

# Interactive Concept Learning for Uncovering Latent Themes in Large Text Collections

Maria Leonor Pacheco<sup>1,2</sup> Tunazzina Islam<sup>3</sup> Lyle Ungar<sup>4</sup> Ming Yin<sup>3</sup> Dan Goldwasser<sup>3</sup>

<sup>1</sup>Microsoft Research <sup>2</sup>University of Colorado Boulder

<sup>3</sup>Purdue University <sup>4</sup>University of Pennsylvania

maria.pacheco@colorado.edu ungar@cis.upenn.edu

{islam32, mingyin, dgoldwas}@purdue.edu

## Abstract

Experts across diverse disciplines are often interested in making sense of large text collections. Traditionally, this challenge is approached either by noisy unsupervised techniques such as topic models, or by following a manual theme discovery process. In this paper, we expand the definition of a theme to account for more than just a word distribution, and include generalized concepts deemed relevant by domain experts. Then, we propose an interactive framework that receives and encodes expert feedback at different levels of abstraction. Our framework strikes a balance between automation and manual coding, allowing experts to maintain control of their study while reducing the manual effort required.

## 1 Introduction

Researchers and practitioners across diverse disciplines are often interested making sense of large text collections. Thematic analysis is one of the most common qualitative research methods used to approach this challenge, and it can be understood as a form of pattern recognition in which the themes (or codes) that emerge from the data become the categories for analysis (Braun and Clarke, 2012; Roberts et al., 2019). In standard practice, researchers bring their own objectives or questions and identify the relevant themes or patterns recognized while analyzing the data, potentially grounding them in a relevant theory or framework. Themes in thematic analysis are broadly defined as “patterned responses or meaning” derived from the data, which inform the research question.

With the explosion of data and the rapid development of automated techniques, disciplines that traditionally relied on qualitative methods for the analysis of textual content are turning to computational methods (Brady, 2019; Hilbert et al., 2019). Topic modeling has long been the go-to NLP technique to identify emerging themes from text collections (Blei et al., 2003; Boyd-Graber et al., 2017;

Baden et al., 2022). Despite its wide adoption, topic modeling does not afford the same flexibility and representation power of qualitative techniques. For this reason, many efforts have been dedicated to understanding the ways in which topic models can be flawed (Mimno et al., 2011), and evaluating their coherence and quality (Stevens et al., 2012; Lau et al., 2014; Röder et al., 2015). More recently, Hoyle et al. (2021) showed that human judgements and accepted metrics of topic quality and coherence do not always agree. Given the noisy landscape surrounding topic modeling, manual qualitative methods are still prevalent across fields for analyzing nuanced and verbally complex data (Rose and Lennerholt, 2017; Lauer et al., 2018; Antons et al., 2020).

Human-in-the-loop topic modeling approaches aim to address these issues by allowing experts to correct and influence the output of topic models. Given that topics in topic models are defined as distributions over words, the feedback received using these approaches is usually limited to identifying representative words and imposing constraints between words (Hu et al., 2011; Lund et al., 2017; Smith et al., 2018). In this paper, we argue that themes emerging from a document collection should not just be defined as a word distribution (similar to a topic model), but as a distribution over generalized concepts that can help us explain them. We build on the definition put forward by Braun and Clarke (2012), where themes are latent patterned meanings that emerge from the data, and supporting concepts serve as a way to explain themes using theoretical frameworks that are deemed relevant by domain experts. For example, emerging themes in a dataset about Covid-19 can be characterized by the strength of their relationship to stances about the covid vaccine and the moral framing of relevant entities (e.g. The theme “*Government distrust*” is strongly correlated to an *anti-vax* stance and frames *Dr. Fauci* as an entity

enabling *cheating*). This representation of a theme aligns more closely with qualitative practices, as experts can introduce their pre-existing knowledge about the domain. Moreover, higher-level abstractions expand the capabilities of experts to correct and influence theme discovery, as it allows them to formulate concepts to generalize from observations to new examples (Rogers and McClelland, 2004), and to deductively draw inferences via conceptual rules and statements (Johnson, 1988).

Following this rationale, we suggest a new computational approach to support and enhance standard qualitative practices for content analysis. We approach both inductive thematic analysis (i.e. identifying the relevant themes that emerge from the data and developing the code-book), and deductive thematic analysis (i.e. identifying the instances where a known theme is observed). To support this process, we allow researchers to shape the space of themes given machine generated candidates. Then, we allow them to provide feedback over machine judgments that map text to themes using relevant conceptual frameworks.

To showcase our approach, we look at the task of characterizing social media discussions around topics of interest to the computational social science community. Namely, we consider two distinct case studies: The covid-19 vaccine debate in the US, and the immigration debate in the US, the UK and Europe. For each case, the qualitative researchers use different theories to ground theme discovery, each associated with a different set of concepts. For the covid-19 vaccine debate, the theme discovery process is grounded using vaccination stances and morality frames (Roy et al., 2021; Pacheco et al., 2022). For the immigration debates, the theme discovery is grounded using three framing typologies: narrative frames (Iyengar, 1991), policy frames (Card et al., 2015) and immigration frames (Benson, 2013; Hovden and Mjelde, 2019). All of these choices build on previous work and were validated by the qualitative researchers. From a machine learning perspective, these two case studies could be regarded as completely different tasks and have been approached independently in previous work. The reason for this is the data, the context, and the target labels (both the emerging themes and the supporting concepts) are different for each scenario.

To aid experts in theme discovery, we propose an iterative two-stage machine-in-the-loop framework. In the first stage, we provide experts with

an automated partition of the data, ranked example instances, and visualizations of the concept distribution. Then, we have a group of experts work together to explore the partitions, code emerging patterns and identify coherent themes. Once themes are identified, we have the experts select representative examples, write down additional examples and explanatory phrases, and explain themes using the set of available concepts. In the second stage, we incorporate the expert feedback using a neuro-symbolic mapping procedure. The *symbolic* part allows us to explicitly model the dependencies between concepts and the emerging themes using weighted logical formulae (e.g.  $w : \text{policy\_frame}(\text{economic}) \Rightarrow \text{theme}(\text{economic\_migrants})$ ). These rules can be interpreted as soft constraints whose weights are learned from the feedback provided by the experts. The *neural* part allows us to maintain a distributed representation of the data points and themes, which facilitates the live exploration of the data based on distances and similarities, and provides a feature representation for learning the rule weights. After the mapping stage concludes, some instances will be assigned to the identified themes, and the remaining instances will be re-partitioned for a consecutive discovery stage.

We conducted extensive evaluations of the different components, design choices, and stages in our methodology. We showed that our framework allows experts to uncover a set of themes that cover a large portion of the data, and that the resulting mapping from tweets to themes is fairly accurate with respect to human judgements. While we focused on polarized discussions, our framework generalizes to any content analysis study where the space of relevant themes is not known in advance.

## 2 Related Work

This paper suggests a novel approach for identifying themes emerging from text collections. The notion of a theme presented in this work is strongly related to topic models (Blei et al., 2003). However, unlike latent topics that are defined as word distributions, our goal is to provide a richer representation that more strongly resembles qualitative practices by connecting the themes to general concepts that help explain them. For example, when identifying themes emerging from polarized discussions in social media, we look at conceptual frameworks such as moral foundations theory (Haidt and Graham,

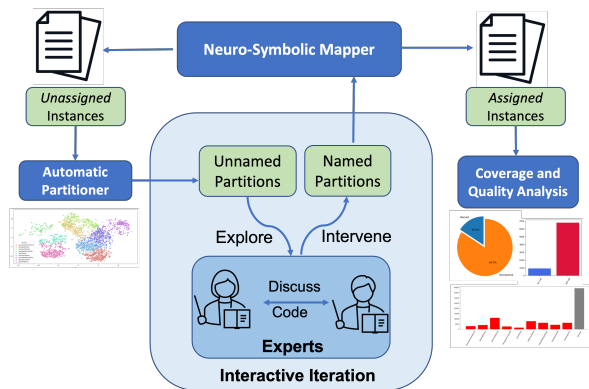


Figure 1: Framework Overview

2007; Amin et al., 2017; Chan, 2021) and framing theory (Entman, 1993; Chong and Druckman, 2007; Morstatter et al., 2018).

