Dissertations

Graduate College

4-1-2023

# Use of Reddit for Social Science Research: A Review of Current Use, Exploration of Potential Sampling Error, and Practical Demonstration Using Reddit to Study Post-Pandemic Teacher Resignation

Caryn Davidson
*Western Michigan University*

Follow this and additional works at: https://scholarworks.wmich.edu/dissertations

Part of the Social Statistics Commons

## Recommended Citation

USE OF REDDIT FOR SOCIAL SCIENCE RESEARCH: A REVIEW OF CURRENT USE, EXPLORATION OF POTENTIAL SAMPLING ERROR, AND PRACTICAL DEMONSTRATION USING REDDIT TO STUDY POST-PANDEMIC TEACHER RESIGNATION

Caryn Davidson, Ph.D.

Western Michigan University, 2023

This dissertation is comprised of three separate studies related to using the social media platform Reddit for social science research. The first paper provides an overview of current social science research that uses Reddit. The second paper explores the impact of sampling choices and potential sources of error when selecting Reddit data to study. The third paper maps the topics of a sample of posts tagged with the #Resignation flair in the r/Teachers subreddit from three months across school year 2021-2022 including reasons why teachers are leaving the post-pandemic classroom.

The practice of using Reddit for social science research is increasing, yet many social science researchers still know little about the possibilities available to them. This first paper reviewed 169 social science articles published in 2021 that used Reddit in their research. Reddit was primarily used for data collection and subject recruitment. The majority of researchers using Reddit for data collection are using qualitative methods for analysis. The review showed that the anonymity of Reddit makes it an especially appropriate platform choice for conducting exploratory research, especially on hard-to-reach populations and/or stigmatizing topics. Reddit was also shown to be successful as part of a broader recruitment strategy, especially to find participants from hard-to-reach populations.

The second paper explores trace selection error by examining the results of an inductive content analysis performed on samples from the (i) top most upvoted and (ii) top most

commented-on posts. These samples were compared with the universe of all responses. Differences between the samples and the universe were more apparent in *what* people said than *how* people were saying it. The sample of top most upvoted posts provided a closer representation of the topic map of the universe than the top most commented-on sample, though both methods produced holes in the topic map showing preliminary indications that these methods may not work well for achieving full topic coverage. Wide variability in how researchers report their methodologies indicates a need for reporting norms for articles utilizing Reddit; this would acknowledge and help ensure proper attention to potential sources of error in research designs. A checklist to begin creating such guidelines is presented.

The third paper utilized a hybrid approach to the Big Data Process. Data science skills were used to scrape #Resignation flair posts from the r/Teachers subreddit and qualitative methods were used to analyze the data. It was found that teachers weigh the costs and benefits of teaching and reach a breaking point where they decide to put themselves first before they quit. The reasons teachers gave for leaving fell into three main categories: Tolls of the Job, Bad Environment, and Issues with the Profession. The most frequently cited reasons for leaving were issues with students, bad administrators, and poor/declining mental health.

This dissertation represents a hybrid approach to the Big Data Process which uses data science skills for data collection and social science skills for analysis and lays the foundation for other social science researchers to do the same in their fields.

Keywords: Reddit, Data Science, Qualitative Methods, Hard-to-Reach Populations, Stigmatizing Topics, Exploratory Research, Social Media Data, Participant Recruitment, Teacher Resignation

USE OF REDDIT FOR SOCIAL SCIENCE RESEARCH: A REVIEW OF CURRENT USE,
EXPLORATION OF POTENTIAL SAMPLING ERROR, AND PRACTICAL
DEMONSTRATION USING REDDIT TO STUDY POST-PANDEMIC
TEACHER RESIGNATION

by

Caryn Davidson

A dissertation submitted to the Graduate College
in partial fulfillment of the requirements
for the degree of Doctor of Philosophy
Evaluation, Measurement, and Research
Western Michigan University
April 2023

Doctoral Committee:

Gary Miron, Ph.D., Chair
Ya Zhang, Ph.D.
June Gothberg, Ph.D.

ACKNOWLEDGEMENTS

Above all else, I would like to acknowledge my family for their support. First and foremost, my husband, Mike Davidson, for believing in me, pushing me, supporting me, and holding me to a higher standard in all things than I could ever hold myself. You are my rock and the true love of my life and none of this would be possible if it wasn't for you. You are also the reason I joined Reddit in the first place so thanks for that as well. Judy Davidson for making it possible for me to go to graduate school, and work, and have children. You welcomed me into your family and your home, and provided endless hours of free babysitting that came with the peace of mind that my children were in good hands. I would have trusted my newborn baby boy with no one else while I attended stats class. A woman can have it all, but only if she has help. My parents, Walter and Joetta Senour, for providing a solid foundation of support, instilling in me a solid work ethic, and always believing that I could do anything I set my mind to. My sister, Kaylan Wagner, for helping me to solve the same problems on this dissertation over and over again. And last but not least, my sons, Jackson and Michael Davidson, for your endless flexibility and independence while Mommy worked and for bringing the greatest joys to my life. I gave up teaching so I could be there for you and it was worth every minute.

I'd also like to acknowledge my university community starting with my committee. Thank you to my advisor, Dr. Gary Miron, for making space in your pandemic sabbotical life to guide me through this process and for never giving up on me. Thank you to Dr. Ya Zhang for your willingness to help and offer practical advice. Thank you to Dr. June Gothberg for demonstrating the joy and rigor of qualitative research methodology, and never saying no. Thank you to the late Dr. Ming Li, former Dean of the College of Education and Human Development (CEHD) and Nicole Leffler, Marketing Director for the CEHD for their support throughout my entire program. Thank you to my fellow

students that navigated coursework, research, comps and dissertation alongside me. Finally, thank you to Dr. Kevin Lee for introducing me to Python and the possibilities.

I'd like to acknowledge my work community that supported me through this journey, starting with my many supervisors over the years: first and foremost, Amanda Ellis, then Jennifer Fraker, Thomas Clouse, and Micki Ray for their faith in me and unending support. My colleague, Maggie Doyle, and her husband, Shawn Doyle, for their willingness to entertain a novice programmer in need, and all my colleagues who read and provided feedback on Paper 3.

Lastly, I acknowledge Brene Brown for showing me that qualitative research is essential even if underappreciated, and Ann Patchett for her comforting and insightful prose on life and writing. Her advice in the story "Getaway Car" published in the book, *This is the Story of a Happy Marriage*, gave me the inspiration for how to chip away at this dissertation at a time when I really needed it.


Caryn Davidson

TABLE OF CONTENTS

## LIST OF TABLES

LIST OF FIGURES

CHAPTER I

INTRODUCTION

Over ten years ago, I left an outstanding group of 5th grade students at one of the best school systems in the country (in the middle of the year, no less) and shut the door on my teaching career. After five years of teaching in Title I schools I was burned out and figured if the life of a teacher felt unsustainable to me then, even in a dream school, it was only going to get worse as my husband and I began a family. The fact that the reality of teaching didn't live up to my dream is not unique to me. I entered the PhD program in Evaluation, Measurement and Research at Western Michigan University with a goal of better understanding why teachers were leaving the classroom and what could be done about the problem. I knew that working conditions surveys like the Teaching, Empowering, Leading, and Learning (TELL) survey, provided a window into the problem and decided to focus my research on teacher working conditions and why teachers leave teaching. As I waded through articles on teacher turnover, attrition, and retention I was overwhelmed by the web of related topics and dismayed by the lack of research into the opinions of former teachers – those that had actually left teaching. The national Teacher Follow-Up Survey provided some insight but given Schaefer et al.'s (2014) finding that teachers often tell cover stories about why they leave, the data gained from questionnaires seemed to only scratch the surface of the problem in a way that did not feel helpful at all.

When I started using the social media platform Reddit, an anonymous, topic-based community, in 2018, I joined the r/teachers subreddit. Subreddits are smaller communities within the platform where moderators enforce rules for community members to engage around a chosen topic. The r/teachers subreddit has been found to function as a virtual teachers' lounge (Carpenter & Willet, 2021) and teachers post about all kinds of things related to the teaching job. I noticed that some teachers were posting stories about quitting teaching and as I read their posts it occurred to me that the anonymous nature of the platform was allowing teachers to speak authentically about their reasons for

1

leaving. The discovery felt like a gold mine of the exact kind of data I'd been trying to figure out how to collect and my thoughts turned from using social media to recruit former teachers to take a questionnaire to using it to leverage this newfound abundance of data.

**Background of the Problem**

Designing studies that effectively capture the authentic opinions of hard-to-reach populations, like former teachers whom have already left the classroom, and/or capture authentic opinions on stigmatizing topics, like why teachers leave teaching, can be difficult given the real-world constraints social science researchers face in executing those designs. In addition to the traditional constraints of labor, time, and cost associated with administering questionnaires, researchers utilizing these methods with hard-to-reach populations and/or stigmatizing topics face the added difficulty of finding appropriate respondents. Furthermore, rapid changes in communication habits have made it more difficult to reach all segments of the general population through traditional survey methods like address-based sampling and random digit dialing causing the ability of surveys to continue providing insights that can generalize from a random sample to the population to be called into question (Reveilhac et al. 2022). Leveraging social media data and data science skills to access the opinions of hard-to-reach populations and/or opinions on stigmatizing topics has been proposed as a potential solution, however the skills and paradigm of data scientists have not yet been fully integrated with the skills of social science researchers (Guber, 2021) causing social science researchers to be skeptical of the work of data scientists. There is a need for social science researchers to become more familiar with the methods used by data scientists so they can be more fully integrated into current social science research methodologies and consequently leverage the potential of social media data to provide valid and reliable information about not just hard-to-reach populations and/or opinions on stigmatizing topics, but any relevant topic, to policy makers in a more timely manner and at a reduced cost.

**Literature Review**

Marpsat and Razafindratsima (2010, p. 4) define several types of difficulty that can define a population as hard-to-reach: the population of interest has very few members, the members are hard to identify, there is no sampling frame for the population, members of the population do not wish to be identified as members due to stigma or other personal tolls of being identified as such, the behaviors of members are not well known thus making them difficult to find. When working with hard-to-reach populations, it is necessary to use non-probability sampling methods which have become increasingly popular due to their ability to overcome real-world constraints (Baker et al., 2013). The inability to randomly sample participants from a sampling frame unavoidably decreases the rigor of any survey design used with hard-to-reach populations.

In 2013, an American Association for Public Opinion Research (AAPOR) task force published a report on non-probability sampling to address the rise in use of these methodologies and the corresponding concerns about their rigor (Baker et al., 2013). The biggest issue with non-probability sampling they identified is the risk of drawing a sample that isn't actually representative of the population. Therefore, statistical adjustments must be made to manage that risk using a range of available methods like sample matching, network sampling, and estimation and weight adjustment methods which they found to range in their application in rigor. In considering when it is appropriate to utilize non-probability methods, they suggest it is helpful to consider the purpose of the survey, whether it is to describe a population or if it is to model relationships between variables. Non-probability sampling may be more appropriate for researchers looking to model concepts.

Only two years later, another AAPOR task force published a report on big data in survey research due to the rise in *its* use. As the report points out, big data has its origins in the physical sciences (think data from instruments taking collections of space), but more recently has come from

varied sources such as online price data, traffic monitoring systems, and social media messages (Japec et al., 2015). They cite Laney (2001), who identified the most prominent characteristics of big data: volume, or the large size of the data; variety, or the complexity of data coming from different systems; and velocity, or the speed at which the data is produced. As Japec et al. discuss, big data is now typically secondary data, or data not intended primarily for research use which causes both ethical and statistical concerns for research use. Data deficiencies, noise accumulation, spurious correlations and incidental endogeneity are all potential problems inherent to big data which are only exacerbated by nonsampling errors (850). They suggested the necessity of a modification of the Total Survey Error (TSE) framework to account for the sources of error inherent to big data processes and provided the diagram in Figure 1, which depicts this process in three stages. In the first, data from different sources are independently generated, in the next they are extracted, transformed, and loaded, commonly referred to as ETL in data science, and in the third they are analyzed which includes the step of first filtering, or sampling from the data, before they are finally analyzed. They explain the errors that can arise from each stage in the process and advise that transparency in these processes can help provide a check on quality throughout these processes.

**Figure 1**

*The Big Data Process Map (Japec et al., 2015)*

They suggest there are advantages and disadvantages to both survey research and big data and recommend using a blended strategy that maximizes "the ability to develop rigorous evidence for the questions of interest for an appropriate investment of resources" (863 Japec et al). Survey research offers control and the ability to generalize to populations, but is costly to execute. Big data are available in large quantities, are relatively easy to collect (for those who know what they are doing), and are available in real-time and though they can provide insight, they cannot promise to provide inference to larger populations. They suggest using the two to complement one another creating the ability to ask and answer new questions that weren't possible before.

Schober et al. (2016) narrowed in on the social media aspect of big data to show how social media analyses and survey methods align and diverge with the intention of providing a base for future research into when the use of social media and survey data might provide the basis for similar or differing conclusions and when it might be more appropriate to use one or the other. In their review of literature comparing the results of social media analyses to survey research results, they found three general points. The first, "when social media and survey results align with one another, they do so through radically different mechanisms" (184), introduces distinction between topic coverage versus population coverage. In survey research, topic coverage is achieved through population coverage through carefully constructed questions that are asked of a representative sample of a population. In social media analyses, topic coverage is achieved through the unique features of the social media platforms in an as yet not understood way. They speculate this can be achieved because the networks within the world of the social media platform are connected to the networks in the social world beyond the social media platform and it is that connection which potentially allows the nonrepresentative sample from social media to be representative of the population. They also found that lexical analysis was the predominant method of transforming social media data into data comparable to survey data

and that there is no clear reason why social media analysis and survey results diverge when they do

diverge. The provide three tables of side-by-side comparisons between different features of surveys

and social media. The first is about how participants understand the activity of responding versus

posting, the second is about the nature of the data, and the last is about practical and ethical

considerations. They conclude that social media analyses can only replace survey methods in cases

where the external gold standard for analyses is not a survey and can only supplement surveys in cases

where the external gold standard is a survey. Finally, they acknowledge the ability of social media

analyses to provide quick and affordable access to phenomena that both have and have not yet been

measured through surveys and call for greater collaboration between data scientists and survey

researchers.

     In 2020 Amaya et al. answered the call to adapt the TSE for big data with the Total Error

Framework (TEF) using the steps in the big data process to enumerate how their error components

play out in big data analysis. The error components they identify are coverage error, sampling error,

specification error, nonresponse/missing data error, measurement/content error, processing error,

modeling/estimation error, and analytic error. In 2021, Sen et al. adapted the TEF to create the Total

Error Framework for Digital Traces of Human Behavior on Online Platforms (TED-On), which can be

seen in Figure 2, to create a common vocabulary to document, communicate and compare research,

and encourage researchers to consider and document these sources of error in their designs and

communications. They introduce a distinction between errors of measurement and errors of

representation. They also introduce the terms "traces" and "users." The first applies to user-generated

content or records of online activity, and the second to the ones who generate these traces. The error

components under measurement are construct based. They are validity, platform affordances error,

trace selection error, trace augmentation error, trace reduction error, and trace measurement error. The

error components under representation are target population based and they are platform coverage error, user selection error, user augmentation error, user reduction error, and adjustment error. Exploring these sources of error more in depth will be important work in the coming years.

**Figure 2**

*Total Error Framework for Digital Traces of Human Behavior on Online Platforms (TED-On) (Sen et al. 2021)*

The predominance of participation in social media platforms over the last two decades has opened up spaces wherein hard-to-reach populations can be found. Instead of having to rely on finding the physical locations where these populations exist, social media platforms create a virtual space where members of hard-to-reach populations can be found. When the target population is not one that is stigmatized, but simply one that does not have a clear sampling frame, researchers have found Facebook can serve as virtual space to find participants. Baltar and Brunet (2011), for example, utilized Facebook to identify Argentinean entrepreneurs living in Spain. Their traditional methods of snowball sampling yielded 80 responses while their Facebook enhanced method, where they searched affinity groups and private messaged users that appeared to fit their inclusion criteria, yielded 1,023 responses. Using Facebook clearly expanded their ability to seek and find participants. Similarly, Cowles et al. (2018) found a significant increase in study enrollment after advertising their study through Facebook ads. Their original methods yielded 123 participants and the Facebook Ads allowed them to recruit an additional 1,138 participants. The ability to purchase reach through the targeting algorithms of Facebook Ads can help researchers locate members of hard-to-reach populations in a more anonymous way through since the ads are only seen by the users and no one in the person's social network needs to know they decided to participate. Grow et al. (2022) conducted a large-scale cross-national online survey to compare the demographics of Facebook users to their self-reported demographics and found that Facebook's advertising platform is a valid way of targeting specific populations, but also suggest that researchers utilize a pre-survey with demographic questions that can help further identify if the participants meet inclusion criteria. Being able to link users to their demographics is a key feature of Facebook that can be helpful in weighting analysis results to more accurately generalize to a population, however, the potential of using Facebook to reach hard-to-reach populations is limited by what people are willing to reveal about themselves. Individuals that don't

8

want to be identified as members of stigmatized groups may keep these aspects of their personas secret in the online space as they do in their everyday lives since the social network of Facebook mirrors a person's true social network. Reddit, on the other hand, is an anonymous social media platform, which provides users the ability to reveal aspects of themselves they may not feel comfortable revealing in their social networks. This affords users the ability to congregate around aspects of their identity and topics that are important to them without fear of stigma and recrimination, and it affords researchers access to members of hard-to-reach populations and their conversations.

*What is Reddit?*

Reddit is a social news aggregation, content rating and discussion website that as of March 2022 is the "9[th]-most-visited website in the world and the 6[th] most-visited website in the United States (Wikipedia, 2022). According to Reddit, the platform has over 50 million unique users active daily, over 100,000 active communities, and includes over 13 billion posts and comments (2022a). Users of Reddit, called "redditors," anonymously participate in affinity groups called "subreddits." Subreddits are dedicated to specific topics and adhere to rules set by moderators. Users "upvote" and "downvote" posts and those receiving the most upvotes are elevated by the site's algorithms so that more users see them with the top posts from all subreddits being presented on the main page of the site in an ever-shifting manner.

The screenshot in Figure 3 shows a post from the r/Teachers subreddit that is tagged with the #Resignation flair. This post has an overall score of 203, meaning the balance of upvotes and downvotes left a net score of 203 upvotes at the time this screenshot was taken - January 30, 2023. The post had 74 comments at the time the screenshot was taken. The username has been crossed out to help protect this user's anonymity. "Upvotes" are similar to "likes" on Facebook, but have a slightly different intended use. According to "Reddiquette," or the "informal expression of the values of many

redditors, as written by redditors themselves," users are encouraged to vote on content according to the following guidelines: "If you think something contributes to conversation, upvote it. If you think it does not contribute to the subreddit it is posted in or is off-topic in a particular community, downvote it," (Reddit, 2022b). Essentially, upvoting and downvoting is supposed to be based upon how well a post contributes to the intended conversation of the subreddit rather than whether or not someone likes or agrees with the opinions expressed in the post.

**Figure 3**

*Example Post Tagged with the #Resignation Flair*



Graham and Rodriguez (2021) found that "voting on Reddit is not a simple, objective rubric to rank content – on the contrary, it is a material-discursive practice that performs localized cultures and meaning-making on the site." Of the thirty topics they identified in the conversations of redditors on several high-volume subreddits in relation to upvoting and downvoting on Reddit, only one of them actually had to do with the platform's intent of the system. They organized the topics into a conceptual framework with four main themes: 1) platform culture, 2) prescriptive device, 3) materialization of

value, and 4) ontology of self. The variety of ways redditors use upvoting and downvoting calls into question whether the popular method of sampling from top posts provides the best data for social science researchers to analyze.

Gaudette et al. (2021) found that upvoting and downvoting created a sort of echo chamber within the subreddit r/The_Donald which reinforced the extreme views of the subreddit's "in-group." Through a comparison of the 1000 most upvoted posts or comments and a random sample of posts and comments, their findings suggest that:

> Reddit's upvoting and downvoting features played a central role in facilitating collective identity formation among those who post extreme right-wing content on r/The_Donald. Reddit's upvoting feature functioned to promote and normalize otherwise unacceptable views against the out-groups to produce a one-sided narrative that serves to reinforce members' extremist views, thereby strengthening bonds between members of the in-group. On the other hand, Reddit's downvoting feature functioned to ensure that members were not exposed to content that challenged their extreme right-wing beliefs, which in turn functioned as an echo chamber for hate" (3503).

Though the cultural norms of different subreddits may vary, it is true across the platform that voting affects what content is seen on Reddit by how many users. Individual user feeds include highly upvoted content from the subreddits they are a member of and if a user specifically visits a particular subreddit, the posts are organized by default to the "top" posts of the day. Posts can also be sorted by "hot," "new," "controversial," and "rising."

**Overall Gap in the Literature, Purpose, and Significance**

Several studies have broadly reviewed research using Reddit. Amaya et al. (2019) provide descriptive information, tips for using Reddit data, and suggestions for conducting surveys via Reddit

and merging survey data with data from Reddit. Medvedev et al. (2019) provide a few examples of research directions that focus on either posts or Reddit users and find a great diversity and richness in terms of the questions and methods being utilized to research with Reddit. Proferes et al. (2021) provided a systematic analysis of 727 manuscripts that used Reddit as a data source over a ten-year period. They found the number of published articles using Reddit for research has increased over time and that researchers are both using data from Reddit for the purpose of studying Reddit-specific phenomena and to study social phenomena more broadly. They also raise ethical issues about utilizing Reddit as a data source. No papers have yet catalogued the most recent social science research utilizing Reddit and few have specifically sought to understand the differences between data sets sampled using different sampling methods. Little research has been conducted on how sampling from social media data affects final analysis and no results were found specifically exploring trace selection error. Finally, current surveys of teacher opinions about leaving the classroom have been limited to current teachers and no current surveys exist exploring the opinions of teachers whom have already left the classroom.

This dissertation seeks to help bridge the gap between data science and social science research so that Reddit data can be leveraged to further social science research goals in several ways. First, the dissertation provides a review of current social science research leveraging Reddit to bring greater knowledge of what is possible using data science methodologies in the social science context. This knowledge will help social science researchers ask and answer new and exciting questions in their fields. Secondly, this paper will take a closer look at the impact of trace selection error, an error component of Sen et al.'s (2021) Total Error Framework for Digital Traces of Human Behavior on Online Platforms (TED-On). Exploring the potential consequences of two common sampling

techniques help bring greater scrutiny to the choices researchers make when selecting digital human traces of human behavior to analyze.

This dissertation will demonstrate how data science methods can be combined with social science methods, creating a variation on the Big Data Process (Japec et al. 2015) that inserts qualitative analysis into the analysis stage. This dissertation also suggests a novel mixed-methods approach for investigating topics of policy interest. The approach begins with the qualitative, exploratory investigation demonstrated in Paper 3, utilizes the results of that exploration to build a questionnaire, and utilizes the recruiting methods discussed in Paper 1 to recruit respondents, hence generating quantitative data that would provide further insight. Finally, this dissertation demonstrates how Reddit data can be leveraged to explore opinions on a stigmatized topic as expressed by a hard-to-reach population, thus providing timely and valid insight into a current issue that policy makers are eager to learn more about.

Other fields of research have already begun to leverage extant Reddit data to reach hard-to-reach populations and/or opinions on stigmatizing topics, but educational research has barely scratched the surface of this resource and certainly has not used it to provide insight to education leaders and policy makers. The window the r/teachers subreddit provides into the day-to-day experiences of teachers allows educational leaders and policy makers a view of what's happening in classrooms that wouldn't be afforded to them during an in-person visit where educators may be motivated to tell cover stories to hide the reality of their situation – a variation of what educators commonly refer to as the "dog and pony show." In the absence of the ability to generate gold-standard survey data about why teachers are leaving the classroom, this dissertation identifies a path forward for better understanding the critical question of why teachers are leaving the classroom.

**Dissertation Format and Related Purposes of the Three Studies**

The three papers in this three-paper dissertation work together to explore how Reddit can be leveraged for social science research. They build from surveying current research methods for leveraging Reddit data in social science research, to contributing to the literature on trace selection error, to demonstrating how Reddit data can be leveraged to address a current issue that is of interest to policy makers today. The first paper in this dissertation will provide a review of over 150 articles published during 2021 that leveraged Reddit in some way for social science research. The second paper will compare two methods of sampling Reddit data - sampling from top most-upvoted posts and sampling from top most-commented-on posts to understand the differences between these two sampling methods as they compare to the results of analyzing the universe of posts. The third paper will use Reddit data to analyze teachers' conversations around leaving teaching which will provide topic coverage to policy makers and provide a foundation for questionnaire creation so that teachers can be surveyed regarding their thoughts on the sustainability of the teaching profession.

*Paper 1: A Survey of Current Social Science Research Methodology Utilizing the Social Media Platform Reddit*

**Purpose.** Though the practice of leveraging Reddit for social science research is increasing, many social science researchers still know little about the possibilities available to them through these methodologies. This paper will provide an overview of current social science research that is leveraging Reddit, broadening the methodological knowledge base and potentially opening the possibility to ask and answer new and exciting questions which hadn't been answerable before. This paper poses that Reddit is an especially appropriate platform choice for conducting exploratory research, especially on hard-to-reach populations and/or stigmatizing topics; and to recruit participants for studies, especially participants from hard-to-reach populations.

**Research Questions.** This paper explores two research questions:

1. How are social science researchers utilizing Reddit to further their research goals?

2. What should a researcher interested in using Reddit to further their research goals take into consideration when formulating a research plan?

**Methodology.** This paper summarizes social science researcher's use of Reddit in research published during 2021. Social science research databases were searched to build a corpus of 169 articles meeting inclusion criteria. Each article was reviewed, and a series of variables were collected. The variables varied slightly depending on whether the article used Reddit for recruitment or in some other way. Variables for recruitment articles included the journal name and field, the country of origin of the authors, whether Reddit was the sole method used for recruitment, whether an incentive was offered, whether the recruitment message was included in the paper, the eventual method utilized once participants were recruited, arguments for using Reddit, and limitations of using Reddit. Variables for all other articles included the journal name and field, the country of origin of the authors, the level of Reddit data analyzed, the number of units analyzed, the method of data scraping, the sampling method used, the method of analysis, whether or not IRB approval was obtained, arguments for using Reddit, limitations of using Reddit. The results of the review were then tabulated.

**Contribution to Research.** The unique features of the Reddit platform afford researchers the ability to reach hard-to-reach populations and/or opinions on stigmatizing topics. This paper provides social science researchers with an understanding of what's possible when it comes to leveraging this social media platform. Having this understanding will help social science researchers ask and answer new questions that may not have been possible before, save time and money on research designs, provide insight into emerging, rapidly evolving, and stigmatizing topics, and leverage the platform for

the recruitment of research participants. This paper also begins to build guidelines for reporting research that leverages Reddit.

***Paper 2: How Sampling Method Affects the Results of Inductive Content Analysis: A Case Study Utilizing Data from the Social Media Platform Reddit***

**Purpose.** This paper seeks to understand how analysis results derived from a sample of the most upvoted posts and the most commented on posts from a given time period about a certain topic compare to analysis of the universe of posts from which the samples were taken. Although researchers are increasingly turning to Reddit as a valid source of research data, little has been published about best practices for sampling Reddit data. This paper begins to build researchers' understanding of how sampling choices affect analyses of Reddit data.

**Research Questions.** This paper will explore three research questions:

1. Is there a difference in the results of an inductive content analysis performed on the top fifty most upvoted posts from a given time period, the top fifty most commented on posts from a given time period, and the universe of posts from that same time period?

2. Is there a difference in the results of sentiment analysis performed on the top 50 most upvoted posts from a given time period, the top 50 most commented on posts from a given time period, and the universe of posts from that same time period?

3. If differences exist, how might these differences be interpreted to inform research design when social science researchers are sampling data from Reddit?

**Methodology.** A Python script provided by Rare Loot (2018) was adapted to scrape the desired data from the JSON coding behind the Pushift.io dataset, a mirrored copy of the Reddit platform and a sample of the fifty most upvoted and fifty most commented on posts were taken. A mixed-methods approach to analysis was conducted to capture both *what* was being said in the posts in

16

each dataset and *how* it was being said in each sample. To measure what was being said in the posts, an inductive content analysis adapted from the procedure outlined by Elo & Kyngäs (2007) was performed on each group. The results the inductive content analysis for each sample was compared with the universe of posts from which they were taken. To measure how original posters were saying things, statistical analysis was performed on several variables calculated using Linguistic Inquiry and Word Count (LIWC)-22.

**Contribution to Research.** This paper begins to build the literature around what Sen et al. (2021) call trace selection error and builds researchers' understanding of how sampling choices affect analyses of Reddit data. By comparing analysis results performed on samples of data to the universe from which they were taken, researchers gain insight into the consequences of two different sampling methods. The nuances of the different methods may be more beneficial for some research questions than others and researchers will be able to apply the insights to their own research situation and make informed choices that make the most sense for them. Building this area of the literature also helps encourage other such experimentation on the effects of sampling methods to be performed.

## Paper 3: Why Teachers Are Leaving the Post-Pandemic Classroom: A Content Analysis of #Resignation Flair Posts in the r/Teachers Subreddit

**Purpose.** It is imperative to understand why teachers are leaving and what can be done to make teaching a more sustainable career, however reaching former teachers can be difficult because they are not easily identifiable and there is no sampling frame for the population. Furthermore, because teachers have been found to tell cover stories about why they leave (Schaefer et al., 2014), traditional survey methods which *have* been undertaken fall short of providing actionable data. This paper maps the topics of a sample of posts tagged with the #Resignation flair in the r/Teachers subreddit from three months across school year 2021-2022 filling the gap in the literature around why

17

teachers are leaving the post-pandemic classroom. Understanding how teachers talk about resigning and the reasons they give for leaving provides a foundation for developing questionnaires that can be used to provide more accurate data for local contexts. This paper provides the first step towards helping policy and decision makers understand what factors matter to teachers when they are making their decisions to stay or go.

**Research Questions.**

1. What topics do teachers discuss in #Resignation flair posts on the subreddit r/Teachers?

2. What reasons do teachers give in #Resignation flair posts on the subreddit r/Teachers for leaving a teaching job or the teaching profession entirely?

**Methodology.** A Python script provided by Rare Loot (2018) was adapted to scrape the desired data from the JSON coding behind the Pushift.io dataset. Data was imported to Microsoft Excel and processed for analysis. An inductive approach to qualitative content analysis based on the procedure outlined by Elo and Kyngäs (2008) was taken. The post title and post body for every post in the dataset was copied into a table in Microsoft Word and the posts were printed in landscape format on 8 ½ x 11 paper. After immersing in the data, codes were recorded in the margins to note either the reasons given for teaching or the "gist" of the post if no reasons were present. During this process it was determined that the posts were falling into several distinct buckets that might provide a useful framework for analyzing the codes and these buckets were recorded for each post as well. Posts were then organized by bucket, and the abstraction process began as the codes recorded in the margins were cut apart and grouped into categories. Pictures were taken to record the groupings and once all were complete, the categories and codes from the pictures were transferred into a digital format. During this process refinement of categories took place. A log trail was kept throughout the inductive content analysis to document the process, record insights and reactions to the data reflexively, and think

through and record decisions made. It was not possible to member check due to the anonymous nature of the data.

   **Contribution to Research.** This paper demonstrates a new hybrid approach to the Big Data Process which blends data science skills for data collection with qualitative social science methods for data analysis. This method provide a number of benefits to social science researchers who wish to provide current insights into relevant policy issues. By replacing traditional data science methods of big data analysis with qualitative methods performed on a sample taken from the big data set, social science researchers can begin to integrate data science and social science research skills that can save social science researchers time, money, and effort while providing valid and reliable information about both hard-to-reach populations and stigmatizing topics.

**References**

Amaya, A., Biemer, P. P., & Kinyon, D. (2020). *Journal of Statistics and Methodology, 8*(1), 89-119. https://doi.org/10.1093/jssam/smz056

Amaya, A., Bach, R., Keusch, F., & Kreuter, F. (2019). New data sources in social science research: Things to know before working with Reddit data. *Social Science Computer Review, 39*(5), 943-960. doi: 10.1177/0894439319893305

Baker, R., Brick, J. M., Bates, N. A., Battaglia, M., Couper, M. P., Dever, J. A., Gile, K. J., Tourangeau, R. (2013). Summary report of the AAPOR task force on non-probability sampling. *Journal of Survey Statistics and Methodology, 1*(2), 90-143. https://doi.org/10.1093/jssam/smt008

Baltar, F., & Brunet, I. (2011). Social research 2.0: Virtual snowball sampling method using Facebook. *Internet Research, 22*(1), 57-74. doi: 10.1108/10662241211199960

Carpenter, J. P., & Bret Staudt Willet, K. (2021). The teachers' lounge and the debate hall: Anonymous self-directed learning in two teaching-related subreddits. *Teaching and Teacher Education, 104*. https://doi.org/10.1016/j.tate.2021.103371

Cowles, C., Berk, S., & Siddiqi, B. (2018). Using Facebook ads to recruit clinical study participants. *Applied Clinical Trials, 27*(12), 14-17.

Elo, S. & Kyngäs, H. (2007). The qualitative content analysis process. *Journal of Advanced Nursing 62*(1), 107-115. doi: 10.1111/j.1365-2648.04569.x

Gaudette, T., Scrivens, R., Davies, G., & Frank, R. (2021). Upvoting extremism: Collective identify formation and the extreme right on Reddit. *new media & society*, *23*(12), 3491-3508. doi: 10.1177/1461444820958123

Graham, T. & Rodriguez, A. (2021). Sociomateriality of rating and ranking devices on social media: A

    case study of Reddit's voting practices. *Social Media + Society*, *7*(3),  doi:

    10.1177/20563051211047667

Grow, A., Perrotta, D., Del Fava, E., Cimentada, J., Rampazzo, F., Gil-Clavel, S., Zagheni, E., Flores,

    R. D., Ventura, I., Weber, I. (2022). Is Facebook's advertising data accurate enough for use in

    social science research? Insights from a cross-national online survey. *Journal of the Royal*

    *Statistical Society, Series A, 185*(S2), 343-363. doi: 10.1111/rssa.12948.

