of health-related online search results. *Decision Support Systems*, 57, 454–462.
- Klein, G., Moon, B., & Hoffman, R. R. (2006). Making sense of sensemaking 1: Alternative perspectives. *IEEE Intelligent Systems*, 21(4), 70–73.
- Knauff, M., & Wolf, A. G. (2010). *Complex cognition: the science of human reasoning, problem-solving, and decision-making*. Springer.
- Kwai, I. (2019). Samoa Lifts State of Emergency After Deadly Measles Epidemic. *New York Times*.

- Liu, B. (2015). *Sentiment Analysis. Sentiment Analysis: Mining Opinions, Sentiments, and Emotions*. <https://doi.org/10.1017/CBO9781139084789>
- Liu, S., Foster, I., Savage, S., Voelker, G. M., & Saul, L. K. (2015). Who is. com? Learning to parse WHOIS records. In *Proceedings of the 2015 Internet Measurement Conference* (pp. 369–380).
- Liu, Z., Nersessian, N., & Stasko, J. (2008). Distributed cognition as a theoretical framework for information visualization. *IEEE Transactions on Visualization and Computer Graphics, 14*(6).
- Lotan, I., Treves, T. a, Roditi, Y., & Djaldetti, R. (2014). Cannabis (medical marijuana) treatment for motor and non-motor symptoms of Parkinson disease: an open-label observational study. *Clinical Neuropharmacology, 37*(2), 41–44. <https://doi.org/10.1097/WNF.0000000000000016>
- Lowe, M. R., & Timko, C. A. (2004). Dieting: really harmful, merely ineffective or actually helpful? *British Journal of Nutrition, 92*(S1), S19–S22.
- Maa, E., & Figi, P. (2014). The case for medical marijuana in epilepsy. *Epilepsia, 55*(6), 783–786.
- Maldonado, R., Berrendero, F., Ozaita, A., & Robledo, P. (2011). Neurochemical basis of cannabis addiction. *Neuroscience, 181*, 1–17.
- Månsson, J. (2014). A dawning demand for a new cannabis policy: A study of Swedish online drug discussions. *International Journal of Drug Policy, 25*(4), 673–681.
- Marshall, C. C., & Bly, S. (2005). Saving and using encountered information: implications for electronic periodicals. In *Proceedings of the Sigchi conference on human factors in computing systems* (pp. 111–120). ACM.
- Masino, S. A., & Rho, J. M. (2019). Metabolism and epilepsy: ketogenic diets as a homeostatic link. *Brain Research, 1703*, 26–30.
- Matarese, L. E., & Harvin, G. K. (2018). The Atkins Diet. In *Clinical Guide to Popular Diets* (pp. 1–14). CRC Press.
- Mavragani, A., & Ochoa, G. (2018). The internet and the anti-vaccine movement: tracking the 2017 EU measles outbreak. *Big Data and Cognitive Computing, 2*(1), 2.



- Mawson, A. R., Ray, B. D., Bhuiyan, A. R., & Jacob, B. (2017). Pilot comparative study on the health of vaccinated and unvaccinated 6-to 12-year-old US children. *J. Transl. Sci*, 3(3), 1–12.
- Mazzi, D. (2018). “The diet is not suitable for all...”: On the British and Irish web-based discourse on the Ketogenic Diet. *Lingue Culture Mediazioni-Languages Cultures Mediation (LCM Journal)*, 5(1), 37–56.
- McCoy, C. G., Nelson, M. L., & Weigle, M. C. (2017). University Twitter engagement: using Twitter followers to rank universities. *ArXiv Preprint ArXiv:1708.05790*.
- Meacham, M. C., Paul, M. J., & Ramo, D. E. (2018). Understanding emerging forms of cannabis use through an online cannabis community: an analysis of relative post volume and subjective highness ratings. *Drug and Alcohol Dependence*, 188, 364–369.
- Miller, L. M. S., & Bell, R. A. (2012). Online health information seeking: the influence of age, information trustworthiness, and search challenges. *Journal of Aging and Health*, 24(3), 525–541.
- Mitchell, J. T., Sweitzer, M. M., Tunno, A. M., Kollins, S. H., & McClernon, F. J. (2016). “I use weed for my ADHD”: a qualitative analysis of online forum discussions on cannabis use and ADHD. *PloS One*, 11(5).
- Moran, M. B., Lucas, M., Everhart, K., Morgan, A., & Prickett, E. (2016). What makes anti-vaccine websites persuasive? A content analysis of techniques used by anti-vaccine websites to engender anti-vaccine sentiment. *Journal of Communication in Healthcare*, 9(3), 151–163.
- Moreno, R., Ozogul, G., & Reisslein, M. (2011). Teaching with concrete and abstract visual representations: Effects on students’ problem solving, problem representations, and learning perceptions. *Journal of Educational Psychology*, 103(1), 32.
- Morphett, K., Herron, L., & Gartner, C. (2019). Protectors or puritans? Responses to media articles about the health effects of e-cigarettes. *Addiction Research & Theory*, 1–8.

- Mosleh, M., Pennycook, G., Arechar, A. A., & Rand, D. (2019). Digital fingerprints of cognitive reflection.
- National Academies of Sciences, Engineering and Medicine. (2017). *The health effects of cannabis and cannabinoids: The current state of evidence and recommendations for research*. National Academies Press.
- Navar, A. M. (2019). Fear-based medical misinformation and disease prevention: from vaccines to statins. *JAMA Cardiology*, 4(8), 723–724.
- Neal, E. G., Chaffe, H., Schwartz, R. H., Lawson, M. S., Edwards, N., Fitzsimmons, G., ... Cross, J. H. (2008). The ketogenic diet for the treatment of childhood epilepsy: a randomised controlled trial. *The Lancet Neurology*, 7(6), 500–506.
- Newman, C. B., Preiss, D., Tobert, J. A., Jacobson, T. A., Page, R. L., Goldstein, L. B., ... Raghuvver, G. (2019). Statin safety and associated adverse events: a scientific statement from the American Heart Association. *Arteriosclerosis, Thrombosis, and Vascular Biology*, 39(2), e38–e81.
- Nicholson, M. S., & Leask, J. (2012). Lessons from an online debate about measles–mumps–rubella (MMR) immunization. *Vaccine*, 30(25), 3806–3812.
- Ninkov, A., & Sedig, K. (2019). VINCENT: A visual analytics system for investigating the online vaccine debate. *Online Journal of Public Health Informatics*, 11(2).
- Ninkov, A., & Sedig, K. (2020). The Online Vaccine Debate: Study of A Visual Analytics System. In *Informatics* (Vol. 7, p. 3). Multidisciplinary Digital Publishing Institute.
- Ninkov, A., & Vaughan, L. (2017). A webometric analysis of the online vaccination debate. *Journal of the Association for Information Science and Technology*, 68(5), 1285–1294. <https://doi.org/10.1002/asi.23758>
- O'Connor, C. (2017). 'Appeals to nature' in marriage equality debates: A content analysis of newspaper and social media discourse. *British Journal of Social Psychology*, 56(3), 493–514.

- Oraby, S., Reed, L., Compton, R., Riloff, E., Walker, M., & Whittaker, S. (2017). And that's a fact: Distinguishing factual and emotional argumentation in online dialogue. *ArXiv Preprint ArXiv:1709.05295*.
- Paoli, A. (2014). Ketogenic diet for obesity: friend or foe? *International Journal of Environmental Research and Public Health*, *11*(2), 2092–2107.
- Pirolli, P., & Card, S. (2005). The sensemaking process and leverage points for analyst technology as identified through cognitive task analysis. In *Proceedings of international conference on intelligence analysis* (Vol. 5, pp. 2–4). McLean, VA, USA.
- Rajanahally, S., Raheem, O., Rogers, M., Brisbane, W., Ostrowski, K., Lendvay, T., & Walsh, T. (2019). The relationship between cannabis and male infertility, sexual health, and neoplasm: A systematic review. *Andrology*, *7*(2), 139–147.
- Redberg, R. F., & Katz, M. H. (2016). Statins for primary prevention: the debate is intense, but the data are weak. *Jama*, *316*(19), 1979–1981.
- Rothenfluh, F., & Schulz, P. J. (2018). Content, quality, and assessment tools of physician-rating websites in 12 countries: quantitative analysis. *Journal of Medical Internet Research*, *20*(6), e212.
- Ruiz, J. B., & Barnett, G. A. (2015). Exploring the presentation of HPV information online: A semantic network analysis of websites. *Vaccine*, *33*(29), 3354–3359.
- Russell, D. M., Stefik, M. J., Pirolli, P., & Card, S. K. (1993). The cost structure of sensemaking. In *Proceedings of the INTERACT'93 and CHI'93 conference on Human factors in computing systems* (pp. 269–276). ACM.
- Russo, E. B., Guy, G. W., & Robson, P. J. (2007). Cannabis, pain, and sleep: lessons from therapeutic clinical trials of Sativex®, a cannabis-based medicine. *Chemistry & Biodiversity*, *4*(8), 1729–1743.
- Salomon, G. (1993). No distribution without individuals' cognition: A dynamic interactional view. *Distributed Cognitions: Psychological and Educational Considerations*, 111–138.