Our work is conceptually similar to recent contributions that characterize themes and issue-specific frames in data, either by manually developing a code-book and annotating the data according to it (Boydston et al., 2014; Mendelsohn et al., 2021), or by using data-driven methods (Demszky et al., 2019; Roy and Goldwasser, 2021). Unlike these approaches, our work relies on interleaved human-machine interaction rounds, in which humans can identify and explain themes from a set of candidates suggested by the model, as well as diagnose and adapt the model’s ability to recognize these themes in documents. This work is part of a growing trend in NLP that studies how human-machine collaboration can help improve automated language analysis (Wang et al., 2021). In that space, two lines of works are most similar to ours. Interactive topic models (Hu et al., 2011; Lund et al., 2017; Smith et al., 2018) allow humans to adapt the identified topics, but the feedback is usually limited to lexical information. Open Framing (Bhatia et al., 2021) allows humans to identify and name frames based on the output of topic models, but lacks our model’s ability for sustained interactions that help shape the theme space, as well as the explanatory power of our neuro-symbolic representation.

### 3 The Framework

We propose an iterative two-stage framework that combines ML/NLP techniques, interactive interfaces and qualitative methods to assist experts in characterizing large textual collections. We define large textual collections as repositories of textual instances (e.g. tweets, posts, documents), where

each instance is potentially associated with a set of annotated or predicted concepts.

In the first stage, our framework automatically proposes an initial partition of the data, such that instances that are thematically similar are clustered together. We provide experts with an interactive interface equipped with a set of *exploratory operations* that allows them to evaluate the quality of the discovered partitions, as well as to further explore and partition the space by inspecting individual examples, finding similar instances, and using open text queries. As experts interact with the data through the interface, they following an inductive thematic analysis approach to identify and code the patterns that emerge within the partitions (Braun and Clarke, 2012). Next, they group the identified patterns into general themes, and instantiate them using the interface. Although intuitively we could expect a single partition to result in a single theme, note that this is not enforced. Experts maintain full freedom as to how many themes they instantiate, if any. Once a theme is created, experts are provided with a set of *intervention operations* to explain the themes using natural language, select good example instances, write down additional examples, and input or correct supporting concepts to characterize the theme assignments. The full set of operations are listed in Tab. 1 and demonstrated in App. A.1.

In the second stage, our framework finds a mapping between the full set of instances and the themes instantiated by the experts. We use the information contributed by the experts in the form of examples and concepts, and learn to map instances to themes using our neuro-symbolic procedure. We allow instances to remain unassigned if there is not a good enough match, and in this case, a consecutive portioning step is done. We refer to instances that are mapped to themes as “named partitions” and unassigned proposed partitions as “unnamed partitions”. Once instances are assigned to themes, experts have access to a comprehensive visual analysis of the state of the system. The main goal of this analysis is to appreciate the trade-off between coverage (how many instances we can account for with the discovered themes) and quality (how good we are at mapping instances to themes). An illustration of the framework can be observed in Fig. 1. Additional details about the coverage and quality analysis are presented in the experimental section.

Below, we discuss the representation of themes and instances, the protocol followed for interaction,

and the mapping and re-partitioning procedures. Our visual interface and our analysis code have been made available to the community<sup>1</sup>.

**Representing Themes and Instances** We represent example instances and explanatory phrases using their S-BERT embedding (Reimers and Gurevych, 2019). To measure the closeness between an instance and a theme, we compute the cosine similarity between the instance and all of the explanatory phrases and examples for the theme, and take the maximum similarity score among them. Our framework is agnostic of the representation used. The underlying embedding objective and the scoring function can easily be replaced.

Operations	Description
Finding Partitions	Experts can find partitions in the space of unassigned instances. We currently support the K-means (Jin and Han, 2010) and Hierarchical Density-Based Clustering (McInnes et al., 2017) algorithms.
Text-based Queries	Experts can type any query in natural language and find instances that are close to the query in the embedding space.
Finding Similar Instances	Experts have the ability to select each instance and find other examples that are close in the embedding space.
Listing Themes and Instances	Experts can browse the current list of themes and their mapped instances. Instances are ranked in order of “goodness”, corresponding to the similarity in the embedding space to the theme representation. They can be listed from closest to most distant, or from most distant to closest.
Visualizing Local Explanations	Experts can visualize aggregated statistics and explanations for each of the themes. To obtain these explanations, we aggregate all instances that have been identified as being associated with a theme. Explanations include wordclouds, frequent entities and their sentiments, and graphs of concept distributions.
Visualizing Global Explanations	Experts can visualize aggregated statistics and explanations for the global state of the system. To do this, we aggregate all instances in the database. Explanations include theme distribution, coverage statistics, and t-sne plots (van der Maaten and Hinton, 2008).

(a) Exploratory Operations

Operations	Description
Adding, Editing and Removing Themes	Experts can create, edit, and remove themes. The only requirement for creating a new theme is to give it a unique name. Similarly, themes can be edited or removed at any point. If any instances are assigned to a theme being removed, they will be moved to the space of unassigned instances.
Adding and Removing Examples	Experts can assign “good” and “bad” examples to existing themes. Good examples are instances that characterize the named theme. Bad examples are instances that could have similar wording to a good example, but that have different meaning. Experts can add examples in two ways: they can mark mapped instances as “good” or “bad”, or they can directly contribute example phrases.
Adding or Correcting Concepts	We allow users to upload additional observed or predicted concepts for each textual instance. For instances and phrases added as “good” and “bad” examples, we allow users to add or edit the values of these concepts. The intuition behind this operation is to collect additional information for learning to map instances to themes.

(b) Intervention Operations

Table 1: Interactive Operations

**Interaction Protocol** We follow a simple protocol where three human coders work together using

<sup>1</sup><https://gitlab.com/mlpacheco/machine-in-the-loop-concepts>

the operations described above to discover themes in large textual corpora. In addition to the three coders, each interactive session is guided by one of the authors of the paper, who makes sure the coders are adhering to the process outlined here.

To initialize the system, the coders will start by using the partitioning operation to find ten initial partitions of roughly the same size. During the first session, the coders will inspect the partitions one by one by looking at the examples closest to the centroid. This will be followed by a discussion phase, in which the coders follow an inductive thematic analysis approach to identify repeating patterns and write them down. If one or more cohesive patterns are identified, the experts will create a new theme, name it, and mark a set of good example instances that help in characterizing the named theme. When a pattern is not obvious, coders will explore similar instances to the different statements found. Whenever the similarity search results in a new pattern, the coders will create a new theme, name it, and mark a set of good example instances that helped in characterizing the named theme.