Guber, D. L. (2021). Public opinion and the classical tradition: Redux in the digital age. *Public*

    *Opinion Quarterly, 85*(4), 1103-1127. doi: 10.1093/poq/nfab053

Japec, L., Kreuter, F., Berg, M., Biemer, P., Decker, P., Lampe, C., Lane, J., O'Neil, C., & Usher, A.

    (2015). Big data in survey research. *Public Opinion Quarterly, 79*(4), 839-880.

    https://doi.org/10.1093/poq/nfv039

Laney, D. (2001). 3-D management: Controlling data volume, velocity, and variety. META Group

    Research Note, February 6.

Marpsat, M., & Razafindratsima, N. (2010). Survey methods for hard-to-reach populations:

    Introduction to the special issue. *Methodolgoical Innovations Online, 5*(2), 3-16.

    https://doi.org/10.4256/mio.2010.0014

Medvedev, A., Lambiotte, R., & Delvenne, J. (2019). The anatomy of Reddit: An overview of

    Academic Research. Dynamics On and Of Complex Networks III, Springer Proceedings in

    Complexity. https://doi.org/10.1007/978-3-030-14683-2_9

Proferes, N., Jones, N., Gilbert, S., Fiesler, C., & Zimmer, M. (2021) Studying Reddit: a systematic

    overview of disciplines, approaches, methods, and ethics. *Social Media + Society*.

    https://doi.org/10.1177/20563051211019004

Rare Loot (2018). *Using Pushshift's API to extract Reddit submissions*.

https://rareloot.medium.com/using-pushshifts-api-to-extract-reddit-submissions-fb517b286563

*Reddit*. (2022, November 13). In *Wikipedia*. https://en.wikipedia.org/wiki/Reddit

Reddit (2022, November 13). *About Reddit.* Reddit. https://www.redditinc.com/

Reddit. (2022, November 13b). *Reddiquette.* Reddit. https://www.reddithelp.com/hc/en-

us/articles/205926439

Reveilhac, M., Steinmetz, S., & Morselli, D. (2022). A systematic literature review of how and

whether social media data can complement traditional survey data to study public opinion.

*Multimedia Tools and Applications, 81*. 10107-10142. https://doi.org/10.1007/s11042-022-

12101-0

Schaefer, L., Downey, C. A., & Clandinin, D. J. (2014). Shifting from stories to live by to stories to

leave by: Early career teacher attrition. *Teacher Education Quarterly, 41*(1), 9-27.

Schober, M. F., Pasek, J., Guggenheim, L., Lampe, C., & Conrad, F. G. (2016). Social media analyses

for social measurement. *Public Opinion Quarterly, 80*(1), 180-211.

https://doi.org/10.1093/poq/nfv048

Sen, I., Flöck, F., Weller, K., Weib, B., & Wagner, C. (2021). *Public Opinion Quarterly, 85*(S1), 399-

422. https://doi.org/10.1093/poq/nfab018

CHAPTER II

A SURVEY OF CURRENT SOCIAL SCIENCE RESEARCH UTILIZING THE SOCIAL MEDIA
PLATFORM REDDIT

Social media data analysis has long been practiced by data scientists and is commonly applied

for marketing purposes in the private sector (Japec et al. 2015), but its use in social science research

and the public sector is less widespread. Traditional methods of conducting valid and reliable surveys

comes at substantial costs in time, money, and effort (Schober et al. 2016) and rapid changes in

communication habits have made it more difficult to reach all segments of the general population

through traditional survey methods like address-based sampling and random digit dialing causing the

ability of surveys to continue providing insights that can generalize from a random sample to the

population to be called into question (Reveilhac et al. 2022). Leveraging social media data and data

science skills holds the potential of yielding valid and reliable results at a fraction of the cost, in some

cases, however the skills and paradigm of data scientists have not yet been fully integrated with the

skills of social science researchers (Guber, 2021) causing social science researchers to be skeptical of

the work of data scientists.

Each social media platform has its own unique features and Reddit has some particular

platform affordances that make it especially appealing for several challenging research functions.

First, the anonymous nature of the platform facilitates access to hard-to-reach populations and

discussions about sensitive and stigmatizing topics that would be either extremely difficult to access or

completely inaccessible to social science researchers otherwise. Secondly, Reddit is a topic-based

community which means it self-organizes around topics of mutual interest, and because of its size and

reach, it creates communities around very specific things that people care about. This feature makes it

easier for researchers to filter all the data on the site to find discussions around their research topics of

interest which may help reduce trace selection error when determining what data should be extracted from the platform, especially for novice data scientists.

Though the practice of leveraging Reddit for social science research is increasing, many social science researchers still know little about the possibilities available to them through these methodologies. This paper will provide an overview of current social science research that is leveraging Reddit, broadening the methodological knowledge base and potentially opening the possibility to ask and answer new and exciting questions which hadn't been answerable before. This paper poses that Reddit is an especially appropriate platform choice for conducting exploratory research, especially on hard-to-reach populations and/or stigmatizing topics; and to recruit participants for studies, especially participants from hard-to-reach populations.

**What is Reddit?**

Reddit is a social news aggregation, content rating and discussion website that as of March 2022 is the "9[th]-most-visited website in the world and the 6[th] most-visited website in the United States (Wikipedia, 2022). According to Reddit, the platform has over 50 million unique users active daily, over 100,000 active communities, and includes over 13 billion posts and comments (2022a). Users of Reddit, called "redditors," anonymously participate in affinity groups called "subreddits." Subreddits are dedicated to specific topics and adhere to rules set by moderators. Users "upvote" and "downvote" posts and those receiving the most upvotes are elevated by the site's algorithms so that more users see them with the top posts from all subreddits being presented on the main page of the site in an ever-shifting manner.

The screenshot in Figure 4 shows a post from the r/Teachers subreddit that is tagged with the #Resignation flair. This post has an overall score of 203, meaning the balance of upvotes and downvotes left a net score of 203 upvotes at the time this screenshot was taken - January 30, 2023.

The post had 74 comments at the time the screenshot was taken. The username has been crossed out to help protect this user's anonymity.

**Figure 4**

*Example Post Tagged with the #Resignation Flair*



"Upvotes" are similar to "likes" on Facebook, but have a slightly different intended use. According to "Reddiquette," or the "informal expression of the values of many redditors, as written by redditors themselves," users are encouraged to vote on content according to the following guidelines: "If you think something contributes to conversation, upvote it. If you think it does not contribute to the subreddit it is posted in or is off-topic in a particular community, downvote it," (Reddit, 2022b). Essentially, upvoting and downvoting is supposed to be based upon how well a post contributes to the intended conversation of the subreddit rather than whether or not someone likes or agrees with the opinions expressed in the post.

Graham and Rodriguez (2021) found that "voting on Reddit is not a simple, objective rubric to rank content – on the contrary, it is a material-discursive practice that performs localized cultures and

meaning-making on the site." Of the thirty topics they identified in the conversations of redditors on several high-volume subreddits in relation to upvoting and downvoting on Reddit, only one of them actually had to do with the platform's intent of the system. They organized the topics into a conceptual framework with four main themes: 1) platform culture, 2) prescriptive device, 3) materialization of value, and 4) ontology of self. The variety of ways redditors use upvoting and downvoting calls into question whether the popular method of sampling from top posts provides the best data for social science researchers to analyze.

Gaudette et al. (2021) found that upvoting and downvoting created a sort of echo chamber within the subreddit r/The_Donald which reinforced the extreme views of the subreddit's "in-group." Through a comparison of the 1000 most upvoted posts or comments and a random sample of posts and comments, their findings suggest that:

> Reddit's upvoting and downvoting features played a central role in facilitating collective identity formation among those who post extreme right-wing content on r/The_Donald. Reddit's upvoting feature functioned to promote and normalize otherwise unacceptable views against the out-groups to produce a one-sided narrative that serves to reinforce members' extremist views, thereby strengthening bonds between members of the in-group. On the other hand, Reddit's downvoting feature functioned to ensure that members were not exposed to content that challenged their extreme right-wing beliefs, which in turn functioned as an echo chamber for hate" (3503).

Though the cultural norms of different subreddits may vary, it is true across the platform that voting affects what content is seen on Reddit by how many users. Individual user feeds include highly upvoted content from the subreddits they are a member of and if a user specifically visits a particular

subreddit, the posts are organized by default to the "top" posts of the day. Posts can also be sorted by "hot," "new," "controversial," and "rising."

**Overall Gap in the Literature, Purpose, and Significance**

Several studies have broadly reviewed research using Reddit. Amaya et al. (2019) provide descriptive information, tips for using Reddit data, and suggestions for conducting surveys via Reddit and merging survey data with data from Reddit. Medvedev et al. (2019) provide a few examples of research directions that focus on either posts or Reddit users and find a great diversity and richness in terms of the questions and methods being utilized to research with Reddit. Proferes et al. (2021) provided a systematic analysis of 727 manuscripts that used Reddit as a data source over a ten-year period. They found the number of published articles using Reddit for research has increased over time and that researchers are both using data from Reddit for the purpose of studying Reddit-specific phenomena and to study social phenomena more broadly. They also raise ethical issues about utilizing Reddit as a data source. No papers have yet catalogued the most recent social science research utilizing Reddit and few have specifically sought to understand the differences between data sets sampled using different sampling methods.

**Literature Review**

Marpsat and Razafindratsima (2010, p. 4) define several types of difficulty that can define a population as hard-to-reach: the population of interest has very few members, the members are hard to identify, there is no sampling frame for the population, members of the population do not wish to be identified as members due to stigma or other personal tolls of being identified as such, the behaviors of members are not well known thus making them difficult to find. When working with hard-to-reach populations, it is necessary to use non-probability sampling methods which have become increasingly popular due to their ability to overcome real-world constraints (Baker et al., 2013). The inability to

randomly sample participants from a sampling frame unavoidably decreases the rigor of any survey design used with hard-to-reach populations.

In 2013, an American Association for Public Opinion Research (AAPOR) task force published a report on non-probability sampling to address the rise in use of these methodologies and the corresponding concerns about their rigor (Baker et al., 2013). The biggest issue with non-probability sampling they identified is the risk of drawing a sample that isn't actually representative of the population. Therefore, statistical adjustments must be made to manage that risk using a range of available methods like sample matching, network sampling, and estimation and weight adjustment methods which they found to range in their application in rigor. In considering when it is appropriate to utilize non-probability methods, they suggest it is helpful to consider the purpose of the survey, whether it is to describe a population or if it is to model relationships between variables. Non-probability sampling may be more appropriate for researchers looking to model concepts.

Only two years later, another AAPOR task force published a report on big data in survey research due to the rise in *its* use. As the report points out, big data has its origins in the physical sciences (think data from instruments taking collections of space), but more recently has come from varied sources such as online price data, traffic monitoring systems, and social media messages (Japec et al., 2015). They cite Laney (2001), who identified the most prominent characteristics of big data: volume, or the large size of the data; variety, or the complexity of data coming from different systems; and velocity, or the speed at which the data is produced. As Japec et al. discuss, big data is now typically secondary data, or data not intended primarily for research use which causes both ethical and statistical concerns for research use. Data deficiencies, noise accumulation, spurious correlations and incidental endogeneity are all potential problems inherent to big data which are only exacerbated by nonsampling errors (850). They suggested the necessity of a modification of the Total Survey Error

28

(TSE) framework to account for the sources of error inherent to big data processes and provided the

diagram in Figure 5, which depicts this process in three stages. In the first, data from different sources

are independently generated, in the next they are extracted, transformed, and loaded, commonly

referred to as ETL in data science, and in the third they are analyzed which includes the step of first

filtering, or sampling from the data, before they are finally analyzed. They explain the errors that can

arise from each stage in the process and advise that transparency in these processes can help provide a

check on quality throughout these processes.

**Figure 5**

*The Big Data Process Map (Japec et al., 2015)*



They suggest there are advantages and disadvantages to both survey research and big data and

recommend using a blended strategy that maximizes "the ability to develop rigorous evidence for the

questions of interest for an appropriate investment of resources" (863 Japec et al). Survey research

offers control and the ability to generalize to populations, but is costly to execute. Big data are

available in large quantities, are relatively easy to collect (for those who know what they are doing),

and are available in real-time and though they can provide insight, they cannot promise to provide inference to larger populations. They suggest using the two to complement one another creating the ability to ask and answer new questions that weren't possible before.

Schober et al. (2016) narrowed in on the social media aspect of big data to show how social media analyses and survey methods align and diverge with the intention of providing a base for future research into when the use of social media and survey data might provide the basis for similar or differing conclusions and when it might be more appropriate to use one or the other. In their review of literature comparing the results of social media analyses to survey research results, they found three general points. The first, "when social media and survey results align with one another, they do so through radically different mechanisms" (184), introduces distinction between topic coverage versus population coverage. In survey research, topic coverage is achieved through population coverage through carefully constructed questions that are asked of a representative sample of a population. In social media analyses, topic coverage is achieved through the unique features of the social media platforms in an as yet not understood way. They speculate this can be achieved because the networks within the world of the social media platform are connected to the networks in the social world beyond the social media platform and it is that connection which potentially allows the nonrepresentative sample from social media to be representative of the population. They also found that lexical analysis was the predominant method of transforming social media data into data comparable to survey data and that there is no clear reason why social media analysis and survey results diverge when they do diverge. The provide three tables of side-by-side comparisons between different features of surveys and social media. The first is about how participants understand the activity of responding versus posting, the second is about the nature of the data, and the last is about practical and ethical considerations. They conclude that social media analyses can only replace survey methods in cases

where the external gold standard for analyses is not a survey and can only supplement surveys in cases where the external gold standard is a survey. Finally, they acknowledge the ability of social media analyses to provide quick and affordable access to phenomena that both have and have not yet been measured through surveys and call for greater collaboration between data scientists and survey researchers.

In 2020 Amaya et al. answered the call to adapt the TSE for big data with the Total Error Framework (TEF) using the steps in the big data process to enumerate how their error components play out in big data analysis. The error components they identify are coverage error, sampling error, specification error, nonresponse/missing data error, measurement/content error, processing error, modeling/estimation error, and analytic error. In 2021, Sen et al. adapted the TEF to create the Total Error Framework for Digital Traces of Human Behavior on Online Platforms (TED-On), which can be seen in Figure 6, to create a common vocabulary to document, communicate and compare research, and encourage researchers to consider and document these sources of error in their designs and communications. They introduce a distinction between errors of measurement and errors of representation. They also introduce the terms "traces" and "users." The first applies to user-generated content or records of online activity, and the second to the ones who generate these traces. The error components under measurement are construct based. They are validity, platform affordances error, trace selection error, trace augmentation error, trace reduction error, and trace measurement error. The error components under representation are target population based and they are platform coverage error, user selection error, user augmentation error, user reduction error, and adjustment error. Exploring these sources of error more in depth will be important work in the coming years.

**Figure 6**

*Total Error Framework for Digital Traces of Human Behavior on Online Platforms (TED-On) (Sen et al. 2021)*



The predominance of participation in social media platforms over the last two decades has opened up spaces wherein hard-to-reach populations can be found. Instead of having to rely on finding

the physical locations where these populations exist, social media platforms create a virtual space where members of hard-to-reach populations can be found. When the target population is not one that is stigmatized, but simply one that does not have a clear sampling frame, researchers have found Facebook can serve as virtual space to find participants. Baltar and Brunet (2011), for example, utilized Facebook to identify Argentinean entrepreneurs living in Spain. Their traditional methods of snowball sampling yielded 80 responses while their Facebook enhanced method, where they searched affinity groups and private messaged users that appeared to fit their inclusion criteria, yielded 1,023 responses. Using Facebook clearly expanded their ability to seek and find participants. Similarly, Cowles et al. (2018) found a significant increase in study enrollment after advertising their study through Facebook ads. Their original methods yielded 123 participants and the Facebook Ads allowed them to recruit an additional 1,138 participants. The ability to purchase reach through the targeting algorithms of Facebook Ads can help researchers locate members of hard-to-reach populations in a more anonymous way through since the ads are only seen by the users and no one in the person's social network needs to know they decided to participate. Grow et al. (2022) conducted a large-scale cross-national online survey to compare the demographics of Facebook users to their self-reported demographics and found that Facebook's advertising platform is a valid way of targeting specific populations, but also suggest that researchers utilize a pre-survey with demographic questions that can help further identify if the participants meet inclusion criteria. Being able to link users to their demographics is a key feature of Facebook that can be helpful in weighting analysis results to more accurately generalize to a population, however, the potential of using Facebook to reach hard-to-reach populations is limited by what people are willing to reveal about themselves. Individuals that don't want to be identified as members of stigmatized groups may keep these aspects of their personas secret in the online space as they do in their everyday lives since the social network of Facebook mirrors a

person's true social network. Reddit, on the other hand, is an anonymous social media platform, which provides users the ability to reveal aspects of themselves they may not feel comfortable revealing in their social networks. This affords users the ability to congregate around aspects of their identity and topics that are important to them without fear of stigma and recrimination, and it affords researchers access to members of hard-to-reach populations and their conversations.

Several researchers have begun to publish methodology articles offering findings that support best practices for recruiting participants online. In 2017, Jamnik & Lane conducted an experiment to compare the survey results between a group of university students recruited through courses (a more traditional method in the field of psychology) and a group recruited through Reddit, specifically the subreddit r/samplesize. They found the Reddit group "provided high-quality data that were inexpensive and comparable to the responses gathered using undergraduate participants," and dubbed the practice to be a "promising tool for the field of psychological assessment, research, and evaluation" (p. 1). During the same year, Shatz (2017) published an article aimed at bringing online recruitment to light as an option for researchers. The advantages and limitations of recruiting through Reddit and guidelines for effective recruitment were provided.

In 2021, Luong & Lomanowska compared the results of surveys taken by a group recruited through Amazon Mechanical Turk (MTurk), a pay for work site, and through Reddit, again through the subreddit r/samplesize. They found the participants recruited through Reddit were demographically diverse, though White participants and participants from the United States were overrepresented. They found the Reddit participants were more internally motivated, though were similar in terms of altruism and motivation to gain self-knowledge and that the reliability and quality of the results were similar across most analyses between the two groups. Jones et al. (2021) recruited participants for a social work survey across four social media platforms: Facebook, LinkedIn, Reddit and Twitter. They found

Facebook and LinkedIn to be the best sources of recruitment for them, providing 66% and 29% of their 2,012 total survey participants. Reddit accounted for 4% and Twitter only 1%, though the authors acknowledge the results from Facebook and LinkedIn benefitted from the researchers existing social networks on those platforms. Richard, Sivo, and Orlowski, et al. (2021) randomly assigned participants to an asynchronous focus group on Reddit and an in-person focus group to determine if there was much of a difference in the number and variety of ideas generated. They found there was not a big difference and that quite a bit of overlap in ideas occurred between the two groups. Richard, Sivo, and Ford, et al. (2021) wrote, "A Guide to Conducting Online Focus Groups via Reddit," to establish protocols for online focus groups. Though the use of Reddit for research is growing exponentially, no papers have yet catalogued the most recent social science research utilizing Reddit. Cataloging the diversity of methods utilizing Reddit for research will help researchers learn more quickly from recent advancements in methodology and further heighten the pace of advancement and the spread of the practice across fields and disciplines.

**Research Questions**

1. How are social science researchers utilizing Reddit to further their research goals?

2. What should a researcher interested in using Reddit to further their research goals take into consideration when formulating a research plan?

**Methodology**

To answer these research questions, a review of published journal articles that utilized Reddit for social science research during the year 2021 was conducted. To gather the data set, an exhaustive search of the social science research databases ACM, EBSCOhost, JStor and SCOPUS for the search term "Reddit" was employed. Each article returned in the search was reviewed to ensure it met the inclusion criteria for this study. Included articles were published in a peer-reviewed academic journal

in 2021, published in the English language, used Reddit in a meaningful way for social science

research, and the full text of the article had to be available. Review articles, meta-analysis articles and

poster abstracts were excluded. The process used to cultivate the set of articles is shown in Figure 7.

**Figure 7**

*Flowchart Illustrating Process of Article Data Set Cultivation*



The final dataset was separated into two groups for analysis: articles that utilized Reddit for

recruitment and all other articles. The articles that utilize Reddit for recruitment represent a clear

subgroup of articles that utilize Reddit for research. Each article was then reviewed, and a list of

variables was recorded for each. The list of variables was generated based on the work of Proferes et

al. (2021) and through a constant comparative generation of new variables that arose while the articles were reviewed. Variables for recruitment articles included:

- the journal name and field;

- the country of origin of the authors;

- whether Reddit was the sole method used for recruitment;

- whether an incentive was offered;

- whether the recruitment message was included in the paper;

- the eventual method utilized once participants were recruited; and

- arguments for using Reddit, and limitations of using Reddit.

Variables for all other articles included:

- the journal name and field;

- the country of origin of the authors;

- the level of Reddit data analyzed;

- the number of units analyzed;

- the method of data scraping;

- the sampling method used;

- the method of analysis;

- whether IRB approval was obtained;

- arguments for using Reddit; and

- limitations of using Reddit.

The results of the variable collection were then tabulated.

## Results

### *Recruitment Articles*

Of the 49 articles that used Reddit for recruitment published in 2021, 26 of them came from the field of human-computer interaction. Eight came from the field of health, eight from psychiatry or psychology, three from communication, two were multidisciplinary, and one each from computer science, cultural studies, and education. The percentage of articles for each field can be seen in Figure 8. Overwhelming, authors of the recruitment articles were from the United States (39), followed by Canada (four), and China, Denmark, Finland, Italy, Japan, Korea, Munich, and the Netherlands with one each. Arguments given for using Reddit to recruit included access to stigmatized topics/groups and the assertion that the anonymity of the site allows for participants to be more honest.

**Figure 8**

*Fields Publishing Articles Using Reddit for Recruitment*

**How Reddit Was Used.** Fourteen out of 49 articles used only Reddit for recruiting, 33 out of 49 used Reddit in addition to other recruitment methods like Facebook, Twitter, and snowball sampling, and two did not report. Twenty-five out of 49 reported using an incentive to help recruit participants, 4 out of 49 stated they did not use an incentive and 24 did not report. Most researchers were recruiting survey participants (26/49), though many were recruiting for interview participants as well (15/49), and some (5/49) were recruiting for both – usually conducting surveys initially and choosing some survey participants to follow-up with interviews. Most researchers recruited between 10-50 participants (23/49), with four recruiting under 10 participants, three recruiting 51-100 participants, two recruiting 500-1,000 participants and two recruiting over 1,000. The percentage of articles recruiting these different numbers of participants can be seen in Figure 9.

**Figure 9**

*Number of Participants Recruited*

Most researchers posted a link to their survey or request for interview participants directly to a subreddit, but some sent direct messages to specific users inviting them to take their survey. With this method, it is possible to generate a sample frame of users based on some criteria and then select a random sample to invite for participation though this method did not produce a remarkably high response rate. Triggs et al.(2021) purposively recruited participants for their study using direct messages to Reddit users whose post history showed they met the desired characteristics for the research questions, and Bhuiyan et al. (2021) sent direct messages to the most active members of subreddits of interest. Rajadesingan et al. (2021) also used direct messaging for recruitment and found they had greater success using an established account of one of the authors to send their messages than they did with a newly established account that had been set up expressly for that purpose. Even so, they found their participation rate was under 10% and they ended up using the subreddit r/PaidStudies to find more participants.

**Common Limitations.** A few limitations were commonly cited by authors of recruitment articles. Most cited, was the lack of diversity in the population they recruited from, and the fact that their results may be biased towards the kind of people that are members of the Reddit community and may not be fully representative of the views of others who are not members. Self-selection bias, the non-representativeness of their samples, and the inability to validate reported information were also cited.

### *All Other Articles*

Of the 120 other articles that used Reddit for social science research published in 2021, 23 of them came from the field of health and 23 from the fields of psychiatry or psychology. Twenty-one came from the field of communication, 15 from human-computer interaction, 13 from arts and humanities, six from computer science, four from education, and three each from multidisciplinary and

library and information science. The percentage of articles represented by each field can be seen in

Figure 10. Again, the authors of the articles were overwhelmingly from the United States (67),

followed by the United Kingdom (15), and Canada (12). Australia had six, Germany five, Italy four,

Finland and Switzerland three, India, Poland, and South Africa two, and Brazil, Korea, New Zealand,

Norway, Russia, Singapore, and the Netherlands all one. As can be seen in Figure 11, most researchers

did not report whether they obtained IRB approval (87 articles or 72.5%), 23 articles reported they

were exempt (19.2%), and 10 reported obtaining approval (8.3%). In 17 articles, the authors discussed

the ethical considerations surrounding the use of extant social media data.

**Figure 10**

*Fields Publishing Articles Using Reddit*

**Figure 11**

*IRB Review Status of Articles*



**Arguments for Using Reddit**. The greatest argument for using Reddit is the non-networked

anonymity of the platform. Two articles specifically mention they were able to gain access to hard-to-

reach groups by studying Reddit data and this is evident by the number of articles in the data set that

have to do with topics like substance abuse and domestic violence. Four articles talked about how the

anonymity of Reddit afforded them access to discussion around sensitive and/or stigmatizing topics.

As Lyons and Brewer put it, "the veil of anonymity and shared experiences make it easier for the users

to openly talk about stigmatizing issues that may be more difficult to discuss face-to-face" (2021).

Three articles asserted the anonymity of Reddit yields more honest discussions, and four articles

pointed out that the unstructured nature of Reddit data, the fact that it represents naturally occurring

conversation, means it is free from researcher priorities or assumptions. Several examples of papers

analyzing sensitive topics and/or hard to reach subjects are: Psychedelic substance use Psychedelic

substance use in the reddit psychonaut community: A qualitative study on motives and modalities

(Pestana et al. 2021), "Therapeutic benefit with caveats? Analyzing social media data to understand the complexities of kratom use" (Smith, et al., 2021), and "A content analysis of reddit users' perspectives on reasons for not following through with a suicide attempt" (Mason et al., 2021). Utilizing extant data saves researchers from having to build and distribute surveys or recruit participants for interviews or focus groups. As Currin-McCulloch et al. (2021) point out "focus groups require extensive resources and introduce methodological limitations including experimenter bias." Finally, the organization of Reddit into subreddits, or interest-based topic groups, inherently builds communities around topics that may be of interest to researchers.

**How Reddit Was Used.** Many social science researchers are using Reddit to scrape and analyze data with a large amount of variability in their methods of scraping and analysis. Some of that variability will be explored here including different methods for finding relevant content, extracting data, and analyzing data.

One area where there was a wide range of variability was methods for what Japec et al. call trace selection, the method of querying all available data on the site to access only those units that are germane to the research topic. Basically, how the data is sampled. This step in the big data process is open to what they call trace selection error, or the potential to include units that should not be included or exclude units that should not be included. Research published in 2021 showed a wide range in skills related to trace selection. Some researchers relied on the Reddit platform's search feature to filter content on the site to meet their inclusion criteria. While this method is certainly easy, it delegates the task of filtering to the platform which has its own goals for displaying content to users which may not align with the research goals. Additionally, as Smith, et al. (2021) noted that, the platform puts page limits on the Reddit search feature which limited how far back in time the search results would display. They found some creative ways to narrow their searches and conduct more of them as a work

around to this problem. The platform does make it possible to do some sorting of search returns, but the features are limited, and researchers with more advanced data skills tended to sort data after scraping. The most advanced methods use data science skills to search the platform for relevant content. One such example is the Lexico-Semantic Similarity Filter, a natural language processing method, which searches for desired word embeddings, or the connections between words in a text, across the platform (Gong et al., 2021). This tool, which has the ability to quickly look across the entire site has a great chance of finding all germane topic than using the built in search feature on the platform, however these tools are not frequently available to the average social science researcher.

Reddit's API can be accessed through PRAW, a package in Python, or through RedditExtractoR or rreddit in R. Consent must be gained from Reddit to directly retrieve data using their API. Gaining consent is a straightforward process. Documentation is available online that requires varying levels of understanding to execute and a rudimentary understanding of Python or R should allow a researcher an entry point into data scraping. However, tailoring available scripts to complete more complicated and/or individualized parameters is much more difficult and requires a deeper level of understanding. As an acknowledgement that these data science skills are not readily available to most social science researchers, Hintz and Betts (2022) provide a ready to use script for researchers to scrape data from Reddit within the R environment that requires little to no knowledge of the platform to successfully use. A little over half the articles did not explain how they scraped data from Reddit. Among the most popular methods of those who did explain were the Python Reddit API Wrapper (PRAW), the Pushshift.io dataset, and manual extracting methods like copy and paste. Less frequently stated methods included accessing the Reddit API (with no elaboration on how this was done), RedditExtractoR, Google's BigQueary, a previously published Reddit dataset, BuzzSumo, Synthesio, and the Google Chrome extension Web Scraper. It is also possible to query the Pushshift

Reddit Dataset, which is essentially a mirrored copy of Reddit. The argument for using the Pushshift

API rather than the Reddit API is that it, "makes it much easier for researchers to query and retrieve

historical Reddit data, provides extended functionality by providing full-text search against comments

and submissions, and has larger single query limits," (Baumgartner et al., 2021). For those with the

skills, it is also possible to use post IDs provided by the Reddit API to retrieve raw JSON of all posts

and comments. Finally, this author came across a tool called thread downloader that generates a PDF

of any thread url which may be of interest for those doing qualitative analysis. Indeed, scraping data

from Reddit proves to be the greatest barrier to research use. The methods of data extraction used in

the articles in the dataset can be seen in Figure 12.

**Figure 12**

*Methods of Data Extraction Used*

Once scraped, most researchers are using qualitative techniques (67 articles or 61.5%) to analyze their data as can be seen in Figure 13. Mixed-methods analysis was used in 27 articles (24.8%) and 15 used quantitative analysis methods (13.8%). The top three methods of analysis are content analysis (28 articles), thematic analysis (26 articles), and statistical analysis (16 articles). The complete list of methods can be found in Table 1. Note that some articles used multiple methods, so the total number of articles listed is greater than 120. Most researchers analyzed original posts (34 articles), threads – meaning original posts and their comment tree (11 articles), comments (10 articles), entire subreddits (nine articles), and posts and top comments only (four articles). Of those that studied original posts, 23 articles analyzed between 1-1,000 posts, five analyzed 1,001-10,000 posts, four analyzed 10,001-100,000, four analyzed over 100,000 and five did not report the number they analyzed. Of those that analyzed threads, 11 analyzed 100 or less threads, seven analyzed 501-1000 threads, three analyzed 501-1000, six analyzed 10001-5,000, and seven analyzed over 5,000. Of those that analyzed comments, five analyzed 1 – 1,000 comments, four analyzed 1,001-5,000, and two analyzed over 50,000.

Many tools are now available to aid in the automation of qualitative analysis, should one be interested. Many of these tools allow quantifiable variables to be assigned to text allowing researchers to summarize distinct characteristics of text and make comparisons between groups of texts on these variables. A few of the tools that showed up in the review of these articles are presented here, however this is not an exhaustive list and more tools will certainly become available over time. One is the Valence Aware Dictionary and sEntiment Reasoner (VADER) which measures emotional valence in text. Another is the Lexicoder Sentiment Dictionary which counts the number of positive or negative words (or group of words) from a dictionary containing more than 4,500 words or expressions (Young & Soroka, 2012). Linguistic Inquiry and Word Count (LIWC) is also available with an expanded

**Figure 13**

*Methodological Orientation of the Articles*



**Table 1**

*Analysis Methods of Social Science Articles Utilizing Reddit Published in 2021*

| Method | Number Of Articles |
| --- | --- |
| Content Analysis | 28 |
| Thematic Analysis | 26 |
| Statistical Analysis | 16 |
| Natural Language Processing | 9 |
| Grounded Theory | 8 |
| Time Series Analysis | 7 |
| Digital Ethnography | 3 |
| Network Analysis | 3 |

| | |
|---|---|
| Conversation Analysis | 2 |
| Meaning Extraction Method | 2 |
| Polarization Measurement | 2 |
| Contrapuntal Analysis, Conversation Analysis, Document Analysis, Ethnography, Image Analysis, Linguistic Style Matching, Semantic Network Analysis, Tool Creation, Triangulation | 1 Each |

series of capabilities in the LIWC-22 version. LIWC-22 allows users to ascertain different quantitative scores on the thoughts, feelings and personality displayed through text, generates word clouds, performs topic modeling based on the meaning extraction method, calculates language style matching, and models the art of narrative to show how stories unfold. The Stanford CoreNLP Natural Language Processing Toolkit was designed to make natural language processing more accessible (Manning et al., 2014). The Affective Norms for English Words (ANEW) tool applies affective ratings to text (Bradley & Lang, 1999). Finally, Google's Perspective API (Application Programming Interface) uses machine learning to derive a toxicity score for text, the SentimentR package in R can quantify the sentiment of text, and the Flesch Reading Ease scale measures text complexity (Flesch, 1948) and can be used through the qantenda package in R.

**Common Limitations.** The most cited limitation of using Reddit data (33 articles) is the lack of available demographic data. The single greatest feature of Reddit for researchers, its anonymity, is also its greatest drawback. Without available demographic data it's impossible to know if views expressed on Reddit are representative of the larger population. Seven articles cited the lack of generalizability of their findings as a limitation, six said their findings were only representative of

Reddit users and three of those noted that people who use Reddit have access to the internet thereby automatically excluding from the results the view of anyone without internet access. Eight articles mentioned self-selection bias as a limitation, two lamented the inability to confirm the identity of users, eight pointed out that there could be imposters in subreddit communities geared toward certain group affiliations participating as members, when in truth they are not. Several researchers mentioned limitations of finding the right data on Reddit – two because their data set was limited to the search terms utilized and one because the Reddit algorithm dictated what they scraped. One article pointed out that deleted posts are not included in the dataset and therefore may bias their findings. Finally, four articles discussed their inability to member-check their findings.