- Secades-Villa, R., Garcia-Rodríguez, O., Jin, C. J., Wang, S., & Blanco, C. (2015). Probability and predictors of the cannabis gateway effect: a national study. *International Journal of Drug Policy*, 26(2), 135–142.
- Sedig, K., Klawe, M., & Westrom, M. (2001). Role of interface manipulation style and scaffolding on cognition and concept learning in learnware. *ACM Transactions on Computer-Human Interaction (TOCHI)*, 8(1), 34–59.
- Sedig, K., Naimi, A., & Haggerty, N. (2017). ALIGNING INFORMATION TECHNOLOGIES WITH EVIDENCEBASED HEALTH-CARE ACTIVITIES: A DESIGN AND EVALUATION FRAMEWORK. *Human Technology*, 13(2).
- Sedig, K., & Parsons, P. (2013). Interaction design for complex cognitive activities with visual representations: A pattern-based approach. *AIS Transactions on Human-Computer Interaction*, 5(2), 84–113.
- Sedig, K., & Parsons, P. (2016). *Design of Visualizations for Human-Information Interaction: A Pattern-Based Framework. Synthesis Lectures on Visualization* (Vol. 4). <https://doi.org/10.2200/S00685ED1V01Y201512VIS005>
- Sedig, K., Parsons, P., & Babanski, A. (2012). Towards a Characterization of Interactivity in Visual Analytics. *Journal of Multimedia Processing and Technologies, Special Issue on Theory and Application of Visual Analytics*, 3(1), 12–28. <https://doi.org/10.1145/0000000.0000000>
- Seeman, N., & Rizo, C. (2009). Assessing and responding in real time to online anti-vaccine sentiment during a flu pandemic. *Healthcare Quarterly (Toronto, Ont.)*, 13, 8–15.
- Seymour, B., Getman, R., Saraf, A., Zhang, L. H., & Kalenderian, E. (2015). When advocacy obscures accuracy online: digital pandemics of public health misinformation through an antifuoride case study. *American Journal of Public Health*, 105(3), 517–523.
- Shi, Y., Mast, K., Weber, I., Kellum, A., & Macy, M. (2017). Cultural fault lines and political polarization. In *Proceedings of the 2017 ACM on Web Science Conference* (pp. 213–217).

- Shneiderman, B., Plaisant, C., & Hesse, B. W. (2013). Improving healthcare with interactive visualization. *Computer*, 46(5), 58–66.
- Skeppstedt, M., Kerren, A., & Stede, M. (2018). Vaccine Hesitancy in Discussion Forums: Computer-Assisted Argument Mining with Topic Models. In *MIE* (pp. 366–370).
- Snyder, J. (2014). Visual representation of information as communicative practice. *Journal of the Association for Information Science and Technology*, 65(11), 2233–2247.
- Springmann, M., Godfray, H. C. J., Rayner, M., & Scarborough, P. (2016). Analysis and valuation of the health and climate change cobenefits of dietary change. *Proceedings of the National Academy of Sciences*, 113(15), 4146–4151.
- Sridhar, D., Foulds, J., Huang, B., Getoor, L., & Walker, M. (2015). Joint models of disagreement and stance in online debate. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)* (pp. 116–125).
- Stefanidis, A., Vraga, E., Lamprianidis, G., Radzikowski, J., Delamater, P. L., Jacobsen, K. H., ... Crooks, A. (2017). Zika in Twitter: temporal variations of locations, actors, and concepts. *JMIR Public Health and Surveillance*, 3(2), e22.
- Swar, B., Hameed, T., & Reychav, I. (2017). Information overload, psychological ill-being, and behavioral intention to continue online healthcare information search. *Computers in Human Behavior*, 70, 416–425.
- Thelwall, M. (2004). *Link Analysis: An Information Science Approach*. Emerald Group Publishing Limited. Retrieved from <http://linkanalysis.wlv.ac.uk/index.html>
- Thelwall, M., Sud, P., & Wilkinson, D. (2012). Link and co-inlink network diagrams with URL citations or title mentions. *Journal of the American Society for Information Science and Technology*, 63(4), 805–816.

- Thelwall, M., & Wilkinson, D. (2004). Finding similar academic Web sites with links, bibliometric couplings and colinks. *Information Processing & Management*, 40(3), 515–526.
- Thelwall, M., & Zuccala, A. (2008). A university-centred European Union link analysis. *Scientometrics*, 75(3), 407–420.
- Thomas, J. C., Diamant, J., Martino, J., & Bellamy, R. K. E. (2012). Using the “Physics” of notations to analyze a visual representation of business decision modeling. In *2012 IEEE Symposium on Visual Languages and Human-Centric Computing (VL/HCC)* (pp. 41–44). IEEE.
- Thomas, J. J., & Cook, K. a. (2005). Illuminating the path: The research and development agenda for visual analytics. *IEEE Computer Society*, 54(2), 184.  
<https://doi.org/10.3389/fmicb.2011.00006>
- Triemstra, J. D., Poeppelman, R. S., & Arora, V. M. (2018). Correlations Between Hospitals’ Social Media Presence and Reputation Score and Ranking: Cross-Sectional Analysis. *Journal of Medical Internet Research*, 20(11), e289.
- Truumees, D., Duncan, A., Mayer, E. K., Geck, M., Singh, D., & Truumees, E. (2020). Cross sectional analysis of scoliosis-specific information on the internet: potential for patient confusion and misinformation. *Spine Deformity*, 1–9.
- Tsou, M.-H., Yang, J.-A., Lusher, D., Han, S., Spitzberg, B., Gawron, J. M., ... An, L. (2013). Mapping social activities and concepts with social media (Twitter) and web search engines (Yahoo and Bing): a case study in 2012 US Presidential Election. *Cartography and Geographic Information Science*, 40(4), 337–348.
- Van Panhuis, W. G., Grefenstette, J., Jung, S. Y., Chok, N. S., Cross, A., Eng, H., ... Cummings, D. (2013). Contagious diseases in the United States from 1888 to the present. *The New England Journal of Medicine*, 369(22), 2152.
- Van Wijk, J. J. (2005). The value of visualization. In *VIS 05. IEEE Visualization, 2005.* (pp. 79–86). IEEE.
- Vaughan, L., & Ninkov, A. (2018). A new approach to web co-link analysis. *Journal of the Association for Information Science and Technology*, 69(6), 820–831.

- Vaughan, L., & You, J. (2008). Content assisted web co-link analysis for competitive intelligence. *Scientometrics*, 77(3), 433–444. <https://doi.org/10.1007/s11192-007-1999-y>
- Velardo, S. (2015). The nuances of health literacy, nutrition literacy, and food literacy. *Journal of Nutrition Education and Behavior*, 47(4), 385–389.
- Vilares, D., & He, Y. (2017). Detecting perspectives in political debates. In *Proceedings of the 2017 conference on empirical methods in natural language processing* (pp. 1573–1582).
- Volkow, N. D., Baler, R. D., Compton, W. M., & Weiss, S. R. B. (2014). Adverse health effects of marijuana use. *The New England Journal of Medicine*, 370, 2219–2227. <https://doi.org/10.1056/NEJMra1402309>
- Wakefield, A. J., Murch, S. H., Anthony, A., Linnell, J., Casson, D. M., Malik, M., ... Harvey, P. (1998). RETRACTED: Ileal-lymphoid-nodular hyperplasia, non-specific colitis, and pervasive developmental disorder in children. Elsevier.
- Walker, M. A., Tree, J. E. F., Anand, P., Abbott, R., & King, J. (2012). A Corpus for Research on Deliberation and Debate. In *LREC* (pp. 812–817). Istanbul.
- Waller, L., Hess, K., & Demetrious, K. (2016). Twitter feeders: An analysis of dominant 'voices' and patterns in a local government mosque controversy. *Australian Journalism Review*, 38(2), 47–60.
- Waseem, Z., & Hovy, D. (2016). Hateful symbols or hateful people? predictive features for hate speech detection on twitter. In *Proceedings of the NAACL student research workshop* (pp. 88–93).
- Wenger, T., Moldrich, G., & Furst, S. (2003). Neuromorphological background of cannabis addiction. *Brain Research Bulletin*, 61(2), 125–128.
- Who.int. (2019). Ten health issues WHO will tackle this year. Retrieved February 12, 2019, from <https://www.who.int/emergencies/ten-threats-to-global-health-in-2019>
- Yoon, H., Sohn, M., Choi, M., & Jung, M. (2017). Conflicting online health information and rational decision making: implication for cancer survivors. *The Health Care Manager*, 36(2), 184–191.

- Yu, L.-C., & Ho, C.-Y. (2014). Identifying emotion labels from psychiatric social texts using independent component analysis. In *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers* (pp. 837–847).
- Zheng, H., Aung, H. H., Erdt, M., Peng, T., Sesagiri Raamkumar, A., & Theng, Y. (2019). Social media presence of scholarly journals. *Journal of the Association for Information Science and Technology*, 70(3), 256–270.



## Chapter 6 – Conclusions

### 6.1. Dissertation Summary

This dissertation has examined how visual analytics systems (VASes) can facilitate making sense of online public health debates. In Chapter 1, I introduced the research. In Chapter 2, I presented a background for the relevant concepts and topics matters related to this research. In Chapter 3, I discussed the design and development of a VAS for investigating the online vaccine debate called VINCENT (VISual aNalytiCs systEm for investigating the online vacciNe debaTe). In Chapter 4, I user tested VINCENT and showed that the VAS helped people to make sense of the online vaccine debate. In Chapter 5, based on the promising findings of the user study, I generalized the development of VINCENT by developing a framework called ODIN (Online Debate entlty aNalyzer) and provided examples of how it could be applied to the development of other VASes for online public health debates. Chapter 6 concludes the dissertation and is divided into three sub-sections: Section 6.2. discusses research contributions and conclusions; Section 6.3. provides a general discussion of the research; Section 6.4. describes and discusses limitations of the research; and Section 6.5. proposes areas for future research.

### 6.2. Research Contributions and Conclusions

This research has contributed to the growing volume of literature demonstrating the value of VASes for evaluating online user generated content (Chang, Ku, & Chen, 2017; Chen, Lin, & Yuan, 2017; Haghigati & Sedig, 2020). I have provided a demonstration of the development and testing of a VAS (VINCENT) for making sense of online public health debates that integrated webometrics, NLP, data visualization, and human-data interaction. Such a study had never been attempted in this domain (public health) before, making this research novel. I have shown that it is feasible to develop such a system.