Next, the coders will look at the local theme explanations and have the option to enhance each theme with additional phrases. Note that each theme already contains a small set of representative instances, which are marked as “good” in the previous step. In addition to contributing “good” example phrases, coders will have the option to contribute some “bad” example phrases to push the representation of the theme away from statements that have high lexical overlap with the good examples, but different meaning. Finally, coders will examine each exemplary instance and phrase for the set of symbolic concepts (e.g. stance, moral frames). In cases where the judgement is perceived as wrong, the coders will be allowed to correct it. In this paper, we assume that the textual corpora include a set of relevant concepts for each instance. In future work, we would like to explore the option of letting coders define concepts on the fly.

**Mapping and Re-partitioning** Each interactive session will be followed by a mapping and re-partitioning stage. First, we will perform the mapping step, in which we assign instances to the themes discovered during interaction. We do not assume that experts will have discovered the full space of latent themes. For this reason, we do not try to assign a theme to each and every instance. We expect that the set of themes introduced by the hu-

man experts at each round of interaction will cover a subset of the total instances available. Following this step, we will re-partition all the unassigned instances for a subsequent round of interaction.

We use DRaiL (Pacheco and Goldwasser, 2021), a neuro-symbolic modeling framework to design a mapping procedure. Our main goal is to condition new theme assignments not only on the embedding distance between instances and good/bad examples, but also leverage the additional judgements provided by experts using the ‘‘Adding or Correcting Concepts’’ procedure. For example, when analyzing the corpus about the Covid-19 vaccine, experts could point out that 80% of the good examples for theme *Natural Immunity is Effective* have a clear *anti-vaccine* stance. We could use this information to introduce inductive bias into our mapping procedure, and potentially capture cases where the embedding distance does not provide enough information. DRaiL uses weighted first-order logic rules to express decisions and dependencies between different decisions. We introduce the following rules:

$$\begin{aligned}
 t_0 - t_n &: \text{Inst}(i) \Rightarrow \text{Theme}(i, t) \\
 a_0 - a_m &: \text{Inst}(i) \Rightarrow \text{Concept}(i, c) \\
 c_0 - c_{n*m} &: \text{Inst}(i) \wedge \text{Concept}(i, c) \Rightarrow \text{Theme}(i, t) \\
 c'_0 - c'_{n*n} &: \text{Inst}(i) \wedge \text{Theme}(i, t) \wedge (t \neq t') \\
 &\Rightarrow \neg \text{Theme}(i, t)
 \end{aligned}$$

The first set of rules  $t_0 - t_n$  and  $a_0 - a_m$  map instances to themes and concepts respectively. We create one template for each theme  $t$  and concept  $c$ , and they correspond to binary decisions (e.g. whether instance  $i$  mentions theme  $t$ ). Then, we introduce two sets of soft constraints:  $c_0 - c_{n*m}$  encode the dependencies between each concept and theme assignment (e.g. likelihood of theme *Natural Immunity is Effective* given that instance has concept *anti-vax*). Then,  $c'_0 - c'_{n*n}$  discourages an instance from having more than one theme assignment. For each rule, we will learn a weight that captures the strength of that rule (i.e. its likelihood of being active for a given input). Then, a combinatorial inference procedure will be run to find the most likely global assignment. Each entity and relation in DRaiL is tied to a neural architecture that is used to learn its weights. In this paper, we use a BERT encoder (Devlin et al., 2019) for all rules. To generate data for learning the DRaiL model, we take the  $K = 100$  closest instances for each good/bad example provided by the experts. Good examples will serve as positive training data.

For negative training data, we take the contributed bad examples, as well as good examples for other themes and concepts. Once the weights are learned, we run the inference procedure over the full corpus.

## 4 Case Studies

We explore two case studies involving discussions on social media: (1) The Covid-19 vaccine discourse in the US, and (2) The immigration discourse in the US, the UK and the EU. For the Covid-19 case, we build on the corpus of 85K tweets released by Pacheco et al. (2022). All tweets in this corpus were posted by users located in the US, are uniformly distributed between Jan. and Oct. 2021, and contain predictions for vaccination stance (e.g. pro-vax, anti-vax) and morality frames (e.g. fairness/cheating and their actor/targets.) (Haidt and Graham, 2007; Roy et al., 2021). For the immigration case, we build on the corpus of 2.66M tweets released by Mendelsohn et al. (2021). All tweets in this corpus were posted by users located in the US, the UK and the EU, written between 2018 and 2019, and contain predictions for three different framing typologies: narrative frames (e.g. episodic, thematic) (Iyengar, 1991), generic policy frames (e.g. economic, security and defense, etc.) (Card et al., 2015), and immigration-specific frames (e.g. victim of war, victim of discrimination, etc.) (Ben-son, 2013; Hovden and Mjelde, 2019). Additional details about the datasets and framing typologies can be found the original publications.

Our main goal is to evaluate whether experts can leverage our framework to identify prominent themes in the corpora introduced above. We recruited a group of six experts in Computational Social Science, four male and two female, within the ages of 25 and 45. The group of experts included advanced graduate students, postdoctoral researchers and faculty. Our studies are IRB approved, and we followed their protocols. For each corpus, we performed two consecutive sessions with three experts following the protocol outlined in Sec. 3. To evaluate consistency, we did an additional two sessions with a different group of experts for the Covid-19 dataset. Each session lasted a total of one hour. In App. A.2, A.3 and A.4, we include large tables enumerating the resulting themes, and describing in detail all of the patterns identified and coded by the experts at each step of the process.

**Coverage vs. Mapping Quality** We evaluated the trade-off between coverage (how many tweets

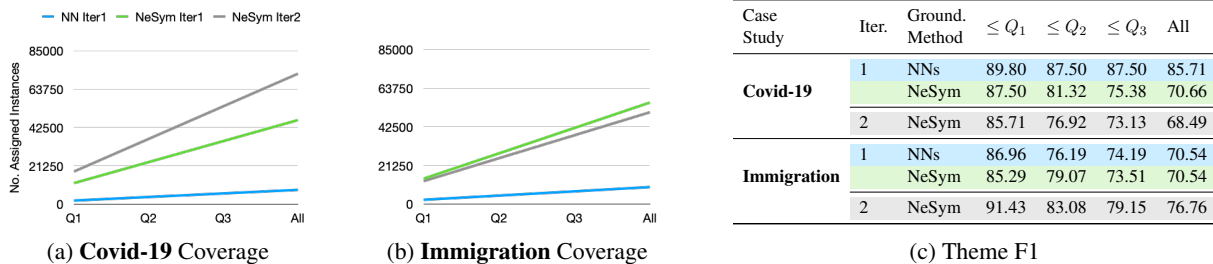


Figure 2: Theme Assignments Where Distance to Theme Centroid  $\leq$  Quartile

we can account for with the discovered themes) and mapping quality (how good we are at mapping tweets to themes). Results are outlined in Fig. 2. To do this evaluation, we sub-sampled a set of 200 mapped tweets for each scenario, uniformly distributed across themes and their proximity to the theme embedding, and validated their assignments manually. The logic behind sampling across different proximities is that we expect mapping performance to degrade the more semantically different the tweets are to the “good” examples and phrases provided by the experts. To achieve this, we look at evaluation metrics at different thresholds using the quartiles with respect to the proximity/similarity distribution. Results for  $Q_1$  correspond to the 25% most similar instances. For  $Q_2$  to the 50% most similar instances, and for  $Q_3$  to the 75% most similar instances. Note that these are continuous ranges and the quartiles serve as thresholds.