## Discussion

Understanding the advantages and potential limitations of using Reddit can help researchers determine whether using Reddit is right for them. Reddit lends itself well to qualitative and exploratory methodology as evidenced by the high number of qualitative analyses represented in the articles reviewed for this paper, though many are finding ways to apply complex statistical analyses to Reddit data as well. Gaining insight into new methodologies for collecting and analyzing data can expand researchers' thinking and open new possibilities for research as they consider data sources and designs that ask and answer new questions. Understanding the ability to reach subjects and topics and depths of conversation that were previously unavailable allows researchers to consider deeper questions. Furthermore, the constant flow of conversation allows for changes to be tracked over time.

Researchers who want to stick their toe in the Reddit waters need to do their homework. It is strongly recommended to set up a Reddit account and become familiar with the platform. Selecting a few subreddits based on personal interests to join is a fun way to start to get to know how things work. Casually following topics related to research interests may spark ideas for research as greater

understanding of what data is available occurs. Similarly, identifying a few subreddits and poking around a little - glancing through the post history, checking out the top posts of the day and the top posts of all time – will give a better sense of the data available over time, though Reddit posts only show for a limited period of back history when scrolling as a user. Historical data is accessible through the pushshift.io dataset and through direct URLs for a post once obtained.

Researchers interested in using Reddit for data collection must consider a few key things. Some subreddits may expressly prohibit data from being taken for research purposes or may require prior permission from moderators to do so. It is important to check the rules of each subreddit from which you want to extract data before laying out a research plan. It is also a good idea to check the most recent terms of service for the Reddit platform and make sure you aren't planning to act in violation of them (Silberman & Record, 2021). As using Reddit for research continues to grow in popularity, it's possible that more subreddits and perhaps even Reddit as a whole may become more aware of researcher activity and take steps to make it more difficult. Reddit may consider adding explicit guidelines regarding the use of Reddit data for research purposes to its Terms of Service and/or to it's Reddit API Terms of Use. And for those interested in the ethical consideration of utilizing research for Reddit, see Proferes et al. (2021) for a thorough discussion. Some researchers take great care to ensure the continued anonymity of Reddit users by not including the name of the subreddit they scraped from, nor any of the usernames, and/or by altering quotes of text to a point where a google search does not reveal their location online (Brewer et al., 2021). Litherland and Mórch (2021) informed members of a subreddit through a post that the study was taking place and sent a direct message to each person who would be quoted in the final paper giving them the ability to withdraw participation. None choose to withdraw. For the work done in Chapter 3 of this dissertation, this author noticed that some of the posts scraped using the Pushshift.io mirror copy of Reddit had

been deleted by the individual users after they'd been copied over and it was decided these should be removed from the dataset to honor the wishes of the author that their content not live online for further public consumption. Considering the ethical implications of utilizing data that users have not created for the purpose of research consumption is a worthwhile endeavor.

Researchers interested in utilizing Reddit for recruitment purposes have a separate set of considerations to keep in mind. To recruit through Reddit, at least one person on the research team will have to have a Reddit user account and manage the group's communication through Reddit. One feature of the Reddit platform, despite its anonymity, is the ability to see how long a user has been on Reddit, how long they have been a member of a particular subreddit, and see the history of their user activity. The Reddit community will be able to determine, and may be less trusting of, an account that is very new. Therefore, establishing a Reddit identity may be a useful endeavor. Keep in mind it is possible to have multiple Reddit accounts so a personal Reddit identity can be kept separate from a professional Reddit identity. It is recommended to seek out the rules of a targeted subreddit before posting, potentially even reaching out to moderators if the rules are unclear about recruitment messages in the subreddit to ensure messages are not removed by moderators. It is also recommended to consider using subreddits dedicated to research recruitment like r/samplesize and r/paidstudies. In addition to multiple ways of recruiting on the Reddit platform, many of the articles surveyed for this paper recommend utilizing multiple methods of recruitment in addition to Reddit, like Facebook or other more specialized social media platforms, to broaden the diversity of the research participants. Finally, for those wanting to member check results, consider submitting results back to the subreddit where recruiting took place to gather feedback (Krsmanovic and Dean, 2021).

Lastly, becoming familiar with the work of others who have utilized Reddit for research is helpful. This paper provides a broad overview of research conducted in 2021 that utilized Reddit for

research. Focusing in on an area of interest and reading some of the papers referenced here or conducting a search using the key word "Reddit" and a desired methodology should return a solid base of articles to being exploring. As time goes on, more and more innovation will come to methods utilizing Reddit for research as researchers become savvier in this area. Researchers may try new ways of guiding the organic conversation on the platform by posting topics of discussion themselves to try and spur conversation in the direction of their research interests. For anyone looking to take such a route, you'll notice that posts tend to get more conversation when the original poster responds to the comments being left to help keep the conversation going. This may be thought of as being similar to moderator a focus group. Taking a more interactive role in guiding the conversation may, however, turn a passive analysis of extant social media data which may not require IRB review into active research with participants which may require IRB review.

At this point, no clear guidelines exist around reporting norms for articles utilizing Reddit, though the growing prevalence of research using this platform necessitates forging one. A short checklist to begin creating such guidelines is presented here. This article proposes that published articles utilizing Reddit in their research should include the following items in their methodology section:

1. Consider the sources of error on the TED-On and transparently report all research decisions that contribute to potential sources of error;

2. State whether IRB review was necessary and/or obtained;

3. Report what ethical considerations were taken throughout the research process to protect the identity of Reddit users whose conversation were used in the study;

4. Report the date data was collected;

5. Report the time frame of the data collected;

6. Report the level of the data collected and analyzed (subreddit, threads, posts, comments, etc.);

7. Report the number of units (subreddits, threads, posts, comments, etc.) that were analyzed;

8. Provide a detailed report of the data extraction method utilized that goes beyond simply stating a software package; and

9. Provide a detailed report of the sampling methods used, if applicable.

A recommended practice is to post any code used for the paper to github, an open-source website community where software developers collaborate (Biester et al., 2021).

**Limitations**

Although an exhaustive search of several of the largest social science databases was conducted, it is possible some articles were missed. Additionally, it is possible the author's interpretation of what articles were using Reddit in a significant way to conduct social science research may have been different from someone else's. This paper did not include articles from a large body of work being done in the machine learning and natural language processing area because these were data science focused, however, understanding the full extent of what is possible in this field will be important for social science researchers in years to come. Business and industry leaders are well ahead of social science researchers in terms of harnessing natural language processing to analyze extant social media data for their purposes.

**Conclusion**

Several features of the Reddit platform afford researchers the ability to ask and answer new research questions. The anonymous nature of the Reddit platform facilitates access to hard-to-reach populations and discussions about sensitive and stigmatizing topics that would be either extremely difficult to access or completely inaccessible to social science researchers otherwise. The fact that Reddit is a topic-based community means the platform is self-organized around specific topics of

53

interest that may make it easier to filter all the data on the site. This paper provided an overview of current social science research that is leveraging Reddit thus familiarizing social science researchers with some data science methods that may be useful to them. These new skills can be integrated with social science research methodologies to ask and answer new and exciting questions which hadn't been answerable before. Reddit is an especially appropriate platform choice for conducting exploratory research, especially on hard-to-reach populations and/or stigmatizing topics; and to recruit participants for studies, especially participants from hard-to-reach populations. New methods that leverage this platform will likely continue to evolve as more social scientist researchers travel down this path.

**Suggestions for Future Research**

As the use of Reddit for social science research continues to grow and Reddit itself continues to grow, it will be important to continue to research novel methodologies by providing surveys of current methodologies such as this paper has provided. Further research should be done on the demographics of Reddit users, including variation of demographics in different subreddits to better understand who is using Reddit. Entwistle et al. found they were able to mine demographic data for 80% of the users in the dataset they analyzed allowing them draw conclusions by demographic group (2021). This method could be applied more broadly across the platform. Additionally, more research is needed into the different sources of error in the TED-On framework, such as how different sampling methods affect the results of analyses of Reddit data, to help researchers make design choices that yield valid and reliable results.

# References

Amaya, A., Biemer, P. P., & Kinyon, D. (2020). *Journal of Statistics and Methodology, 8*(1), 89-119.

    https://doi.org/10.1093/jssam/smz056

Amaya, A., Bach, R., Keusch, F., & Kreuter, F. (2019). New data sources in social science research:

    Things to know before working with Reddit data. *Social Science Computer Review, 39*(5), 943-

    960. doi: 10.1177/0894439319893305

Baker, R., Brick, J. M., Bates, N. A., Battaglia, M., Couper, M. P., Dever, J. A., Gile, K. J.,

    Tourangeau, R. (2013). Summary report of the AAPOR task force on non-probability

    sampling. *Journal of Survey Statistics and Methodology, 1*(2), 90-143.

    https://doi.org/10.1093/jssam/smt008

Baltar, F., & Brunet, I. (2011). Social research 2.0: Virtual snowball sampling method using

    Facebook. *Internet Research, 22*(1), 57-74. doi: 10.1108/10662241211199960

Barthel, M., Stocking, G., Holcomb, J., and Mitchell. A. (2016). Nearly eight-in-ten Reddit users get

    news on the site. Pew Research Center.

    https://www.pewresearch.org/journalism/2016/02/25/seven-in-ten-reddit-users-get-news-on-

    the-site/

Baumgartner, J., Zannettous, S., Keegan, B., Squire, M, & Blackburn, J. (2020). The Pushshift Reddit

    dataset. https://arxiv.org/pdf/2001.08435.pdf

Bhuiyan, M. M., Whitley, H., Horning, M., Lee, S. W., Mitra, T. (2021). Designing transparency cues

    in online news platforms to promote trust: Journalists' & consumers' perspectives.

    *Proceedings of the ACM on Human-Computing Interaction*, *5*(CSCW2), 1-31.

    https://doi.org/10.1145/3479539

Biester, L., Matton, K., Rajendran, J., Provost, E. M., & Mihalcea, R. (2021). Understanding the

impact of COVID-19 on online mental health forums. *ACM Transactions on Management

Information Systems*, 12(4), 1-28. https://doi.org/10.1145/3458770

Bradley, M. M., & Lang, P.J. (1999). *Affective norms for English words (ANEW): Instruction manual

and affective ratings.* Technical Report C-1, The Center for Research in Psychophysiology,

University of Florida.

Brewer, G., Centifanti, L., Castro Caicedo, J., Huxley, G., Peddie, C., Stratton, K., & Lyons, M.

(2021). Experiences of mental distress during COVID-19: Thematic analysis of discussion

forum posts for anxiety, depression, and obsessive-compulsive disorder. *Illness, Crisis & Loss,*

0(0), 1-17. doi:10.1177/10541373211023951

Cowles, C., Berk, S., & Siddiqi, B. (2018). Using Facebook ads to recruit clinical study participants.

*Applied Clinical Trials, 27*(12), 14-17.

Currin-McCulloch, J., Stanton, A., Boyd, R., Neaves, M., & Jones, B. (2021). Understanding breast

cancer survivors' information-seeking behaviours and overall experiences: a comparison of

themes derived from social media posts and focus groups. *Psychology & Health,* 36(7), 810-

827. doi:10.1080/08870446.2020.1792903

Dixon, S. (2022). *Reddit usage reach in the United States 2021, by age group*. Statista.

https://www.statista.com/statistics/261766/share-of-us-internet-users-who-use-reddit-by-age-

group/

Dixon, S. (2022). *Reddit usage reach in the United States 2021, by ethnicity.* Statista.

https://www.statista.com/statistics/261770/share-of-us-internet-users-who-use-reddit-by-

ethnicity/

Dixon, S. (2022). *Reddit usage reach in the United States 2021, by gender.* Statista.

    https://www.statista.com/statistics/261765/share-of-us-internet-users-who-use-reddit-by-

    gender/

Entwistle, C., Horn, A. B., Meier, T., & Boyd, R. L. (2021). Dirty laundry: The nature and substance

    of seeking relationship help from strangers online. *Journal of Personal and Social*

    *Relationships*, 38(12), 3472-3496. https://doi.org/10.1177/02654075211046635

Flesch, R. (1948). A new readability yardstick. *Journal of Applied Psychology, 32*(3), 221–233.

    https://doi.org/10.1037/h0057532

Gaudette, T., Scrivens, R., Davies, G., & Frank, R. (2021). Upvoting extremism: Collective identify

    formation and the extreme right on Reddit. *new media & society*, *23*(12), 3491-3508. doi:

    10.1177/1461444820958123

Gong, C., Saha, K., & Chancellor, S. (2021). "The smartest decision for my future": Social- media

    reveals challenges and stress during post-college life transitions. *ACM Human-Computer*

    *Interaction,* 5(CSCW3), 1-29. doi:10.1145/3476039

Graham, T. & Rodriguez, A. (2021). Sociomateriality of rating and ranking devices on social media: A

    case study of Reddit's voting practices. *Social Media + Society*, *7*(3),  doi:

    10.1177/20563051211047667

Grow, A., Perrotta, D., Del Fava, E., Cimentada, J., Rampazzo, F., Gil-Clavel, S., Zagheni, E., Flores,

    R. D., Ventura, I., Weber, I. (2022). Is Facebook's advertising data accurate enough for use in

    social science research? Insights from a cross-national online survey. *Journal of the Royal*

    *Statistical Society, Series A, 185*(S2), 343-363. doi: 10.1111/rssa.12948.

Guber, D. L. (2021). Public opinion and the classical tradition: Redux in the digital age. *Public*

    *Opinion Quarterly, 85*(4), 1103-1127. doi: 10.1093/poq/nfab053

Hintz, E. A. & Betts, T. (2022). Reddit in communication research: Current status, future directions

    and best practices. Annals of the International Communication Association, 46(2), 116-133.

    doi: 10.1080/23808985.2022.2064325.

Jamnik, M. R. & Lane, D. J. (2017). The use of Reddit as an inexpensive source for high-quality data.

    *Practical Assessment, Research, and Evaluation*. 22(5). https://doi.org/10.7275/j18t-c009

Japec, L., Kreuter, F., Berg, M., Biemer, P., Decker, P., Lampe, C., Lane, J., O'Neil, C., & Usher, A.

    (2015). Big data in survey research. *Public Opinion Quarterly, 79*(4), 839-880.

    https://doi.org/10.1093/poq/nfv039

Jones, A., Walters, J., & Brown, A. (2021). Participant recruitment in social work: a social media

    approach. *Social Work Research*, 44(4), 247-255. doi:10.1093/swr/svaa017

Krsmanovic, A., & Dean, M. (2021). How women suffering from endometriosis disclose about their

    disorder at work. *Health Communication*, 37(8), 992-1003. doi:10.1080/10410236.1880053

Laney, D. (2001). 3-D management: Controlling data volume, velocity, and variety. META Group

    Research Note, February 6.

Litherland, K. T., & Morch, A. I. (2021). Instruction vs. Emergence on r/place: Understanding the

    growth and control of evolving artifacts in mass collaboration. Computers in Human Behavior,

    122. https://doi.org/10.1016/j.chb.2021.106845

Luong, R., & Lomanowska, A., M. (2021). Evaluating Reddit as a crowdsourcing platform for

    psychology research projects. *Teaching of Psychology*, 49(4).

    https://doi.org/10.1177/00986283211020739

Lyons, M., & Brewer, G. (2021). Experiences of intimate partner violence during lockdown and the

    COVID-19 pandemic. *Journal of Family Violence*, 37, 969-077.

    https://doi.org/10.1007/s10896-021-00260-x

Manning, C. D., Surdeanu, M., Bauer, J., Finkel, J., Bethard, S. J., & McClosky, D. (2014). *The Stanford CoreNLP Natural Language Processing Toolkit. In Proceedings of 52nd Annual meeting of the Association for Computational Linguistics: System Demonstrations.* Association for Computational Linguistics. https://nlp.stanford.edu/pubs/StanfordCoreNlp2014.pdf

Marpsat, M., & Razafindratsima, N. (2010). Survey methods for hard-to-reach populations: Introduction to the special issue. *Methodolgoical Innovations Online, 5*(2), 3-16. https://doi.org/10.4256/mio.2010.0014

Mason, A., Jang, K., Morley, K., Scarf, D., Collings, S. C., & Riordan, B. C. (2021). A content analysis of Reddit users' perspectives on reasons for not following through with a suicide attempt. *Cyberpsychology, Behavior, and Social Networking*, *24*(10). https://doi.org/10.1089/cyber.2020.0521

Medvedev, A., Lambiotte, R., & Delvenne, J. (2019). The anatomy of Reddit: An overview of academic research. *Dynamics On and Of Complex Networks III, Springer Proceedings in Complexity*. https://doi.org/10.1007/978-3-030-14683-2_9

Pestana, J., Beccaria, F., & Petrilli, E. (2021) Psychedelic substance use in the Reddit psychonaut community. A qualitative study on motives and modalities. *Drugs and Alcohol Today*, *21*(2), 112-123. doi.10.1108/DAT-03-2020-0016

Proferes, N., Jones, N., Gilbert, S., Fiesler, C., & Zimmer, M. (2021) Studying Reddit: a systematic overview of disciplines, approaches, methods, and ethics. *Social Media + Society*. https://doi.org/10.1177/20563051211019004

Rajadesingan, A., Duran, C., Resnick, P, & Budak, C. (2021) 'Walking into a fire hoping you don't catch': Strategies and designs to facilitate cross-partisan online discussions. *Proceedings of the ACM on Human-Computing Interaction*, *5*(CSCW2), 1-30. https://doi.org/10.1145/3479537

*Reddit*. (2022, November 13). In *Wikipedia*. https://en.wikipedia.org/wiki/Reddit

Reddit (2022, November 13). *About Reddit.* Reddit. https://www.redditinc.com/

Reddit. (2022, November 13b). *Reddiquette.* Reddit. https://www.reddithelp.com/hc/en-us/articles/205926439

Reveilhac, M., Steinmetz, S., & Morselli, D. (2022). A systematic literature review of how and whether social media data can complement traditional survey data to study public opinion. *Multimedia Tools and Applications, 81*. 10107-10142. https://doi.org/10.1007/s11042-022-12101-0

Richard, B., Sivo, S. A., Ford, R. C., Murphy, J., Boote, D. N., Witta, E., & Orlowski, M. (2021). A guide to conducting online focus groups via Reddit. *International Journal of Qualitative Methods*, 20, 1-9. doi:10.1177/16094069211012217.

Richard, B., Sivo, S. A., Orlowski, M., Ford, R.C., Murphy, J., Boote, D. N., & Willa, E. L. (2021). Qualitative research via focus groups: Will going online affect the diversity of your findings? *Cornell Hospitality Quarterly*, 62(1), 32-45. doi:10.1177/1938965520967769.

Schober, M.F., Pasek, J., Guggenheim, L., Lampe, C., Conrad, F.G. (2016). Social media analyses for social measurement. *Public Opinion Quarterly, 80*(1), 180-211. doi: 10.1093/poq/nfv048

Sen, I., Flöck, F., Weller, K., Weib, B., & Wagner, C. (2021). *Public Opinion Quarterly, 85*(S1), 399-422. https://doi.org/10.1093/poq/nfab018

Shatz, I. (2017). Fast, free, and targeted: Reddit as a source for recruiting participants online. *Social Science Computer Review*. 35(4). doi:10.1177/0894439316650163

Silberman, W. R. & Record, R. A. (2021). We post it, u Reddit: Exploring the potential of Reddit for health interventions targeting college populations. *Journal of Health Communication*, 26(6), 381-390. doi: 10.1080/10810730.2021.1949648

Smith, K. E., Rogers, J. M., Schriefer, D., & Grundmann, O. (2021) Therapeutic benefit with caveats?: Analyzing social media data to understand the complexities of kratom use. *Drug and Alcohol Dependence,* 226. https://doi.org/10.1016/j.drugalcdep.2021.108879

Triggs, A. H., Moller, K., Neumayer, C. (2021). Context collapse and anonymity among queer Reddit users. *New Media & Society*, *23*(1), 5-21). https://doi.org/10.1177/1461444819890353

Young, L., & Soroka, S. (2012) Affective news: The automated coding of sentiment in political texts. *Political Communication*, 29(2), 205-231. Doi.10.1080/10584609.2012.671234

CHAPTER III

HOW METHODS OF SAMPLING AFFECT THE RESULTS OF INDUCTIVE CONTENT
ANALYSIS: A CASE STUDY UTILIZING DATA FROM THE SOCIAL MEDIA PLATFORM
REDDIT

Advancements in artificial intelligence (AI) continue to expand the realm of what's possible

for computer-assisted qualitative analysis of Big Data, including social media data such as that found

on the social media platform, Reddit, however, qualitative researchers have doubts about AI's ability

to replace human instruments. Through interviews with fifteen qualitative researchers, Fueston and

Brubaker (2021) found the researchers lacked trust in the ability of AI to capture the same

understanding of nuance as a person and had concerns about a lack of transparency regarding the

biases baked into AI analysis, though they were amenable to incorporating certain tools into their

processes to supplement and support, rather than replace human analysis. Nelson et al. (2021) echo

these findings when they suggest based on the results of a comparison of a hand-coded content

analysis of news articles to three different computer-assisted methods (dictionary, supervised machine

learning, and unsupervised machine learning) that computer-assisted methods should complement but

not go so far as replace traditional approaches.

Though qualitative researchers may not be ready to hand over analysis tasks to AI, they may

benefit from the use of data science methods during the data collection phase and turn to traditional

qualitative methods in the analysis phase – a hybrid approach to the Big Data Process outlined by

Japec et al. (2016) which can be seen in Figure 14. This approach allows qualitative researchers to

access a trove of qualitative data on a full range of topics, including opinions of hard-to-reach

populations and opinions on stigmatizing topics. Human instruments have less processing power than

AI and qualitative researchers must make choices about how to filter, or sample from, all the available

data to create a more manageable dataset for human-coded analysis.

**Figure 14**

*Hybrid Approach to the Big Data Process*



The Total Error Framework for Digital Traces of Human Behavior on Online Platforms (TED-On) provides a map of potential sources of error when applying the big data process to social media data (Sen et al., 2021). This paper will focus on one type of error identified in the framework, trace selection error. As Sen et al. explain, "typically, researchers query the available traces to select those that are broadly relevant to the construct of interest. To the extent these queries fail to capture all relevant posts or include irrelevant posts, they create a type of measurement error we call trace selection error" (410). When social science researchers perform this step, they often must also sample from all the available traces that were netted in these queries, to reduce the dataset to a manageable size for analysis which provides yet another source of error. Based on the findings in Paper 1 that many researchers rely on non-probabilistic methods of sampling from all available units of analysis

which meet their criteria, it becomes necessary to investigate what the effects of these common sampling methods are on the results of qualitative analysis.

This paper seeks to understand how analysis results derived from a sample of the most upvoted posts and the most commented on posts from a given time period about a certain topic compare to analysis of the universe of posts from which the samples were taken. Although researchers are increasingly turning to Reddit as a valid source of research data, little has been published about best practices for sampling Reddit data. This paper begins to build researchers' understanding of how sampling choices affect analyses of Reddit data.

**Background: What is Reddit?**

Reddit is a social news aggregation, content rating and discussion website that as of March 2022 is the "9th-most-visited website in the world and the 6th most-visited website in the United States" (Wikipedia, 2022). According to Reddit, the platform has over 50 million unique users active daily, over 100,000 active communities, and includes over 13 billion posts and comments (2022a). Reddit is an anonymous platform which affords users the ability to share their opinions in an unfiltered way without fear of recrimination. Users of Reddit, called "redditors," participate in affinity groups called "subreddits." Subreddits are dedicated to specific topics and adhere to rules set by moderators and users frequently have conversations about sensitive topics, the content of which would either not normally be available to researchers or would come at the cost of extensive relationship building.

Several surveys have been conducted which draw inferences regarding the primary users of the site. In 2016, Pew Research Center found that about 7% of U.S. adults were using Reddit and that 69% of those users were men, 58% were between the ages of 18-29, 33% were between 30-49, 7% between 50-64, 1% were 65 or older, and 63% were White (Barthel et al.). In 2021, Statista found that 23% of all male adults, 12% of all female adults, 17% of White adults and 17% of Black adults in the U.S. use

Reddit (Dixon, 2022b and c). They found that 36% of users are aged between 18-29, 22% are aged 30-49, 10% are aged 50-64 and 3% are 65 or older (Dixon, 2022a). Though users tend to be younger and male on the aggregate, individual subreddits have their own cultures and may reflect different demographic groups. Seiter and Brophy (2022) found that users on Reddit provided social support more frequently than YouTube and Facebook users and also engaged in less aggressive communication. They believe this is due to both the anonymity of Reddit and the rules for civil discourse, or "Reddiquette," found on the Reddit platform.

**Literature Review**

Users "upvote" and "downvote" posts and those receiving the most upvotes are elevated by the site's algorithms within subreddits and sitewide. "Upvotes" are similar to "likes" on Facebook but have a slightly different intended use. According to "Reddiquette," or the "informal expression of the values of many redditors, as written by redditors themselves," users are encouraged to vote on content according to the following guidelines: "If you think something contributes to conversation, upvote it. If you think it does not contribute to the subreddit it is posted in or is off topic in a particular community, downvote it," (Reddit, 2022b). Essentially, upvoting and downvoting is supposed to be based upon how well a post contributes to the intended conversation of the subreddit rather than whether someone likes or agrees with the opinions expressed in the post.

It appears in practice, however, that the true meaning of upvoting may vary across subreddits. Graham and Rodriguez (2021) found that "voting on Reddit is not a simple, objective rubric to rank content – on the contrary, it is a material-discursive practice that performs localized cultures and meaning-making on the site." Of the thirty topics they identified in the conversations of redditors on several high-volume subreddits in relation to upvoting and downvoting on Reddit, only one of them actually had to do with the platform's intent of the system. They organized the topics of conversation

into a conceptual framework with four main themes: 1) platform culture, 2) prescriptive device, 3) materialization of value, and 4) ontology of self. The variety of ways redditors use upvoting and downvoting calls into question whether this popular method of sampling from top posts provides the best data for social science researchers to analyze. Gaudette et al. (2021) found that upvoting and downvoting created a sort of echo chamber within the subreddit r/The_Donald which reinforced the extreme views of the subreddit's "in-group" and prevented content that might challenge the group's views from being widely seen. Though the cultural norms of different subreddits may vary, it is true across the platform that voting affects what content is seen on Reddit and by how many users. Individual user feeds include highly upvoted content from the subreddits they are a member of and if a user specifically visits a particular subreddit, the posts are organized by default to the "top" posts of the day. Posts can also be sorted by "hot," "new," "controversial," and "rising."

Despite the lack of research on the efficacy of sampling methods, there is plenty of research taking place which utilizes Reddit data for research purposes. Commonly used methods include sampling posts from across the entire platform which meet certain inclusion criteria, sampling everything from a limited number of subreddits meeting inclusion criteria, sampling from the top most upvoted content, sampling from the top most commented on content, random sampling and purposive sample. Computer-assisted methods of sampling are also gaining in popularity, though they are much less accessible to the average social science researcher and may still return a body of text that is beyond the capacity of human analysis when used in a social media environment. Hintz and Betts (2022) have provided a script using the R platform that is intended to make sampling Reddit data more accessible to the average researcher, however researchers must still decide which posts to study after scraping all the relevant data that is returned using their method.

### *Top Most Upvoted*

Many researchers sample from the top posts believing them to be the most important posts in a subreddit. Struik et al. (2021) used the Reddit search feature to sort posts by top posts of the month and scraped the resulting data. Some researchers use this method and also recognize there may be limitations to this method. Davis and Kettrey (2021), for example, sorted posts by upvotes and noticed there was a steep drop off in the number of upvotes a post received beyond the six most upvoted threads and chose to focus their analysis on the top six threads only. They noted their sampling method "poses the limitation of excluding comments on the less visible threads, which may represent less popular voices." Carpenter and Willet (2021) sampled from top posts and noted this as a limitation of their study stating, "by purposefully sampling the top-voted posts and comments, we may have missed important trends in the contributions to these subreddits that were not top-voted. For instance, a study of the most downvoted posts and comments would likely reveal more controversial content" (p. 10).

### *Top Most Commented On*

Researchers also sample from the most commented on posts. Garg et al. (2021) suggest that sampling from the most commented on posts ensures analysis of "posts that resonated with a high number of subscribers (i.e., a higher number of comments by a higher number of distinct redditors) and represented "successful interactions" (p. 6). They do not claim for this sampling method to be one that supplies a representative sample of the universe of posts but do argue it to be a valid way of sampling and sampled the top 4% of most commented-on posts for every month during their selected timeframe. Bhandari and Sun (2021) chose to sample from the most commented on posts after an initial exploration of their data revealed a difference between the kinds of posts receiving a lot of upvotes and the kinds of posts receiving a lot of comments. Additionally, the authors cited a desire to capture "more nuanced community dynamics" when explaining their sampling rationale. Kimiafar et

al. (2021) sampled from the most commented on posts about COVID-19 but did not provide a rationale for doing so.

### *Random Sample*

Some researchers took a random sample without explaining their rationale for doing so and some performed more sophisticated methods of random sampling. Jungherr et al. (2021) took a random sample of 5,000 memes from a larger set of data scraped from a particular subreddit for a previous study to perform manual text analysis. Seraj et al. (2021) chose to randomly sample 1,000 first posts from a particular subreddit. Knittel et al. (2021) randomly sampled from daily discussion threads from two different subreddits, sampling twice as heavily from one subreddit than the other to reflect the fact the one subreddit had twice as many daily discussion threads as the other. Stevens et al. (2021) took what they call a systematic random sample of 50 posts from each week over a five-year period and then further randomly sample five posts from each set of 50 after a manual review of posts.

### *Purposive*

Alaggia and Wang (2020) created inclusion criteria and manually screened posts from a certain time period selecting all those that met their criteria and stopped coding once categorical saturation was reached. Kaufman et al. (2021) similarly purposively sampled comments related to a specific topic though they did not explain how this was done. de Saint Laurent et al. (2021) purposively selected memes from a larger dataset to focus in on a narrower topic. Del Valle and Smit (2021) purposively sampled comments and replies containing certain key words found to reflect mnemonic activity – the focus of their study and Kaufman. They argued that purposive sampling was appropriate for their research because their intention was not to make any inferences to a larger population. Smith et al. (2021) utilized a purposive sampling approach which used criteria comprised of a combination of number of comments, upvotes and downvotes to select posts, so although their method was purposive

it uses some of the key features of other sampling methods. Gray and Chivukula (2021) purposively sampled posts that represented what they call "typical cases" after conducting an initial exploration of a larger dataset.

*Computer Assisted Methods*

Gong et al. (2021) created a filtering mechanism they call the Lexico-Semantic Similarity Filter (LSSF) which leverages word embeddings, a natural language processing technique, that measures how words cluster together and helps identify target clusters of words in posts that are commonly found together in posts of interest. Using this technique allowed them to sort through large quantities of data from ten different subreddits to find posts related to their research interest. Feuston and Brubaker (2021) discuss how one researcher in their study used semantic network analysis, which similarly uses clusters of words that frequently appear together, to identify text for analysis from a larger body of text and another developed a machine learning classifier to identify text. They also highlight several potential limitations of using computer assisted methods and discuss how finding the right balance of collaboration between humans and computer-assisted methods is important. Their discussion further underscores the fact that using these tools requires not only execution skills, but the ability to appropriately apply the skills to have valid and reliable outcomes.

**Research Questions**

This paper will address three research questions:

1. Is there a difference in the results of an inductive content analysis performed on the top fifty most upvoted posts from a given time period, the top fifty most commented on posts from a given time period, and the universe of posts from that same time period?

2. Is there a difference in the results of sentiment analysis performed on the top 50 most upvoted posts from a given time period, the top 50 most commented on posts from a given time period, and the universe of posts from that same time period?

3. If differences exist, how might these differences be interpreted to inform research design when social science researchers are sampling data from Reddit?

**Methodology**

Due to the primary author's research interests, posts tagged with the #resignation flair from the r/teachers subreddit were chosen as the focus of study. It was determined by the Western Michigan University Human Subjects Institutional Review Board that analyzing anonymous, available to the public social media data did not constitute as human subjects research and did not require any review. A Python script provided by Rare Loot (2018) was adapted to scrape the desired data from the JSON coding behind the Pushift.io dataset. The Pushshift dataset acts as a mirrored copy of all posts on the Reddit website and was found to be easier to retrieve data from than using Reddit's API directly. Limitations in coding ability necessitated scraping one day's worth of data at a time.

An initial exploration of data from the r/teachers subreddit was conducted to determine when the #resignation flair was added as an option. All data from the day it was introduced, 8/10/2020, through 10/10/2022 were scraped, filtered to only those posts with the #resignation flair, and combined into one dataset. Data elements scraped for each post included: a unique post id, the post title, the post text, the number of upvotes the post received, the number of comments a post received, the flair attached to the post, the Reddit user id of the original poster, the date and time of the post, the post URL, and the subreddit the post was submitted to. In some cases, the post had been deleted by the original poster right after it was posted, and the mirrored copy did not include post text but rather the message [deleted]. These posts were filtered and removed.

70

The number of upvotes and number of comments were found to be taken from the time a post was copied into the Pushshift.io dataset, which usually occurs soon after the post was made and before the post has had a chance to be upvoted or commented on. Therefore, it was decided to retrieve the number of upvotes and number of comments the post had at the time of archiving, if applicable, or on the date of retrieval if the post hadn't yet been archived by clicking on the public facing URL and manually copying this information into the dataset. Through this process, it was discovered that some posts were deleted from Reddit by the original author after they had been copied into the Pushshift.io dataset. To honor the wishes of the original author to not have their post continue to be publicly available, these posts were removed from the dataset. Additionally, it was discovered that some posts were not actually submitted to the r/teachers subreddit and/or did not actually have the resignation flair attached to them and these were deleted as well. All data were scraped and manually retrieved during the month of October 2022.