This system, VINCENT, is made up of four components: the online presence map, the word cloud, the map of website locations, and the emotion bar charts. The online presence map is a scatterplot of the websites that displays each website in proximity to one another based on their shared online presence. Websites that are plotted closer together share more online presence. In addition, each website's individual online presence is encoded using the size of the circle on the map. The larger a circle, the more inlinks and, therefore, the larger the site's online presence it has. The map of website locations displays a representation of the locations of each website on a world map. The word cloud is a representation of the 25 most common, yet unique, words that are related to the vaccine debate from each website or group of websites. Words are sized based on the frequency of their appearance on the website or group of websites. The user can control which word cloud is displayed by using the website selector. The emotion bar chart represents the positive and negative emotions found in websites' text towards specific vaccines and vaccines in general. The two bar charts represent the negative (red) and positive (green) emotions detected by IBM's Natural Language Understanding API. Each bar is composed of several rectangles that individually refer to specific websites. The width of each of these individual rectangles represents the degree of detected emotion on that specific website. The wider the rectangle, the more emotional the text is when discussing the selected vaccine. The bar charts change in response to the data that is chosen on the vaccine selector.

Moreover, in this dissertation I showed that VASes like VINCENT can improve users' abilities to complete sense-making tasks about online public health debates. Users of VINCENT were more effective at the prescribed tasks, found the tasks easier to complete, and were more confident in their responses than those who did not use the system to complete the tasks. Specifically, the user study of VINCENT showed that the system improved participants' abilities to identify, locate, and compare the debate position of websites, online presence of websites, similarity of websites, emotional tone of websites, geographic locations of websites, and specific focuses of websites. This

demonstrates that it is possible to create such systems and that these systems ultimately help users in making sense of online debates.

This dissertation has also presented ODIN, a framework for developing VASes for making sense of online public health debates. In this framework, I have generalized online public health debates and identified seven attributes (presence, shared presence, geographic location, registrant, age, focus, and sentiments) that make up these debates. I discussed and described these attributes and provided four examples of how the framework could be applied to various online public health debates (vaccines, statins, cannabis, dieting plans). Without frameworks, it is difficult for people to consider how to begin to approach problems like those discussed in this dissertation. ODIN has provided future researchers in the area with a starting point to build such systems. In other words, while VINCENT has demonstrated to future developers of VASes for online public health debates that it is possible to develop systems that improve a users' ability to complete a variety of sense-making tasks, ODIN provides them with precise attributes to consider and methods of analyses that can be implemented.

### **6.3. Discussion**

This dissertation has demonstrated that it is feasible to develop VASes to make sense of online public health debates. The value of such systems for users was demonstrated in a user study of VINCENT, in which users of the system were able to complete sense making tasks about the online vaccine debate more accurately, quicker, and with more confidence than non-system users. The observations ultimately led to continued investigation of online public health debates and the development of a framework for creating these types of systems. In this section, I discuss some important observations made over the course of this research.

Online debates are complex information spaces that are difficult to make sense of on the surface level (Albrecht, 2006; Zummo, 2017). It is not easy for people to find and appropriately evaluate the necessary information from an online public health debate if

given ample time, let alone when asked to do so quickly. Based on the results of this research, without the assistance of VASes, many people struggle to quickly determine some of the most fundamental aspects of a debate, such as the position held by a website or what the focus of a website is. One particular observation from our study of VINCENT was that many non-system users incorrectly identified the position of a vaccine website (i.e., they labelled pro-vaccine websites as anti-vaccine or vice-versa). This has profound implications and demonstrates the necessity for the development of systems like VINCENT. If people struggle to determine things like the vaccine position, it is unlikely that they have the ability to perform the more nuanced evaluations of content such as comparing the focuses or the sentiments of various websites (another finding from our study). VASes have been shown to greatly improve people's abilities to quickly make sense of complex information spaces, and this research has expanded and elaborated on this point.

Visual analytics, webometrics, and/or natural language processing (NLP) are constantly evolving and changing areas, and researchers who focus on them need to be cognizant of this, because it is relevant to every aspect of the development of VASes such as VINCENT. Tools for visualizing and interacting with data have come and gone and are updated constantly. Tableau, the interaction visualization software used to develop VINCENT, saw many changes to its functionality during the period in which this research was carried out. In addition, Tableau was acquired by Salesforce in 2019, sparking fears that it could become proprietary, or could be changed, or potentially disbanded altogether (although currently, this is not the case; it is still accessible). With regard to webometrics and NLP data analysis tools, many changes have predictably occurred since the beginning of this research in 2015. For example, WHOIS (a tool for collecting information on the registrant of a website) has been effectively shut down as a result of changing privacy laws in Europe (Layton & Elaluf-Calderwood, 2019). What this means is that for visual analytics researchers who are interested in developing systems for making sense of online debates, it is important to be able to quickly adapt and respond to the changing environment and not be overly dependent on one specific

resource or approach to data collection. Researchers in this area must be quick to adapt their research strategies so that they can capitalize on the opportunities presented by new tools and resources. Hopefully, over time, the availability of these resources will become more consistent and reliable. Even though this is unlikely to happen, increased consistency in publicly available tools and resources would lead to further growth of this research area and greater adoption of visual analytics in other domains.

Visual analytics researchers need to take time to consider the skillsets of the people who will use the systems they are developing. During the user study of VINCENT, issues emerged regarding the participants' ability to complete tasks. These issues were the result of not enough consideration to what the participants' skillsets would be. For example, one task had the participants examine the concentration of vaccine positions in various geographic regions of North America. Many participants struggled to understand the terminology of geographic regions of North America. These regions included terms such as "Mid-West." Upon further investigation, this unfamiliarity of the general public with geographical regions is not a unique finding (Ottati, 2015). This example highlights the need for geographic regions to be highlighted or explained explicitly in the VAS.

While I approached this research with the motivation of studying online public health debates specifically, it became increasingly clear over the course of my studies that the differences and distinctions between online public health debates and online debates in general are fairly negligible. In other words, online public health debates and online debates both share similar qualities in how they manifest themselves online, and therefore in how they can be analyzed by those with a vested interest in making sense of these debates. It is feasible that the observations and conclusions in this dissertation could apply to other online debates in areas not related to public health, although more research is needed to confirm.

## 6.4. Limitations

There are several limitations to this research that must be noted to place the findings of this dissertation in appropriate context. To begin, there were limitations regarding the resources that were available to perform data analysis, data visualization, and human-data interaction for VINCENT. For the data analytics, I used various tools to meet the needs that the system required, which were described in Chapter 3. The limitations of these data analytics tools are important to consider as they can influence how the users of the system will comprehend the data and the information space. For example, to conduct emotion analysis, IBM's Natural Language Understanding API (Grimes, 2016) was utilized. While this resource was useful for this project, it has limitations. For example, it can only analyze one webpage at a time. The result of this was that only a subset of each website's text was analyzed and therefore available to users of VINCENT. Another limitation of the NLU API is that this tool (like any NLP tool) can only achieve a certain level of accuracy as compared to the "gold standard" (Dale, 2018; Dolianiti et al., 2019). In other words, this means that while tools like the NLU API can evaluate text much quicker than humans are able, there is a ceiling to the accuracy of the tool's analysis when compared to what an in-depth and timely human analysis of the text would indicate.

Another example of a limitation of the data analysis methods used relates to co-link analysis. I used MOZ Link Explorer to collect inlinks, which only provided free inlink data at a domain-level, rather than at the page-level or site-level previous as other resources have previously provided (Thelwall & Wilkinson, 2004; Vaughan & Ninkov, 2018; Vaughan & You, 2008). While the results of the co-link analysis with domain-level inlinks still revealed valuable insights about the relationship between the various websites, it would have been useful to have had data for these other inlink levels as well.

Tableau, which was the software used for development of the data visualization and human-data interactions built into the system, also imposed limitations on the design of VINCENT that had to be accommodated. One important limitation was selector

functionality. While it was possible to select a data point on any of the maps, it was not possible to connect that selection to the word cloud. The word cloud selection had to be performed independently. As I described in Chapter 4, workarounds that made the system work as seamlessly as possible were developed. However, these two selection processes would ideally have been merged into one to limit confusion for users.

With regard to the user study of VINCENT, there are limitations to the design and logistics of the study that must be mentioned. First, I was limited in terms of who the participants in the study could be. As a PhD student I had limited funds, human power, and resources to attract potential participants. As a result, my study had 34 total participants (17 who used the system, 17 who did not). Ideally, I would have been able to have a larger pool of participants. Furthermore, all of these participants were Western University students and not people who necessarily had a background in public health or knowledge of the topic of the debate the study was focused on. While the conclusion that VINCENT helped users from the general public make sense of the online debate is still important, this study was not conducted on a group of participants with extensive expertise in public health. If the participants were all public health experts, the results may have been different or more enlightening in terms of the needs of public health stakeholders. However, for this project it was challenging to recruit sufficient participants who were Western students, let alone public health stakeholders, thus making an expansion of the participant pool to include those with such expertise unfeasible. Again, this is a result of the need for participants to meet the requirements already discussed and be available for participation in the study, which took for each person (in total) anywhere from 1.5 to 2 hours.

There may have been limitations arising from the design of the user-study of VINCENT. Particularly, when asked to measure online presence (in this dissertation, this means the number of inlinks to a website), the participants who user-tested the system had the relevant data readily available to them. The participants without the system, on the other hand, did not have this data and did not necessarily know how online presence should be

measured in this study (they were told that online presence was the amount of attention a website receives, as displayed in Appendix C). Based on some of the feedback from the interviews, this may have been a confusing point for some participants and limited what they possibly could have done in their response. It could also be an issue related to the internal validity of the study. It is likely that participants without the system would not have been able to complete this task even if they had known exactly how online presence should be measured, given the resources they had. However, it is also possible that they would have known what tools were needed to complete the tasks, looked up each website's online presence, and compared them during the experiment. Even though this is less likely since the participants were non-experts in this research area, it is still worth noting that this potential limitation exists in the study. In future research, which will be discussed further in Section 6.5., clearer instruction should be provided to the participants to eliminate this potential issue. It would also be useful to repeat a study such as this but compare participants' abilities to make sense of an online debate with different VASes available to them that would have different designs and features.