To evaluate the impact of our neuro-symbolic mapping procedure (NeSym), we compared it against a nearest neighbors (NNs) approach that does not leverage conceptual frameworks and looks only at the language embedding of the tweets and theme examples and explanatory phrases. For the first iteration of Covid-19, we find that the approximate performance of the NeSym mapping at  $Q_1$  is better (+2 points) than the approximate full mapping for NNs, while increasing coverage x1.5. For immigration, we have an even more drastic result, having an approximate 15 point increase at a similar coverage gain. In both cases, experts were able to increase the number of themes in subsequent iterations<sup>2</sup>. While the coverage increased in the second iteration for Covid, it decreased slightly for Immigration. For Covid, most of the coverage increase can be attributed to a single theme (*Vax Efforts Progression*), which accounts for 20% of the mapped data. In the case of Covid, this large jump

<sup>2</sup>Due to effort required and cost, we only do a subsequent interactive session over the NeSym mapping.

in coverage is accompanied by a slight decrease in mapping performance. In the case of Immigration, we have the opposite effect: as the coverage decreases the performance improves, suggesting that the mapping gets stricter. This confirms the expected trade-off between coverage and quality. Depending on the needs of the final applications, experts could adjust their confidence thresholds.

To perform a fine-grained error analysis, we looked at the errors made by the model using manual validation. In Fig. 3 we show the confusion matrix for the Covid case. We find that the performance varies a lot, with some themes being more accurate than others. In some cases, we are good at capturing the general meaning of the theme but fail at grasping the stance similarities. For example, *Anti Vax Spread Missinfo* gets confused with *Pro Vax Lie*, where the difference is on who is doing the lying. In other cases, we find that themes that are close in meaning have some overlap (e.g. *Alt Treatments* with *Vax Doesn't Work*). We also find that unambiguous, neutral themes like *Vax Appointments*, *Got The Vax* and *Vax Efforts Progression* have the highest performance. Lastly, we observe that for some errors, none of the existing themes are appropriate (Last row: *Other*), suggesting that there are still undiscovered themes. Upon closer inspection, we found that the majority of these tweets are among the most distant from the theme embedding. The full distribution of *Other* per interval can be observed in App. A.6. We include the confusion matrix for immigration in App. A.6.

Given our hypothesis that themes can be characterized by the strength of their relationship to high-level concepts, we consider mappings to be better if they are more cohesive. In the Covid case, we expect themes to have strong relationships to vaccination stance and morality frames. In the Immigration case, we expect themes to have strong relationships to the framing typologies. To measure this, we define a theme purity metric for each

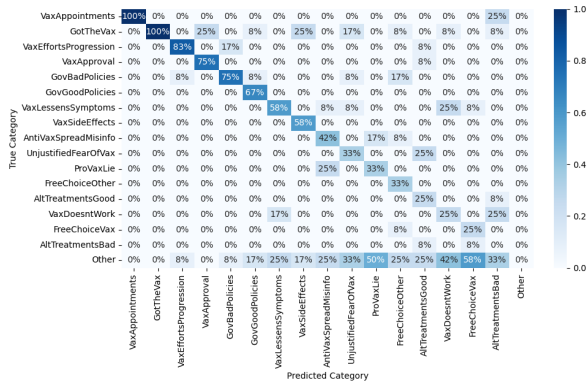


Figure 3: **Confusion matrix for Covid after second iteration.** Values are normalized over the predicted themes (cols), and sorted from best to worst.

Iter.	Ground Method	Covid Vaccine			Immigration		
		# Thm	Cover	Purity	# Thm	Cover	Purity
Baselines	LDA (Var. Bayes)	9	39.8	63.72	13	26.8	57.14
	LDA (Gibbs)	<b>79.8</b>	63.90		55.9	54.86	
	NNs		9.3	68.81	11.1	58.44	
1	NeSym		54.3	<b>69.97</b>	<b>65.8</b>	<b>61.72</b>	
Baselines	LDA (Var. Bayes)	16	26.1	65.02	19	18.3	57.94
	LDA (Gibbs)		73.1	65.14	46.8	<b>59.25</b>	
	NeSym		<b>84.3</b>	<b>65.50</b>	<b>59.6</b>	59.19	

Table 2: **Dataset Coverage and Avg. Concept Purity.** For LDA, we assigned a tweet to its most probable topic if the probability was  $\geq 0.5$ .

concept. For example, for stance this is defined as:  $Purity_{stance} = \frac{1}{N} \sum_{t \in Themes} \max_{s \in Stance} |t \cap s|$ . Namely, we take each theme cluster and count the number of data points from the most common stance value in said cluster (e.g. the number of data points that are *anti-vax*). Then, we take the sum over all theme clusters and divide it by the number of data points. We do this for every concept, and average them to obtain the final averaged concept purity. In Tab. 2 we show the average concept purity for our mappings at each iteration in the interaction. We can see that the NeSym procedure results in higher purity with respect to the NNs procedure, even when significantly increasing coverage. This is unsurprising, as our method is designed to take advantage of the relationship between themes and concepts. Additionally, we include topic modeling baselines that do not involve any interaction, and find that interactive themes generally result in higher purity partitions than topics obtained using LDA. Details about the steps taken to obtain LDA topics can be found in App. A.5.

**Effects of Consecutive Iterations** In Fig. 2 we observed different behaviors in subsequent iterations with respect to coverage and performance.

To further inspect this phenomenon, we looked at the tweets that shifted predictions between the first and second iterations. Fig. 4 shows this analysis for Immigration. Here, we find that a considerable number of the tweets that were assigned to a theme in the first iteration were unmatched (i.e. moved to the *Unknown*) in the second iteration. This behavior explains the decrease in coverage. Upon closer inspection, we found that the majority of these unmatched tweets corresponded to assignments that were in the last and second to last intervals with respect to their similarity to the theme embedding. We also observed a non-trivial movement from the *Unknown* to the new themes (shown in red), as well as some shifts between old themes and new themes that seem reasonable. For example, 1.2% of the total tweets moved from *Role of Western Countries* to *Country of Immigrants*, 1% moved from *Academic Discussions* to *Activism*, and close to 3% of tweets moved from *Trump Policy* and *UK Policy* to *Criticize Anti Immigrant Rhetoric*. This behavior, coupled with the increase in performance observed, suggests that as new themes are added, tweets move to a closer fit. In App. A.7 we include the shift matrix for Covid, as well as the distribution of the unmatched tweets with respect to their semantic similarity to the theme embedding. For Covid, we observe that the increase in coverage is mostly attributed to the addition of the *Vax Efforts Progression* theme, which encompasses all mentions to vaccine development and roll-out. Otherwise, a similar shifting behavior can be appreciated.

### Consistency between Different Expert Groups

To study the subjectivity of experts and its impact on the resulting themes, we performed two parallel studies on the Covid corpus. For each study, a different group of experts performed two rounds of interaction following the protocol outlined on Sec. 3. The side-by-side comparisons of the two studies can be observed in Tab. 3. We find that the second group of experts is able to obtain higher coverage and higher concept purity with a slightly reduced number of themes. To further inspect this phenomenon, as well as the similarities and differences between the two sets of themes, we plot the overlap coefficients between the theme-to-tweet mappings in Fig. 5. We use the Szymkiewicz–Simpson coefficient, which measures the overlap between two finite sets and is defined as:  $overlap(X, Y) = \frac{|X \cap Y|}{\min(|X|, |Y|)}$ .