A graph was then generated to show how many #resignation flair posts were made in the r/teachers subreddit by month for the 2020-2021 and 2021-2022 school years (see Figure 15) to help determine which data to utilize for the study. The 2021-2022 school year was selected due to the higher number of posts in that school year in addition to the increased distance from the COVID-19 pandemic. The months of September, January and May were selected to provide a cross-section of the school year from months with a high volume of posts. All posts from these three months were combined into one final dataset and the find and replace feature in excel was utilized to clean the post title and post text of the special characters that were inserted during the scraping process.

**Figure 15**

*Total Number of Resignation Flair Posts in r/teaches by Month*



To create the final datasets for this study, the data were sorted in ascending order first by number of upvotes and the top 50 were tagged, then by the number of comments and the top 50 were tagged. Different configurations of data were used for the qualitative and quantitative analysis to meet assumptions of statistical tests. Qualitative analyses were performed on three datasets: The top 50 most-upvoted or "MUV" for short, the top 50 most-commented-on, or "MCO" for short, and the entire universe of posts. For the quantitative analysis, posts were separated into four datasets, those that appeared in both the MUV and MCO datasets, or "Both" for short, those that were unique to MUV, those that were unique to MCO, and everything else leftover in the universe, or "EEU" for short.

A mixed-methods approach to analysis was conducted to capture both *what* was being said in the posts in each dataset and *how* it was being said. To measure what was being said in the posts, an inductive content analysis adapted from the procedure outlined by Elo & Kyngäs (2007) was performed on each group. To measure how original posters were saying things, statistical analysis was performed on several variables calculated using Linguistic Inquiry and Word Count (LIWC)-22.

*Inductive Content Analysis*

To answer the first research question, an inductive approach to qualitative content analysis based on the procedure outlined by Elo and Kyngäs (2008) was taken. The procedure was adapted to allow for the simultaneous analysis of multiple datasets with the goal of comparing how the inclusion of additional codes changed the grouping and categorization process. To help narrow the focus of the content analysis, the research question, "What reasons do teachers give for wanting to leave or leaving their current teaching job or the profession?" was used and the unit of analysis was specific reasons. Not all posts included reasons and in these cases the goal became to capture the main "gist" of a post, or what the post was mostly about. Finally, it is important to note that the datasets for the qualitative analysis vary slightly from those of the quantitative datasets. A log trail was kept throughout the qualitative content analysis to document the process, record insights and reactions to the data in a reflexive manner, and think through and record decisions that were made. Being a former teacher and parent of school-age children brought life experiences and perspectives to the analysis that deepened the researcher's understanding of the topic though epoche was used to refrain from allowing those experiences to bias analysis.

The post title and post body for every post in each dataset was copied into a table in Microsoft Word and the posts were printed in landscape format on 8 ½ x 11 paper. Beginning with the Both dataset, posts were read through several times to immerse into the data and open codes were recorded in the margins to represent either reasons or the gist of the post. Posts from the MUV and MCO datasets were read next, in that order, and open codes were recorded in the margins. During this process it was determined that the posts were falling into several distinct buckets that might provide a useful framework for analyzing the codes and these categories were recorded for each post as well.

73

The next step was to immerse in the Universe data set, open code in the margins, and categorize each post. Throughout this process, several additional categories emerged, and a final framework of categories began to be defined. At this point, it was unclear whether or not analysis of codes should be structured around the categories since they had changed with the addition of data from the Universe set so the difference was noted and analysis of codes in the first three datasets moved forward using the original buckets that had emerged during open coding of those sets.

For the sake of ease, codes from the Both dataset were hand-written in purple on blank 8 ½ x 11 pieces of paper. Then codes from the MUV dataset were written in blue either on a separate sheet of paper or at the top of the Both paper depending on how many codes there were for a given category. Finally, posts from the MCO dataset were written in green either on a separate sheet of paper or at the bottom of the Both paper, again, depending on how many codes there were for a given category. During the transfer process, codes were reworded when necessary to more accurately and more succinctly reflect the intent of the text from which they derived and reasons that were already given were starred for each additional occurrence to record frequency. Pictures of the layout were taken to document the structure of the codes at this stage.

The handwritten codes from the Both dataset were then cut apart for the two buckets that did not relate to a teacher's position on a continuum of leaving. Groupings were made for each of the buckets and pictures were taken to document the structure. The handwritten codes from the MUV dataset for the same buckets were then cut apart and added to the groupings, rearranging when necessary, and pictures were taken again. The MUV codes were then removed and the MCO codes were added, regrouping when necessary and pictured. At this point, the decision was made to sort the posts in the Universe dataset into buckets so that the codes from these two buckets could be added to the groupings since they were the smallest categories. This acted as both a pilot for the rest of the

74

process and an important source of information gathering. The posts having to do with the continuum of leaving were sorted based on the refined framework that had emerged during the process of open coding the Universe dataset and since the framework was not finalized until the end of the process, the placement in the correct bucket was revisited for each post as they were cut apart. Final pictures were taken of the groupings once all codes were included.

Next, the three buckets relating to a teacher's position on a continuum of leaving were analyzed. During the grouping of the first bucket, several categories emerged which provided structure for the groupings of the other two buckets. The same procedure described earlier of grouping codes from the Both dataset and documenting with pictures, adding codes from the MUV dataset and documenting with pictures, then removing the MUV codes, adding the MCO codes and documenting with pictures was followed for each bucket. The codes were then set aside, kept in baggies organized by category.

Codes for each bucket in the Universe set were then cut apart and sorted into the categories identified during the grouping of the first three datasets. They were then placed in baggies along with a colored piece of paper to mark their bucket membership to keep them organized until all codes were sorted. Once sorted, the baggies for each category in all four datasets were laid out along the continuum of leaving by bucket as in the array pictured in Figure 16 as a staging ground. The codes from each category were then sorted into groups along a continuum arranged adjacent to the staging ground, dubbed the analysis arena. Related categories like emotional toll, physical toll and personal life toll were sorted at the same time, each in their own respective place along the continuum. Pictures were taken to record the groupings. Once all the groupings were completed, the abstraction process began based on the pictures taken to record groupings. This consisted of creating a digital map of the categories in each configuration which were then compared and analyzed.

75

**Figure 16**

*Staging Ground Array of Data*

| | Thinking About Quitting | Planning to Quit | In the Process of Quitting | Officially Quit |
|---|---|---|---|---|
| Unique to the Point on the Continuum | Both, MCO MUV, Universe | Both, MCO MUV, Universe | Both, MCO MUV, Universe | Both, MCO MUV, Universe |
| Emotional Toll | Both, MCO MUV, Universe | Both, MCO MUV, Universe | Both, MCO MUV, Universe | Both, MCO MUV, Universe |
| Physical Toll | Both, MCO MUV, Universe | Both, MCO MUV, Universe | Both, MCO MUV, Universe | Both, MCO MUV, Universe |
| Personal Life Toll | Both, MCO MUV, Universe | Both, MCO MUV, Universe | Both, MCO MUV, Universe | Both, MCO MUV, Universe |
| Bad Environment | Both, MCO MUV, Universe | Both, MCO MUV, Universe | Both, MCO MUV, Universe | Both, MCO MUV, Universe |
| Students | Both, MCO MUV, Universe | Both, MCO MUV, Universe | Both, MCO MUV, Universe | Both, MCO MUV, Universe |
| Parents | Both, MCO MUV, Universe | Both, MCO MUV, Universe | Both, MCO MUV, Universe | Both, MCO MUV, Universe |
| Admin | Both, MCO MUV, Universe | Both, MCO MUV, Universe | Both, MCO MUV, Universe | Both, MCO MUV, Universe |
| Profession | Both, MCO MUV, Universe | Both, MCO MUV, Universe | Both, MCO MUV, Universe | Both, MCO MUV, Universe |
| Logistics | Both, MCO MUV, Universe | Both, MCO MUV, Universe | Both, MCO MUV, Universe | Both, MCO MUV, Universe |

### *Statistical Analysis of LIWC Variables*

To answer the second research question, the four LIWC summary variables: Analytical Thinking, Clout, Authenticity and Emotional Tone were utilized to measure how things shared in posts were being said. These are considered style variables, or those that convey how people are communicating rather than what they are communicating (Tausczik & Pennebaker, 2010). Analytical thinking "captures the degree to which people use words that suggest formal, logical, and hierarchical thinking patterns," clout "refers to the relative social status, confidence, or leadership that people display through their writing or talking," authenticity is "a reflection of the degree to which a person is self-monitoring," and emotional tone combines positive and negative tone into one variable with

numbers above 50 suggesting a more positive tone and numbers below 50 suggesting a more negative one (LIWC, n.d.). The software essentially counts the number of words in a post that fall into these categories and determines a final score for the category based on the counts. Word counts for each posts were also calculated and compared.

One-way between-subjects analysis of variance and Kruskal-Wallis tests were performed as appropriate for each of the five LIWC variables. Assumptions for these models were checked, followed by an omnibus F-test of whether the average score for each of the variables was significantly different among the datasets. The alpha level for the omnibus F-test was set to .05. With significant results of this F-test, post-hoc Tukey HSD pairwise comparisons for groups with unequal n's were made to determine the pattern of differences. The family-wise alpha level was set to be .05 for the post-hoc analysis.

**Results**

A total of 336 posts tagged with the #resignation flair were scraped from September 2021, January 2022, and May 2022. The descriptive statistics for the number of posts, the range of the number of upvotes and comments, and the mean and median of the number of upvotes and comments for each of datasets used in this paper can be seen in Table 2. Once posts reach a certain number of upvotes, the actual number appears to be rounded to the nearest 100 on the Reddit platform and is listed as, for example, "2.8 k", which was transcribed as 2,800 into the dataset. As one would expect, the average number of upvotes are highest in the Both and MUV datasets and the average number of comments are highest in the Both and MCO datasets. Averages for the universe are much lower than in the Both, MUV, and MCO datasets.

**Table 2**

*Descriptive Statistics for Upvotes and Comments from All Datasets*

| | Total # of Posts | Range of # of Upvotes | Mean # of Upvotes | Median # of Upvotes | Range of # of Comments | Mean # of Comments | Median # of Comments |
|---|---|---|---|---|---|---|---|
| Both Top 50 Most Upvoted and Top 50 Most Commented On (Both) | 36 | 184-2,800 | 747 | 448 | 42-522 | 146 | 97 |
| Unique Top 50 Most upvoted (MUV) | 14 | 177-375 | 238 | 223 | 15-29 | 25 | 24 |
| Unique Top 50 Most Commented On (MCO) | 14 | 8-163 | 84 | 83 | 40-169 | 66 | 59 |
| Everything Else in the Universe (EEU) | 272 | 0-174 | 30 | 15 | 0-40 | 10 | 7 |
| Universe | 336 | 0-2,800 | 118 | 22 | 0-522 | 27 | 9 |

### *Results of LIWC Analysis*

**Descriptive Statistics for LIWC Variables.** The word count varied across the entire universe of posts from one word to 2,339, with the averages being similar across the Both, MUV, EEU and Universe datasets and being about sixty words lower for the MCO dataset. Average analytic thinking scores were lowest for the MUV dataset and higest for the Both and MCO datasets suggesting that perhaps posts that have higher levels of analytic thinking get more overall engagement – both upvotes and comments, and that posts with lower levels of analytic thinking may get upvoted but not

commented on. The average clout score of posts was highest in the MCO dataset (mean = 40, median = 24) with twenty-seven percent of posts in the universe having a clout score of one. Details on clout scores for each dataset can be seen in Table 5. The range of emotional tone scores shows that both positive and negative posts were made about resigning and the averages across all datasets show the majority of posts about resigning were negative in tone with the most negative posts being found in the MUV dataset. The range of authenticity scores show that both inauthentic and extremely authentic posts were made about resigning, however the averages across all datasets in this case show that the majority of posts were authentic with the most authentic posts being found in the MUV dataset. The MUV dataset is both the most negative and the most authentic. Details on for each variable across datasets can be seen in Tables 3-7.

**Table 3**

*Descriptive Statistics for Word Count from All Datasets*

|  | *Mean* | *Standard Deviation* | *Min* | *Max* |
| --- | --- | --- | --- | --- |
| Both Top 50 Most Upvoted and Top 50 Most Commented On (Both) | 189 | 153 | 28 | 723 |
| Unique Top 50 Most upvoted (MUV) | 183 | 159 | 38 | 634 |
| Unique Top 50 Most Commented On (MCO) | 126 | 125 | 18 | 400 |
| Everything Else in the Universe (EEU) | 197 | 233 | 1 | 2,339 |
| Universe | 193 | 219 | 1 | 2,339 |

**Table 4**

*Descriptive Statistics for Analytic Thinking Scores from All Datasets*

| | *Mean* | *Standard Deviation* | *Min* | *Max* |
|---|---|---|---|---|
| Both Top 50 Most Upvoted and Top 50 Most Commented On (Both) | 35.44 | 21.31 | 1 | 85.33 |
| Unique Top 50 Most upvoted (MUV) | 31.11 | 24.54 | 2.69 | 79.89 |
| Unique Top 50 Most Commented On (MCO) | 37.27 | 29.15 | 4.85 | 86.73 |
| Everything Else in the Universe (EEU) | 32.49 | 23.34 | 1 | 99 |
| Universe | 32.95 | 23.37 | 1 | 99 |

**Table 5**

*Descriptive Statistics for Clout Scores from All Datasets*

| | *Mean* | *Standard Deviation* | *Min* | *Max* |
|---|---|---|---|---|
| Both Top 50 Most Upvoted and Top 50 Most Commented On (Both) | 23.42 | 34.79 | 1 | 99 |
| Unique Top 50 Most upvoted (MUV) | 15.77 | 37.73 | 1 | 64.43 |
| Unique Top 50 Most Commented On (MCO) | 40.12 | 37.74 | 1 | 94.01 |
| Everything Else in the Universe (EEU) | 14.01 | 22.64 | 1 | 99 |
| Universe | 16.18 | 25.38 | 1 | 99 |

**Table 6**

*Descriptive Statistics for Emotional Tone Scores from All Datasets*

|  | *Mean* | *Standard Deviation* | *Min* | *Max* |
|---|---|---|---|---|
| Both Top 50 Most Upvoted and Top 50 Most Commented On (Both) | 32.56 | 26.88 | 1 | 82.53 |
| Unique Top 50 Most upvoted (MUV) | 25.64 | 32.68 | 2.41 | 98.96 |
| Unique Top 50 Most Commented On (MCO) | 30.36 | 21.95 | 1 | 81.12 |
| Everything Else in the Universe (EEU) | 33.85 | 27.95 | 1 | 99 |
| Universe | 33.23 | 27.76 | 1 | 99 |

**Table 7**

*Descriptive Statistics for Authenticity Scores from All Datasets*

|  | *Mean* | *Standard Deviation* | *Min* | *Max* |
|---|---|---|---|---|
| Both Top 50 Most Upvoted and Top 50 Most Commented On (Both) | 81.44 | 29.91 | 3.12 | 99 |
| Unique Top 50 Most upvoted (MUV) | 87.71 | 19.58 | 25.16 | 99 |
| Unique Top 50 Most Commented On (MCO) | 78.91 | 26.66 | 1.17 | 99 |
| Everything Else in the Universe (EEU) | 85.26 | 22.82 | 1 | 99 |
| Universe | 84.69 | 23.67 | 1 | 99 |

**Comparison of Group Means.** The assumptions of normality and homogeneity of variance were satisfied for the Word Count and Analytic variables. The assumption of normality was met for Clout, but not the assumption of homogeneity of variance. A one-way between-subjects analysis of variance was performed for the Word Count, Analytic and Clout variables and the Brown-Forsythe test was used for the Clout variable to account for the violation of homogeneity of variance. The assumption of normality was not met for Authenticity or Tone and a Kruskal-Wallis Test was performed for these variables.

There was not a statistically significant difference between groups for the Word Count or the Analytic variables as determined by one-way ANOVA. The ANOVA summary for Word Count can be seen in Table 8 and for Analytic Thinking can be seen in Table 9. There was a statistically significant difference between groups for the Clout variable as determined by one-way ANOVA ($F(3,332) = 6.055$, $p < .001$) (Table 10) and post hoc Tukey HSD family wise comparisons showed the differences to be statistically significant between the MCO and EEU groups ($p = .048$) and the MUV and EEU groups ($p < .001$). There was not a statistically significant difference between groups for the Emotional Tone and Authenticity variables as determined by Kruskal-Wallis tests.

**Table 8**

*The ANOVA Summary for Word Count*

| Source | SS | Df | MS | F |
|---|---|---|---|---|
| Between Groups | 69205.864 | 3 | 23068.621 | .698 |
| Within Groups | 16043928.1 | 332 | 48325.085 | |
| Total | 16113134.0 | 335 | | |

**Table 9**

*The Summary ANOVA Table Showing the Difference in Mean Analytic Thinking Score Across the*

*Four Datasets*

| Source | SS | Df | MS | F |
|---|---|---|---|---|
| Between Groups | 589.567 | 3 | 196.522 | .783 |
| Within Groups | 182327.860 | 332 | 549.180 | |
| Total | 182917.427 | 335 | | |

**Table 10**

*The Summary ANOVA Table Showing the Difference in Mean Clout Score Across the Four Datasets*

| Source | SS | Df | MS | F |
|---|---|---|---|---|
| Between Groups | 11196.411 | 3 | 3732.137 | 6.055* |
| Within Groups | 204647.746 | 332 | 616.409 | |
| Total | 215844.157 | 335 | | |

* $p < .001$

### *Results of Inductive Content Analysis*

**Type of Post Descriptives for all Datasets.** Since the data so readily sorted into some distinct

buckets early on, it made sense to compare how these buckets were represented across the datasets

which can be seen in Table 11. "Commentary & News" posts were overrepresented in both the MUV

(14%) and MCO (24%) datasets when compared to the Universe (8.9%), with the MCO dataset have

almost three times as many of these posts. "Thinking About Quitting" posts were underrepresented in

the MUV dataset where they made up only 10% of posts compared to 17.9% of posts in the Universe,

whereas posts about people who had officially quit were overrepresented in the MUV dataset,

accounting for 40% of posts whereas they made up only 25% of posts in the Universe. "Leaving Not

by Choice" posts were not represented at all in either the MUV or MCO datasets and made up 3.6% of posts in the Universe. Finally, "Flipside" posts were overrepresented in both the MUV and MCO datasets with 14% and 12% respectively compared to only 5.1% in the Universe.

**Table 11**

*Number and Percentage of Type of Posts by Dataset*

| Category | #(%) in the MUV Dataset | #(%) in the MCO Dataset | #(%) in the Universe |
|---|---|---|---|
| Commentary & News | 7 (14%) | 12 (24%) | 30 (8.9%) |
| Thinking About Quitting | 5 (10%) | 8 (16%) | 60 (17.9%) |
| Planning to Quit | 2 (4%) | 2 (4%) | 69 (20.5%) |
| In the Process of Quitting | 5 (10%) | 5 (10%) | 37 (11%) |
| Officially Quit | 20 (40%) | 13 (26%) | 84 (25%) |
| Leaving Not by Choice | 0 (0%) | 0 (0%) | 12 (3.6%) |
| Flipside | 7 (14%) | 6 (12%) | 17 (5.1%) |
| Support | 2 (4%) | 3 (6%) | 6 (1.8%) |
| Excluded | 2 (4%) | 1 (2%) | 21 (6.3%) |

**Category Structure Comparison.** The categories each bucket was organized into were also compared to determine differences in the analysis of the three data sets. Comparisons of the buckets that were not part of the Continuum of Quitting: "Flipside", "Commentary & News", "Support" and

"Leaving Not by Choice", will be made first. Next, each bucket along the Continuum of Quitting will be compared.

The "Flipside" bucket existed in the MUV and Universe datasets but not in the MCO dataset. Sampling from the most commented on posts would have masked this important bucket of posts that showed how people felt after they had moved on. Similarly, the "Not by Choice" bucket did not appear in either the MUV or the MCO datasets and only revealed itself during analysis of the Universe. In looking more closely at the details of "Flipside" bucket in the MUV and Universe datasets, several categorical shifts occurred as more posts within the universe were added. The organization of categories for the "Flipside" buckets can be seen in detail in Figures 17 and 18. The category of "Personal Benefits" was present in the MUV dataset but not the Universe as these items were shifted into areas that made more sense. "Emotional Benefits" was renamed to "Emotional Health Benefits" and became a subset of a broader category called "Health" which also included "Physical Health Benefits" and "Social Health Benefits." "Regrets - Not Fulfilled" moved under a new category called "Dealing with the Aftermath;" a "Workload" category was created; and "Ability to Work from Home" was moved under "Freedom."

**Figure 17**

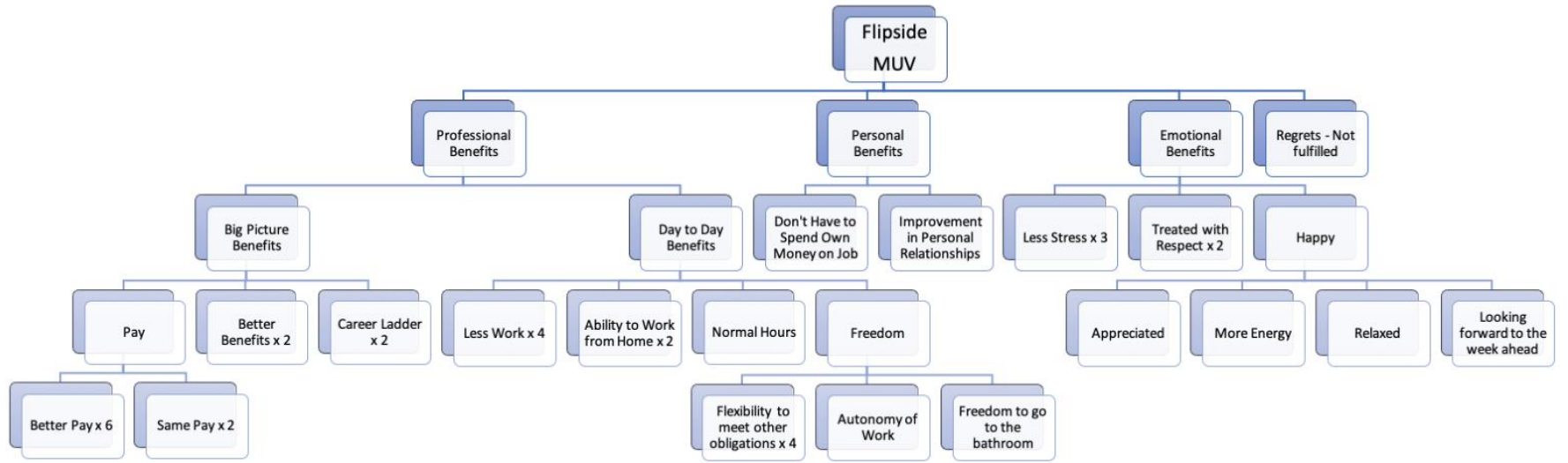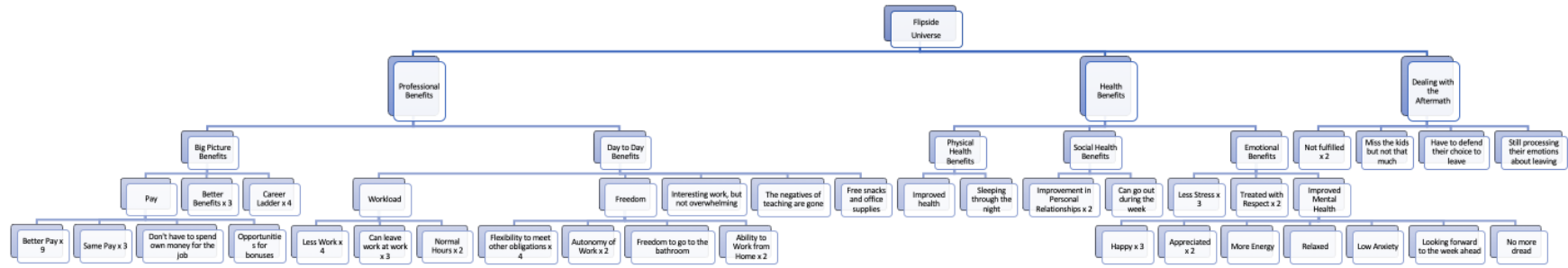*Categories in the Flipside Bucket of the MUV Dataset*

**Figure 12**

*Categories in the Flipside Bucket of the Universe Dataset*

The "Support" bucket did not have much complexity in the MUV and MCO datasets and became clearly split into two subcategories as more details were added from the Universe. Details for the "Support" buckets can be found in Figures A1-A3 in Appendix A. In the "Commentary & News" bucket, the category of "Low Pay" was not present in the MUV dataset and "News Stories" only existed as a category in the MCO dataset showing that posts sharing news articles received a lot of comments. With the addition of other posts in the universe dataset, it made more sense to categorize the news stories based on their topics. Similarly, the relatively small "Admin Commentary" category disappeared when analyzing the Universe dataset and the items were subsumed into a larger "Random Comments" category. Finally, although the category "Polls" existed in all three datasets, one kind of poll emerged as being especially popular in the analysis of the Universe dataset leading to the creation of a sub-category to accommodate. The organization of categories for the "Commentary & News" buckets can be found in Figures A4-A6 in Appendix A.

Comparisons of the buckets along the Continuum of Quitting will now be discussed. Perhaps most notable is the fact that analysis of the MUV and MCO datasets led to a three-bucket continuum: "Thinking About Quitting," "I'm Quitting," and "I Quit." Analysis of the Universe refined the continuum into four buckets that better capture the points along the continuum: "Thinking About Quitting," "Planning to Quit," "In the Process of Quitting," and "Officially Quit." Either sampling method, in this case, flattened the continuum and masked the presence of unique circumstances people discussed while they were planning to quit versus when they were in the process of quitting. Similarly, some posts that had been included in the original "I Quit" bucket were moved into the "In the Process of Quitting" bucket and "I Quit" was renamed "Officially Quit." As a better sense of the continuum took shape, it was clear there was a difference in experience between having put in notice of

resignation but still have teaching duties to finish before leaving and being completely done with teaching. A closer look at each bucket will now be taken.

For the "Thinking About Quitting" bucket, the MUV and MCO had the same main categories, but the MCO analyses had more codes falling under the categories. MCO added "Thankless work" and "Lack of support and resources" under "Profession;" "Toxic school culture," "Students don't want to learn" and "Too depleted to hold kids accountable" under "Bad environment;" "Mental health suffering from pandemic" under "Emotional Toll;" and "Have to be so healthy just to get by" and "I'm at high-risk for COVID" under "Physical Toll." Details for the "Thinking About Quitting buckets can be found in Figures A7-A12 in Appendix A. In the Universe, the "Tolls of the Job" category went from being divided into two subcategories into being divided into three with the addition of "Personal" to "Emotional," and "Physical." The "Emotional" category was built out even more with two categories, "COVID" and "Dealing with health issues" covering a broader range of codes. "Students" became a category with six subcategories, each with multiple instances, whereas previously there had just been two student-related codes under the "Bad Environment" category. The "Profession" category started to show workload as a significant issue with seven related codes whereas there was only one code related to workload ("Too many hours with kids") in the MCO dataset and no mention of workload in the MUV dataset. The "Logistics" category went from having one code to six. The most prevalent tolls of the job in the MUV and MCO datasets were "emotionally drained" and "physically drained" and in the Universe they were "Burned out shell of a human" with eight related instances, followed by "emotionally drained" with five for emotional tolls and "Exhausted" with nine related instances for physical which replaced "Physically drained" in the other datasets. "Bad admin" remained the most prevalent category under "Bad Environment" in all three datasets and "Low pay" in

the MUV and MCO datasets was replaced by "Requires too many hours of work" in the Universe as the most prevalent code under "Profession."

At this point in the continuum the buckets diverge as the two buckets of "I'm Quitting" and "I Quit" in the MUV and MCO datasets split into three in the Universe dataset: "Planning to Quit," "In the Process of Quitting," and "Officially Quit." This change at the macro level is extremely significant in and of itself, and also makes it more difficult to continue comparing the subcategories between the MUV, MCO and Universe datasets at the micro level so the "I'm Quitting" buckets from the MUV and MCO datasets will be compared to the "Planning to Quit" and "In the Process of Quitting" buckets. Details for each of these buckets can be found in Figures A13-A20 in Appendix A. In the MUV and MCO datasets, only a few logistic codes were present (two for MUV and three for MCO), but in the Universe there were a total of 58 individual codes related to logistics. This shows that logistical questions are a huge concern for the average poster, but that these kinds of posts don't necessarily get highly upvoted or commented on. Overall, the MUV dataset contained more codes and therefore details than the MCO dataset. Codes related to workload in the "Profession" category and codes related to students in the "Bad Environment" category were prevalent enough to warrant their own categories in the MUV dataset whereas they were just one of several issues in the MCO dataset. Therefore, the MUV dataset better matched the Universe dataset in this regard. "Workload" was the most prevalent code in the "Profession" category in the MUV dataset with four codes and the most prevalent code in both the "Planning to Quit" (four also) and "In the Process of Quitting" (ten) buckets of the Universe dataset whereas there was no most prevalent code for the "Profession" category in the MCO dataset. The "Tolls of the Job" category in the MCO dataset was further divided into "Emotional Tolls" and "Physical Tolls" in the MUV dataset.

The "I Quit" buckets of the MUV and MCO datasets will now be compared to the "Officially Quit" bucket of the Universe. The first and most obvious difference is the change in name in the Universe dataset to give a more accurate semblance after its refinement. The "Unique to I Quit" category was largely the same between the MUV and MCO datasets and all three datasets showed "Time to put myself first" as the most prevalent category, though it was much more fleshed out with codes in the Universe. The "Poor Benefits" category was more fleshed out in MCO dataset than the MUV, however all three datasets showed "Low Pay" to be the most prevalent issue with the profession. "Bad admin" was the most prevalent code in the "Bad Environment" category in the MUV and dataset but was a close second to "Students" in the MUV and Universe datasets. The "Tolls of the Job" category was more fleshed out in the Universe than in the MUV and MCO datasets and also made up a larger percentage of all codes represented in the dataset. Finally, general comments about mental and physical health suffering were more prevalent in the Universe dataset though comments about specific mental and physical issues were present in all three. Details can be found in Figures A21-A25 in Appendix A.

**Discussion**

It has been established that differences between the three datasets do exist and in answer to the third research question, the potential implications of those differences will now be discussed. Perhaps not surprisingly, analyzing the entire Universe of data gave a more detailed picture than taking a sample regardless of sampling method. However, the data in the two different samples provided a fairly good window into the universe, though the pictures were slightly skewed in a few places and left some holes. The differences were more apparent within the content of what people were talking about than how they were talking about it.

There were no statistically significant differences between the datasets on four out of the five style variables tested: word count, analytic thinking, emotional tone, and authenticity. The way people were talking about resigning and how many words they used to talk about it did not vary greatly between the sampled data and the universe of data. There was a statistically significant difference in the Clout variable between the MUV and EEU and the MCO and EEU datasets showing that posts where people spoke with more authority were more likely to become the top most-upvoted and top most-commented on posts. Therefore, researchers selecting one of these sampling methods may be missing out on the opinions and experience of posters who speak with less authority on their topic which could potentially be either desirable or undesirable based on a researcher's interests. Further research could be done in other topic areas to see if this finding holds true.

The difference in representation of the kinds of posts across the three datasets varied enough to warrant this as a significant consideration for researchers hoping to sample from a universe of data. The over-representation of certain kinds of posts in a sample may give researchers a skewed view of what topics matter most to people. In this particular case, members of the r/teachers subreddit may be providing support to people by upvoting posts where people talk about getting a new job which skews a top most-upvoted sample towards a certain kind of post rather than a full spectrum of experience around a topic. A researcher may be interested in what kinds of posts specifically get the most positive responses from other users in which case this sampling method would be very appropriate. However, if the goal is to gather information about a wide range of experience, this method is not going to be the best choice. Both samples, and especially the MCO sample, were skewed more heavily towards posts about commentary and news than the universe. A sample of most-commented on posts is less representative of the universe of posts in this regard, but if a researcher is interested in studying conversations around a topic this sampling method might give them the most amount and best data to

analyze. Since Reddit was originally established as a news sharing social media platform, it makes

sense that these kinds of posts might generate a lot of conversation on the platform and the

conversation might be extremely valuable for certain researchers to analyze based on their research

interests. There were a lot of polls within the "Commentary & News" bucket and the analysis of

responses to these poll questions could be extremely interesting to researchers. Sampling from top-

most commented on posts could help researchers filter posts to identify polls that would give them a

deeper look into the community.

The lack of representation of certain kinds of posts demonstrates the potential of entire

viewpoints to be missed when sampling. The "Flipside" bucket of posts was present in the MUV

dataset, but not the MCO dataset showing that people showed support for those that had found new

careers or teaching jobs in a better environment but didn't have much to say about it. Sampling from

top most commented-on posts would have excluded this perspective. "Leaving Not by Choice" posts

were not present in either sample and though they made up only a small percentage of posts in the

Universe (3.6%) this bucket of posts was an unexpected finding in the dataset despite the author's

familiarity with teacher resignation. Other such unexpected findings may be overlooked by researchers

sampling using these methods.

Perhaps the most significant difference between the analysis of the three datasets is the shift in

the Continuum of Quitting that occurred as the buckets were more fleshed out during analysis of the

Universe dataset. Inclusion of every post in the universe provided a better understanding of the

continuum as subtler differences between different points along the continuum were explored. Posts

were re-categorized as the buckets became better defined. Since this researcher is ultimately interested

in the content of the posts and gaining a better understanding of teacher resignation, these details are

extremely important to the research interest. However, it is possible that if a larger sample had been

taken the fleshed-out continuum may have still taken shape so it could be a question of sample size rather than sampling method.