The framework presented in this dissertation also has limitations based on the scope in which it was developed. First, I looked at four specific online public health debates in the development of the framework (vaccines, cannabis, statins, dieting plans). While it appears that online debates tend to follow a similar structure, looking at more debates and ones in areas other than public health may have expanded or given new perspectives on the needs for the framework. For example, even with the four debates examined, the issue of polarity varied. In the vaccine debate, there was extreme polarity between the anti-vaccine and pro-vaccine groups. However, in the dieting plan debate, there appeared to be less polarity, and more of the debate was about which dieting plan is best.

Furthermore, I developed this framework by examining website and Twitter data and did not include other popular social media platforms such as Facebook, Instagram, or Reddit. Each social media platform has unique qualities, and analyses of these other platforms could involve expansion of the framework. For example, on Reddit, users can have followers, similar to Twitter. However, Reddit also has additional functionality in that



users can follow domain specific sub-Reddits (e.g., “books”, “hockey”, or “politics”). Incorporating function such as these into future studies could also reveal more about the user and/or their relationships with other users.

## 6.5. Future Research

There is enormous potential for growth in this research area of VASes for making sense of online public health debates, or of online debates in general. First, this dissertation has examined online debates primarily in terms of how they manifest themselves on websites and on Twitter. However, there are numerous other online domains (i.e., Reddit, Facebook, Instagram) that can and should be considered as more work in this area is developed or as new resources are created. Social media platforms each have unique qualities (as I discussed in Section 6.4) that are worthy of exploration in their own right (Alhabash & Ma, 2017; Bossetta, 2018; Shane-Simpson, Manago, Gaggi, & Gillespie-Lynch, 2018).

This dissertation has exclusively focused on what I have labeled “macro-level” debates (e.g., Getman et al., 2018; Ninkov & Vaughan, 2017), which occur between websites and social media profiles and manifest themselves as contradicting information provided on a topic. Micro-level debates (e.g., (Herring, Job-Sluder, Scheckler, & Barab, 2002; Nicholson & Leask, 2012; Oraby et al., 2017)), which occur as specific instances of debated topics between specific social media profiles on a website or web forum, are also an important area of research for understanding online debates. Developing VASes or frameworks for making sense of micro-level debates that help public health stakeholders or the general public investigate this level of online debates could also be a fruitful research area. These debates are complex, hard to discern, and different from macro-level debates for a number of reasons. One reason is that micro-level debates are conversational, meaning that there is a back and forth between the opposed parties. As a result, there is more opportunity for irrelevant or parallel topics to enter the debate. This can be challenging for researchers if they are interested in maintaining focus on a specific

topic of debate. Another challenge of doing research on micro-level debates is the vastness of information. The research in this dissertation took a subset of relevant macro-level debate entities, which was manageable since there is a finite number of macro-level debate entities that are publicly available. There are many more cases to explore for micro-level debates, occurring on various platforms. Not all of these cases are publicly available. Sampling of micro-level debates for analysis would be challenging but could yield important findings about how individuals discuss these topics online.

This research has been conducted with consideration of four online public health debates (vaccines, statins, cannabis, dieting plans), with the majority of the focus placed on the online vaccine debate. Other online public health debates, or online debates in general, may also be important to consider (Fox & Duggan, 2013; Miller & Bell, 2012). Some examples of debates include gun control, universal health care, electronic cigarettes or climate change. Each of these online debates are found to present contradicting and or conflicting information (Truumees et al., 2020; Yoon, Sohn, Choi, & Jung, 2017), but each also has unique qualities that are important to consider. As the growing body of research on online debates increases, the nature of these debates may become clearer, and the ways in which VASes can help facilitate sense-making of these debates could be further improved.

VASes could be developed for making sense of online public health debates with features different to those that were presented in VINCENT or the other systems based on ODIN that I described in Chapter 5. An example of such a feature that could be implemented was uncovered from the study of VINCENT, where some participants reported having difficulty seeing small details in the visualizations (that was specifically noted in relation to the word cloud). Giving users the ability to make the visualizations bigger or full screen could be useful, as it might let them focus on a particular visualization when the task they are assigned requires them to do this. In fact, when I was developing the emotion bar charts of VINCENT, an alternate idea that was considered was to use tree maps (Johnson & Shneiderman, 1999; Long, Hui, Fook, & Zainon, 2017) that,

collectively, would have been sized based on the total emotion detected. This would have been useful because the user could have quickly compared the total emotions based on the size of the entire tree map but also may have been able to view the individual emotion scores more easily than with the bar charts. Sizing tree maps, however, was not possible in Tableau when I was developing this tool. Examining emerging systems that use these different methods could be fruitful.

Furthermore, it would also be important for future research to conduct user studies that compare different types of systems to each other to identify which are the best for making sense of online debates. In this dissertation, I compared the performance of people who used the system to complete sense-making tasks with people who did not use the system. By comparing users of one system to users of another system, it would be possible to pinpoint various design features that allow sense-making tasks to be completed more effectively or quickly.

Finally, another important area of future research would be one that takes a more user-driven approach to the study and development of VASes for online public health debates as opposed to the information-driven approach taken by this research. The information-driven approach to this research topic involved centering the work on identifying the possible attributes of online debate entities, determining how to measure and analyze them, and considering how VASes can facilitate making sense of this information space. A more user-driven approach to this research topic, on the other hand, would involve understanding what the user wants to know about the information space and what they prefer to do with regard to tools that help them engage with this information. Taking a user-driven approach could explore important questions such as: 1) What does the user base his or her decisions about this information on?; 2) What are the differences between the general public and public health stakeholders in terms of their needs as they try to make sense of these debates?; 3) How do people make sense of online debates instinctually? By approaching this research area from a user-driven perspective, new

insights as to what users think about these debates and how VASes can be developed to meet or improve these perceived requirements can be developed.

## 6.6. References

- Albrecht, S. (2006). Whose voice is heard in online deliberation?: A study of participation and representation in political debates on the internet. *Information, Community and Society*, 9(1), 62–82.
- Alhabash, S., & Ma, M. (2017). A tale of four platforms: Motivations and uses of Facebook, Twitter, Instagram, and Snapchat among college students? *Social Media+ Society*, 3(1), 2056305117691544.
- Bossetta, M. (2018). The digital architectures of social media: Comparing political campaigning on Facebook, Twitter, Instagram, and Snapchat in the 2016 US election. *Journalism & Mass Communication Quarterly*, 95(2), 471–496.
- Chang, Y. C., Ku, C. H., & Chen, C. H. (2019). Social media analytics: Extracting and visualizing Hilton hotel ratings and reviews from TripAdvisor. *International Journal of Information Management*, 48, 263-279.
- Chen, S., Lin, L., & Yuan, X. (2017). Social media visual analytics. In *Computer Graphics Forum* (Vol. 36, pp. 563–587). Wiley Online Library.
- Dale, R. (2018). Text analytics apis, part 1: The bigger players. *Natural Language Engineering*, 24(2), 317–324.
- Dolianiti, F. S., Iakovakis, D., Dias, S. B., Hadjileontiadou, S. J., Diniz, J. A., Natsiou, G., Hadjileontiadis, L. J. (2019). Sentiment analysis on educational datasets: a comparative evaluation of commercial tools. *Educational Journal of the University of Patras UNESCO Chair*. <https://doi.org/10.26220/UNE.2987>
- Fox, S., & Duggan, M. (2013). Health online 2013. *Health*, 2013, 1-55.
- Getman, R., Helmi, M., Roberts, H., Yansane, A., Cutler, D., & Seymour, B. (2018). Vaccine hesitancy and online information: The influence of digital networks. *Health Education & Behavior*, 45(4), 599–606.
- Grimes, S. (2016). Sentiment, emotion, attitude, and personality, via Natural Language Processing. Retrieved January 20, 2019, from <https://www.ibm.com/blogs/watson/2016/07/sentiment-emotion-attitude-personality-via-natural-language-processing/>

- Haghighati, A., & Sedig, K. (2020). VARTTA: A Visual Analytics System for Making Sense of Real-Time Twitter Data. *Data*, 5(1), 20.
- Herring, S., Job-Sluder, K., Scheckler, R., & Barab, S. (2002). Searching for safety online: Managing" trolling" in a feminist forum. *The Information Society*, 18(5), 371–384.
- Johnson, B., & Shneiderman, B. (1999). Tree-Maps: A Space-Filling Approach to the Visualization of Hierarchical. *Readings in Information Visualization: Using Vision to Think*, 152–159.
- Layton, R., & Elaluf-Calderwood, S. (2019). A Social Economic Analysis of the Impact of GDPR on Security and Privacy Practices. In *2019 12th CMI Conference on Cybersecurity and Privacy (CMI)* (pp. 1–6). IEEE.
- Long, L. K., Hui, L. C., Fook, G. Y., & Zainon, W. M. N. W. (2017). A Study on the Effectiveness of Tree-Maps as Tree Visualization Techniques. *Procedia Computer Science*, 124, 108–115.
- Miller, L. M. S., & Bell, R. A. (2012). Online health information seeking: the influence of age, information trustworthiness, and search challenges. *Journal of Aging and Health*, 24(3), 525–541.
- Nicholson, M. S., & Leask, J. (2012). Lessons from an online debate about measles–mumps–rubella (MMR) immunization. *Vaccine*, 30(25), 3806–3812.
- Ninkov, A., & Vaughan, L. (2017). A webometric analysis of the online vaccination debate. *Journal of the Association for Information Science and Technology*, 68(5), 1285–1294. <https://doi.org/10.1002/asi.23758>
- Oraby, S., Reed, L., Compton, R., Riloff, E., Walker, M., & Whittaker, S. (2017). And that's a fact: Distinguishing factual and emotional argumentation in online dialogue. *ArXiv Preprint ArXiv:1709.05295*.
- Ottati, D. F. (2015). Geographical literacy, attitudes, and experiences of freshman students: A qualitative study at Florida International University. FIU Electronic Theses and Dissertations. 1851.