In cases where we observe high overlap between

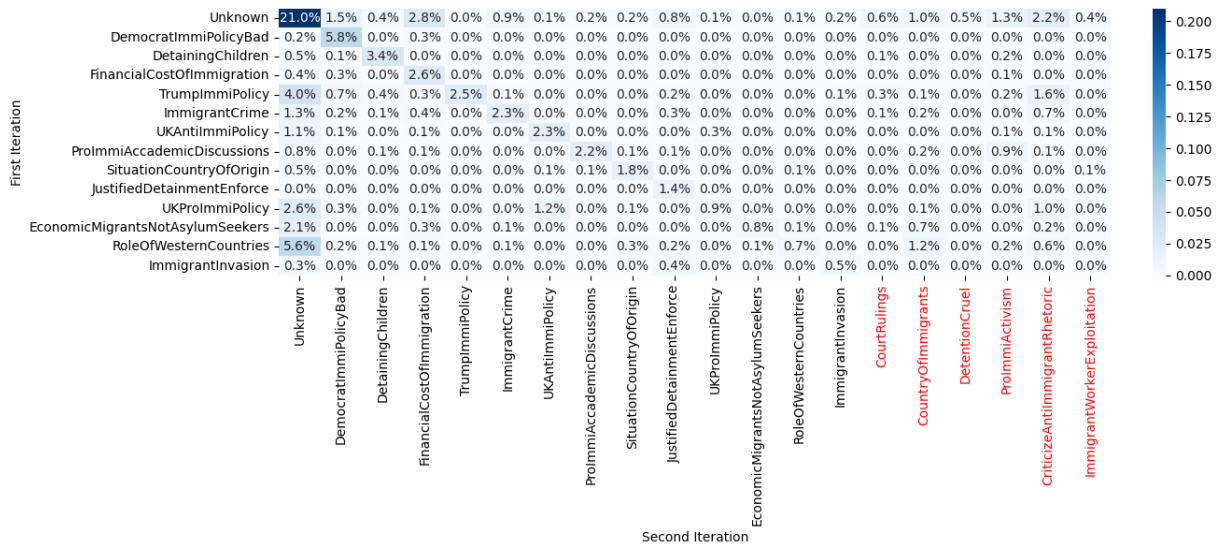


Figure 4: **Shifting predictions for Immigration.** Themes added during second iteration are shown in red, and values are normalized over the full population.

Iter.	Metric	Group 1	Group 2
1	Num Themes	9	8
	Coverage	54.30	61.80
	Stance Purity	83.18	87.43
	Moral Frame Purity	56.75	65.52
	2	Num Themes	16
Coverage		84.30	85.90
Stance Purity		80.12	84.31
Moral Frame Purity		50.88	52.17

Table 3: Two Different Groups of Experts on Covid

the two groups, we find that there is essentially a word-for-word match between the two discovered themes. For example, *Vax Lessens Symptoms*, which was surprisingly named the same by the two groups, as well as *Vax Availability* vs. *Vax Appointments*, *Got The Vax* vs. *I Got My Vax*, and *Vax Side Effects* vs. *Post Vax Symptoms*. In other cases, we find that different groups came up with themes that have some conceptual (and literal) overlap, but that span different sub-segments of the data. For example, we see that the theme *Reasons the US Lags On Vax* defined by the second group, has overlap with different related themes in the first group, such as: *Gov. Bad Policies*, *Vax Efforts Progression*, and *Unjustified Fear of Vax*. Similarly, while the second group defined a single theme *Vax Personal Choice*, the first group attempted to break down references to personal choices between those directly related to taking the vaccine (*Free Choice Vax*), and those that use the vaccine as analogies for other topics, like abortion (*Free Choice Other*). While some themes are clearly present in the data

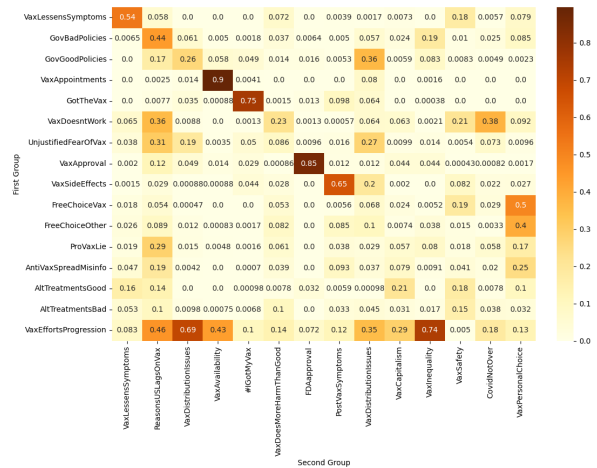


Figure 5: Theme Overlap Coefficient Heatmap between Different Groups of Experts

and identified by the two groups, we see that subjective decisions can influence the results. The first group was inclined to finer grained themes (with the exception of *Vax Efforts Progression*), while the second group seemed to prefer more general themes. In future work, we would like to study how the variation observed with our approach compares to the variation encountered when experts follow fully manual procedures, as well the impact of the crowd vs. experts working alone.

**Abstract Themes vs. Word-level Topics** To get more insight into the differences between topics based on word distributions and our themes, we looked at the overlap coefficients between topics obtained using LDA and our themes. Fig. 6 shows



the coefficients for Immigration. While some overlap exists, the coefficients are never too high (a max. of 0.35). One interesting finding is that most of our themes span multiple related topics. For example, we find that *Trump Policy* has similar overlap with *undocumented\_ice\_workers\_trump*, *migrants\_migrant\_trump\_border*, and *children\_parent\_kids\_trump*. While all of these topics discuss Trump policies, they make reference to different aspects: workers, the border and families. This supports our hypothesis that our themes are more abstract in nature, and that they capture conceptual similarities beyond word distributions. Overlap coefficients for Gibbs sampling, Covid, and subsequent iterations can be seen in App. A.8.

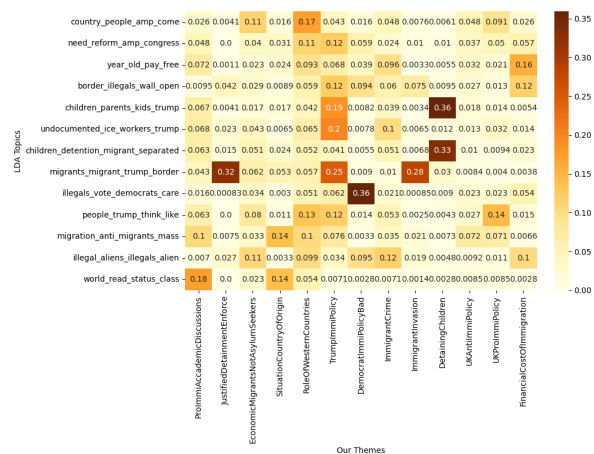


Figure 6: Overlap Coefficients between LDA Var. Bayes and our Themes (First Iter. Immigration).

## 5 Limitations

The study presented in this paper has three main limitations. (1) While the design of the framework does not prohibit the utilization of longer textual forms, the two case studies presented deal with short texts. When dealing with longer text forms, we need to consider the cognitive load of having experts look at groups of instances. In our ongoing work, we employ strategies such as summarization, highlighting and other visualization techniques to deal with these challenges. (2) In the studies presented, qualitative researchers worked in groups to identify themes. Our goal in comparing two independent groups of researchers was to evaluate the degree of subjectivity by observing if the themes identified by the two groups would diverge. This setup might not always be realistic, as a lot of times qualitative researchers work independently or asynchronously. In the future, we will explore the effect

of the crowd in minimizing subjectivity, as well as the role that the computational tools play in more challenging settings. (3) Finally, we did not include a comprehensive user study to gather input from the experts about their experience with our framework. We consider this to be an important next step and we are actively working in this direction.