Beyond the significance of this structural difference, numerous such differences existed between each sample and the universe of data at the subcategory level with subtler structural differences taking shape as more details were added. Categories became further subdivided providing more detailed understanding, though in some cases the most prevalent details remained the same despite adding more details or did not shift much. The more significant differences were cases where details were completely missing from the MUV and MCO datasets which comes back again to the purpose of the research being conducted – if the research is exploratory with a goal of understanding as much as possible about a topic, or what Schober et al. (2016) call topic coverage, the two sampling methods studied in this paper are going to inherently leave holes in understanding. However, given real-world constraints, researchers may decide the time saved in analyzing less data may be worth the loss in detail of the final analysis.

An unexpected methodological finding of this paper was the value of exploring a universe of data through a sample of the top 50 most-upvoted and top 50 most-commented on posts. Analyzing a subset of the data first allowed the researcher to begin to develop a categorical framework that was then applied to the larger set of data. Though the entire content analysis process was inductive, the initial categorical framework discovered during analysis of the samples was used as a starting point for adding in additional codes from the universe rather than starting over from scratch. The framework was modified as necessary and aided in managing a larger set of data suggesting this method may provide a good entry point into analyzing larger quantities of social media data qualitatively.

**Limitations**

Though this paper begins to build the body of literature around the consequences of different sampling methods, it does have some limitations. The data for this study was scraped from one subreddit and may not generalize to data scraped from other subreddits or social media platforms. Additionally, decisions made when building the universe of data for this study may have affected the results. Data selected from a different year, from different months, or with a different flair may have yielded different results. Filtering available data by flair was an easy way to limit the universe of data and ensure posts were relevant to the topic of research interest, however, there are other ways to filter available data that may have provided a more comprehensive dataset on the same topic. Because human coded qualitative analysis was a part of the study design it was necessary to create a manageably sized dataset. The script used to scrape data yielded some inaccuracies in filtering by the established criteria and though posts that should not have been included were excluded from the dataset (posts deleted by users and posts that weren't actually from the subreddit of interest, for example) there is no way of knowing if the converse happened - whether posts that should have been scraped were not.

Another limitation of this paper is that analyses for the third paper, which focused on the universe of data, were conducted simultaneously with the qualitative analyses of the universe for this paper and took place after quantitative analyses had been completed. Twenty-one posts that had been included in the quantitative analyses were excluded from the qualitative analyses after further review revealed they did not meet inclusion criteria for the third paper, leaving a slight difference in the posts analyzed quantitatively versus qualitatively. Another limitation of this paper is that small size of the samples. It is possible some of the differences between the samples and the universe are due to too

small of a sample size rather than the sampling method. Finally, this study focused on two methods of sampling when more sampling methods are available and could have been incorporated into the study.

**Directions for Future Research**

The increasing prevalence of Reddit as a source of extant research data necessitates that more research be done on the consequences of sampling methods on this social media platform. Replication of this study design across other subreddits on the Reddit platform and on other social media platforms to better understand how context affects the differences found in this study. This study design should also be extended to include other sampling methods like random sampling, purposive sampling, and computer-assisted sampling methods. The finding that sampling from top posts clearly leaves out some perspectives begs the question of whether a random or purposive sample might do a better job of fully representing all perspectives. Further research should also be done that varies the number of units being sampled. This paper only sampled the Top 50 most-upvoted and most commented-on posts and it's possible that sampling the Top 100 may have given a better representation of the Universe. Finally, this study design could be modified to show if the sampling methods used to generate the dataset for this study affected the results. For example, using a key word search to filter scraped data rather than a flair to see if this decision has any consequences, or taking a sample from more time periods or different time periods to investigate the consequences of those decisions.

**Conclusion**

Overall, there is no way to say that the MUV or MCO sampling method is better than the other. Every research topic is going to have nuance unique to that topic and the differences noted in this article will not hold true in the exact same way for all topics. This paper does show, however, that sampling from top posts from a universe of posts whether they are top most-upvoted or top most-commented-on, has the potential to provide a skewed view of the universe. What is certain, is that

96

some perspectives that are less popular or occur less frequently are likely to be excluded when using the MUV or MCO sampling methods and researchers need to be aware of this likelihood when making sampling decisions. Researchers should consider their stage of research when choosing a sampling method. If a researcher is at the exploratory phase, they may want to learn as much as they can about all aspects of a topic and sampling from top posts may not benefit them as much as taking a random sample might. However, if much is known about a topic already or if a researcher is most interested in what the most popular opinions on a topic may be, focusing on the top most-upvoted posts may make the most sense for them. Likewise, if a researcher is most interested in conversation dynamics, sampling the top most-commented on posts will provide them with richer conversation data to analyze. The lessons learned in this paper can be applied by researchers to their own context and can inform their decisions as they weigh whether the limitations presented here are of significant concern for them given their research goals.

# References

Alaggia, R. & Wang, S. (2020). "I never told anyone until the #metoo movement": What can we learn from sexual abuse and sexual assault disclosures made through social media? *Child Abuse & Neglect, 103*, 1-10. https://doi.org/10.1016/j.chiabu.2019.104312

Bandari, A., & Sun, B. (2021). An online home for the homeless: A content analysis of the subreddit r/homeless. *new media & society*, *00*(0), 1-18. doi: 10.1177/14614448211048615

Barthel, M., Stocking, G., Holcomb, J., and Mitchell. A. (2016). *Nearly eight-in-ten Reddit users get news on the site*. Pew Research Center. https://www.pewresearch.org/journalism/2016/02/25/seven-in-ten-reddit-users-get-news-on-the-site/

Carpenter, J. P. & Willet, K. B. S. (2021). The teachers' lounge and the debate hall: Anonymous self-directed learning in two teaching-related subreddits. *Teaching and Teacher Education, 104*. https://doi.org/10.1016/j.tate.2021.103371

Davis, A. J. & Kettrey, H. H. (2021). Clear and omnipresent danger: Digital age culture wars and reactions to drag queen story hour across diverse subreddit communities. *Social Currents 0*(0), 1-20. doi: 10.1177723294965211050019

del Valle, M. E., & Smit, R. (2021). Moonwalking together: Tracing Redditors' digital memory work on Michael Jackson. *Convergence: The International Journal of Research into New Media Technologies*, *27*(6), 1811-1832. doi: 10.1177/13548665211003878

de Saint Laurent, C., Glaveanu, V. P., & Literat, I. (2021). Internet memes as partial stories: Identifying political narratives in coronavirus memes. *Social Media + Society*, *7*(1), 1-13. doi: 10.1177/2056305121988932

Dixon, S. (2022a). *Reddit usage reach in the United States 2021, by age group*. Statista.

https://www.statista.com/statistics/261766/share-of-us-internet-users-who-use-reddit-by-age-group/

Dixon, S. (2022b). *Reddit usage reach in the United States 2021, by ethnicity.* Statista. https://www.statista.com/statistics/261770/share-of-us-internet-users-who-use-reddit-by-ethnicity/

Dixon, S. (2022c). *Reddit usage reach in the United States 2021, by gender.* Statista. https://www.statista.com/statistics/261765/share-of-us-internet-users-who-use-reddit-by-gender/

Elo, S. & Kyngäs. (2008) The qualitative content analysis process. *Journal of Advanced Nursing*, *62*(1), 107-115. doi: 10.1111/j.1365-2648.2007.04569.x

Feuston, J. L. & Brubaker, J. R. (2021). Putting tools in their place: The role of time and perspective in human-AI collaboration for qualitative analysis. *Proceedings of the ACM on Human-Computer Interaction*, *5*(CSCW2), 1-25. https://doi.org/10.1145/3479856

Garg, R., Kapadia, Y., & Sengupta, S. (2021). Using the lenses of emotion and support to understand unemployment discourse on Reddit. *Proceedings of the ACM on Human-Computer Interaction*, *5*(CSCW1), 1-24. https://doi.org/10.1145/3449088

Gaudette, T., Scrivens, R., Davies, G., & Frank, R. (2021). Upvoting extremism: Collective identify formation and the extreme right on Reddit. *new media & society*, *23*(12), 3491-3508. doi: 10.1177/1461444820958123

Gong, C., Saha, K., Chancellor, S. (2021). "The smartest decision for my future": Social media reveals challenges and stress during post-college life transitions. *Proceedings of the ACM on Human-Computer Interaction*, *5*(CSCW2), 1-29. https://doi.org/10.1145/3476039

Graham, T. & Rodriguez, A. (2021). Sociomateriality of rating and ranking devices on social media: A

    case study of Reddit's voting practices. *Social Media + Society*, *7*(3). doi:

    10.1177/20563051211047667

Gray, C. M., & Chivukula, S. S. (2021). "That's dastardly ingenious": Ethical argumentation strategies

    on Reddit. *Proceedings of the ACM on Human-Computing Interaction*, *5*(CSCW1), 1-25.

    https://doi.org/10.1145/3449144.

Hintz, E. A. & Betts, T. (2022). Reddit in communication research: Current status, future directions

    and best practices. *Annals of the International Communication Association*, *46*(2), 116-133.

    doi: 10.1080/23808985.2022.2064325.

Japec, L., Kreuter, F., Berg, M., Biemer, P., Decker, P., Lampe, C., Lane, J., O'Neil, C., & Usher, A.

    (2015). Big data in survey research. *Public Opinion Quarterly, 79*(4), 839-880.

    https://doi.org/10.1093/poq/nfv039

Jungherr, A., Posegga, Ol, & An, J. (2021). Populist supporters on Reddit: A comparison of content

    and behavioral patterns within publics of supporters of Donald Trump and Hillary Clinton.

    *Social Science Computer Review*, *40*(3), 1-22. doi: 10.1177/0894439321996130

Kaufman, M. R., Bazell, A. T., Collaco, A. & Sedoc, J. (2021). "This show hits really close to home

    on so many levels": An analysis of *Reddit* comments about HBO's *Euphoria* to understand

    viewers' experiences of and reactions to substance use and mental illness. *Drug and Alcohol*

    *Dependence*, *220*. https://doi.org/10.1016/j.drugalcdep.2020.108468

Kimiafar, K. Dadkhah, M., Sarbaz, M., & Mehraeen, M. (2021). An analysis on top commented posts

    in reddit social network about COVID-19. *Journal of Medical Signals & Sensors*, *11*(1), 62-65.

    doi: 10.4103/jmss.JMSS_36_2

Knittel, M., Kollig, F., Mason, A., & Wash, R. (2021). "Anyone else have this experience?": Sharing the emotional labor of tracking data about me. *Proceedings of the ACM on Human-Computing Interaction*, *5*(CSCW1), 1-30. https://doi.org/10.1145/3449153.

LIWC. (n.d.). *LIWC Analysis.* https://www.liwc.app/help/liwc

Nelson, L.K., Burk, D., Knudsen, M., & McCall, L. (2021). The future of coding: A comparison of hand-coding and three types of computer-assisted text analysis methods. *Sociological Methods & Research*, *50*(1), 202-237. doi: 101.1177/049124118769114.

Rare Loot (2018). *Using Pushshift's API to extract Reddit submissions*. https://rareloot.medium.com/using-pushshifts-api-to-extract-reddit-submissions-fb517b286563

*Reddit*. (2022, November 13). In *Wikipedia*. https://en.wikipedia.org/wiki/Reddit

Reddit (2022, November 13a). *About Reddit.* Reddit. https://www.redditinc.com/

Reddit. (2022, November 13b). *Reddiquette.* Reddit. https://www.reddithelp.com/hc/en-us/articles/205926439

Schober, M. F., Pasek, J., Guggenheim, L., Lampe, C., & Conrad, F. G. (2016). Social media analyses for social measurement. *Public Opinion Quarterly, 80*(1), 180-211. https://doi.org/10.1093/poq/nfv048

Seiter, C. R. & Brophy, N. S. (2022). Social support and aggressive communication on social network sites during the COVID-19 pandemic. *Health Communication*, *37*(10), 1295-1304. https://doi.org/10.1080/104110236.2021.1886399

Sen, I., Flöck, F., Weller, K., Weib, B., & Wagner, C. (2021). A total error framework for digital traces of human behavior on online platforms. *Public Opinion Quarterly, 85*(S1), 399-422. https://doi.org/10.1093/poq/nfab018

Seraj, S., Blackburn, K. G, & Pennebaker, J. W. (2021). Language left behind on social media exposes the emotional and cognitive costs of a romantic breakup. *PNAS*, *118*(7). https://doi.org/10.1073/pnas.2017154118

Smith, K. E., Rogers, J. M., Schriefer, D., & Grundmann, O. (2021) Therapeutic benefit with caveats?: Analyzing social media data to understand the complexities of kratom use. *Drug and Alcohol Dependence, 226*. https://doi.org/10.1016/j.drugalcdep.2021.108879

Stevens, H. R., Acic, I., & Taylor, L. D. (2021) Uncivil reactions to sexual assault online: Linguistic features of news reports predict discourse incivility. *Cyberpsychology, Behavior and Social Networking*, *24*(12), 815-821. https://doi.org/10.1089/cyber.2021.0075

Struik, L., & Yang, Y. (2021). e-cigarette cessation: Content analysis of a quit vaping community on Reddit. *Journal of Medical Internet Research*, *23*(10). doi: 10.2196/28303

Tausczik, Y. R., & Pennebaker, J. W. (2010). The psychological meaning of words: LIWC and computerized text analysis methods. *Journal of Language and Social Psychology*, *29*(1), 24-54. doi: 10.1177/0261927X09351676

CHAPTER IV

WHY TEACHERS ARE LEAVING THE POST-PANDEMIC CLASSROOM: A CONTENT ANALYSIS OF #RESIGNATION FLAIR POSTS IN THE R/TEACHERS SUBREDDIT

Several decades of research has shown effective teachers to be the number one factor in determining student achievement (Ferguson, 1991; Nye et al., 2004). The importance of effective teachers matters more now than ever as schools struggle to address the impacts of interrupted learning due to the COVID-19 pandemic. Before the pandemic, our nation was already at risk of experiencing a teacher shortage crisis, especially in high need subject areas and in high-poverty schools (Sutcher et al., 2016). The most recent data from the Teacher Follow-up Survey estimated that 7.7% of teachers left the profession after school year 2012-13 and 8.1% of teachers changed schools for a combined turnover rate of 15.8%. In schools where 75% or more students qualified for free or reduced-price lunch, teacher attrition rates were estimated to be at 9.8% with an additional 12.2% changing schools for an overall turnover rate of 22% (Goldring et al., 2014). For comparison, the rates for Singapore, Finland and Ontario, Canada, some of the highest-performing countries in the world, are between 3-4% (Sutcher et al.). After the pandemic, these trends have only intensified. According to the results of a survey of National Education Association members that was released early in 2022, 55% of teachers surveyed were thinking about leaving the profession earlier than they had planned (GBAO Strategies). In a 2021 Rand Corporation survey of teachers, nearly a quarter of teachers reported they were likely to leave their jobs at the end of the year, a figure which was up from one-sixth of teachers prior to the pandemic (Steiner & Woo).

High rates of teacher turnover perpetuate poor working conditions that make it difficult to build the capacity needed to improve student achievement (Allensworth et al., 2009; Ronfeldt et al., 2013). Furthermore, estimates of the overall economic costs of teacher attrition have ranged from $2.1 billion annually in Texas alone, to between $4.9 to $7.3 billion nationally (Synar & Maiden, 2012). It

is imperative to understand why teachers are leaving and what can be done to make teaching a more sustainable career, however teacher turnover is a multifaceted problem and several deficiencies in the literature exist. First of all, the literature has struggled to both define the cause of teacher turnover and to identify potential solutions. The landscape is complicated by the long list of related topic areas: teacher attrition, teacher burnout, teacher retention, job satisfaction, work and life balance, working conditions and the numerous subtopics found within each. The literature on working conditions has become more robust over time, however there is a lack of agreement amongst researchers studying working conditions on how the construct is defined (Merrill, 2021) which hinders the ability to draw clear conclusions across the literature.

Research into teacher opinion about working conditions and leaving the classroom has relied on traditional survey methodology. As Reveilhac et al. (2021) explain, "survey measures are limited by the questions covered in a questionnaire and don't leave much space for "spontaneous expressions of opinion" (p. 10110). Furthermore, Schaefer et al. (2014) found that ex-teachers tell safe "cover stories" about why they leave teaching. They argue, "that while cover stories may make it easier for teachers to leave teaching and, ironically, for researchers to precisely determine why teachers leave teaching, cover stories may ultimately lead to policies or "fixes" based on teachers' alibis for leaving, instead of their more complex and harder truths" (p. 25). In light of this, it is important to understand the true reasons why teachers are leaving and more authentic sources of teacher voice must be leveraged than what is provided through standard surveys of teacher satisfaction and working conditions.

This paper maps the topics of a sample of posts tagged with the #Resignation flair in the r/Teachers subreddit from three months across school year 2021-2022 filling the gap in the literature around why teachers are leaving the post-pandemic classroom. The window the r/teachers subreddit

provides into the day-to-day experiences of teachers allows educational leaders and policy makers a view of what's happening in classrooms that wouldn't be afforded to them during an in-person visit where educators may be motivated to tell cover stories to hide the reality of their situation – a variation of what educators commonly refer to as the "dog and pony show." Understanding how teachers talk about resigning and the reasons they give for leaving provides a foundation for developing questionnaires that can be used to provide more accurate data for local contexts. This paper provides the first step towards helping policy and decision makers understand what factors matter to teachers when they are making their decisions to stay or go.

**Background: What is Reddit?**

Reddit is a social news aggregation, content rating and discussion website that as of March 2022 is the "9th-most-visited website in the world and the 6th most-visited website in the United States" (Wikipedia, 2022). According to Reddit, the platform has over 50 million unique users active daily, over 100,000 active communities, and includes over 13 billion posts and comments (2023a). Users of Reddit, called "redditors," participate in affinity groups called "subreddits." Subreddits are dedicated to specific topics and adhere to rules set by moderators. Reddit is an anonymous platform which affords users the ability to share their opinions in an unfiltered way without fear of recrimination. Therefore, it is common for users to have conversations about sensitive topics, such as experiences with illicit drugs, or domestic violence, or the real reasons for resigning from teaching, the content of which would either not normally be available to researchers or would come at the cost of extensive relationship building.

Users "upvote" and "downvote" posts and those receiving the most upvotes are elevated by the site's algorithms within subreddits and sitewide. "Upvotes" are similar to "likes" on Facebook but have a slightly different intended use. According to "Reddiquette," or the "informal expression of the

values of many redditors, as written by redditors themselves," users are encouraged to vote on content according to the following guidelines: "If you think something contributes to conversation, upvote it. If you think it does not contribute to the subreddit it is posted in or is off topic in a particular community, downvote it," (Reddit, 2022b). Upvoting and downvoting are supposed to be based upon how well a post contributes to the intended conversation of the subreddit rather than whether someone likes or agrees with the opinions expressed in the post.

The screenshot in Figure 19 shows a post from the r/Teachers subreddit that is tagged with the #Resignation flair. This post has an overall score of 203, meaning the balance of upvotes and downvotes left a net score of 203 upvotes at the time this screenshot was taken (January 30, 2023). The post had 74 comments at the time the screenshot was taken. The username has been crossed out to help protect this user's anonymity.

**Figure 19**

*Example Post Tagged with the #Resignation Flair*

Several surveys have been conducted which draw inferences regarding the primary users of the site. In 2016, Pew Research Center found that about 7% of U.S. adults were using Reddit and that 69% of those users were men, 58% were between the ages of 18-29, 33% were between 30-49, 7% between 50-64, 1% were 65 or older, and 63% were White (Barthel, et al.). In 2021, Statista found that 23% of all male adults, 12% of all female adults, 17% of White adults and 17% of Black adults in the U.S. use Reddit (Dixon, 2022b and c). They found that 36% of users are aged between 18-29, 22% are aged 30-49, 10% are aged 50-64 and 3% are 65 or older (Dixon, 2022a). Though users tend to be younger and male on the aggregate, individual subreddits have their own cultures and may reflect different demographic groups. Seiter and Brophy (2022) found that users on Reddit provided social support more frequently than YouTube and Facebook users and engaged in less aggressive communication. They believe this is due to both the anonymity of Reddit and the rules for civil discourse, or "Reddiquette," found on the Reddit platform.

The r/Teachers subreddit began on December 23, 2008, is self-described as "a sub for all things teacher related!" and has over 400,000 members (Reddit, 2023c). Carpenter and Willet (2021) performed an in-depth study of r/Teachers compared to the subreddit r/Education and found that r/Teachers "can be described as a virtual teachers' lounge: a conversational, chatty space, with a personal and affective flavor" (p. 7) where teachers post personal stories intended for an audience of other teachers.

**Literature Review**

The literature has struggled to define the cause of teacher turnover and identify potential solutions. The issue is further complicated by the lengthy list of related topic areas: attrition, burnout, retention, job satisfaction, working conditions, work and life balance and the numerous subtopics found therein. Rinke shows that "traditional approaches to understanding the problem of teacher

retention typically involve analysis of either individual or workplace characteristics that impact career decisions" (2006, p. 1) and posits that teachers' own perspectives are a necessary addition to the literature. The problem is clearly multifaceted, and many pieces of the problem have been described, but no one solution has yet been found. Furthermore, much of the research has taken place in an international context. Researchers and policymakers in the United States can learn from both international research on teacher turnover and international exemplars of a potential better way. Several international studies have looked at better understanding why teachers leave the teaching profession (Buchanan, 2009; El Helou et a., 2016), but there is a need to better understand why teachers in the United States are leaving their schools.

### *Organizational Characteristics that Contribute to Turnover*

Ingersoll (2001) was the first to find that working conditions contribute to higher rates of teacher turnover and shift the research away from student demographics as a predictor of turnover (for a review see Simon & Johnson, 2015). This study prompted a new round of research aimed at verifying these results and determining which organizational characteristics, or which working conditions, were the most important to teachers when making their decisions to stay or leave (e.g., Allensworth et al., 2009; Bascia & Rottman, 2011; Hughes, 2012; Johnson et al., 2012; Ladd, 2009). Collegial relationships, principal leadership and school culture have been shown to have the greatest effect on reducing teacher turnover in several studies (Johnson et al. 2012; Ladd, 2009). With an understanding of the role a principal can play in determining teacher working conditions, Sterrett et. al (2018) used results from the TELL (Teaching Empowering Learning and Leading) survey in North Carolina to identify a lack of agreement with the items in the construct of time and created a survey and follow up interviews to understand how principals manage and could better manage the time aspect of teachers' duties to address this shortcoming of the teacher work environment. Pay incentives

rewarding effective teachers have not been found to improve student achievement or increase teacher retention (Shifrer et al., 2017). Finally, teacher induction has been shown to mitigate teacher turnover amongst first year teachers (Ronfeldt & McQueen, 2017). The literature on working conditions has become more robust over time, however there is a lack of agreement amongst researchers studying working conditions on how the construct is defined which hinders the ability to draw clear conclusions across the literature (Merrill, 2021).

*Teacher Workload*

Teaching is a demanding job that requires considerable time commitment. In the United States, a 1999 diary study of the time teachers spend working revealed that, on average, teachers worked 10.28 hours per day when including what the authors call "work invasiveness," work issues entering the teachers' personal lives, and an average of eight and a half of those hours were spent at school. The average contractual hours for teachers in the study was six and a half, four hours less than the average number of hours worked per day by the teachers (Drago, et al. 1999). Goldring et al. found that 51% of teachers that left during 2012-2013 stated their workload was more manageable in their new position and 60.8% found it easier to balance work and personal life (2014). Buchanan (2009) recommends reducing class size or face-to-face teaching load based on his findings that one of the five main differences between ex-teachers' current profession and teaching is their workload and responsibility. The Australian ex-teachers he interviewed now work less outside of the workday, experience less workload creep, and have more relaxed work responsibilities. Other international studies explore issues of work intensification - how the requirements of teaching have become more demanding, (Butt & Lance, 2000; Easthope & Easthope, 2000) and how workload and satisfaction are related (Ballet & Kelchtermans, 2008).

In a recent book that studies characteristics of effective international education systems including those of Australia (Victoria and New South Wales), Canada (Alberta and Ontario), Finland, Shanghai, and Singapore, Darling-Hammond et al. found that teacher collaboration and the structural supports to allow for it within the school day were key features of these systems.

> Teachers in high-performing countries spend a great deal of their learning time in collaboration with peers. This is possible because, in many of these countries, teachers spend less of their working day directly in front of students than do teachers in the United States. According to the TALIS [Teaching and Learning International Survey], US teachers, for example, spend about 27 hours a week teaching students directly, about 50% more than the international average of about 19 hours. By contrast, teachers in Singapore spend about 17 hours a week teaching (OECD, 2014d). In Shanghai it is about 15 hours." (2017, p. 113)

Kraft and Papay (2014) found the rate at which teachers grow varies drastically from school to school and over a ten-year period significant differences in teacher effectiveness can result. They point to several factors that contribute to these within school differences including time for teacher collaboration, collegial relationships, and positive school leadership.

### Teacher Characteristics that Contribute to Turnover

Numerous studies exist which focus on the personal characteristics of individual teachers, and a few are mentioned here. A study in Norway investigated how discipline problems, time pressure, low student motivation and value dissonance (personal values not aligned with school values) contributed to the three dimensions of burnout (emotional exhaustion, depersonalization, and reduced personal accomplishment). All four potential stressors were associated with emotional exhaustion and time pressure had the strongest correlation. Discipline problems most strongly predicted depersonalization and personal accomplishment (Skaalvik & Skaalvik, 2017). Another study found

that teachers who left the classroom had weaker self-efficacy and held higher standards for themselves which may have accelerated their burnout whereas teachers who stayed had higher self-efficacy, were more likely to seek out support, and were better at setting boundaries for themselves (Hong, 2012). Finally, the Haberman Star Teacher Interview has been developed to try and select teacher candidates with the right values and belief systems to teach in urban schools and has been cautiously recommended as a method of selecting candidates that will persist (Waddel & Marszalek 2018).

*Narrative Inquiry into Teacher Turnover*

Narrative inquiry has attempted to take a broader look at the individual characteristics that affect teacher turnover. Schaefer et al. (2014) found that ex-teachers tell safe "cover stories" about why they leave teaching. They argue:

> ...while cover stories may make it easier for teachers to leave teaching and, ironically, for researchers to precisely determine why teachers leave teaching, cover stories may ultimately lead to policies or "fixes" based on teachers' alibis for leaving, instead of their more complex and harder truths. (p. 25)

As the truths hidden beneath these cover stories unfolded, Schaefer et al. found a tension between the teachers' personal and professional landscapes; a tension between their imagined selves as professional adults and their actual lives as teachers which ultimately led them to leave teaching. Downey et al., (2014) further explored this idea in a full-length book. They previously had thought teachers moved back and forth between their personal lives and their professional lives, but revised this concept when they realized that teachers never leave their personal selves. The professional self layers over the personal self, and if the two are not able to mesh, either the professional self must be abandoned or the person experiences difficulty.

In a 2015 study, Clandinin et al. interviewed forty beginning teachers about their experience of coming to teaching. Analysis of their interviews uncovered seven themes and three of those seven (balance, trying not to let teaching consume them, and can I keep doing this? Is this teaching?) had to do with the tension between personal and professional landscapes. Clandinin et al. (2015) found that many of the beginning teachers struggled with composing a life as a teacher:

> In this study, the participants helped us to see the importance of framing a teacher's life as much more than who they are in schools. They encouraged us to see them as people living and composing their lives both in school and out of school. Listening to participants, we were struck by the discontinuities and disconnects between their lives in and out of school. Participants told us about how difficult it was for them to compose and live a sustainable life both in school and out of school contexts. (p. 13)

This idea of a sustainable life both in school and out of school contexts is acutely missing from the research and the overall discourse around teacher turnover. Schaefer et al. (2012) assert there is a "need to shift the conversation from one focused only on retaining teachers, toward a conversation about sustaining teachers" (p. 106).

**Reflection on My Identity**

I am an early career teacher leaver. I taught in high-poverty schools for five of my five and a half years of teaching. I found the working conditions and work and life balance to be extremely challenging. I became a teacher leader both in my school and in the larger education community and although this work was extremely rewarding, it was work in addition to my teaching duties and contributed to my own burnout.

Many of my colleagues also left high-poverty schools and the teaching profession altogether. Some left because it did not make sense to have children in day care and teach, some left to pursue

112

graduate degrees, some left simply because they could not keep up with the pace of teaching, and others left for a combination of these reasons. As someone who only recently paid off the debt I accrued for a teaching degree I am not using, someone who wanted to be a teacher, but did not find the career to be sustainable, I have put a lot of thought and now research into what could solve the problem of teacher turnover.

I have wondered about my own personality traits and how they contributed to my burnout. For example, I am an introvert and being "on" all day long, five days a week, was very taxing for me. Is this why teaching did not work well for me? I have wondered whether I did not put up good enough boundaries between my work and personal life. Was I overly dedicated to my students and did not take good enough care of myself? Is this why teaching did not work well for me? I have wondered about the schools where I taught. Did I just not find the right school culture to be a part of? Or is there something unsustainable about being a teacher, especially in a high-poverty school? As I learned more about teachers that had left the classroom before me and looked forward to the time when my husband and I would begin our family, I could not imagine how I could possibly balance teaching with starting a family and began pursuing opportunities that would diversify my skill set. When a change in my husband's job forced us to relocate, I used the opportunity as a safe cover story for officially leaving teaching.

This research stems from a belief that teacher voice needs to be more present in conversations about how to make the teaching profession more sustainable. My own experiences with teaching, including witnessing others, have caused me to develop opinions about what needs to change about the teaching profession to make it more sustainable. As I travel down this path of research, I am very cognizant of this fact and am consciously walking a line between bracketing my own experiences with

teaching and using them to guide the focus of my research and interpret the experiences of others as shared through their posts on Reddit.

**Research Questions**

This paper will explore two research questions to better understand why teachers are leaving the classroom:
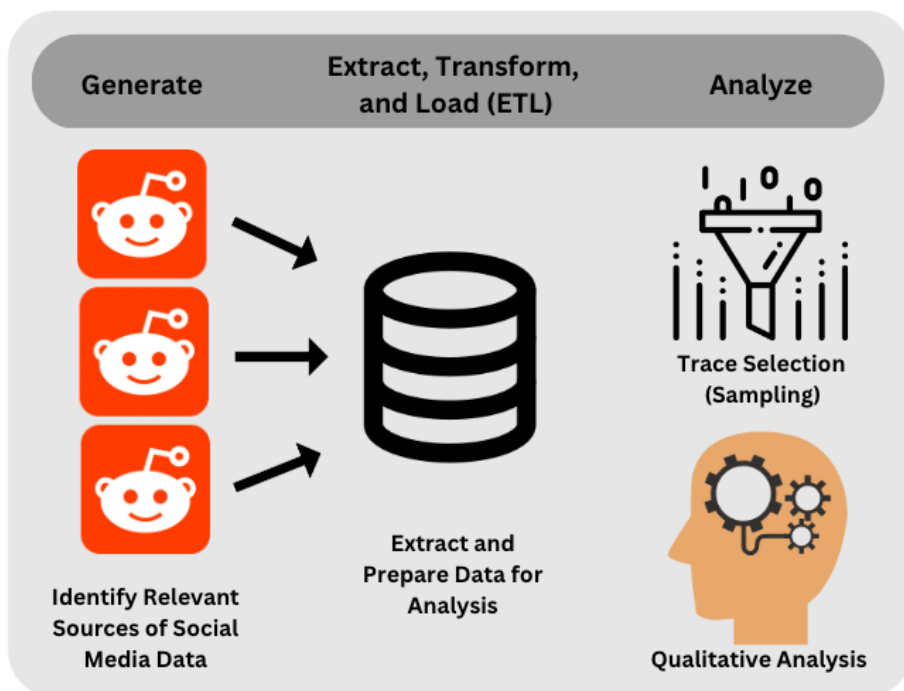
1. What topics do teachers discuss in #Resignation flair posts on the subreddit r/Teachers?

2. What reasons do teachers give in #Resignation flair posts on the subreddit r/Teachers for leaving a teaching job or the teaching profession entirely?

**Methodology**

Previous quantitative research focused on why teachers leave schools is only reporting on the cover stories teachers tell and limits the possible reasons teachers could give by framing the topic through the researcher's lens. This paper leverages a new hybrid approach to the Big Data Process which blends data science skills for data collection with qualitative inductive content analysis for data analysis which can be seen in Figure 20. By utilizing publicly available, yet anonymously shared posts where teachers discuss resignation, there is a greater likelihood of studying the authentic reasons why teachers resign. A qualitative approach is necessary to dig deeper into the phenomenon of teacher resignation to better understand what is unsustainable about the profession and causes teachers to resign. Qualitative research is founded on the constructivist belief that each person has their own interpretation of reality and "these meanings are varied and multiple, leading the researcher to look for the complexity of views" (Creswell & Poth, 2018, p. 24). It is the complexity inherent in the reasons teachers leave which makes qualitative research appropriate for this study.