- Shane-Simpson, C., Manago, A., Gaggi, N., & Gillespie-Lynch, K. (2018). Why do college students prefer Facebook, Twitter, or Instagram? Site affordances, tensions between privacy and self-expression, and implications for social capital. *Computers in Human Behavior, 86*, 276–288.
- Thelwall, M., & Wilkinson, D. (2004). Finding similar academic Web sites with links, bibliometric couplings and colinks. *Information Processing & Management, 40*(3), 515–526.
- Truumees, D., Duncan, A., Mayer, E. K., Geck, M., Singh, D., & Truumees, E. (2020). Cross sectional analysis of scoliosis-specific information on the internet: potential for patient confusion and misinformation. *Spine Deformity*, 1–9.
- Vaughan, L., & Ninkov, A. (2018). A new approach to web co-link analysis. *Journal of the Association for Information Science and Technology, 69*(6), 820–831.
- Vaughan, L., & You, J. (2008). Content assisted web co-link analysis for competitive intelligence. *Scientometrics, 77*(3), 433–444. <https://doi.org/10.1007/s11192-007-1999-y>
- Yoon, H., Sohn, M., Choi, M., & Jung, M. (2017). Conflicting online health information and rational decision making: implication for cancer survivors. *The Health Care Manager, 36*(2), 184–191.
- Zummo, M. (2017). A linguistic analysis of the online debate on vaccines and use of fora as information stations and confirmation niche. *International Journal of Society, Culture & Language, 5*(1), 44–57.

## Appendices

### Appendix A - Set of Websites

Name	Domain
Adult Vaccination	<a href="http://www.adultvaccination.org/">http://www.adultvaccination.org/</a>
Age of Autism	<a href="http://www.ageofautism.com/">http://www.ageofautism.com/</a>
Australian Vaccination-risks Network	<a href="http://avn.org.au/">http://avn.org.au/</a>
Experimental Vaccines	<a href="http://experimentalvaccines.org/">http://experimentalvaccines.org/</a>
Families Fighting Flu	<a href="http://www.familiesfightingflu.org/">http://www.familiesfightingflu.org/</a>
Gavi The Vaccine Alliance	<a href="http://www.gavi.org/">http://www.gavi.org/</a>
History of Vaccines	<a href="http://www.historyofvaccines.org/">http://www.historyofvaccines.org/</a>
Immunization Action Coalition	<a href="http://www.immunize.org/">http://www.immunize.org/</a>
Immunize BC	<a href="http://www.immunizebc.ca/">http://www.immunizebc.ca/</a>
Immunize Canada	<a href="http://immunize.ca">http://immunize.ca</a>
Institute for Vaccine Safety	<a href="http://www.vaccinesafety.edu/">http://www.vaccinesafety.edu/</a>
National Vaccine Information Center	<a href="http://www.nvic.org/">http://www.nvic.org/</a>
Parents Requesting Open Vaccine Education	<a href="http://vaccineinfo.net/">http://vaccineinfo.net/</a>
Prevent Childhood Influenza	<a href="http://www.preventchildhoodinfluenza.org/">http://www.preventchildhoodinfluenza.org/</a>
Sabin Vaccine Institute	<a href="http://www.sabin.org/">http://www.sabin.org/</a>
Safe Minds	<a href="http://www.safeminds.org/">http://www.safeminds.org/</a>
SaneVax	<a href="http://sanevax.org/">http://sanevax.org/</a>
Shots of Prevention	<a href="http://shotofprevention.com/">http://shotofprevention.com/</a>
The Immunization Partnership	<a href="http://www.immunizeusa.org/">http://www.immunizeusa.org/</a>
The Informed Parent	<a href="http://www.informedparent.co.uk/">http://www.informedparent.co.uk/</a>
The Thinking Moms Revolution	<a href="http://thinkingmomsrevolution.com/">http://thinkingmomsrevolution.com/</a>
Think Twice Global Vaccine Institute	<a href="http://thinktwice.com/">http://thinktwice.com/</a>
Vaccinate Your Family	<a href="https://www.vaccinateyourfamily.org/">https://www.vaccinateyourfamily.org/</a>
Vaccination Information Network	<a href="http://www.vaccinationinformationnetwork.com/">http://www.vaccinationinformationnetwork.com/</a>
Vaccination Liberation	<a href="http://vaclib.org/">http://vaclib.org/</a>
Vaccination News	<a href="http://www.vaccinationnews.org/">http://www.vaccinationnews.org/</a>
Vaccine Choice Canada	<a href="http://vaccinechoicecanada.com">http://vaccinechoicecanada.com</a>
Vaccine Injury Help Center	<a href="http://www.vaccineinjuryhelpcenter.com/">http://www.vaccineinjuryhelpcenter.com/</a>
Vaccine Injury Info	<a href="http://www.vaccineinjury.info/">http://www.vaccineinjury.info/</a>
Vaccine Liberation Army	<a href="http://vaccineliberationarmy.com/">http://vaccineliberationarmy.com/</a>
Vaccine Resistance Movement	<a href="http://vaccineresistancemovement.org/">http://vaccineresistancemovement.org/</a>
Vaccine Truth	<a href="http://vaccinetruth.org/">http://vaccinetruth.org/</a>
Vaccines Today	<a href="http://www.vaccinestoday.eu/">http://www.vaccinestoday.eu/</a>
Vaccines.gov	<a href="http://www.vaccines.gov/">http://www.vaccines.gov/</a>
Vaxxter	<a href="http://vaxxter.com">http://vaxxter.com</a>
Voices for Vaccines	<a href="http://www.voicesforvaccines.org/">http://www.voicesforvaccines.org/</a>
World Association for Vaccine Education	<a href="http://novaccine.com/">http://novaccine.com/</a>



## Appendix B - Tasks

Task 1.1 Looking at all given websites, how are the pro-vaccine and anti-vaccine websites distributed?

- More pro-vaccine websites
- More anti-vaccine websites
- Equal pro- and anti-vaccine websites

Task 1.2 How many of each vaccine position are there?

Post-Task 1 I found the previous task easy to complete.

- Strongly Agree
- Agree
- Somewhat agree
- Neither agree nor disagree
- Somewhat disagree
- Disagree
- Strongly disagree

Task 2.1 Looking at the pro-vaccine websites, identify the website with the strongest Presence.

Task 2.2 Looking at the anti-vaccine websites, identify the website with the strongest Presence.

Post-Task 2 I found the previous task easy to complete.

- Strongly Agree
- Agree
- Somewhat agree
- Neither agree nor disagree
- Somewhat disagree
- Disagree
- Strongly disagree

Task 3.1 Identify and list the name of every website located outside of North America.

Task 3.2 For each of these websites, identify the country in which the website is located.

Task 3.3 For each of these websites, identify the vaccine position it takes.

Post-Task 3 I found the previous task easy to complete.

- Strongly Agree
- Agree
- Somewhat agree
- Neither agree nor disagree
- Somewhat disagree
- Disagree
- Strongly disagree

Task 4 Given the following 4 pairs of websites, compare and assign a Similarity rating. Explain each response.

4.1 Vaccination Liberation and Age of Autism

4.2 Prevent Childhood Influenza and Sabin Vaccine Institute

4.3 Institute for Vaccine Safety and World Association for Vaccine Education

4.4 Think Twice Global Vaccine Institute and Vaccination Liberation

Post-Task 4 I found the previous task easy to complete.

Strongly Agree

Agree

Somewhat agree

Neither agree nor disagree

Somewhat disagree

Disagree

Strongly disagree

Task 5 Identify if each of the following 4 words have a strong Focus amongst: (1) the pro-vaccine websites or (2) the anti-vaccine websites.

5.1 Cancer

5.2 Pregnancy

5.3 Mumps

5.4 Virus

Post-Task 5 I found the previous task easy to complete.

Strongly Agree

Agree

Somewhat agree

Neither agree nor disagree

Somewhat disagree

Disagree

Strongly disagree

Task 6 Identify if the following websites' Focus on *autism* should be regarded as: (1) strong; (2) weak or (3) none.

6.1 Shots of Prevention

6.2 The Thinking Mom's Revolution

6.3 Families Fighting the Flu

Post-Task 6 I found the previous task easy to complete.

- Strongly Agree
- Agree
- Somewhat agree
- Neither agree nor disagree
- Somewhat disagree
- Disagree
- Strongly disagree

Task 7 With regard to the HPV vaccine, identify the websites that have stronger negative than positive Emotion (i.e., are more negative about the HPV vaccine than positive).

- Vaccination News
- Vaxxter
- Voices for Vaccines
- Vaccines.gov

Post-Task 7 I found the previous task easy to complete.

- Strongly Agree
- Agree
- Somewhat agree
- Neither agree nor disagree
- Somewhat disagree
- Disagree
- Strongly disagree

Task 8 Looking at the Polio vaccine, identify which website has the strongest positive Emotion. Explain your response.

Post-Task 8 I found the previous task easy to complete.

- Strongly Agree
- Agree
- Somewhat agree
- Neither agree nor disagree
- Somewhat disagree
- Disagree
- Strongly disagree

Task 9 Looking at the anti-vaccine websites, identify which of the following vaccines has the strongest negative Emotion. Explain your response.

- HPV
- Flu
- Polio
- Measles

Post-Task 9 I found the previous task easy to complete.

- Strongly Agree
- Agree
- Somewhat agree
- Neither agree nor disagree
- Somewhat disagree
- Disagree
- Strongly disagree

Task 10.1 Identify which of the following areas has the highest concentration of pro-vaccine websites.

- Western North America
- North Eastern North America
- Europe
- Midwestern USA

Task 10.2 Identify which of the following areas has the highest concentration of anti-vaccine websites.

- Western North America
- North Eastern North America
- Europe
- Midwestern USA

Post-Task 10 I found the previous task easy to complete.

- Strongly Agree
- Agree
- Somewhat agree
- Neither agree nor disagree
- Somewhat disagree
- Disagree
- Strongly disagree

## Appendix C - Defined Terms

**Emotion**—The feelings and attitudes that are connected to specified words/phrases found on a website. Values include:

**Positive Emotion**—Emotions that encompass feelings such as joy, enjoyment, satisfaction, and pleasure. Can be invoked by a sense of well-being, inner peace, love, safety, or contentment. This emotion ranges from weakest positive to strongest positive.