## 6 Summary

We presented a concept-driven framework for uncovering latent themes in text collections. Our framework expands the definitions of a theme to account for theoretically informed concepts that generalize beyond word co-occurrence patterns. We suggest an interactive protocol that allows domain experts to interact with the data and provide feedback at different levels of abstraction. We performed an exhaustive evaluation using two case studies and different groups of experts. Additionally, we contrasted the extracted themes against the output of traditional topic models, and showed that they are better at capturing conceptual similarities that go beyond word distributions.

## Acknowledgements

We thank the anonymous reviewers of this paper for all of their feedback. This work was partially supported by an NSF CAREER award IIS-2048001.

## References

Avnika B Amin, Robert A Bednarczyk, Cara E Ray, Kala J Melchiori, Jesse Graham, Jeffrey R Huntsinger, and Saad B Omer. 2017. Association of moral values with vaccine hesitancy. *Nature Human Behaviour*, 1(12):873–880.

David Antons, Eduard Grünwald, Patrick Cichy, and Oliver Salge. 2020. [The application of text mining methods in innovation research: current state, evolution patterns, and development priorities](#). *RD Management*, 50.

Christian Baden, Christian Pipal, Martijn Schoonvelde, and Mariken A. C. G van der Velden. 2022. [Three gaps in computational text analysis methods for social sciences: A research agenda](#). *Communication Methods and Measures*, 16(1):1–18.

Rodney Benson. 2013. *Shaping Immigration News: A French-American Comparison*. Communication, Society and Politics. Cambridge University Press.

Vibhu Bhatia, Vidya Prasad Akavoor, Sejin Paik, Lei Guo, Mona Jalal, Alyssa Smith, David Assefa Tofu, Edward Edberg Halim, Yimeng Sun, Margrit Betke,

- Prakash Ishwar, and Derry Tanti Wijaya. 2021. [Open-Framing: Open-sourced tool for computational framing analysis of multilingual data](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 242–250, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- David M. Blei, Andrew Y. Ng, and Michael I. Jordan. 2003. Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3(null):993–1022.
- Jordan Boyd-Graber, Yuening Hu, and David Minmo. 2017. *Applications of Topic Models*.
- Amber Boydston, Dallas Card, Justin H. Gross, Philip Resnik, and Noah A. Smith. 2014. Tracking the development of media frames within and across policy issues.
- Henry E. Brady. 2019. [The challenge of big data and data science](#). *Annual Review of Political Science*, 22(1):297–323.
- Virginia Braun and Victoria Clarke. 2012. *Thematic analysis.*, pages 57–71.
- Dallas Card, Amber E. Boydston, Justin H. Gross, Philip Resnik, and Noah A. Smith. 2015. [The media frames corpus: Annotations of frames across issues](#). In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 438–444, Beijing, China. Association for Computational Linguistics.
- Eugene Y Chan. 2021. Moral foundations underlying behavioral compliance during the covid-19 pandemic. *Personality and individual differences*, 171:110463.
- Dennis Chong and James N Druckman. 2007. Framing theory. *Annu. Rev. Polit. Sci.*, 10:103–126.
- Dorottya Demszky, Nikhil Garg, Rob Voigt, James Zou, Jesse Shapiro, Matthew Gentzkow, and Dan Jurafsky. 2019. Analyzing polarization in social media: Method and application to tweets on 21 mass shootings. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2970–3005.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Robert M Entman. 1993. Framing: Toward clarification of a fractured paradigm. *Journal of communication*, 43(4):51–58.
- Jonathan Haidt and Jesse Graham. 2007. When morality opposes justice: Conservatives have moral intuitions that liberals may not recognize. *Social Justice Research*, 20(1):98–116.
- Martin Hilbert, George Barnett, Joshua Blumenstock, Noshir Contractor, Jana Diesner, Seth Frey, Sandra González-Bailón, PJ Lamberson, Jennifer Pan, Tai-Quan Peng, Cuihua (Cindy) Shen, Paul E. Smaldino, Wouter van Atteveldt, Annie Waldherr, Jingwen Zhang, and Jonathan J. H. Zhu. 2019. [Computational communication science computational communication science: A methodological catalyzer for a maturing discipline](#). *International Journal of Communication*, 13(0).
- Jan Fredrik Hovden and Hilmar Mjelde. 2019. [Increasingly controversial, cultural, and political: The immigration debate in scandinavian newspapers 1970–2016](#). *Javnost - The Public*, 26(2):138–157.
- Alexander Hoyle, Pranav Goel, Andrew Hian-Cheong, Denis Peskov, Jordan Boyd-Graber, and Philip Resnik. 2021. [Is automated topic model evaluation broken? the incoherence of coherence](#). In *Advances in Neural Information Processing Systems*, volume 34, pages 2018–2033. Curran Associates, Inc.
- Yuening Hu, Jordan Boyd-Graber, and Brianna Satinoff. 2011. [Interactive topic modeling](#). In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 248–257, Portland, Oregon, USA. Association for Computational Linguistics.
- Shanto. Iyengar. 1991. *Is anyone responsible? : how television frames political issues*. American politics and political economy series. University of Chicago Press, Chicago.
- Xin Jin and Jiawei Han. 2010. *K-Means Clustering*, pages 563–564. Springer US, Boston, MA.
- Ralph H. Johnson. 1988. [Gilbert harman change in view: Principles of reasoning](#) (cambridge, ma: Mit press 1986). pp. ix 147. *Canadian Journal of Philosophy*, 18(1):163–178.
- Jey Han Lau, David Newman, and Timothy Baldwin. 2014. [Machine reading tea leaves: Automatically evaluating topic coherence and topic model quality](#). In *Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics*, pages 530–539, Gothenburg, Sweden. Association for Computational Linguistics.
- Claire Lauer, Eva Brumberger, and Aaron Beveridge. 2018. [Hand collecting and coding versus data-driven methods in technical and professional communication research](#). *IEEE Transactions on Professional Communication*, 61(4):389–408.
- Jeffrey Lund, Connor Cook, Kevin Seppi, and Jordan Boyd-Graber. 2017. [Tandem anchoring: a multiword anchor approach for interactive topic modeling](#). In

- Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 896–905, Vancouver, Canada. Association for Computational Linguistics.
- Andrew Kachites McCallum. 2002. Mallet: A machine learning for language toolkit. [Http://mallet.cs.umass.edu](http://mallet.cs.umass.edu).
- Leland McInnes, John Healy, and Steve Astels. 2017. hdbscan: Hierarchical density based clustering. *The Journal of Open Source Software*, 2(11):205.
- Julia Mendelsohn, Ceren Budak, and David Jurgens. 2021. [Modeling framing in immigration discourse on social media](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2219–2263, Online. Association for Computational Linguistics.
- David Mimno, Hanna Wallach, Edmund Talley, Miriam Leenders, and Andrew McCallum. 2011. [Optimizing semantic coherence in topic models](#). In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing*, pages 262–272, Edinburgh, Scotland, UK. Association for Computational Linguistics.
- Fred Morstatter, Liang Wu, Uraz Yavanoglu, Stephen R. Corman, and Huan Liu. 2018. [Identifying framing bias in online news](#). *Trans. Soc. Comput.*, 1(2):5:1–5:18.
- Maria Leonor Pacheco and Dan Goldwasser. 2021. [Modeling content and context with deep relational learning](#). *Transactions of the Association for Computational Linguistics*, 9:100–119.
- Maria Leonor Pacheco, Tunazzina Islam, Monal Mahajan, Andrey Shor, Ming Yin, Lyle Ungar, and Dan Goldwasser. 2022. [A holistic framework for analyzing the COVID-19 vaccine debate](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5821–5839, Seattle, United States. Association for Computational Linguistics.
- Radim Rehurek and Petr Sojka. 2011. Gensim–python framework for vector space modelling. *NLP Centre, Faculty of Informatics, Masaryk University, Brno, Czech Republic*, 3(2).
- Nils Reimers and Iryna Gurevych. 2019. [Sentence-BERT: Sentence embeddings using Siamese BERT-networks](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.
- Kate Roberts, Anthony Dowell, and Jing-Bao Nie. 2019. Attempting rigour and replicability in thematic analysis of qualitative research data; a case study of codebook development. *BMC Medical Research Methodology*, 19.
- Michael Röder, Andreas Both, and Alexander Hinneburg. 2015. [Exploring the space of topic coherence measures](#). In *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining, WSDM '15*, page 399–408, New York, NY, USA. Association for Computing Machinery.
- T. Rogers and James L. McClelland. 2004. Semantic cognition: A parallel distributed processing approach.
- Jeremy Rose and Christian Lennerholt. 2017. Low cost text mining as a strategy for qualitative researchers. *Electronic Journal on Business Research Methods*, forthcoming.
- Shamik Roy and Dan Goldwasser. 2021. [Analysis of nuanced stances and sentiment towards entities of US politicians through the lens of moral foundation theory](#). In *Proceedings of the Ninth International Workshop on Natural Language Processing for Social Media*, pages 1–13, Online. Association for Computational Linguistics.
- Shamik Roy, Maria Leonor Pacheco, and Dan Goldwasser. 2021. [Identifying morality frames in political tweets using relational learning](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 9939–9958, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Alison Smith, Varun Kumar, Jordan Boyd-Graber, Kevin Seppi, and Leah Findlater. 2018. [Closing the loop: User-centered design and evaluation of a human-in-the-loop topic modeling system](#). In *23rd International Conference on Intelligent User Interfaces, IUI '18*, page 293–304, New York, NY, USA. Association for Computing Machinery.
- Keith Stevens, Philip Kegelmeyer, David Andrzejewski, and David Buttler. 2012. [Exploring topic coherence over many models and many topics](#). In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 952–961, Jeju Island, Korea. Association for Computational Linguistics.
- Laurens van der Maaten and Geoffrey Hinton. 2008. [Visualizing data using t-SNE](#). *Journal of Machine Learning Research*, 9:2579–2605.
- Zijie J. Wang, Dongjin Choi, Shenyu Xu, and Diyi Yang. 2021. [Putting humans in the natural language processing loop: A survey](#). In *Proceedings of the First Workshop on Bridging Human–Computer Interaction and Natural Language Processing*, pages 47–52, Online. Association for Computational Linguistics.

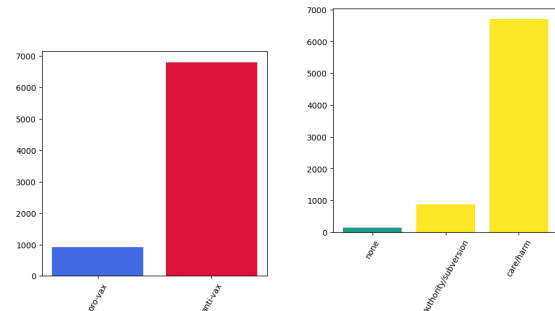
# A Appendix

## A.1 Tool Screenshots

### A.1.1 Exploratory Operations

Method  
  
 K (# initial clusters, only needed if using K-means)

Figure 7: Cluster Instances



(a) Stance (b) Moral Foundation

Figure 13: Visualizing Local Explanations: Attribute Distribution for *The Vaccine Doesn't Work*

Query by theme OR Write a text query  
 Theme:

Figure 8: Text-based Queries

Showing tweets similar to:  
 Thank you for your leadership on this critical issue, @GovSisolak. <https://t.co/luYNNX1Df>

id	tweet_id	text	stance	distance	good	morality	mf	theme_id	select
74343	74342	Thank you for your leadership on this critical issue, @GovSisolak. <a href="https://t.co/luYNNX1Df">https://t.co/luYNNX1Df</a>	pro-vax	0.13269954919815063	True	moral	authority/subversion	13	<input type="checkbox"/>
878	877	We know you care about this issue as much as we do. @POTUS @JoeBiden @FLOTUS @docsinpolitics <a href="https://t.co/7bp9wWICY">https://t.co/7bp9wWICY</a> <a href="https://t.co/uvvmfzPfg">https://t.co/uvvmfzPfg</a>	pro-vax	0.18669486045837402	True	moral	authority/subversion	13	<input type="checkbox"/>
2983	2982	Thank You @POTUS! So productive having REAL leadership from the @WhiteHouse!! #Biden #BuildBackBetter #COVID19 #COVID #vaccine <a href="https://t.co/mo50EiNesh">https://t.co/mo50EiNesh</a>	pro-vax	0.17249584197998047	True	moral	authority/subversion	13	<input type="checkbox"/>

Figure 9: Finding Similar Tweets

Query by theme OR Write a text query  
 Theme:

Show 5 entries

id	tweet_id	text	stance	distance	good	morality	mf	theme_id	select
74343	74342	Thank you for your leadership on this critical issue, @GovSisolak. <a href="https://t.co/luYNNX1Df">https://t.co/luYNNX1Df</a>	pro-vax	0.13269954919815063	True	moral	authority/subversion	13	<input type="checkbox"/>
2983	2982	Thank You @POTUS! So productive having REAL leadership from the @WhiteHouse!! #Biden #BuildBackBetter #COVID19 #COVID #vaccine <a href="https://t.co/mo50EiNesh">https://t.co/mo50EiNesh</a>	pro-vax	0.17249584197998047	True	moral	authority/subversion	13	<input type="checkbox"/>

Figure 10: Listing Arguments and Examples

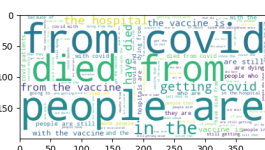


Figure 11: Visualizing Local Explanations: Word Cloud Example for *The Vaccine Doesn't Work*

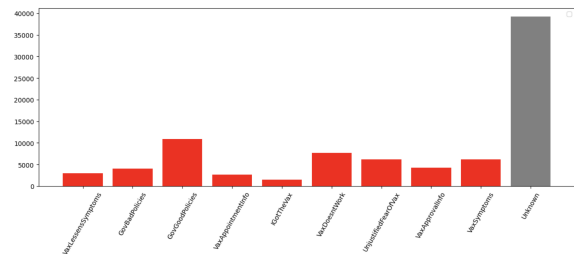


Figure 14: Visualizing Global Explanations: Theme Distribution

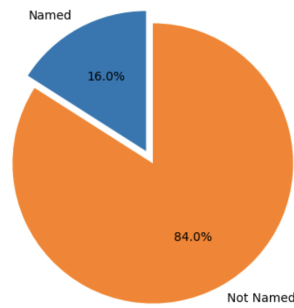


Figure 15: Visualizing Global Explanations: Coverage

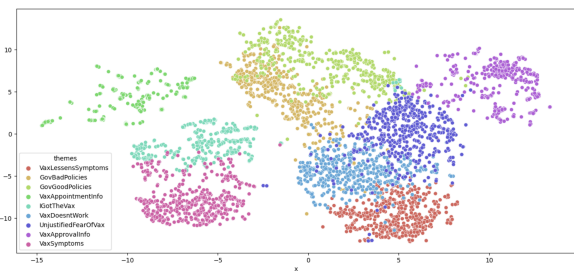


Figure 16: Visualizing Global Explanations: 2D t-SNE

### Top 10 Positive Entities

entity
vaccine
a comprehensive school response
student academic and mental health recovery plans
the model