**Figure 20**

*Hybrid Approach to the Big Data Process*



### *Data Collection*

Posts tagged with the #Resignation flair from the r/Teachers subreddit on the Reddit social media platform were chosen as the unit of study. It was determined by the Western Michigan University Human Subjects Institutional Review Board that analyzing anonymous, publicly available social media data did not constitute as human subjects research and did not require any review. A Python script provided by Rare Loot (2018) was adapted to scrape the desired data from the JSON coding behind the Pushift.io dataset. The Pushshift dataset acts as a mirrored copy of all posts on the Reddit website and was found to be easier to retrieve data from than using Reddit's API (Application Programming Interface) directly. Limitations in coding ability necessitated scraping one day's worth of data at a time.
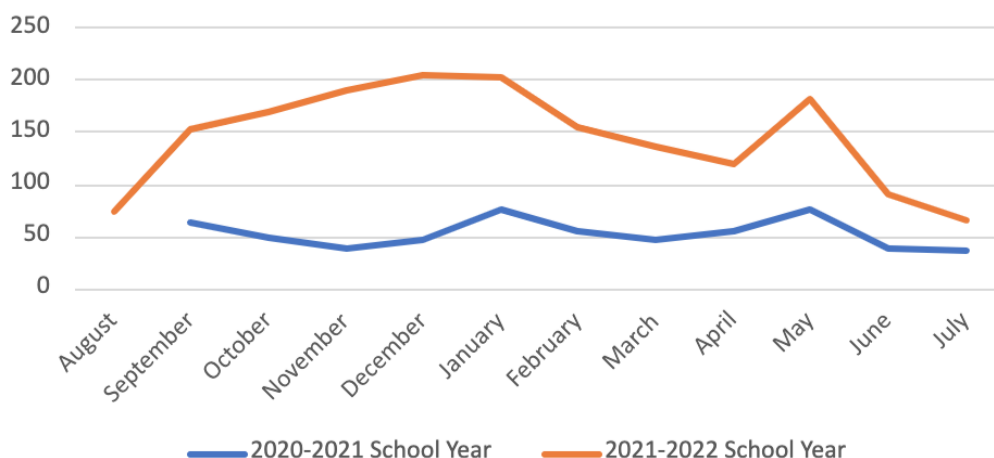
An initial exploration of data from the r/Teachers subreddit was conducted to determine when the #Resignation flair was added. All posts in the r/Teachers subreddit from the day it was introduced, 8/10/2020, through 10/10/2022 were scraped and imported in Microsoft Excel, filtered to only those posts with the #Resignation flair, and combined into one dataset. Data elements scraped for each post included: a unique post id, the post title, the post text, the number of upvotes the post received, the number of comments a post received, the flair attached to the post, the Reddit user id of the original poster, the date and time of the post, the post URL, and the subreddit the post was submitted to. In some cases, the post had been deleted by the original poster right after it was posted, and the mirrored copy did not include post text but rather the message "[deleted]." These posts were filtered and removed.

The number of upvotes and number of comments were found to be taken from the time a post was copied into the Pushshift.io dataset, which usually occurs soon after the post was made and before the post has had a chance to be upvoted or commented on. Therefore, it was decided to retrieve the number of upvotes and number of comments the post had at the time of archiving, if applicable, or on the date of retrieval if the post had not yet been archived by clicking on the public facing URL and manually copying this information into the dataset. Through this process, it was discovered that some posts were deleted from Reddit by the original author after they had been copied into the Pushshift.io dataset. To honor the wishes of the original author to not have their post continue to be publicly available, these posts were removed from the dataset. Additionally, it was discovered that some posts were not actually submitted to the r/Teachers subreddit and/or did not actually have the resignation flair attached to them and these were deleted as well. All data were scraped and manually retrieved in October 2022.

A graph was then generated to show how many #Resignation flair posts were made in the r/Teachers subreddit by month for the 2020-2021 and 2021-2022 school years (see Figure 21) to help determine which data to utilize for the study. The 2021-2022 school year was selected due to the higher number of posts in that year in addition to the increased distance from the COVID-19 pandemic. The months of September, January and May were selected to provide a cross-section of the school year from months with a high volume of posts. All posts from these three months were combined into one final dataset and the find and replace feature in excel was utilized to clean the post titles and post texts of the special characters that were inserted during the scraping process.

**Figure 21**

*Total Number of Resignation Flair Posts in r/teaches by Month*



*Inductive Content Analysis*

An inductive approach to qualitative content analysis based on the procedure outlined by Elo and Kyngäs (2008) was taken. The post title and post body for every post in the dataset was copied into a table in Microsoft Word and the posts were printed in landscape format on 8 ½ x 11 paper. After immersing in the data, codes were recorded in the margins to note either the reasons given for teaching or the "gist" of the post if no reasons were present. During this process it was determined that the posts

were falling into several distinct buckets that might provide a useful framework for analyzing the codes and these buckets were recorded for each post as well. Posts were then organized by bucket, and the abstraction process began as the codes recorded in the margins were cut apart and grouped into categories. Pictures were taken to record the groupings and once all were complete, the categories and codes from the pictures were transferred into a digital format. During this process refinement of categories took place. A log trail was kept throughout the inductive content analysis to document the process, record insights and reactions to the data reflexively, and think through and record decisions made. It was not possible to member check due to the anonymous nature of the data.

**Results**

The final dataset consisted of 315 posts. Inclusion criteria consisted of posts discussing the experience of American, public school (including public charter school), K-12 teachers. Due to the anonymity of Reddit, it is not possible to verify inclusion criteria have been met for each post, however, any posts that clearly did not meet the criteria were excluded. Examples include posts from preschool or day care center teachers, posts from principals or HR representatives, and posts discussing teaching in another country. The upvotes and comments of these posts varied, with averages skewing towards the lower end of engagement. Details can be seen in Table 1.

**Table 12**

*Descriptive Statistics for Upvotes and Comments*

|  | Total # of Posts | Range of # of Upvotes | Mean # of Upvotes | Median # of Upvotes | Range of # of Comments | Mean # of Comments | Median # of Comments |
|---|---|---|---|---|---|---|---|
| Universe | 315 | 0-2,800 | 118 | 22 | 0-522 | 27 | 9 |

*Buckets Identified in the Data*

        To answer the first research question, "What topics do teachers discuss in #Resignation flair posts on the subreddit r/Teachers?" the topics within the main buckets that posts were sorted into were explored. The buckets are: "Commentary, Polls & News," "Leaving Not by Choice," "Support for Leaving," "Thinking About Quitting," "Planning to Quit," "In the Process of Quitting," "Officially Quit," and "Flipside." These eight buckets can be split into two groups, those representing stages along "The Continuum of Quitting," and those that are not on the continuum, or non-continuum buckets. Most posts with the #Resingation flair discussed a teacher's experience at various stages of quitting, with the largest bucket of posts being from teachers that had "Officially Quit" (26.7% of all posts). Additional details about the number and percentage of posts by bucket can be seen in Table 2.

**Table 13**

*Number and Percentage of Posts by Bucket*

| Category | # (%) in the Universe |
|---|---|
| Commentary, Polls & News | 30 (9.5%) |
| Leaving Not by Choice | 12 (3.8%) |
| Support for Leaving | 6 (1.9%) |
| Thinking About Quitting | 60 (19.1%) |
| Planning to Quit | 69 (21.9%) |
| In the Process of Quitting | 37 (11.8%) |
| Officially Quit | 84 (26.7%) |
| Flipside | 17 (5.4%) |

        **Buckets Along "The Continuum of Quitting."** A significant finding of this content analysis is "The Continuum of Quitting" which is derived from the identified buckets. This continuum describes the stages a teacher goes through as they move from "Thinking About Quitting," through

"Planning to Quit," to being "In the Process of Quitting," to being officially done ("Officially Quit"), to having moved on to the next chapter of their life ("Flipside"). Many teachers spend no time on this continuum. Some may only think briefly about quitting once a year when they must declare their intent for the next year to their school, and some may spend years on the continuum agonizing over whether they should move from thinking about quitting to planning to quit. The continuum can operate on two levels: a teacher's relationship with their current job and their relationship with the profession, meaning a teacher may think about quitting their current teaching job to find another one at a different school, or they may be thinking about leaving the profession all together, or they may be considering both at the same time. The placement of posts into these categories was determined by the content of the post, or what the poster was talking about. The main characteristics of each bucket will now be explained - an example post will be given and the codes that fell under logistical concerns along with those that were unique to each bucket will be detailed.

*Thinking About Quitting.* The sixty posts in this bucket represent teachers that are thinking about quitting but have not yet committed to doing so. These posts often contained solicitations for advice on logistical matters, as in this post:

> For those of you that quit teaching how did you do it? I'm six years in and I just can't do this anymore. I am mentally and emotionally exhausted every day. What steps did you take to get out? Any advice would be appreciated.

Five teachers sought advice on other career options, three sought general advice about how to quit, three expressed concerns about breaking contracts, and one person each expressed concerns about leaving mid-year, asked whether you have to give a full two weeks' notice when you quit, and wondered if it looks bad to leave your school because everyone else is. Several non-logistical characteristics were found to be unique to this bucket. Five teachers expressed feelings of conflict

about quitting, four expressed the feeling that they "can't go on like this," three expressed feeling trapped in the teaching profession, and one person expressed the idea of putting themself first as they decided to take mental health leave. These codes can be seen in Figure 22.

*Planning to Quit.* This bucket contained sixty-nine posts representing teachers that have made the decision to quit but have not yet turned in a letter of resignation. They may be just beginning to plan or have a plan fully in place and are just working out the final details of execution. Wherever they are, the main feature is that they have not quit yet, but they are committed to doing so. The unique and logistics codes found in this bucket can be seen in Appendix B, Figure B1. A lot of these posts include solicitations for advice on logistical matters as in this post: "How long did it take for you to find another job? I've probably filled out 100/150 applications and heard back ONCE. I. Need. Out." And this post:

> I have decided that this will be my last year teaching. It is my third year and I need to leave for my own sake. I have three questions: 1. When/how should I tell my admin that I plan on leaving? 2. If I let my teaching license expire, what happens if I want to teach again later? (Illinois license) 3. Those of you who have resigned before, what is some advice you wish you were given before you resigned? Thanks!

Fifteen teachers at this stage were searching for advice around other career options, seven were seeking logistical advice about the quitting process, four expressed concerns about breaking contracts, three sought advice about how to highlight their teacher skills on a resume, one expressed worries about losing insurance if they quit, and one was seeking advice about retirement. Unique to this bucket were two teachers that shared potential opportunities for a new job and asked for advice about whether the opportunity was worth it or not, two that talked about how they had reached their breaking point, and one expressed that they could not wait to quit.
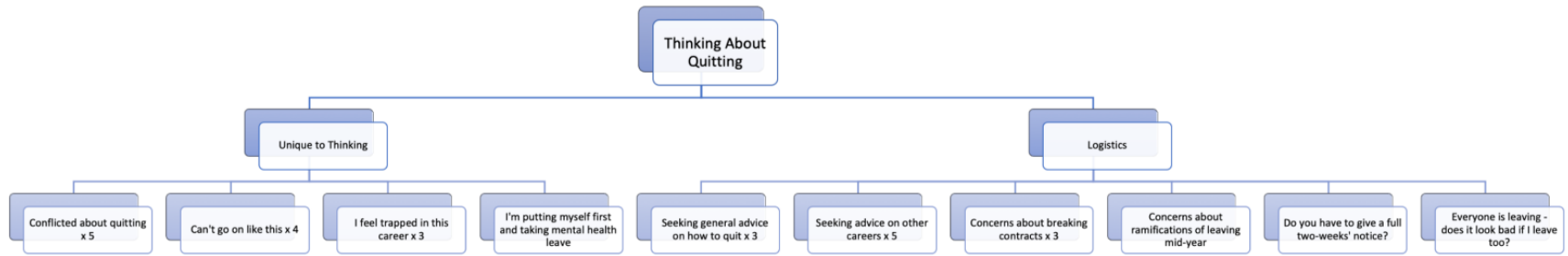
121

***In the Process of Quitting.*** The thirty-seven posts in this bucket represent teachers that turned

in a letter of resignation and are still finishing a period of teaching. Sometimes they are finishing a

two-week or 30-day period and sometimes they are finishing out the rest of the school year.

Sometimes the teacher has found a new job and is trying to understand how to best meet the

expectations of two different contracts, as in this post:

> Contract start date is next week. Resigned today due to spouse's health and accepting a job
>
> closer to home. They stated they would need to fill the spot before they release me. New
>
> district is willing to wait. They can hold me even though I resigned prior to my contract start
>
> date?

Five teachers had questions about contracts at this stage, four had gotten new jobs and were seeking

advice on how to inform their current job they were leaving, three were seeking advice on how to tell

their students they were leaving, three wanted to know if they could use up their sick days before they

left, three were having issues with compensation, three were seeking advice on how to provide

feedback to their district about why they were leaving, one shared their district was refusing to accept

**Figure 22**

*Codes in the "Unique to Thinking About Quitting" and "Logistics" Categories of the "Thinking About Quitting" Bucket*

their resignation, one was having trouble with an administrator over sub plans, and one had a legal question. Unique to this stage in the continuum was that more people (four) were expressing the idea that it is time to put themself first, two told stories about administrators being unprofessional about their quitting, and two expressed that they were "tapped out" until they were done. The following were also expressed by one teacher each: that they were happy to be quitting, that people tried to talk them out of quitting, that annoying things about their job were so much more annoying now that they were almost done, that they were quitting without a plan, and that they were not sure if they made the right choice. The unique and logistics codes for this bucket can be seen in Appendix B, Figure B2.

*Officially Quit.* This bucket was the largest with eighty-four posts. The posts in this bucket represent teachers that have resigned and finished their teaching responsibilities or are just sharing the news that they have quit and are not discussing their process of quitting. They may be going to teach in another school, or they may be done teaching forever. Here is an example:

> I cried in front of my principal - But I did it. Almost had a panic attack as I was approaching their office. They were understanding as I explained how hard this year has been behavior wise while simultaneously dealing with a chronic illness. I've been deeply depressed since August and I finally feel like I can breathe.

This bucket only had one request for logistical support – a question about how to get their retirement contributions back after leaving. The theme of it being time to put themself first showed up strongly in this bucket with fifteen teachers expressing some variation on this theme: five stated it outright in those words, three added "I feel bad, but...," two expressed that they hate to walk away from their dream, but..., two stated their health is more important to them, one that you should not have to be on meds to teach, one that they deserve better, and one that they just want a decent life. Five teachers posted just to share their joy that they'd finally quit, five teachers shared they had decided to stay

home with a baby, three expressed that it felt good to quit, three shared they had quit without a plan, one shared that they'd broken down when they quit, one that it wasn't fair to the kids for them to be checked out, one that they used to love teaching, one that the good of teaching doesn't outweigh the stress and anxiety, and one expressed frustration that teachers are there for students but no one is there for them. The codes unique to this bucket can be found in Appendix B, Figure B3.

*Flipside.* The seventeen posts from this bucket express the vantage point of former teachers that have moved on from teaching and have a new job. These posts are mostly from people that have left teaching all together, such as this one:

> I just popped by to say hello from the afterlife of teaching, in case anyone is wondering what it's like. I resigned at the end of Q3 and started a job in corporate training. I got a 64% pay raise when I switched and I work from home with an awesome computer setup. I have 15 days vacation, 5 days personal leave, additional sick leave, and health insurance is stupid cheap, plus I get about a dozen paid holidays. I've done less in my first few weeks on the job than I did in a few days at school, and everyone I interact with treats me like a human being. I'm not exhausted every day by 3pm, I like my own children again, and I actually look forward to waking up in the morning...

The categories in this bucket followed a different structure than the other buckets along the "Continuum of Quitting" and can be seen in Appendix B, Figure B4. Teachers talked about health benefits, professional benefits, and dealing with the aftermath of leaving teaching in their posts. The most frequently cited benefit of being on the flipside of teaching was improved mental health. Three former teachers talked about being happy now and three talked about having less stress. Having more energy, being relaxed, having low anxiety, no more dread, and looking forward to the week ahead were all cited by one former teacher each. Two former teachers reported they are now treated with

125

respect in their jobs and two more reported they are now appreciated. Two former teachers shared that their personal relationships had improved since leaving teaching and one reported they can now go out and socialize during the week. One teacher reported improved health and one that they could now sleep through the night.

On the professional benefits side, there were big picture and day-to-day benefits. Twelve teachers talked about the pay, nine saying their pay was better and three saying it was the same. One person cited not having to spend their own money on their job and another talked about having the opportunity for bonuses. Four former teachers talked about the presence of a career ladder in their new field and three talked about having better benefits. On a day-to-day level former teachers saw differences in their workload and in their level of freedom and flexibility. Four former teachers found they had less work than before, three found they could leave their work at work, and two talked about having normal hours. Four former teachers talked about having the flexibility to meet other obligations, two talked about having autonomy in their work, two talked about having the ability to work from home, and one talked about having the freedom to go to the bathroom whenever they need to. Despite the overall positive tone of these posts, there were a few posts that discussed the negative aspects of dealing with the aftermath of leaving teaching. Two former teachers reported that they no longer feel fulfilled, one that they miss the kids (but not that much), one that they have to defend their choice to leave the profession to others who do not understand, and one that they are still processing their emotions about leaving.

**Non-Continuum Buckets.** The posts from these buckets did not contain reasons for quitting teaching or thinking about quitting teaching but were still related to the topic of teacher resignation.

*Commentary, Polls, & News.* The thirty posts in this bucket represented people stating their own opinions about teacher resignation, sharing media about teacher resignation, or posing polls to the

group such as this one: "I was just thinking about how many teachers are already considering leaving the profession next year. My question is: 1) Does this apply to you? And 2) What are the main factors that are playing into your decision?" Polling the group to see who is leaving at different points in time was the most popular poll with four different iterations. Two people expressed that they were not seeing lots of teachers resigning where they were and wanted to hear from others if they were. Other poll questions, along with the other categories and codes for this bucket, can be seen in Appendix B, Figure B5.

Three news stories were shared about the number of teachers resigning, one person shared that their district had over 1,000 teacher positions posted on their job website, and one shared that half the teachers at their school had left. Two people talked about teacher pay being too low and one chided Biden for not raising salaries like he said he would. Two people shared links to a video and a blog post from former teachers explaining why they had quit. One person suggested, "let's all quit and see what they do," one thanked resigning teachers for their service, one complained about good teachers being pushed out by administrators, one wondered why principals are surprised when teachers leave, one expressed that the amount of work required to be a teacher is untenable and one teacher declared the subreddit to be a collective funeral mourning the profession.

***Support for Leaving.*** The six posts in this bucket represent unsolicited support and encouragement for leaving teaching that fell into two categories, "Put Yourself First" and "Support for Finding a New Job." Two people gave public service announcements about being burned out, one of which tried to point out the kind of unhealthy rationalizations teachers can make to continue working under extremely stressful conditions, like this one:

'...Going to therapy because of the toll your job is taking is totally fine – there are so many

other stressful careers, right?' If this sounds too familiar – know that it's the environment that's

broken and not you...It's not until you experience a healthy workplace you realize how much

toxicity you have rationalized as normal.

One person declared that teaching is not worth it and encouraged everyone to take care of themselves

and one said people need to recognize when to put themself first and stated it is okay to quit even if

you love teaching. Three people shared suggestions for job opportunities, one suggested to check a

school's turnover rate before applying, and two people offered encouragement for finding a new job

sharing that they had done it and "you can do it too." The categories and codes found in the "Support"

bucket can be seen in Appendix B, Figure B6.

    ***Not by Choice.*** Twelve posts discussed experiences of having to leave a teaching position not

by choice. Most of the codes in this bucket had to do with non-renewals. The categories and codes for

this bucket can be seen in Appendix B, Figure B7. Four teachers told specific stories about being non-

renewed, one of which was a teacher in their tenure year who was being pushed out before reaching

tenure. Three were being offered the opportunity to resign before being officially non-renewed and

were seeking advice about which scenario is better, two said they were already looking for something

else, and two got non-renewed unexpectedly. Here is an example of a post demonstrating several of

these characteristics:

    Hello all. I am a second year teacher out of Oklahoma. Yesterday I was informed that my

    teaching contract will not be renewed for next year. I still have a week until the board meeting.

    In my spare time, I have been pursuing a business degree as I wanted to leave the profession in

    a year or so anyways. I do not plan on coming back to education. My question is, should I

    accept the non-renewal so I have access to unemployment? Or should I resign?

Two posters expressed self-consolations like, "Anyway. At least now I can actually take care of my

family," after sharing about their situation.

*Reasons for Quitting*

To answer the second research question, "What reasons do teachers give in #Resignation flair posts on the subreddit r/Teachers for leaving a teaching job or the teaching profession entirely?" the paper will now turn to the results of inductive content analysis on reasons coded from the dataset. The reasons teachers gave for quitting or thinking about quitting fell into three main categories: "Tolls of the Job," "Bad Environment," and "Issues with the Profession." These main categories and their subcategories can be seen in Appendix B, Figure B8. Note that main categories and subcategories are listed in the order of most frequently appearing to least frequently appearing within the dataset. Each category will now be explored in detail.

**Tolls of the Job.** Teachers talked a lot about the aspects of their job that took a toll on their well-being. The tolls fell into four main categories, "Mental Health Tolls," "Physical Health Tolls," "Feelings of Burn Out," and "Personal Tolls." These subcategories, along with the codes that fall beneath them and their frequencies, can be seen in Appendix B, Figure B9.

*Mental Health Tolls.* The most frequently discussed toll of teaching was poor or declining mental health. These general references to poor or declining mental health occurred thirty-five times. Here is one such reference:

This job has ruined my mental health. I'm back seeing my therapist regularly. I've started getting ulcers in my mouth. I can't sleep at night and I've lost weight from the stress. I can't even take time off or sick days because there's such a terrible substitute shortage.

Depression and anxiety were specifically listed fourteen times each as being a reason for quitting as in this post: "Mentally, I'm done. Depression and anxiety have kicked my ass to the point where my stomach churns before heading to school, to the point where I had days with back to back panic attacks, to the point where I need to love myself more than I care about my students." Teachers

expressed dread about being in the classroom (ten references), shared they were experiencing panic attacks about going in to work (eight references), and talked about having low self-efficacy (five), being emotionally exhausted (three), having suicidal thoughts (two), and generally being tired of being abused (one).

***Physical Health Tolls.*** Seventeen posts made references to poor or declining physical health and seven more talked about being physically drained or running themself ragged. Six posts discussed the impacts of either getting COVID or being at considerable risk for COVID on their decision to leave such as in this post:

> In February of this year, COVID came very close to killing me. Since then, I have really been focused on my health and anxiety levels. I know I am a good teacher, but I also know that this just isn't the job for me anymore.

Six posts referenced issues around being pregnant or post-partum, three referenced issues with chronic illness due to stress, three talked about being physically exhausted, two about having trouble sleeping, and one talked about the pressure of having to be so healthy just to get by.

***Feelings of Burn Out.*** Burn out can be both an emotional and a physical feeling so codes discussing burn out didn't fit clearly into one category or the other. Fourteen posts, talked about feelings of burn out, as in this one:

> I'm burnt out, every weekend I spend just barely coping with the stress of facing another school week. With it's long hours, little personal time, and so much noise. When I get home at the end of the day, all I can do is vaguely exist for a few hours before going to sleep for the coming day.

Thirteen talked about the impact of stress, five talked about being exhausted and three talked about how they were medicating or self-medicating to get by:

Teaching has sucked every single last bit of energy out of me. I'm never happy. Even when I

do go out and do things, I feel overwhelmed and distracted because I'm thinking about work.

I'm in my 20s. I shouldn't be feeling like this. I go to therapy, I take meds. Nothing helps. I feel

like a complete shell of a person.

*Personal Tolls.* Nine people, including this person, discussed the impact teaching was having

on their ability to have a personal life outside of school, "I feel like I never get to do any of the things

that other young professionals my age get to do because I'm so destroyed from work." Seven

specifically mentioned not having enough time for their family and four talked about the difficulty of

balancing personal stress and the stress of teaching, as in this post, "So many terrible things have

happened in my personal life and the added stress of teaching has not helped...I know I need to put

myself and health first but I feel terrible."

**Bad Environment.** Working conditions have long been known to be contributing factors to

teachers leaving teaching. This study found three main categories that mattered to teachers: "Human

Factors," "Physical Environment Factors," and issues with the "Broader Education Context," which

can be seen, along with the codes that fall beneath them and their frequencies, in Appendix B, Figure

B10.

*Human Factors.* The most frequently cited human factor, by far, was bad administrators with

twenty-five references, such as this one, "The administrators made the environment so toxic, and I felt

like I had to walk on eggshells constantly." Eight people talked about how administrators did not

provide an adequate discipline structure, as in this post:

I spent an hour in the office yesterday seeing through the discipline of two students who were

sexually harassing me in class, talking about my body and whether or not they'd "smash" me

right in front of my face. Neither student was suspended or given more than a slap on the wrist.

I have to teach them again today after my admin promised to "fix" the issue.

Six talked about a general lack of support from administrators and two talked about frustrations with the feedback they were getting on their teaching. Fifteen people talked about bad colleagues with problems ranging from being bullied by their colleagues to simply feeling alone because everyone else was so busy trying to keep their own head above water. Eleven people talked about bad school cultures and eleven talked about issues with parents being rude, disrespectful and/or unsupportive. Issues with students also played a very prominent role amongst the reasons teachers gave for leaving. Fourteen talked about students being disrespectful, thirteen generally mentioned difficult student behaviors, eight referenced student violence and fights, seven talked about students not wanting to learn, six talked about students stealing personal property from them, four talked about students damaging or destroying teacher or school property, three talked about being physically assaulted, three about being sexually harassed, three about being sexually assaulted, and two talked about students making school shooting jokes. Here is an example from a post that addresses several of these issues with students:

> What I can't deal with is the non-stop fighting, constant talking and yelling over me, complete disrespect, apathy, complete refusal to listen or follow directions, weekly and now almost daily fist fights in my classroom, students are sawing through chairs and breaking things constantly (deliberately), they steal stuff out of my desk, they curse at me, they tell me how much they hate my class (I'm actually super nice and have pretty fun, upbeat classes).

***Physical Environment Factors.*** Concerns about COVID may be less prominent during the 2021-2022 school year than they were during the 2020-2021 school year, but COVID was still a significant factor for some teachers when making the decision to leave. Eleven expressed concerns

about their safety in a school due to COVID, like this one: "What do you actually put in a letter of resignation? Do I put that I'm leaving due to the school having zero safety policies that are actually followed?" and six expressed frustrations with the effects of COVID on the educational system. Teachers also talked about the rigid school schedule including two references each to having very few breaks, no time to eat, and no time to go to the bathroom, and one reference to having no time to plan.

*Broader Education Context.* Concerns that the education system is failing and not wanting to be a part of it and issues with a district were the most prominent factors discussed by teachers with five references each. These were followed by having no budget for supplies and feeling the effects of politics and the culture wars with four references each, issues with the Board of Education with three references, and finally issues with the state department of education and how difficult it is to take maternity leave with two each. As one person put it, "The last two years in particular have been hell, with the board forcing its way into day to day decisions despite having no educational background."

**Issues with the Profession.** Finally, teachers expressed their frustrations with "Labor Issues," and the fact that teaching is a "High Demand, Low Reward Job." The categories, along with the codes that fall beneath them and their frequencies, are detailed in Appendix B, Figure B11.

*Labor Issues.* The most frequently cited labor issue was low pay. One teacher explained in a very straightforward manner why she was quitting:

> I am a Second year teacher, I teach SPED RESOURCE 9-12. I have only taught during COVID. I turned in my resignation. My admin was great, the students were great and I loved my fellow staff. The problem: I work as a waitress 3 nights a week and make more from that than I do teaching. What's wrong with this world?

Seven teachers also expressed frustrations with having too much work to do as in this post, "Then there's all the work. SO MUCH WORK. Not only professional work, like lesson planning and

grading, but emotional work!" and six cited the fact that it is necessary for them to work overtime. Four talked about the extra work they have been forced to take on as other teachers have left, two talked about having too much paperwork, and two simply stated they were required to spend too many hours with students. Five teachers talked about having poor benefits and two talked about the lack of a career ladder where they could advance in the profession.

*High Demand, Low Reward Job.* Four teachers talked about how teachers are not appreciated, four about how teachers do not receive support, and three talked about how teachers are not respected. Two teachers talked about how they felt like they were good at their jobs but it just is not enough, two thought the demands of the teaching job are not well understood by those outside the profession, two talked about what a physically demanding job it is, and one talked about how teachers are expected to be heroes.

**Summary.** In total, the "Bad Environment" category had 197 codes, the "Tolls of the Job" category had 172, and "Issues with the Profession" had 62 as can be seen in Table 3. When all codes at the subcategory level are tallied up, the most prominent subcategory is "Human Factors" with 145, followed by "Mental Health Tolls" with 92. The tallies for the rest of the subcategories can be seen in Table 4. Looking one level deeper still to the next level of subcategories, again tallied up to this level, the most cited reasons for leaving the classroom were issues with students with 67 references, bad administrators with 41, and poor/declining mental health with 35. The entire "Top Ten" list can be found in Table 5. It is surprising that issues with students were the most frequently cited reasons for leaving as this is not something often found in the literature. The issues with students teachers cited in their posts were frequently behavior issues which are related to issues with administrators not creating a solid discipline structure.

**Table 14**

*Total Number of Codes for Each Main Category*

| Main Category | # of Codes |
|---|---|
| Bad Environment | 197 |
| Tolls of the Job | 172 |
| Issues with the Profession | 62 |

**Table 15**

*Total Number of Codes for Each Subcategory*

| Subcategory | # of Codes |
|---|---|
| Human Factors | 145 |
| Mental Health Tolls | 92 |
| Labor Issues | 45 |
| Physical Health Tolls | 45 |
| Feelings of Burn Out | 35 |
| Broader Education Context | 28 |
| Physical Environment Factors | 24 |
| Personal Tolls | 20 |
| High Demand, Low Reward Job | 17 |

**Table 16**

*Top Ten Most Frequently Cited Reasons for Leaving*

| Reason | # of References |
|---|---|
| Issues with Students | 67 |
| Bad Administrators | 41 |
| Poor/Declining Mental Health | 35 |
| Too Much Work, Not Enough Time | 22 |
| Poor/Declining Physical Health | 17 |

| | |
|---|---|
| Low Pay | 17 |
| Depression | 14 |
| Anxiety | 14 |
| Burned Out Shell of a Human | 14 |
| Stress | 13 |

**Discussion**

Mapping the topics found in #Resignation flair posts in the r/Teachers subreddit shows what issues surrounding resignation are important to teachers. Teachers want to discuss with one another the phenomenon of teacher resignation; they want to seek advice about and support one another through the logistics of the resignation process; they want to share their stories about resignation and connect around their experiences. Reddit provides space for interested teachers to build a community where they can do all these things.

Media related to teacher resignation ranging from blog posts, to podcasts, to YouTube videos, to major news stories are aggregated here under the "Commentary, Polls & News" category creating a repository of all things teacher resignation related that have made their way into the broader culture at large. With so many teachers talking about wanting to leave teaching, it was surprising to find teachers sharing stories of being non-renewed. The fact that so many teachers were being offered the option to resign rather than be non-renewed reveals that statistics regarding teacher resignations may be slightly inflated by teachers who were forced to resign to avoid having a non-renewal on their record.

The posts from the "Flipside" bucket echoed the results found by Buchanan (2009) and Goldring et al. (2014). Teachers who shared stories of leaving spoke of the greater benefits and fewer costs in their new jobs. Teachers cited same or better pay, better benefits, greater flexibility and freedom, a more manageable workload, being valued in their work, and reported better well-being emotionally, physically, and personally. Studying the differences between teaching and former

teachers' new jobs helps shine a light on what needs to change about the teaching profession. Especially as more teachers leave the profession and share their stories of experiencing greater benefits with fewer costs, the profession will continue to hemorrhage teachers while simultaneously disincentivizing people from entering the field.

Posts from the "Continuum of Quitting" show the process teachers move through as they consider, plan to, and eventually quit. A frequent topic of conversation across the first three stages of the "Continuum of Quitting" was questions around the logistics of resigning. As teachers connect across states and districts to discuss these logistics, they are building a new etiquette of resignation as the number of teachers leaving pushes the boundaries on how and when they can, should, and do resign.

The idea of putting yourself first was a prominent theme that showed up primarily in the "In the Process of Quitting" and "Officially Quit" buckets. There was only one instance of putting yourself first in the "Thinking About Quitting" bucket – someone who had decided to put themself first and take mental health leave so they could sort out whether they wanted to quit. There were no references to putting yourself first in the posts from the "Planning to Quit" bucket, there were four in the "In the Process of Quitting" bucket, and there were fifteen in the "Officially Quit" bucket. Teaching is a caregiving profession and teachers often put the needs of others before themselves. In some cases, poor and declining mental and physical health eventually forced teachers to choose to put themselves despite their guilt for having to do so. The act of putting themself first was clearly difficult for many teachers.

Some teachers spoke of a breaking point, a moment when they knew they needed to quit. For some, the moment came after weeks of panic attacks on the way to school in the morning, others during an interaction with a student, or when something happened in their personal life that pushed

them over the edge. Many teachers were able to make the decision and make plans to leave at the end of the school year, but some reached their breaking point and needed to quit as soon as possible. Leaving mid-year has always been taboo in teaching but based on these posts it seems to be happening more frequently as more teachers reach their breaking points and realize they need to put themselves first despite the consequences for students, schools, and communities.

The findings of previous research around teacher turnover are echoed in the results in this paper - references to poor working conditions, especially the impact of bad administrators and unmanageably high workloads, and feelings of burnout are all present. What this study reveals that previous research has missed is how much some teachers' mental and physical health declines due to their jobs. Being expected to perform a high demand, low reward job in a bad environment takes a toll on teachers' emotional, physical, and personal well-being. Teachers cope with this situation for as long as they can until they reach a breaking point – until the tolls of the job outweigh the benefits and they make the difficult decision to leave. Each person has a cost/benefit equation they calculate for their job, and teachers are no different. There have always been teachers who don't see the cost/benefit equation for teaching working in their favor but the number of teachers reaching their breaking point over the last few years seems to be increasing. If we were experiencing a teacher shortage crisis before the pandemic, what are we experiencing now? A full-blown teacher shortage disaster? Based on the frequency of reasons cited by teachers, it seems that human factors are the most prominent factors leading teachers to leave with issues with students and bad administrators topping the list. However, it cannot be assumed that these are the most prominent factors in every school. Leaders and policy makers need a way to determine which factors are tipping the balance in the cost/benefit equation towards teachers leaving at the local level.