**Negative Emotion**—Emotions that encompass feelings such as sadness, anger, fear, and disgust. Can be invoked by a sense of conflict, injustice, betrayal, danger, loss, and disadvantage. This emotion ranges from weakest negative to strongest negative.

**Focus**—The topics that are discussed on a website. Examples of some values include: autism, diseases, research or injury. Values range from weak focus to strong focus.

**Online Presence**—The online attention that a website receives. Values for this range from weakest presence to strongest presence.

**Shared Online Presence**—The degree to which multiple websites have online attention directed to them from the same sources. Values for this range from weak shared online presence to strong shared online presence.

**Similarity**—The degree of closeness of websites with regard to their vaccine position, shared online presence, and focus. Values include:

**None**—Websites that have different vaccine positions.

**Low**—Websites that have the same vaccine position, but differ in focus and/or shared online presence.

**High**—Websites that have the same vaccine position, shared focus and online presence.

**Vaccine Position**—The view a website takes on vaccines. Values include:

**Anti-Vaccine**—Discouraging people from vaccinating, advocating for the right to choose not to

vaccinate, and/or linking vaccinations to other health issues.

**Pro-Vaccine**—Encouraging people to vaccinate, spreading scientific information about vaccinations, and/or refuting false claims made by anti-vaccination groups.

## Appendix D - No System Post-Tasks Questionnaire

1. I was confident about the answers I provided for the given tasks.
  - Strongly Agree
  - Agree
  - Somewhat agree
  - Neither agree nor disagree
  - Somewhat disagree
  - Disagree
  - Strongly disagree
2. I found it easy to complete all the given tasks.
  - Strongly Agree
  - Agree
  - Somewhat agree
  - Neither agree nor disagree
  - Somewhat disagree
  - Disagree
  - Strongly disagree
3. Any additional comments about your experience completing these tasks.
4. Would you like to participate in an interview session? Your answer determines whether you might be asked, not whether you will be asked.
  - Yes, I want to be an interview candidate. I realize that I may not be invited.
  - Do not contact me at all about the interview session.

## Appendix E - System Post-Tasks Questionnaire

1. I was confident about the answers I provided for the given tasks.
  - Strongly Agree
  - Agree
  - Somewhat agree
  - Neither agree nor disagree
  - Somewhat disagree
  - Disagree
  - Strongly disagree
  
2. I found it easy to complete all the given tasks.
  - Strongly Agree
  - Agree
  - Somewhat agree
  - Neither agree nor disagree
  - Somewhat disagree
  - Disagree
  - Strongly disagree
  
3. I was able to clearly understand and evaluate the various individual website's online presence using the Online Presence map.
  - Strongly Agree
  - Agree
  - Somewhat agree
  - Neither agree nor disagree
  - Somewhat disagree
  - Disagree
  - Strongly disagree
  
4. I was able to clearly understand and evaluate the shared online presence of multiple websites using the Online Presence map.
  - Strongly Agree
  - Agree
  - Somewhat agree
  - Neither agree nor disagree
  - Somewhat disagree
  - Disagree
  - Strongly disagree

5. I was able to clearly understand and evaluate the focus of individual websites and groups of websites using the Word Cloud.
  - Strongly Agree
  - Agree
  - Somewhat agree
  - Neither agree nor disagree
  - Somewhat disagree
  - Disagree
  - Strongly disagree
  
6. I was able to clearly understand and evaluate the geographic dispersion of the various websites on the Map of Website Locations.
  - Strongly Agree
  - Agree
  - Somewhat agree
  - Neither agree nor disagree
  - Somewhat disagree
  - Disagree
  - Strongly disagree
  
7. I was able to clearly understand and evaluate the distribution of emotion that several websites collectively had towards various vaccines.
  - Strongly Agree
  - Agree
  - Somewhat agree
  - Neither agree nor disagree
  - Somewhat disagree
  - Disagree
  - Strongly disagree
  
8. I was able to clearly understand and evaluate the distribution of emotion that a single website had towards various vaccines.
  - Strongly Agree
  - Agree
  - Somewhat agree
  - Neither agree nor disagree
  - Somewhat disagree
  - Disagree
  - Strongly disagree



9. I found it easy to connect the information across the various visualizations together.
- Strongly Agree
  - Agree
  - Somewhat agree
  - Neither agree nor disagree
  - Somewhat disagree
  - Disagree
  - Strongly disagree
10. I found it easy to control the visualizations to see what I wanted to know.
- Strongly Agree
  - Agree
  - Somewhat agree
  - Neither agree nor disagree
  - Somewhat disagree
  - Disagree
  - Strongly disagree
11. I found it helpful to have the various types of information integrated into one system that I controlled.
- Strongly Agree
  - Agree
  - Somewhat agree
  - Neither agree nor disagree
  - Somewhat disagree
  - Disagree
  - Strongly disagree
12. What did you find worked well within the system?
13. What did you find confusing within the system?
14. Any additional comments about the visualizations and system?
- Would you like to participate in an interview session? Your answer determines whether you might be asked, not whether you will be asked.
- Yes, I want to be an interview candidate. I realize that I may not be invited.
  - Do not contact me at all about the interview session.

## Appendix F - Interview Questions

### Interview Questions—Review of Tasks

*[Note: For each participant who engages in an interview, his/her task responses and post-task questionnaire will be examined and certain responses will be chosen as desirable for gathering further information. Those responses will be used below.]*

1. How did you go about completing the tasks?
2. On the questionnaire you said <response>. Could you explain this in more detail?

### Interview Questions—Control Group Participants

*[Note: I will be spending 5 min showing the control group the introduction video to OVE. After the video, the interview questions will re-commence. The system will be available for us to look at and use as during the discussion.]*

1. What are your initial thoughts about this visual analytics system?
2. How would your approach to completing the tasks have been different if you had had the system to use?
3. Your previous response to *[insert various tasks]* easiness was *[insert their responses]*. How do you think you would have felt using the system? Could you elaborate?
4. Do you have any other comments or questions about the system?

### Interview Questions—Treatment Group Participants

*[Note: The system will be available for us to use as during the discussion.]*

1. After seeing the system again, what are your initial thoughts about it?
2. How would your approach to completing the tasks have been different if you had not had the system to use?
3. Your previous response to *[insert various tasks]* easiness was *[insert their responses]*. How do you think you would have felt without using the system? Could you elaborate?
4. Do you have any other comments or questions about the system?

### Interview Questions—Final

1. Do you have any other comments you would like to make?

## Appendix G – Ethics Approval for Study



Date: 13 March 2019

To: Dr. Kamran Sedig

Project ID: 111310

Study Title: Integrating Webometrics, Natural Language Processing, Visualization, and Human-Data Interaction into a Visual Analytics System

Application Type: NMREB Initial Application

Review Type: Delegated

Full Board Reporting Date: April 5 2019

Date Approval Issued: 13/Mar/2019

REB Approval Expiry Date: 13/Mar/2020

Dear Dr. Kamran Sedig

The Western University Non-Medical Research Ethics Board (NMREB) has reviewed and approved the WREM application form for the above mentioned study, as of the date noted above. NMREB approval for this study remains valid until the expiry date noted above, conditional to timely submission and acceptance of NMREB Continuing Ethics Review.

This research study is to be conducted by the investigator noted above. All other required institutional approvals must also be obtained prior to the conduct of the study.

#### Documents Approved:

Document Name	Document Type	Document Date	Document Version
consent_form_1	Written Consent/Assent	22/Feb/2019	2
consent_form_2	Written Consent/Assent	07/Mar/2019	3
defined_terms	Other Data Collection Instruments	21/Feb/2019	
Demographics_ Questionnaire	Online Survey	21/Feb/2019	
Flyer	Recruitment Materials	05/Feb/2019	
interview_invite	Recruitment Materials	27/Jun/2019	
interview_script	Interview Guide	05/Feb/2019	
No-System_Post-Task_ Questionnaire	Online Survey	21/Feb/2019	
participation_number	Paper Survey	06/Feb/2019	
Poster	Recruitment Materials	05/Feb/2019	
recruitment_email	Recruitment Materials	27/Jun/2019	
System_Post-Task_ Questionnaire	Online Survey	21/Feb/2019	
Tasks	Online Survey	21/Feb/2019	
url_list	Other Data Collection Instruments	21/Feb/2019	
verbal_recruitment_script	Oral Script	07/Mar/2019	3

#### Documents Acknowledged:

Document Name	Document Type	Document Date	Document Version
Figure_1_OVE_System	Supplementary Tables/Figures		
Figure_2_OVE_Online_Presence_Map	Supplementary Tables/Figures		
Figure_3_OVE_Word_Cloud	Supplementary Tables/Figures		
Figure_4_OVE_Geo_Map	Supplementary Tables/Figures		
Figure_5_OVE_Emotion_Bar_Chart	Supplementary Tables/Figures		

Page 1 of 2

Figure\_6\_Chrome\_Browser\_Control\_Group Supplementary Tables/Figures

No deviations from, or changes to the protocol should be initiated without prior written approval from the NMREB, except when necessary to eliminate immediate hazard(s) to study participants or when the change(s) involves only administrative or logistical aspects of the trial.

The Western University NMREB operates in compliance with the Tri-Council Policy Statement Ethical Conduct for Research Involving Humans (TCPS2), the Ontario Personal Health Information Protection Act (PHIPA, 2004), and the applicable laws and regulations of Ontario. Members of the NMREB who are named as Investigators in research studies do not participate in discussions related to, nor vote on such studies when they are presented to the REB. The NMREB is registered with the U.S. Department of Health & Human Services under the IRB registration number IRB 00000941.

Please do not hesitate to contact us if you have any questions.

Sincerely,

Kelly Patterson, Research Ethics Officer on behalf of Dr. Randal Graham, NMREB Chair

Note: This correspondence includes an electronic signature (validation and approval via an online system that is compliant with all regulations).