*Towards an Evaluation Framework for Teacher Sustainability*

Based on the findings in this paper, an Evaluation Framework for Teacher Sustainability is proposed. The proposed framework can be seen in Appendix B, Figure B12. The framework is based upon the categories found in the "Reasons Why Teachers Leave" section of this paper. Some of the categories have been renamed and/or moved to provide framework categories free from a negative connotation. For example, the category "Bad Environment" becomes "Environment" in the evaluation framework. The "Profession" category found in the data becomes simply "Labor Issues" in this framework and the category "High Demand, Low Reward Job" becomes "Cost/Benefit Check" and moves to a new location. Understanding that each school represents a unique context, this framework, along with the details found in this paper for each category, can serve as the starting point for survey or discussion questions that can be used with teachers to determine which factors need attention at specific schools and districts. Schools are dynamic places and contexts change from year to year. It is necessary to keep a pulse on which factors need to be addressed in which schools to help balance individual teachers' cost/benefit equations towards the direction of teaching being worth it. Teachers deserve not to have to choose between their mental and physical health or their jobs. At all levels, leaders and policy makers need to get creative about finding solutions to the teacher resignation problem. Looking for ways to reduce the burden on teachers from things as small as providing lunch and recess coverage so this duty does not fall on teachers to as grand as doubling the number of teachers and cutting the hours of teaching time in half should be explored. What could teachers accomplish if they were provided with the time, space, and emotional, physical, and personal well-being they deserve to do their jobs well?

**Limitations**

This paper has several limitations. First, the data analyzed represent only three-months of posts tagged with the #Resignation flair in the r/Teachers subreddit. It is possible that including posts from other times of the year may have added categories that are not currently present. It is also possible that posts not tagged with the #Resignation flair included discussion of the reasons why teachers leave that may have been relevant to answering the second research question but were not included in the dataset. Another limitation of this paper is that the findings cannot be assumed to represent the experience of all teachers who resign, they can only be assumed to represent the teachers who posted about their experiences resigning in the r/Teachers subreddit during the time period sampled. It is possible there are differences between someone who posts about their experiences resigning in a social media platform and someone who does not. Additionally, the findings of this paper will note generalize to teachers at all levels, teachers of all subjects, or teachers from all demographic areas. Finally, it was not possible to member check these results due to the anonymous nature of the data.

**Directions for Future Research**

This study constitutes exploratory research and the categories which took shape need to undergo refinement. Further research should be conducted analyzing posts from across a greater time period to help provide this refinement. The prominence of categories could also be analyzed over time to see if different factors are more critical at different times of the school year or if the prominence of categories were different during COVID-19 or if they will continue to change as we move fully past the pandemic. As of January 1, 2023, the moderators of the r/Teachers subreddit have determined, based on feedback from their members, that discussions about resigning from teaching are no longer appropriate for the subreddit. They have retired the #Resignation flair and are redirecting people to the subreddit r/TeachersInTransition to have these conversations. Therefore, future research should

incorporate this subreddit into the data collection process. It is also possible that conversations about teacher resignation are happening in other subreddits, and a more exhaustive search could be conducted to capture more conversations from across Reddit. As the categories in this paper are refined, the Evaluation Framework for Teaching Sustainability can also be refined, and questions aligned to this framework should be created to help local leaders assess areas of concern in their local contexts.

Having a map of topics around teacher resignation allows researchers to understand what further research is possible in this area. The non-continuum buckets provide insight into teachers leaving not by choice, how teachers offer support to others around burnout and resignation, and how teachers discuss the phenomenon of teacher resignation. Researchers interested in teacher conversation around resignation could focus on comments on the kinds of posts found in the "Commentary, Polls & News" bucket. Aggregating the results of the multiple polls questions found within the subreddit could give flashpoints of teacher opinion over time – especially interesting since the same questions seemed to be asked frequently. This bucket also provides an aggregated source of media around teacher resignation, including lesser-known media published across various platforms which might be difficult to come across in other ways. This sample of posts identified several news articles from prominent news sources in addition to a teacher blog, a teacher created YouTube video and a teacher created podcast called The Great Teacher Resignation. Researchers interested in how people make difficult life decisions and the kind of support required to do so could study how people seek and receive support in this anonymous, online space. Finally, each stage of the Continuum of Teaching showed unique characteristics and future research could focus on any of these stages to better understand teacher experience at that stage.

**Conclusion**

There is no one reason teachers leave teaching – there are an infinite number of combinations of variables, both personal and professional, that teachers weigh. Unfortunately, the costs of teaching are outweighing the benefits of teaching for more teachers, and they are choosing to leave classrooms. To reverse this trend, education leaders and policy makers need to address all aspects of the teaching cost/benefit equation. The benefits of teaching need to increase to help offset the costs, and the costs need to be reduced. Leaders and policy makers need to listen at the school and district level to understand which factors might have the most impact in shifting this balance in their local context. Some solutions may seem obvious, like raising teacher pay, but this alone may not be enough. Better pay may encourage teachers to stay in bad situations longer, but it does not change the fact that they are experiencing bad work environments and suffering from the emotional, physical, and personal tolls of the job. Schools should be a place where teachers feel balanced, supported, and like they have enough time to meet the demands of their job – for their own sake, and for the sake of their students. The findings of this paper around the reasons why teachers leave the classroom provides policy makers with a foundation for questionnaire creation so that teachers can be surveyed regarding their thoughts on the sustainability of the teaching profession.

Teachers are the number one determining factor in student success – it is time to treat them in a way that honors this fact. It is time to step back and reconsider the structure of schools to address the high workload teachers experience to bring more balance to their lives. Creative changes to schedules and staffing to reduce the number of hours teacher spend with students, freeing them up to collaborate with peers, plan lessons, improve their craft, and build relationships with families during school hours so they do not have to do these things on their own time would be a great place to start.

## References

Allensworth, E., Ponisciak, S., & Mazzeo, C. (2009). The schools teachers leave: Teacher mobility in

    Chicago public schools. Chicago, IL: Consortium on Chicago School Research.

Ballet, K., & Kelchtermans, G. (2009). Struggling with workload: Primary teachers' experience of

    intensification. *Teaching and Teacher Education, 25*, 1150-1157.

Barthel, M., Stocking, G., Holcomb, J., and Mitchell. A. (2016). *Nearly eight-in-ten Reddit users get*

    *news on the site*. Pew Research Center.

    https://www.pewresearch.org/journalism/2016/02/25/seven-in-ten-reddit-users-get-news-on-

    the-site/

Bascia, N., & Rottman, C. (2011). What's so important about teachers' working conditions? The fatal

    flaw in North American education reform. *Journal of Education Policy*, *26*, 787-802.

    doi:10.1080/02680939.2010.543156

Buchanan, J. (2009). Where are they now? Ex-teachers tell their life-work stories. *Issues in*

    *Educational Research, 19*, 1-13.

    https://journals.sagepub.com/doi/10.1177/000494411205600207

Butt, G., & Lance, A. (2005). Secondary teacher workload and job satisfaction: Do successful

    strategies for change exist? *Educational Management Administration & Leadership, 33*, 401-

    422. doi: 10.1177/1741143205056304

Carpenter, J. P., & Bret Staudt Willet, K. (2021). The teachers' lounge and the debate hall:

    Anonymous self-directed learning in two teaching-related subreddits. Teaching and Teacher

    Education (104). https://doi.org/10.1016/j.tate.2021.103371

Clandinin, D. J., Long, J., Schaefer, L., Downey, C. A., Steeves, P., Pinnegar, E., McKenzie Robblee,
S., & Wnuk, S. (2015). Early career attrition: Intentions of teachers beginning. *Teaching
Education, 26*(1), 1-16. doi:10.1080/10476210.2014.996746

Creswell, J. W., Poth, C. N. (2018). *Qualitative inquiry & research design: Choosing among five
approaches*. Thousand Oaks, CA: SAGE Publications, Inc.

Darling-Hammond, L., Burns, D., Campbell, C., Goodwin, A. L., Hammerness, K., Low, E. L.,
McIntyre, A., Sato, M., & Zeichner, K. (2017). *Empowered educators: How high-performing
systems shape teaching quality around the world*. San Francisco, CA: Jossey-Bass.

Dixon, S. (2022a). *Reddit usage reach in the United States 2021, by age group*. Statista.
https://www.statista.com/statistics/261766/share-of-us-internet-users-who-use-reddit-by-age-
group/

Dixon, S. (2022b). *Reddit usage reach in the United States 2021, by ethnicity.* Statista.
https://www.statista.com/statistics/261770/share-of-us-internet-users-who-use-reddit-by-
ethnicity/

Dixon, S. (2022c). *Reddit usage reach in the United States 2021, by gender.* Statista.
https://www.statista.com/statistics/261765/share-of-us-internet-users-who-use-reddit-by-
gender/

Downey, C. A., Schaefer, L., & Clandinin, D. J. (2014). *Narrative conceptions of knowledge: Towards
understanding teacher attrition*. Bingley, United Kingdom: Emerald Publishing Limited. doi:
10.1108/S1479-368720140000023003

Drago, R., Caplan, R., Costanza, D., Brubaker, T., Cloud, D., Harris, N., Kashian, R., & Riggs, T. L.
(1999). New estimates of working time for elementary school teachers. *Monthly Labor Review,
122*(4), 31-40.

Easthope, C., & Easthope, G. (2000). Intensification, extension and complexity of teachers' workload. *British Journal of Sociology of Education, 2,* 43-58.

El Helou, M., Nabhani, M., & Bahous, R. (2016). Teachers' views on causes leading to their burnout. *School Leadership & Management, 36*, 551-567. doi: 10.1080/13632434.2016.1247051

Elo, S. & Kyngäs, H. (2007). The qualitative content analysis process. Journal of Advanced Nursing 62(1), 107-115. doi: 10.1111/j.1365-2648.04569.x

Ferguson, R. (1991). Paying for public education: New evidence on how and why money matters. *Harvard Journal on Legislation, 28*, 465-498.

GBAO Strategies (2022). Poll Results: Stress and burnout pose threat of educator shortages. Washington, DC. https://www.nea.org/sites/default/files/2022-02/NEA%20Member%20COVID-19%20Survey%20Summary.pdf

Goldring, R., Sheyla, T., & Riddles, M. (2014). Teacher attrition and mobility: Results from the 2012-2013 Teacher Follow-Up Survey. National Center for Education Statistics, U.S. Department of Education. https://nces.ed.gov/pubs2014/2014077.pdf

Hirsch, E. & Emerick, S. (2006). Teaching working conditions are student learning conditions: A report on the 2006 North Carolina teacher working conditions survey. Center for Teaching Quality. https://files.eric.ed.gov/fulltext/ED498770.pdf

Hong, J. Y. (2012). Why do some beginning teachers leave the school, and others stay? Understanding teacher resilience through psychological lenses. *Teachers and Teaching, 18*(4), 417-440. doi: 10.1080/13540602.2012.696044

Hughes, G. (2012). Teacher retention: Teacher characteristics, school characteristics, organizational characteristics, and teacher efficacy. *The Journal of Educational Research*, *105*, 245-255. doi: 10.1080/00220671.2011.584922

Ingersoll, R. (2001). Teacher turnover and teacher shortages: An organizational analysis. *American Educational Research Journal*, *39*, 499-534. doi: 10.3102/00028312038003499

Johnson, S., Kraft, M. A., & Papay, J. P. (2012). How context matters in high-need schools: The effects of teachers' working conditions on their professional satisfaction and their students' achievement. Teachers College Record 114(100306).

Kraft, M. A., & Papay, J. P. (2014). Can professional environments in schools promote teacher development? Explaining heterogeneity in returns to teaching experience. *Educational Evaluation and Policy Analysis, 36*(4), 476-500. doi:10.3102/0162373713519496

Ladd, H. (2009). Teachers' perceptions of their working conditions: How predictive of policy-relevant outcomes? *Educational Evaluation and Policy Analysis*, *33*, 235-26. doi: 10.3102/0162373711398128

Merrill, B. C. (2021). Configuring a construct definition of teacher working conditions in the United States: A systematic narrative review of researcher concepts. Review of Educational Research 91(2). Doi: 10.3102/0034654320985611

Nye, B., Konstantopoulos, S., & Hedges, L. (2004). How large are teacher effects? *Educational Evaluation and Policy Analysis, 26*, 237-257.

Rare Loot (2018). *Using Pushshift's API to extract Reddit submissions*. https://rareloot.medium.com/using-pushshifts-api-to-extract-reddit-submissions-fb517b286563

*Reddit*. (2022, November 13). In *Wikipedia*. https://en.wikipedia.org/wiki/Reddit

Reddit (2022, November 13a). *About Reddit.* Reddit. https://www.redditinc.com/

Reddit. (2022, November 13b). *Reddiquette.* Reddit. https://www.reddithelp.com/hc/en-us/articles/205926439

Reddit. (2023, January 24c). *r/Teachers*. https://www.reddit.com/r/Teachers/

Reveilhac, M., Steinmetz, S., & Morselli, D. (2022). A systematic literature review of how and

    whether social media data can complement traditional survey data to study public opinion.

    *Multimedia Tools and Applications, 81*. 10107-10142. https://doi.org/10.1007/s11042-022-

    12101-0

Rinke, C. R. (2006). Understanding teachers' careers: Linking professional life to professional path.

    Educational Research Review, 3(1), 1-13. https://doi.org/10.1016/j.edurev.2007.10.001

Ronfeldt, M., Loeb, S., & Wyckoff, J. (2013). How teacher turnover harms student achievement.

    *American Educational Research Journal, 50*, 4-36. doi: 10.3102/0002831212463813

Ronfeldt, M., & McQueen, K. (2017). Does new teacher induction really improve retention? *Journal

    of Teacher Education, 68*(4), 394-410. doi: 10.1177/0022487117702583

Schaefer, L., Downey, C. A., & Clandinin, D. J. (2014). Shifting from stories to live by to stories to

    leave by: Early career teacher attrition. *Teacher Education Quarterly, 41*(1), 9-27.

Schaefer, L., Long, J., & Clandinin, D. J. (2012). Questioning the research on early career teacher

    attrition and retention. *Alberta Journal of Educational Research, 58*, 106-121.

Seiter, C. R. & Brophy, N. S. (2022). Social support and aggressive communication on social network

    sites during the COVID-19 pandemic. *Health Communication*, *37*(10), 1295-1304.

    https://doi.org/10.1080/104110236.2021.1886399

Shifrer, D., Turley, R. L., & Heard, H. (2017). Do teacher financial awards improve teacher retention

    and student achievement in an urban disadvantaged school district? *American Educational

    Research Journal, 54*(6), 1117-1153. doi: 10.3102/0002831217716540

Simon, N., & Johnson, S. (2015). Teacher turnover in high-poverty schools: what we know and can

    do. *Teachers College Record*, 117(3), 1-36.

Skaalvik, E. M, & Skaalvik, S. (2017). Dimensions of teacher burnout: relations with potential

    stressors at school. *Social Psychology of Education, 20*(4), 775-790. doi: 10.1007/s11218-017-

    9391-0

Steiner, E. D., & Woo, A. (2021). Job-related stress threatens the teacher supply: Key findings from

    the 2021 State of the U.S. Teacher Survey. Rand Corporation.

Sterret, W. L., Parker, M. A., & Mitzner, K. (2018). Maximizing teacher time: The collaborative

    leadership role of the principal. Journal of Organizational and Educational Leadership 3(2).

Sutcher, L., Darling-Hammond, L., & Carver-Thomas, D. (2016). *A coming crisis in teaching?*

    *Teacher supply, demand, and shortages in the U.S*. Palo Alto, CA: Learning Policy Institute.

Synar, E. & Maiden, J. 2012. A comprehensive model for estimating the financial impact of teacher

    turnover. *Journal of Education Finance 38*(2), 130-144.

Waddel, J. H & Marszalek, J. M. (2018). Haberman Star Teacher Interview as a predictor of success in

    urban teacher preparation. *Educational Policy Analysis Archives, 26(*35), 1-33. doi:

    10.14507/epaa.26.2808

CHAPTER V

CONCLUSION

This chapter concludes this dissertation by summarizing key findings from each paper and drawing conclusions about one way data science skills can be integrated with social science research methodology to provide valid insight into policy relevant issues at a reduced cost in time, money and effort. Limitations and future directions of research will also be discussed.

**Review of Main Findings**

The three papers in this three-paper dissertation worked together to explore how Reddit can be leveraged for social science research. They built from surveying current research methods for leveraging Reddit data in social science research, to contributing to the literature on trace selection error, to demonstrating how Reddit data can be leveraged to address a current issue that is of interest to policy makers today. The first paper in this dissertation provided a review of over 150 articles published during 2021 that leveraged Reddit in some way for social science research. Forty-nine of the articles leveraged Reddit for recruitment and 26 of those came from the field of human-centered computing. The two main arguments for using Reddit for recruitment were the access it provides to stigmatized topics and groups and the anonymity of the platform which allows for honesty. The majority of researchers used Reddit as one of several methods for recruitment. Commonly cited limitations of using Reddit for recruitment were the lack of diversity in the population of Reddit users and that it caused the participant pool to be biased towards the kind of person that uses social media. The other 120 articles also most commonly cited the anonymity of the platform and its ability to afford researchers access to hard-to-reach populations and opinions on stigmatizing topics as arguments for leveraging Reddit. Saving time and money and the ability to access organic data that is free from reearcher bias were also cited. The most commonly cited limitations were the lack of demographic data available for users and the inability to member check understanding and results. Reddit data was

most frequently analyzed using qualitative methods (67 articles) and mixed-methods (27 articles). A wide range of data science skills were represented and great variability of transparency in reporting was noted.

The second paper compared two methods of sampling Reddit data - sampling from top most-upvoted posts and sampling from top most-commented-on posts to understand the differences between these two sampling methods as they compare to the results of analyzing the universe of posts. Overall, differences between the samples and the universe were more apparent in what people were saying rather than how people were saying it. There were no statistically significant differences between the datasets on four out of the five style variables tested: word, count, analytic thinking, emotional tone, and authenticity. There was a statistically significant difference between the most-commented on sample and everything else in the universe on clout. Several categories of posts were over- or under-represented in the samples showing there may be biases towards certain topics in samples of top most-commented on and top most-upvoted posts. There were also categories in the universe which didn't appear in the samples at all showing that researchers interested in achieving topic coverage should not rely on these sampling methods to achieve it. Overall, sampling from most-commented on posts proved to good way to find posts with lots of polls in them, if that is something that is of interest to a researcher, but did not provide as robust of a topic map as the sample of most upvoted posts. Both samples failed to provide full topic coverage, but the map made from the sample of most upvoted posts included more detail.

The third paper utilized Reddit data to map the topics of teachers' conversations around leaving teaching and provide insight into the reasons why teachers leave. Teachers commented on teacher resignation, posed poll questions, shared news stories about teacher resignation, provided support for leaving teaching and talked about leaving not by choice in addition to a range of topics

150

related to the Continuum of Quitting. The Continuum of Quitting consists of Thinking About Quitting, Planning to Quit, In the Process of Quitting, Officially Quit, and the Flipside. The majority of posts, over a quarter of the total, were about being officially quit. The reasons teachers gave for leaving fell into three main categories: Tolls of the Job, Bad Environment, and Issues with the Profession. The most frequently cited reasons for leaving were issues with students (67 citations), bad administrators (41), and poor/declining mental health (35).

**Conclusion**

Working with Big Data is not new to the private sector and the past few decades have seen interest in leveraging data science skills for social science research, especially to mine public opinion data from social media platforms (Schober et al. 2016). The consensus is that big data analysis can complement traditional survey research but will not replace it due to the inherent inability to generalize knowledge to populations in the same way that traditional survey research can (Japec et al., 2015; Schober et al., 2016; Sen et al. 2021; Reveilhac et al., 2022). Leveraging social media data for data collection can save social science researchers time and money on research designs, provide insight into emerging, rapidly evolving, and stigmatizing topics, and leverage the platform for the recruitment of research participants

Each social media platform has its own unique features and Reddit has some particular platform affordances that make it especially appealing for several challenging research functions. First, the anonymous nature of the platform facilitates access to hard-to-reach populations and discussions about sensitive and stigmatizing topics that would be either extremely difficult to access or completely inaccessible to social science researchers otherwise. Secondly, Reddit is a topic-based community which means it self-organizes around topics of mutual interest, and because of its size and reach, it creates communities around very specific things that people care about. This feature makes it

both more likely that stigmatizing and obscure topics are being discussed on the platform and easier for researchers to filter all the data on the site to find discussions around their research topics of interest, which may help reduce trace selection error when determining what data should be extracted from the platform, especially for novice data scientists. Because of these two unique features, Reddit is an especially appropriate platform choice for conducting exploratory research, especially on hard-to-reach populations and/or stigmatizing topics; and to recruit participants for studies, especially participants from hard-to-reach populations.

This dissertation helps bridge the gap between data science and social science research so that Reddit data can be leveraged to further social science research goals in several ways. First, the dissertation provided a review of current social science research leveraging Reddit to bring greater knowledge of what is possible using data science methodologies in the social science context. This knowledge will help social science researchers ask and answer new and exciting questions in their fields. Though qualitative researchers may not be ready to hand over analysis tasks to Artificial Intelligence, they may benefit from the use of data science methods during the data collection phase and turn to traditional qualitative methods in the analysis phase – a hybrid approach to the Big Data Process outlined by Japec et al. (2016) which can be seen in Figure 17.

This approach allows qualitative researchers to access a trove of qualitative data on a full range of topics, including opinions of hard-to-reach populations and opinions on stigmatizing topics. However, human instruments have less processing power than AI and qualitative researchers must make choices about how to filter, or sample from, all the available data to create a more manageable dataset for human-coded analysis. Social science researchers must become more familiar with all sources of error in the Total Error Framework for Digital Traces of Human Behavior on Online Platforms (TED-On)

**Figure 17**

*Hybrid Approach to the Big Data Process*



(Sen et al., 2021). This dissertation took a closer look at the impact of trace selection error, one error component in the TED-On to help researchers consider potential sources of error that may occur when selecting data from social media platforms to analyze. The findings of Paper 1 showed a wide range of trace selection methods being utilized in current research. Because the practice of leveraging social media data for social science research is still relatively new, it will be imperative that clearer norms for reporting and developed and adhered to moving forward to ensure rigorous design choices that will produce valid and reliable results.

At this point, no clear guidelines exist around reporting norms for articles utilizing Reddit. A short checklist to begin creating such guidelines was presented in Paper 1 of this dissertation. This checklist can help researchers report their findings and can also help editors and peer reviewers ensure high levels of transparency around methodological choices appear in their journals. Published articles utilizing Reddit in their research should include the following items in their methodology section:

1. Consider the sources of error on the TED-On and transparently report all research decisions that contribute to potential sources of error;

2. State whether IRB review was necessary and/or obtained;

3. Report what ethical considerations were taken throughout the research process to protect the identity of Reddit users whose conversation were used in the study;

4. Report the date data was collected;

5. Report the time frame of the data collected;

6. Report the level of the data collected and analyzed (subreddit, threads, posts, comments, etc.);

7. Report the number of units (subreddits, threads, posts, comments, etc.) that were analyzed;

8. Provide a detailed report of the data extraction method utilized that goes beyond simply stating a software package; and

9. Provide a detailed report of the sampling methods used, if applicable.

A recommended practice is to post any code used for the paper to github, an open-source website community where software developers collaborate (Biester et al., 2021).

It will also be critical for social scientists to forge partnerships with data scientists moving forward to further integrate the two fields. This dissertation proposed one method of combining skills from each field and others will surely be developed as time goes on. Japec et al. Suggest a minimum of four roles required to successfully work with big data: a domain expert, a researcher, a computer scientist and a system administrator (859). Ensuring each of these roles are represented on a team utilizing social media data for research will further help ensure sources of error are kept to a minimum.

Finally, this dissertation suggests a novel mixed-methods design for investigating opinions of hard-to-reach populations and/or stigmatizing topics which can be used to inform decision makers on issues of policy interest. The design can be seen in Figure 18. The approach begins with the hybrid

approach to the Big Data Process, where data science methods are used during the data collection phase and traditional qualitative methods are used in the analysis phase, which was demonstrated in Paper 3. The results of qualitative analysis can be utilized to build a questionnaire, which can then be distributed through social media platforms to help recruit survey participants. This approach provides two levels of insight into policy relevant issues by first providing a map of a relevant topic and then providing data that is more representative of a population's opinions on that topic.

Other fields of research have already begun to leverage extant Reddit data to reach hard-to-reach populations and/or opinions on stigmatizing topics, but educational research has barely scratched the surface of this resource and certainly has not used it to provide insight to education leaders and policy makers. The window the r/teachers subreddit provides into the day-to-day experiences of teachers allows educational leaders and policy makers a view of what's happening in classrooms that wouldn't be afforded to them during an in-person visit where educators may be motivated to tell cover stories to hide the reality of their situation – a variation of what educators commonly refer to as the "dog and pony show." In the absence of the ability to generate gold-standard survey data about why teachers are leaving the classroom, this dissertation identifies a path forward for better understanding the critical question of why teachers are leaving the classroom. This path can be used to study why people leave other professions or can be applied to any policy issue about which current survey results are not available. Social media data, especially anonymous social media data like that found on Reddit, can provide timely insight to policy makers on a range of topics that may be emerging, rapidly evolving, and/or not freely discussed in other arenas, and provide the basis for further quantitative research on these topics.

**Figure 18**

*Mixed-Methods Design for Investigating Opinions of Hard-to-Reach Populations and/or Stigmatizing*

*Topics*

It is important to keep in mind that the social media landscape is constantly evolving and that what works today may not continue to work tomorrow (Japec et al., p. 850). What happens when users of Reddit or members of particular subreddits become aware they are being observed and their opinions are being used for research and policy decision purposes? Will the fact of mere observation render these online spaces less safe and thus hinder the sharing of authentic opinion which makes the platform so appealing for research use? Will imposters be motivated to invade what have heretofore been spaces where anonymity has afforded users space to be honest and change the culture through misrepresentation or trolling? Whether these things happen or not, Reddit could change their terms of service limiting the availability and use of data within the platform. Change is unavoidable and researchers will have to stay up to date with the most current trends in leveraging social media data for social media research.

**Limitations**

Though this dissertation makes several contributions to research methodology, it does have some limitations. An exhaustive search of several of the largest social science databases was conducted for the first paper, however it is possible some articles were missed. Additionally, it is possible the author's interpretation of what articles were using Reddit in a significant way to conduct social science research may have been different from someone else's. This paper did not include articles from a large body of work being done in the machine learning and natural language processing area because these were data science focused. Future research could focus on understanding the full extent of what is possible using data science methods of analysis. The data for papers 2 and 3 were scraped from one subreddit and may not generalize to data scraped from other subreddits or social media platforms and cannot be generalized to the general population of teachers. Additionally, decisions made when building the universe of data for this study may have affected the results. Data

157

selected from a different year, from different months, or with a different flair may have yielded different results. Filtering available data by flair was an easy way to limit the universe of data and ensure posts were relevant to the topic of research interest, however, there are other ways to filter available data that may have provided a more comprehensive dataset on the same topic. Furthermore, the script used to scrape data yielded some inaccuracies in filtering by the established criteria and though posts that should not have been included were excluded from the dataset (posts deleted by users and posts that weren't actually from the subreddit of interest, for example) there is no way of knowing if the converse happened - whether posts that should have been scraped were not.

Another limitation of this dissertation is that analyses for the third paper, which focused on the universe of data, were conducted simultaneously with the qualitative analyses of the universe for this paper and took place after quantitative analyses had been completed. Twenty-one posts that had been included in the quantitative analyses were excluded from the qualitative analyses after further review revealed they did not meet inclusion criteria for the third paper, leaving a slight difference in the posts analyzed quantitatively versus qualitatively. Another limitation of Paper 2 is that small size of the samples. It is possible some of the differences between the samples and the universe are due to too small of a sample size rather than the sampling method. Finally, this study focused on two methods of sampling when more sampling methods are available and could have been incorporated into the study. The second paper focused on two methods of sampling when more sampling methods are available and could have been incorporated into the study. Further research is needed to help build the literature on trace selection error, including replication of the findings in Paper 2, and exploration into the consequences of other sampling methods for trace reduction. Additionally, the findings Paper 3 will note generalize to teachers at all levels, teachers of all subjects, or teachers from all demographic

areas. Finally, it was not possible to member check these results due to the anonymous nature of the data.

**Directions for Future Research**

As the use of Reddit for social science research continues to grow and Reddit itself continues to grow, it will be important to continue to research novel methodologies by providing surveys of current methodologies such as this paper has provided. Further research should be done on the demographics of Reddit users, including variation of demographics in different subreddits to better understand who is using Reddit. Entwistle et al. found they were able to mine demographic data for 80% of the users in the dataset they analyzed allowing them draw conclusions by demographic group (2021). This method could be applied more broadly across the platform. Additionally, more research is needed into the different sources of error in the TED-On framework, such as how different sampling methods affect the results of analyses of Reddit data, to help researchers make design choices that yield valid and reliable results.

The increasing prevalence of Reddit as a source of extant research data necessitates that more research be done on the consequences of sampling methods on this social media platform. Replication of the study design from the second paper should occur across other subreddits on the Reddit platform and on other social media platforms to better understand how context affects the differences found in this study. This study design should also be extended to include other sampling methods like random sampling, purposive sampling, and computer-assisted sampling methods. The finding that sampling from top posts clearly leaves out some perspectives begs the question of whether a random or purposive sample might do a better job of fully representing all perspectives. Further research should also be done that varies the number of units being sampled. This paper only sampled the Top 50 most-upvoted and most commented-on posts and it's possible that sampling the Top 100 may have given a

better representation of the Universe. Finally, this study design could be modified to show if the sampling methods used to generate the dataset for this study affected the results. For example, using a key word search to filter scraped data rather than a flair to see if this decision has any consequences, or taking a sample from more time periods or different time periods to investigate the consequences of those decisions.

The study in the third paper constitutes exploratory research and the categories which took shape need to undergo refinement. Further research should be conducted analyzing posts from across a greater time period to help provide this refinement. The prominence of categories could also be analyzed over time to see if different factors are more critical at different times of the school year or if the prominence of categories were different during COVID-19 or if they will continue to change as we move fully past the pandemic. As of January 1, 2023, the moderators of the r/Teachers subreddit have determined, based on feedback from their members, that discussions about resigning from teaching are no longer appropriate for the subreddit. They have retired the #Resignation flair and are redirecting people to the subreddit r/TeachersInTransition to have these conversations. Therefore, future research should incorporate this subreddit into the data collection process. It is also possible that conversations about teacher resignation are happening in other subreddits, and a more exhaustive search could be conducted to capture more conversations from across Reddit. As the categories in this paper are refined, the Evaluation Framework for Teaching Sustainability can also be refined, and questions aligned to this framework should be created to help local leaders assess areas of concern in their local contexts.

Having a map of topics around teacher resignation allows researchers to understand what further research is possible in this area. The non-continuum buckets provide insight into teachers leaving not by choice, how teachers offer support to others around burnout and resignation, and how

160

teachers discuss the phenomenon of teacher resignation. Researchers interested in teacher conversation around resignation could focus on comments on the kinds of posts found in the "Commentary, Polls & News" bucket. Aggregating the results of the multiple polls questions found within the subreddit could give flashpoints of teacher opinion over time – especially interesting since the same questions seemed to be asked frequently. This bucket also provides an aggregated source of media around teacher resignation, including lesser-known media published across various platforms which might be difficult to come across in other ways. This sample of posts identified several news articles from prominent news sources in addition to a teacher blog, a teacher created YouTube video and a teacher created podcast called The Great Teacher Resignation. Researchers interested in how people make difficult life decisions and the kind of support required to do so could study how people seek and receive support in this anonymous, online space. Finally, each stage of the Continuum of Teaching showed unique characteristics and future research could focus on any of these stages to better understand teacher experience at that stage.

**Contribution to Evaluation, Measurement and Research**

This dissertation makes several critical contributions to the field of Evaluation, Measurement and Research. This paper provides social science researchers with an understanding of what's possible when it comes to leveraging Reddit for social science research which will help social science researchers ask and answer new questions, save time and money on research designs, provide insight into emerging, rapidly evolving, and stigmatizing topics, and leverage the platform for the recruitment of research participants. This paper also begins to build guidelines for reporting on research that leverages Reddit.

This dissertation makes a contribution to measurement by beginning to build the literature around what Sen et al. (2021) call trace selection error and builds researchers' understanding of how

161

sampling choices affect analyses of Reddit data. By comparing analysis results performed on samples of data to the universe from which they were taken, researchers gain insight into the consequences of two different sampling methods. The nuances of the different methods may be more beneficial for some research questions than others and researchers will be able to apply the insights to their own research situation and make informed choices that make the most sense for them. Building this area of the literature helps encourage other such experimentation on the effects of sampling methods to be performed.

This dissertation makes a contribution to research practice by proposing and demonstrating a new hybrid approach to the Big Data Process which blends data science skills for data collection with qualitative social science research methods for data analysis. The results of this exploratory research can then be used to develop questionnaires which can be distributed through the Reddit platform to recruit survey respondents.