## Appendix H – Flyer for Study



### **PARTICIPANTS NEEDED FOR VISUAL ANALYTICS RESEARCH**

We are looking for volunteers to take part in a study that explores how individuals examine a set of vaccine websites with or without a visual analytics tool. To participate you must be a registered student at Western, at least 18 years of age, and be able to use a mouse or track pad, keyboard, and computer without any assistance.

If you are interested and agree to participate you would be asked to complete a series of tasks exploring a set of websites, and then provide feedback regarding your experiences completing the tasks

Your initial participation would involve 1 session. This session will be about 60 minutes long. If you agree to participate in a follow-up interview, there may be an additional audio-recorded session lasting a maximum of 30 minutes.

In appreciation for your time, you will receive a \$10 [Tim Horton's gift card](#).

For more information about this study, or to volunteer for this study, please contact:  
Anton Ninkov, Faculty of Information and Media Studies

## Appendix I – Poster for Study



### **PARTICIPANTS NEEDED FOR VISUAL ANALYTICS RESEARCH**

We are looking for volunteers to take part in a study that explores how individuals examine a set of vaccine websites with or without a visual analytics tool. To participate you must be a registered student at Western, at least 18 years of age, and be able to use a mouse or track pad, keyboard, and computer without any assistance.

If you are interested and agree to participate you would be asked to:

Complete a series of tasks exploring a set of websites,

Provide feedback regarding your experiences completing the tasks

Your initial participation would involve 1 session. This session will be about 60 minutes long.

If you agree to participate in a follow-up interview, there may be an additional audio-recorded session lasting a maximum of 30 minutes.

In appreciation for your time, you will receive a Tim Horton's gift card.

This study is being conducted under the supervision of Dr. Kamran Sedig, director of the INSIGHT LAB here at Western.

For more information about this study, or to volunteer for this study,

please contact:

Anton Ninkov

Faculty of Information and Media Studies

## Appendix J – In-Class Verbal Recruitment Script



### In-class recruitment verbal script

Hello, my name is Anton Ninkov and I am a doctoral candidate in the Faculty of Information and Media Studies. I work under the supervision of Dr. Kamran Sedig, who is a professor in the Faculty of Information and Media Studies and the director of the INSIGHT LAB. In our lab, we are studying how visual analytics tools can support users to explore a set of websites and I'm recruiting participants who meet the following inclusion criteria. To participate, you do not need to use, like, or know anything about visualizations. You just need to be a student at Western University, at least 18 years of age, and be able to use a mouse or track pad, keyboard, and computer without any assistance. This study is voluntary, confidential, and will have no effect on your academic standing. It is completely independent of this course.

This research will hopefully lead to a deeper understanding of how visual analytics tools can facilitate learning.

If you volunteer as a participant in this study, you will be asked to complete multiple tasks requiring you to explore a set of websites. You will also be asked to express your opinions about the difficulty of completing the tasks. Some participants may be invited to an audio-recorded interview, although you can opt out of this part.

The exploration session should take approximately 60 minutes of your time, while the interview session should take approximately 30 minutes.

Participants will be compensated with a Tim Horton's gift card upon completion of the exploration session.

If you are interested in participating, please contact me via email. If you have any questions, feel free to ask me now.

Thank you.

**Appendix K – Participation Number Form**Participation Number

Participation Number \_\_\_\_\_

Group A or B \_\_\_\_\_

Name \_\_\_\_\_

Email address \_\_\_\_\_

## Appendix L – Letter of Information and Consent for Exploration Session



### Letter of Information and Consent

**Project Title:**

Integrating Webometrics, Natural Language Processing, Data Visualization, and Human-Data Interaction into a Visual Analytic System

**Document Title:**

Letter of Information and Consent – Exploration Session

**Principal Investigator + Contact:**

Dr. Kamran Sedig, Professor, Faculty of Information and Media Studies  
Western University

**Co-Investigator + Contact:**

Anton Ninkov, PhD Candidate, Library and Information Science  
Western University

**1. Invitation to Participate**

You are being invited to participate in this research study about visual analytics systems and their value in exploring sets of websites. You were invited to participate in this study because you responded to a call for participants either on a poster or directly from Anton.

**2. Why is this study being done?**

Visual analytics is an important and active research area that integrates data visualization, human-data interaction and data analytics. The world is bombarded with vast amounts of data, the result of the growth and widespread adoption of the Internet in the general public's daily life, and there is therefore a need for information scientists to develop new ways to aid people in thinking about and understanding the information

This study is being conducted to examine and compare how people explore a group of websites with or without the assistance of a visual analytics tool. We are looking to



see what the potential value of a visual analytic system is and whether it can have aid people in an activity such as this.

### 3. How long will you be in this study?

This study will last between 45 minutes – 1 hour.

### 4. What will happen during this study?

If you decide to participate then you will be "randomized" into one of the two groups described below. Randomization means that you are put into a group by chance (like flipping a coin). There is no way to predict which group you will be assigned to. You will have 50/50 chance of being placed in either/any group. Neither you nor the researchers can choose what group you will be in.

*Control Group* – participants who will complete tasks without the use of a visual analytics tool.

*Treatment Group* – participants who will complete tasks with the use of a visual analytics tool.

Inclusion criteria for participation in this study include being a registered student at Western, at least 18 years of age, and ability to use a mouse or track pad, keyboard, and computer without any assistance.

### 5. What are the study procedures?

If you agree to participate, the study procedures will include the following steps:

- 1) Demographic questionnaire (<5 minutes). In this step, you will be asked to answer questions about your background, familiarity with visual interfaces and vaccines.
- 2) Familiarization period with computer (10 minutes) – Here you will be given time to familiarize yourself with the tools and computer layout provided to you.
- 3) Tasks (30 minutes) – In this step, you will be competing a set of tasks involving exploring a set of vaccine websites.
- 4) Post-Task questionnaire (<5 minutes) – Finally, you will be asked to do a Post-Task Questionnaire that will ask you about your experiences while completing the Tasks.

Note – During Step 3 (Tasks), your activity on the computer will be recorded via screen recording to help the researchers understand what exactly you did and keep track of how long each task took. This recording is mandatory for participation in the study.

The study will take place at an office at University of Western Ontario at Middlesex College.

**6. What are the risks and harms of participating in this study?**

There are no known or anticipated risks or discomforts associated with participating in this study.

**7. What are the benefits of participating in this study?**

There is unlikely to be any benefit to you by participating in this study (other than the compensation discussed subsequently). However, by participating in this study you will be part of a larger research project that could help society in general and future visual analytics research.

**8. Can participants choose to leave the study?**

Yes. If you decide to withdraw from the study, you have the right to request (e.g., by phone, in writing, etc.) withdrawal of information collected about you. If you wish to have your information removed please let the researcher know and your information will be destroyed from our records. Once the study has been published we will not be able to withdraw your information

**9. How will participants' information be kept confidential?**

Your name will only be recorded on the consent form and participation number allocation form, which will be stored for 7 years in a locked safe located in Middlesex College at Western University, separate from our study files. You will be described using your participation code, not your name, in the results.

Your responses will be collected through a secure online survey platform called Qualtrics. Qualtrics uses encryption technology and restricted access authorizations to protect all data collected. In addition, Western's Qualtrics server is in Ireland, where privacy standards are maintained under the European Union safe harbour framework. The data will then be exported from Qualtrics and securely stored on Western University's server.

**10. Are participants compensated to be in this study?**

Yes, participants will be compensated with a \$10 Tim Horton's gift card. You will receive your compensation once your participation in the study has ended.

**11. What are the rights of participants?**

Your participation in this study is voluntary. You may decide not to be in this study. Even if you consent to participate you have the right to not answer individual questions or to withdraw from the study at any time. If you choose not to participate or to leave the study at any time it will have no effect on your academic standing. You do not waive any legal right by consenting to this study.

**12. Whom do participants contact for questions?**

The participants should contact the principal investigator of the study, Kamran Sedig, for any follow up information or further questions.

If you have any questions about your rights as a research participant or the conduct of this study, you may contact The Office of Human Research Ethics. This office oversees the ethical conduct of research studies and is not part of the study team. Everything that you discuss will be kept confidential.

**This letter is yours to keep for future reference.**

**13. Consent****Project Title:**

Integrating Webometrics, Natural Language Processing, Data Visualization, and Human-Data Interaction into a Visual Analytic System

**Document Title:**

Exploration Session Consent Form – Participant copy

**Principal Investigator + Contact:**

Dr. Kamran Sedig, Professor, Faculty of Information and Media Studies  
Western University

I have read the Letter of Information, have had the nature of the study explained to me and I agree to participate. All questions have been answered to my satisfaction.

\_\_\_\_\_  
Print Name of Participant

\_\_\_\_\_  
Signature

\_\_\_\_\_  
*Date (DD-MMM- YYYY)*

My signature means that I have explained the study to the participant named above. I have answered all questions.

\_\_\_\_\_  
Print Name of Person  
Consent

\_\_\_\_\_  
Signature

\_\_\_\_\_  
*Date (DD-MMM- YYYY) Obtaining*

## Appendix M – Letter of Information and Consent for Interview Session



### Letter of Information and Consent

**Project Title:**

Integrating Webometrics, Natural Language Processing, Data Visualization, and Human-Data Interaction into a Visual Analytic System

**Document Title:**

Letter of Information and Consent – Interview Session

**Principal Investigator + Contact:**

Dr. Kamran Sedig, Professor, Faculty of Information and Media Studies  
Western University

**Co-Investigator + Contact:**

Anton Ninkov, PhD Candidate, Library and Information Science  
Western University

**1. Invitation to Participate**

You are being invited to participate in this research study about visual analytics systems and their value in exploring sets of websites. You were invited to participate in this study because you responded that you would be willing to participate in a follow up interview in the Post-Task Questionnaire during the Exploration Session.