**References**

Biester, L., Matton, K., Rajendran, J., Provost, E. M., & Mihalcea, R. (2021). Understanding the impact of COVID-19 on online mental health forums. *ACM Transactions on Management Information Systems*, 12(4), 1-28. https://doi.org/10.1145/3458770

Japec, L., Kreuter, F., Berg, M., Biemer, P., Decker, P., Lampe, C., Lane, J., O'Neil, C., & Usher, A. (2015). Big data in survey research. *Public Opinion Quarterly, 79*(4), 839-880. https://doi.org/10.1093/poq/nfv039

Reveilhac, M., Steinmetz, S., & Morselli, D. (2022). A systematic literature review of how and whether social media data can complement traditional survey data to study public opinion. *Multimedia Tools and Applications, 81*. 10107-10142. https://doi.org/10.1007/s11042-022-12101-0

Schober, M.F., Pasek, J., Guggenheim, L., Lampe, C., Conrad, F.G. (2016). Social media analyses for social measurement. *Public Opinion Quarterly, 80*(1), 180-211. doi: 10.1093/poq/nfv048

Sen, I., Flöck, F., Weller, K., Weib, B., & Wagner, C. (2021). *Public Opinion Quarterly, 85*(S1), 399-422. https://doi.org/10.1093/poq/nfab018

A. Figures Showing Categories for All Buckets Across All Datasets

**Figure A1**

*Categories in the Support Bucket of the MUV Dataset*

**Figure A2**

*Categories in the Support Bucket of the MCO Dataset*

**Figure A3**

*Categories in the Support Bucket of the Universe Dataset*

**Figure A4**

*Categories in the Commentary & News Bucket of the MUV Dataset*

**Figure A5**

*Categories in the Commentary & News Bucket of the MCO Dataset*

**Figure A6**

*Categories in the Commentary & News Bucket of the Universe Dataset*

**Figure A7**

*Categories in the Not by Choice Bucket of the Universe Dataset*

**Figure A8**

*Categories in the Thinking About Quitting Bucket of the MUV Dataset*



**Figure A9**

*Categories in the Thinking About Quitting Bucket of the MCO Dataset*

**Figure A10**

*Categories in the Thinking About Quitting Bucket of the Universe Dataset*



**Figure A11**

*Detail of the Unique and Logistics Categories in the Thinking About Quitting Bucket of the Universe Dataset*

**Figure A12**

*Details of the Reasons Why Teachers Leave Categories in the Thinking About Quitting Bucket of the Universe Dataset*
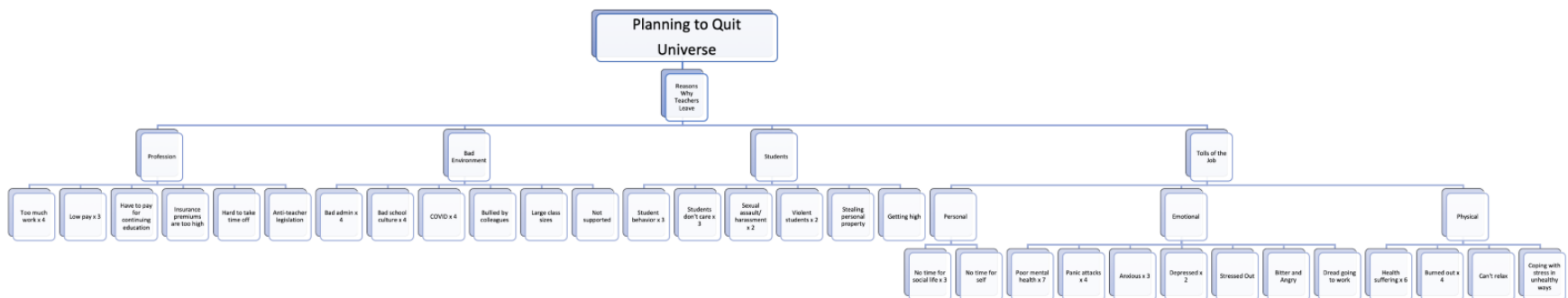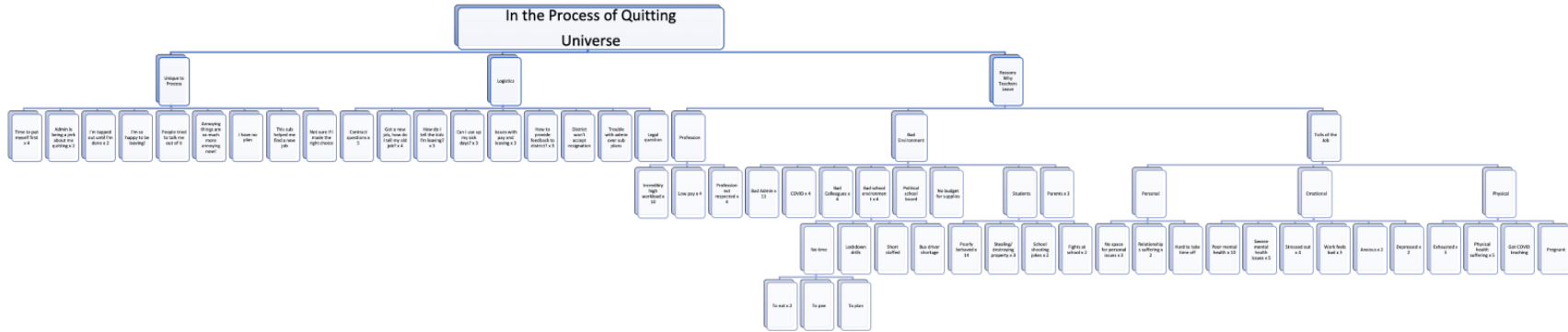
**Figure A13**

*Categories in the I'm Quitting Bucket of the MUV Dataset*

**Figure A14**

*Categories in the I'm Quitting Bucket of the MCO Dataset*



**Figure A15**

*Categories in the Planning to Quit Bucket of the Universe Dataset*

**Figure A16**

*Detail of the Unique and Logistics Categories in the Planning to Quit Bucket of the Universe Dataset*

**Figure A17**

*Details of the Reasons Why Teachers Leave Categories in the Planning to Quit Bucket of the Universe Dataset*

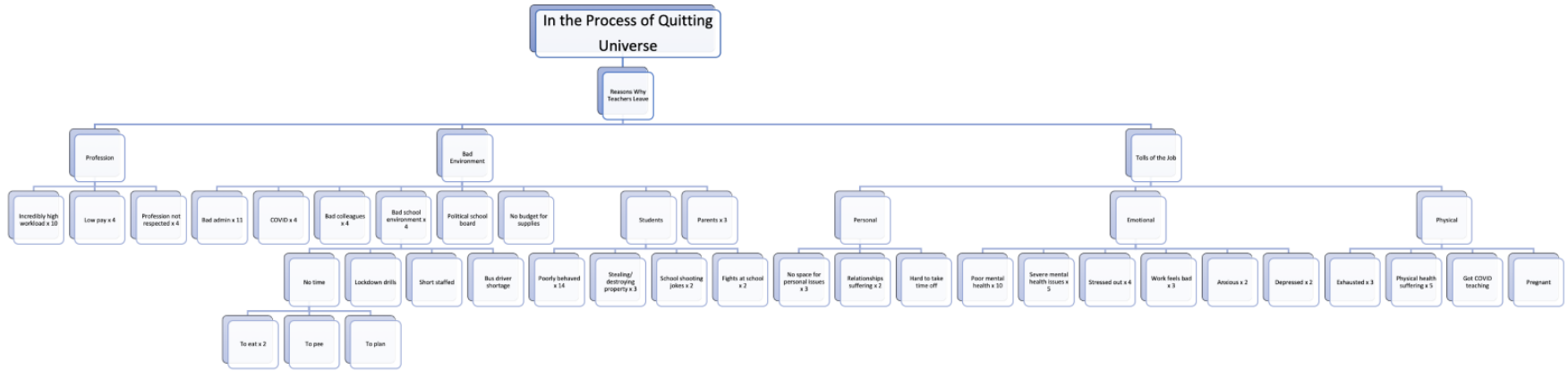**Figure A18**

*Categories in the In the Process of Quitting Bucket of the Universe Dataset*



**Figure A19**

*Detail of the Unique and Logistics Categories in the In the Process of Quitting Bucket of the Universe Dataset*

**Figure A20**

*Details of the Reasons Why Teachers Leave Categories in the In the Process of Quitting Bucket of the Universe Dataset*

**Figure A21**
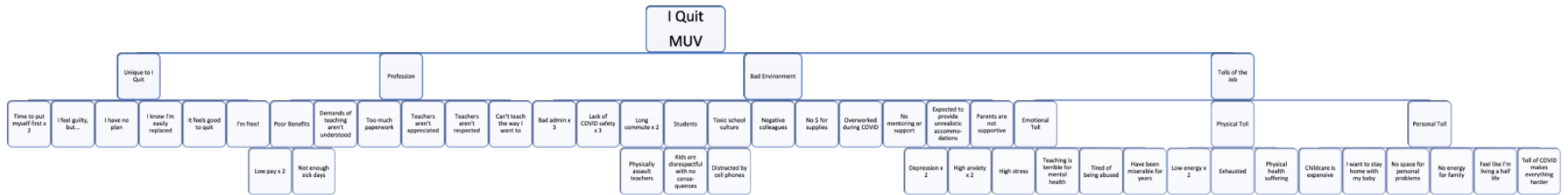
*Categories in the I Quit Bucket of the MUV Dataset*

**Figure A22**

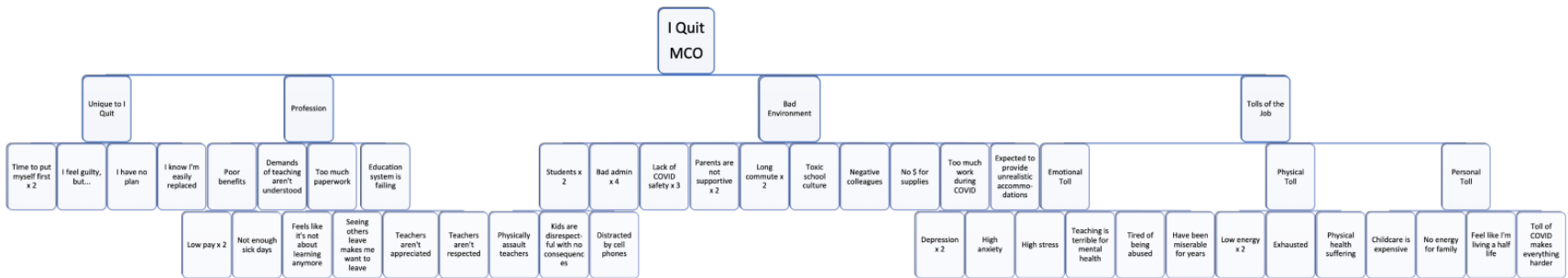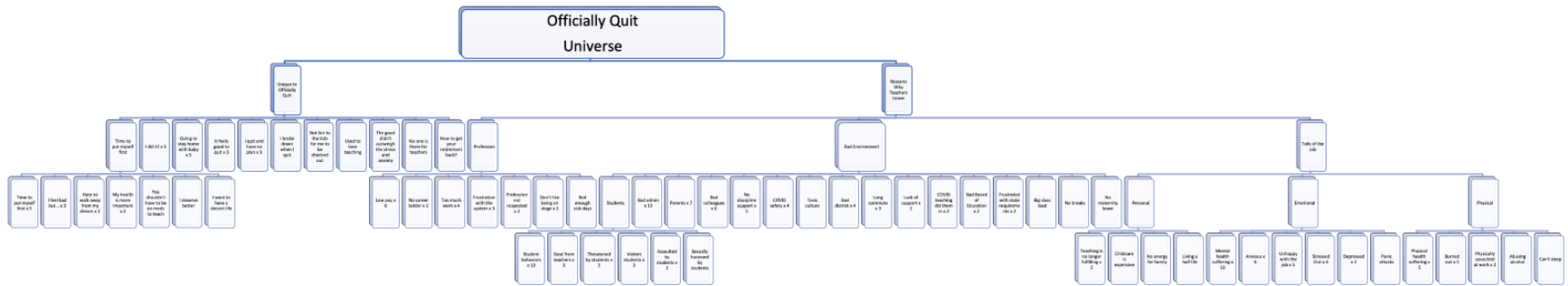*Categories in the I Quit Bucket of the MCO Dataset*

**Figure A23**

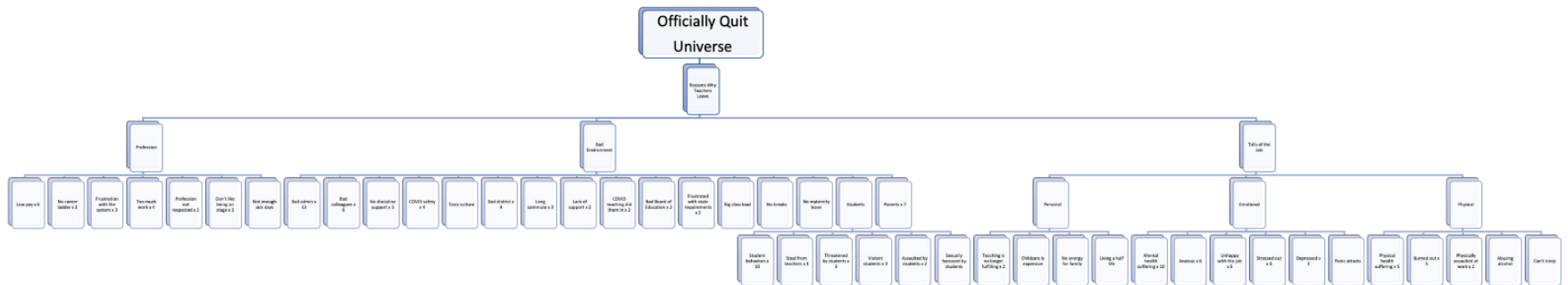*Categories in the Officially Quit Bucket of the Universe Dataset*

**Figure A24**

*Detail of the Unique and Logistics Categories in the Officially Quit Bucket of the Universe Dataset*

**Figure A25**

*Details of the Reasons Why Teachers Leave Categories in the Officially Quit Bucket of the Universe Dataset*

# B. Figures Showing Categories and Codes

**Figure B1**

*Codes in the "Unique to Planning to Quit" and "Logistics" Categories of the "Planning to Quit" Bucket*
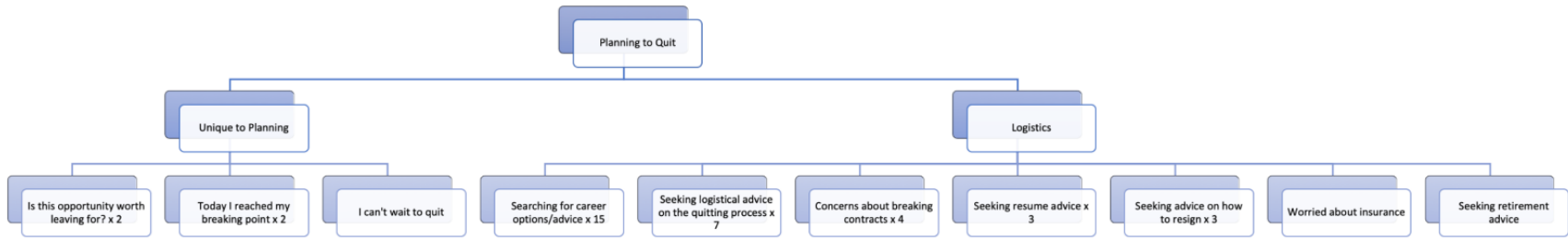
**Figure B2**

*Codes in the "Unique to In the Process of Quitting" and "Logistics" Categories of the "In the Process of Quitting" Bucket*
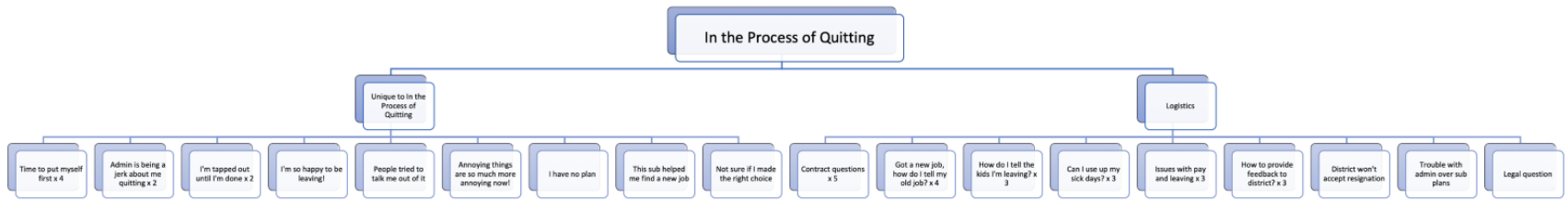
**Figure B3**

*Codes in the "Unique to Officially Quit" Category of the "Officially Quit" Bucket*

**Figure B4**

*Categories and Codes in the "Flipside" Bucket*

**Figure B5**

*Categories and Codes in the "Commentary, Polls & News" Bucket*

```
                              ┌─────────────────────┐
                              │ Commentary, Polls & │
                              │       News          │
                              └─────────────────────┘
        ┌──────────────┬──────────────────┬──────────────────┐
     ┌──────┐      ┌──────────┐       ┌─────────┐      ┌───────────┐
     │ Polls│      │ So Many  │       │ Low Pay │      │  Random   │
     └──────┘      │Teachers  │       └─────────┘      │ Comments  │
                   │are Leaving│                       └───────────┘
                   └──────────┘
```

**Polls:**
- Who's Leaving?
  - Are you leaving teaching? Why? x 2
  - How many are leaving this December?
  - Who's leaving this summer?
- I'm not seeing it - do you? x 2
- What careers have you gone to?
- If you've quit, where are you from?
- What admin moves made you leave?
- How bad will it get?
- Where do licenses get taken away if you leave mid-year?

**So Many Teachers are Leaving:**
- Many teachers are leaving the profession; burnout x 3
- There are 1,000 job openings in my district
- Many teachers are leaving my school; too much stress

**Low Pay:**
- Pay is too low x 2
- Biden never raised our pay like he said he would

**Random Comments:**
- Link to teacher explaining why they left x 2
- Let's all quit and see what they do
- Thank you for your service resigning teachers
- Good teachers get pushed out
- Why are principals surprised when teachers leave?
- The amount of work required to teach is untenable
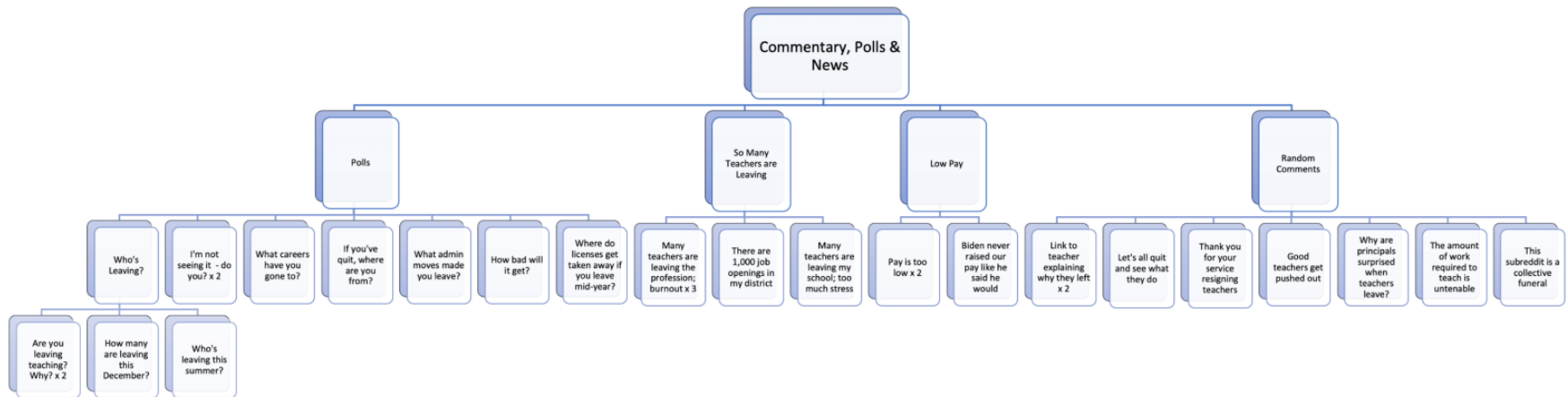- This subreddit is a collective funeral

**Figure B6**

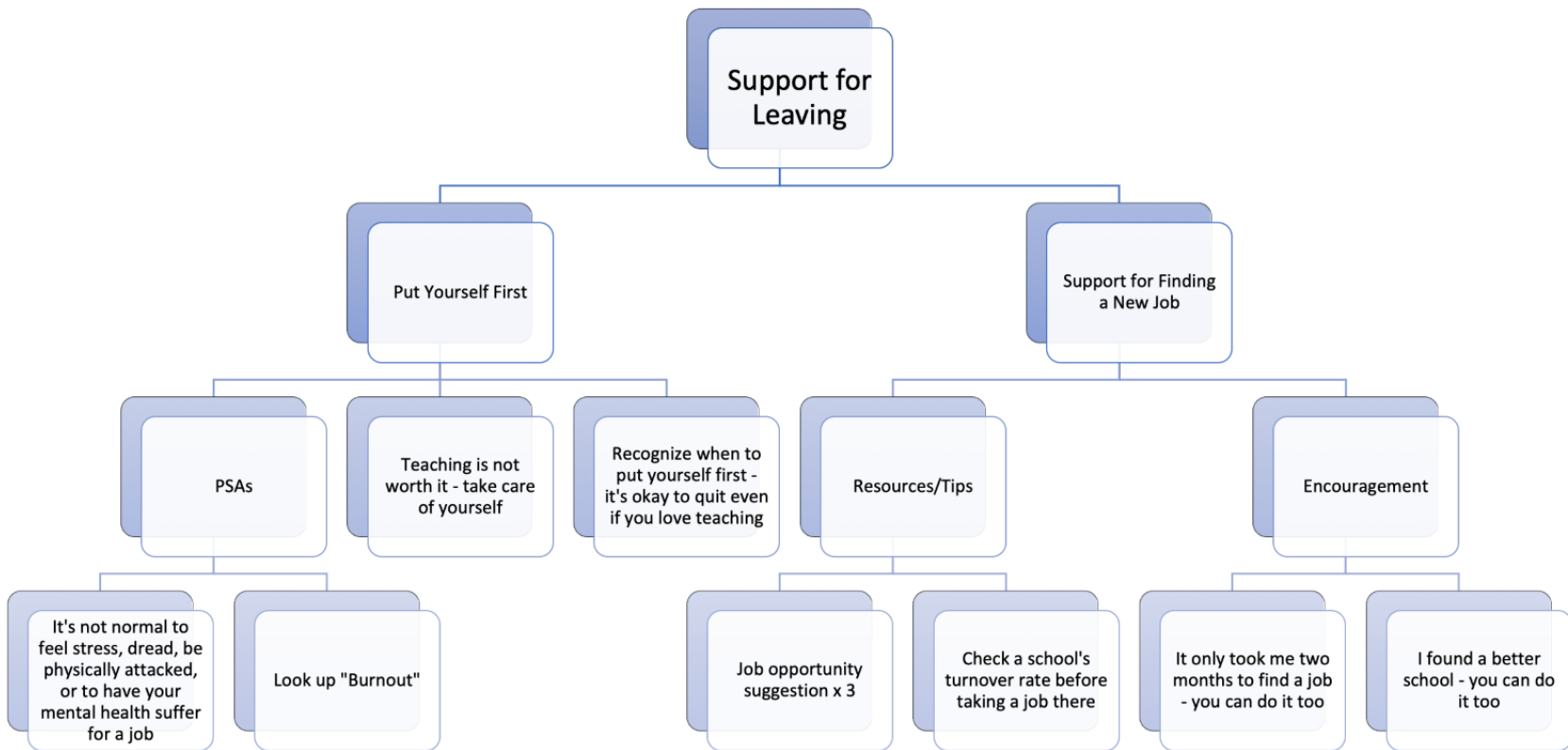*Categories and Codes in the "Support for Leaving" Bucket*

**Figure B7**

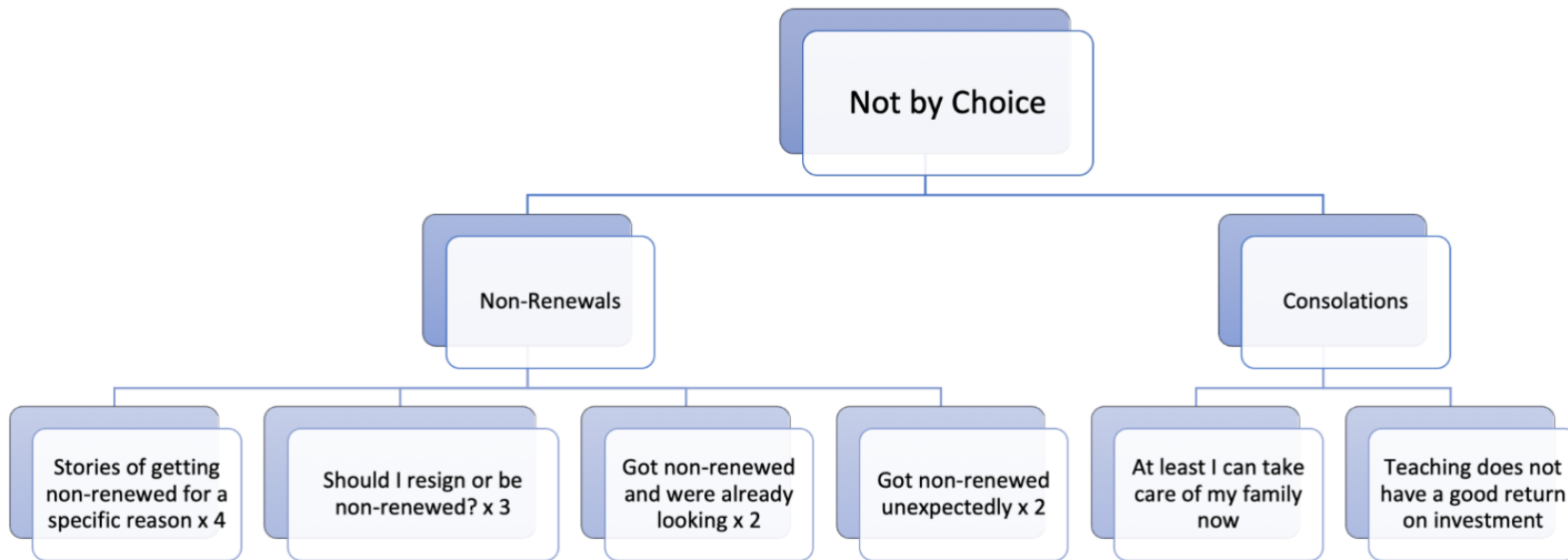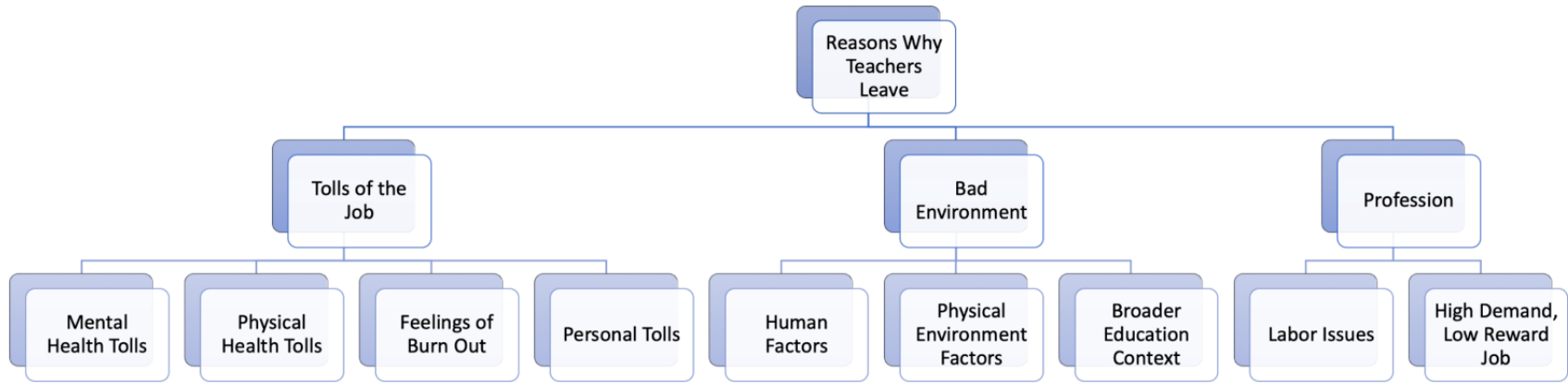*Categories and Codes in the "Not by Choice" Bucket*

**Figure B8**

*Reasons Why Teachers Leave Main Categories and Subcategories*



**Figure B9**

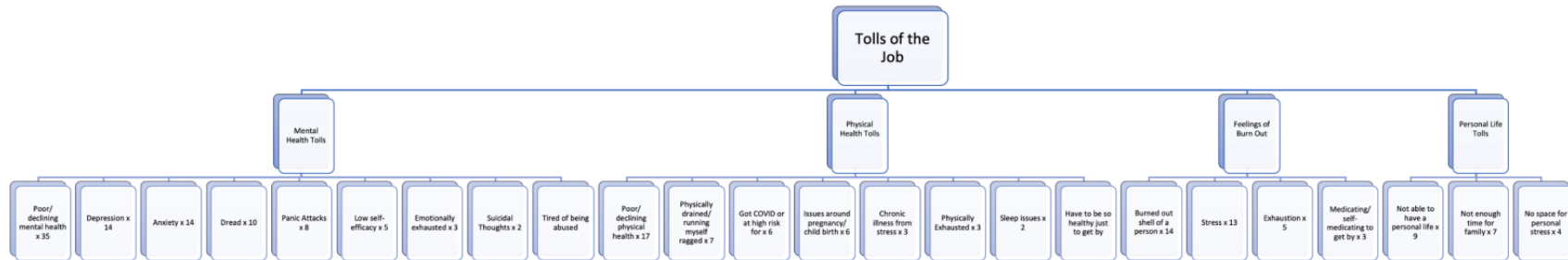*Tolls of the Job Codes and Frequencies*

**Figure B10**
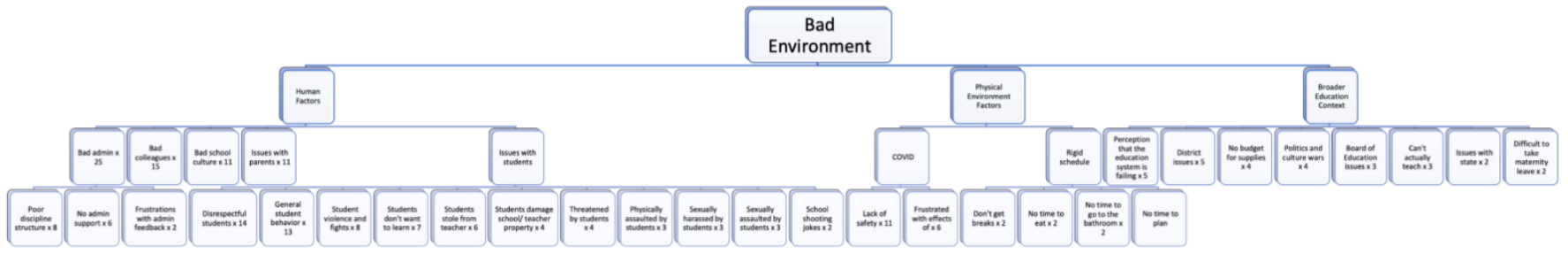
*Bad Environment Codes and Frequencies*

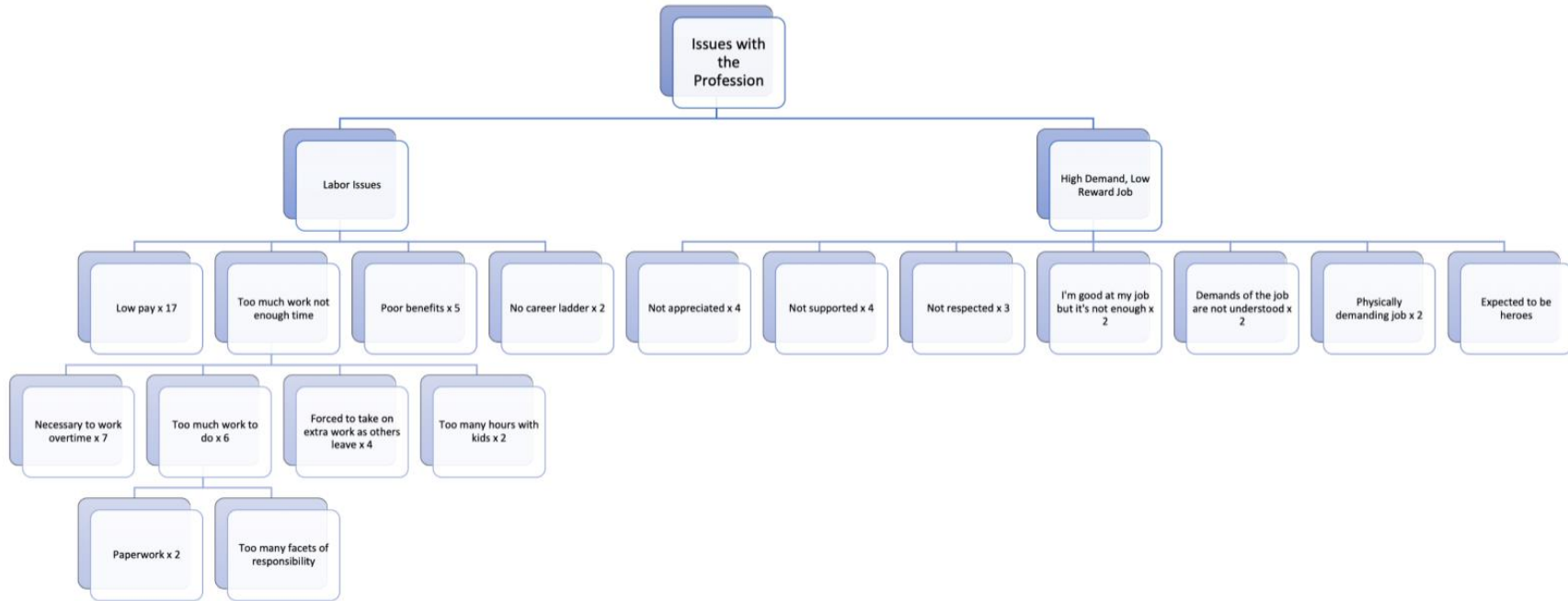**Figure B11**

*Issues with the Profession Codes and Frequencies*

**Figure B12**

*Evaluation Framework for Teacher Sustainability*