**2. Why is this study being done?**

Visual analytics is an important and active research area that integrates data visualization, human-data interaction and data analytics. The world is bombarded with vast amounts of data, the result of the growth and widespread adoption of the Internet in the general public's daily life, and there is therefore a need for information scientists to develop new ways to aid people in thinking about and understanding the information

This study is being conducted to examine and compare how people explore a group of websites with or without the assistance of a visual analytics tool. We are looking to see what the potential value of a visual analytic system is and whether it can have aid people in an activity such as this.

**3. How long will you be in this study?**

This study will last less than 30 minutes.

**4. What will happen during this study?**

Inclusion criteria for participation in this study include being a registered student at Western, at least 18 years of age, and ability to use a mouse or track pad, keyboard, and computer without any assistance.

**5. What are the study procedures?**

If you agree to participate, the study procedure will be one Interview Session (30 minutes). In this interview, you would be asked to answer a series of questions in which you would elaborate on your experiences during the Exploration Session. The interview would be audio recorded (no video recording) so that the researcher can go back, transcribe, and thoughtfully consider the responses of the participant. Audio recording of the interview is mandatory for participation in the study.

The study will take place at an office at University of Western Ontario at Middlesex College.

**6. What are the risks and harms of participating in this study?**

There are no known or anticipated risks or discomforts associated with participating in this study.

**7. What are the benefits of participating in this study?**

There is unlikely to be any benefit to you by participating in this study (other than the compensation discussed subsequently). However, by participating in this study you will be part of a larger research project that could help society in general and future visual analytics research.

**8. Can participants choose to leave the study?**

Yes. If you decide to withdraw from the study, you have the right to request (e.g., by phone, in writing, etc.) withdrawal of information collected about you. If you wish to have your information removed please let the researcher know and your information will be destroyed from our records. Once the study has been published we will not be able to withdraw your information

**9. How will participants' information be kept confidential?**

Your name will only be recorded on the consent form and participation number allocation form, which will be stored for 7 years in a locked safe located in Middlesex College at Western University, separate from our study files. You will be described using your participation code, not your name, in the results.

The audio recording of the interview will be recorded using a tape recorder. The tapes containing the audio of the interviews will be securely stored in a locked safe separately from any identifiable information for 7 years at the co-investigator Anton Ninkov's residence in Ottawa, Ontario.

**10. Are participants compensated to be in this study?**

Yes, participants will be compensated with a \$10 Tim Horton's gift card. You will receive your compensation once your participation in the study has ended.

**11. What are the rights of participants?**

Your participation in this study is voluntary. You may decide not to be in this study. Even if you consent to participate you have the right to not answer individual questions or to withdraw from the study at any time. If you choose not to participate or to leave the study at any time it will have no effect on your academic standing. You do not waive any legal right by consenting to this study.

**12. Whom do participants contact for questions?**

The participants should contact the principal investigator of the study, Kamran Sedig, for any follow up information or further questions.

If you have any questions about your rights as a research participant or the conduct of this study, you may contact The Office of Human Research Ethics email:. This office oversees the ethical conduct of research studies and is not part of the study team. Everything that you discuss will be kept confidential.

**This letter is yours to keep for future reference.**

**13. Consent****Project Title:**

Integrating Webometrics, Natural Language Processing, Data Visualization, and Human-Data Interaction into a Visual Analytic System

**Document Title:**

Interview Session Consent Form – Participant copy

**Principal Investigator + Contact:**

Dr. Kamran Sedig, Professor, Faculty of Information and Media Studies  
Western University

I have read the Letter of Information, have had the nature of the study explained to me and I agree to participate. All questions have been answered to my satisfaction.

\_\_\_\_\_  
Print Name of Participant

\_\_\_\_\_  
Signature

\_\_\_\_\_  
*Date (DD-MMM- YYYY)*

I consent to the use of unidentified quotes obtained during the study in the dissemination of this research.

YES  NO

My signature means that I have explained the study to the participant named above. I have answered all questions.

\_\_\_\_\_  
Print Name of Person  
Consent

\_\_\_\_\_  
Signature

\_\_\_\_\_  
*Date (DD-MMM- YYYY)* Obtaining



## Curriculum Vitae

### Anton Boudreau Ninkov

#### Education

**PhD** (expected October 2020) – Library and Information Science  
September 2014 – Present

**London, Ontario**

**Western University**

**Faculty of Information and Media Studies**

**Thesis:** *Making Sense of Online Public Health Debates with Visual Analytics Systems*

**Supervisor:** Dr. Kamran Sedig

**Master of Science** – Print Media

September 2011 – May 2013

**Rochester, New York**

**Rochester Institute of Technology**

**School of Media Sciences**

**Thesis:** *A Multivariate Analysis of the Human Factors and Preferences Towards Digital Publishing Platforms for the iPad*

**Supervisor:** Professor Chris Bondy

**Bachelor of Art** – Joint Honours BA Communication and Sociology

September 2007 – May 2011

**Ottawa, Ontario**

**University of Ottawa**

**International Communication and French Immersion**

June 2010 – July 2010

**Lille, France**

**L'Université Catholique de Lille**

**Series of Web-Design Courses**

May 2011 – August 2011

**Ottawa, Ontario**

**Algonquin College**

**1<sup>st</sup> year course in Statistics**

September 2006 – June 2007

**Syracuse, New York**

**Syracuse University**

Proficient with Tableau, Adobe Package (including Photoshop, InDesign, Premier Pro, Illustrator), SAS/SPSS, HTML/CSS/JavaScript/PHP programming languages, Microsoft Office Suite, Final Cut Pro, ResourceSpace, and experience with R, Python & D3.js

## **Publications**

Ninkov, A. & Sedig, K. (2020), The Online Vaccine Debate: Study of a Visual Analytics System. *Informatics, 7(1)*.

Ninkov, A. & Sedig, K. (2019), VINCENT: A visual analytics system for investigating the online vaccine debate. *Online Journal of Public Health Informatics, 11(2)*.

Vaughan, L., & Ninkov, A. (2018). A new approach to web co-link analysis. *Journal of the Association for Information Science and Technology, 69(6)*, 820-831.

Ninkov, A. & Vaughan, L. (2017), A webometric analysis of the online vaccination debate. *Journal of the Association for Information Science and Technology, 68(5)*, 1285-1294.

## **Scholarships**

Ontario Graduate Scholarship, University of Western Ontario, September 2017 – August 2018

Western Graduate Research Scholarship, PhD LIS, University of Western Ontario, September 2014 - August 2018

## **Conferences & Presentations**

Anton Ninkov. Insights in Tableau: An Applied Data Visualization Webinar-based Course. Population Data BC. Presentation Title: “VINCENT: A visual analytics systems developed with tableau for analyzing the online Vaccine debate.” July, 2020.

Anton Ninkov. Graduate Research Day. Faculty of Information and Media Studies. Western University. Presentation Title: “Aspect Based Sentiment Analysis.” April, 2016

Anton Ninkov & Liwen Vaughan. Brown Bag Lecture Series. Faculty of Information and Media Studies. Western University. Presentation Title: “The online vaccination debate: Comparing online presence of anti and pro vaccination information.” October, 2015

## **Experiences**

### **Faculty of Information & Media Studies – Western University**

#### *PhD Candidate*

- Research in Visual Analytics Systems (VASes) for online public health debates. Developed a VAS to investigate the online vaccine debate that integrates webometrics and natural language processing data analysis with data visualization and human-data interaction called VINCENT (Visual aNalytiCs system for investigating the online vacciNe debaTe).
- Collected, transformed, cleaned, and managed data for VINCENT from a variety of sources including: MOZ Link Explorer, Alexa SEO, Google Analytics, Topsy, Internet Archives, IBM Watson, and custom-built data collection tools.
- SPSS & SAS were used for statistical analysis of data.
- Built interactive data visualizations for VINCENT using Tableau as well as D3.
- Tested VINCENT with human subjects. Results demonstrated that there was an extreme benefit for using VINCENT to explore the data and some considerations for future design of VASes of online debates were developed from the responses.
- Examined future applications of VASes in other online public health debates, highlighting the isomorphic structure these online debates take and developed a framework called ODIN (Online Debate entlty aNalyzer) to help with future development of VASes for these online debates.

### **Faculty of Information & Media Studies - Western University**

#### *Research Assistant for Paul Benedetti*

- Worked with Professor Paul Benedetti on a webometrics based online public health project.
- Examined over 800 chiropractor websites for inclusion of various terms (i.e. “Autism”, “Vaccines”, and “Subluxation”).
- Results of study were referenced in an investigation of chiropractors in The Globe and Mail (<https://www.theglobeandmail.com/canada/article-chiropractors-at-a-crossroads-the-fight-for-evidence-based-treatment/>).

### **Faculty of Information & Media Studies - Western University**

#### *Teaching Assistant (TA)*

- Taught 6 courses of 25-30 first and second year undergraduate students in the Faculty of Information and Media Studies.
- Prepared lectures of one – two hours on course curriculum including: media studies, essay writing, and communication history.
- Graded and evaluated students in numerous ways including exams, papers, presentations, and participation.
- Created a handbook to provide TAs a resource to ensure consistency between semesters/tutorials and help TAs with preparing for their weekly 2-hour tutorials.

### **Database Publishing Consultants Incorporated (DPCI)**

#### *Applications Specialist*

- Analyzed, implemented, researched, and supported cross-media publishing workflows and digital asset management systems (i.e. Adobe DPS, ResourceSpace, K4, among others).
- Worked directly with clients to assess their needs and devise workflow solutions.
- Analyzed, designed, and implemented Content Management Systems.
- Developed, designed, and analyzed graphic user interfaces.
- Developed HTML5 technologies.
- Developed and designed mobile applications.
- Worked on projects with Adobe Digital Publishing Suite and customized digital publications.

### **RIT Press**

#### *Graduate Publishing Assistant*

- Responsible for converting a variety of printed publications to E-Book format.
- Managed the uploading of content to the subscription only online publication *The Haydn Journal*.
- Researched information on various digital publishing software, apps, and workflows to meet RIT Press' needs.
- Conducted marketing research on a variety of potential publications (some of which are currently being produced).
- Helped with editing content for multiple publications.