(a) Top Positive Entities

### Top 10 Negative Entities

entity
the vaccine
covid
biden
trump

(b) Top Negative Entities

Figure 12: Visualizing Local Explanations: Most Frequent Positive and Negative Entities for *Bad Governmental Policies*

## A.1.2 Intervention Operations

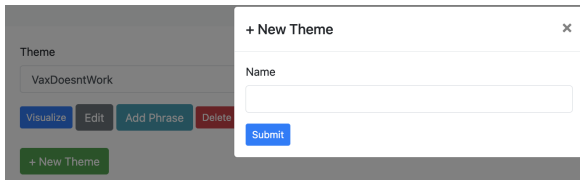


Figure 17: Adding New Themes



Figure 18: Marking Instances as *Good*

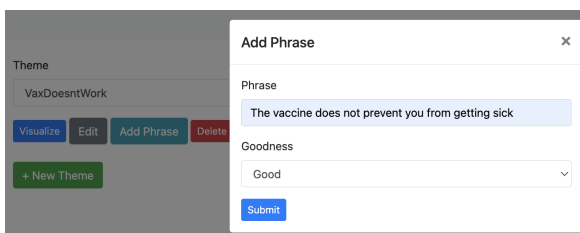


Figure 19: Adding *Good* Examples

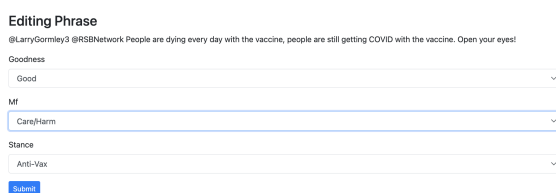


Figure 20: Correcting Attributes - Stances and Moral Foundations

## A.2 Interactive Sessions for Covid: First Group of Experts

Table 4 and 5 outline the patterns discovered by the the first group of experts on the first a second iteration, respectively.

## A.3 Interactive Sessions for Covid: Second Group of Experts

Table 6 and 7 outline the patterns discovered by the second group of experts on the first a second iteration, respectively.

## A.4 Interactive Sessions for Immigration

Table 8 and 9 outline the patterns discovered by the experts for immigration.

## A.5 Topic Modeling Details

To obtain LDA topics with Variational Bayes sampling we use the Gensim implementation (Rehurek and Sojka, 2011). To obtain LDA topics with Gibbs sampling we use the MALLET implementation (McCallum, 2002). In both cases, we follow all the preprocessing steps suggested by Hoyle et al. (2021), with the addition of the words covid, vaccin\* and immigra\* to the list of stopwords.

## A.6 Fine-Grained Results

The confusion matrix for Immigration can be seen in Fig. 21. Distribution of errors that do not match any existing theme, according to their similarity interval can be seen in Fig 22.

Cluster	Experts Rationale	New Named Themes
K-Means 0	Discusses what the vaccine can and cannot do. Emphasis in reducing COVID-19 symptoms in case of infection (“like a bad cold”). Contains tweets with both stances.	VaxLessensSymptoms
K-Means 1	A lot of mentions to political entities. Politicians get in the way of public safety	GovBadPolicies
K-Means 2	A lot of tweets with mentions and links. Not a lot of textual context. Some examples thanking and praising governmental policies. <b>Theme added upon inspecting similar tweets</b>	GovGoodPolicies
K-Means 3	Overarching theme related to vaccine rollout. Mentions to pharmacies that can distribute, distribution in certain states, places with unfulfilled vax appointments. <b>Too broad to create a theme</b>	-
K-Means 4	Broadcast of vaccine appointments. Which places you can get vaccine appointments at.	VaxAppointments
K-Means 5	“I got my vaccine” type tweets	GotTheVax
K-Means 6	Mixed cluster, not a clear theme in centroid. Two prominent flavors: the vaccine not working and people complaining about those who are scared of vaccine.	VaxDoesntWork UnjustifiedFearOfVax
K-Means 7	Tweets look the same as K-Means 5	-
K-Means 8	Tweets about development and approval of vaccines	VaxApproval
K-Means 9	Tweets related to common vaccine side-effects	VaxSideEffects

Table 4: **First Iteration:** Patterns Identified in Initial Clusters and Resulting Themes

Cluster	Experts Rationale	New Named Themes
K-Means 0	Tweets weighting health benefits/risks, but different arguments. (e.g. it works, doesn’t work, makes things worse...) <b>Too broad to create a theme.</b>	-
K-Means 1	Messy cluster, relies on link for information.	-
K-Means 2	Relies on link for information.	-
K-Means 3	A lot of mentions to government lying and misinformation. “misinformation” is used when blaming antivax people. “experts and government are lying” is used on the other side. References to alt-treatments on both sides. <b>Text lookup “give “us the real meds”, “covid meds”</b>	AntiVaxSpreadMisinfo ProVaxLie AltTreatmentsGood AltTreatmentsBad
K-Means 4	Some examples are a good fit for old theme, VaxDoesntWork. <b>Other than that no coherent theme.</b>	-
K-Means 5	Tweets about free will and choice. <b>Text lookup “big gov”, “free choice”, “my body my choice”</b> Case “my body my choice” - a lot of mentions to abortion People using covid as a metaphor for other issues.	FreeChoiceVax FreeChoiceOther
K-Means 6	Almost exclusively mentions to stories and news.	-
K-Means 7	Availability of the vaccine, policy. Not judgement of good or bad, but of how well it progresses.	VaxEffortsProgression
K-Means 8	Assign to previous theme GotTheVax	-
K-Means 9	Vaccine side effects. Assign to previous theme, VaxSymptoms	-

Table 5: **Second Iteration:** Patterns Identified in Subsequent Clusters and Resulting Themes

Cluster	Experts Rationale	New Named Themes
K-Means 0	People asking people to get vaccinated. Some skeptical but acknowledge it reduces symptoms. It works but it has limitations. More specifically, it lessens the symptoms.	VaxLessensSymptoms
K-Means 1	Republicans have hurt the vax rate in the US. Finding someone (or some party) to blame. Politicians are hurting people with policy. Vaccine in the US is behind, trying to explain why	ReasonsUSLagsOnVax
K-Means 2	A lot of them are just replies. Cluster is for links and usernames.	-
K-Means 3	Availability and distribution of the vaccine. How stances of people in different states affect it. Vaccine distribution issues due to local policy.	VaxDistributionIssuesDueToLocalPolicy
K-Means 4	Clear cluster. Vaccine info, availability info.	VaxAvailabilityInfo
K-Means 5	Testimonials, #IGotMyVax	#IGotMyVax
K-Means 6	Some themes match the vaccine lessens symptoms. Other theme: no need to get the vaccine, it doesn't work. Vaccine does more harm than good.	VaxDoesMoreHarmThanGood
K-Means 7	Same as K-means 5	-
K-Means 8	About covid vaccine updates. FDA approval. In other cases it depends on the content on the link. So you can't really tell.	FDAApproval
K-Means 9	Obvious. Vaccine symptoms, vaccine effects. Post vaccination symptoms.	PostVaxSymptoms

Table 6: **Second Group's First Iteration:** Patterns Identified in Initial Clusters and Resulting Themes

Cluster	Experts Rationale	New Named Themes
K-Means 0	Links and promotions	-
K-Means 1	Looks like previous theme IGotMyVax, assign them.	-
K-Means 2	Very short tweets with links, and no context. Could be availability but not sure. Decided against adding theme	-
K-Means 3	Two themes observed. One old one, regarding VaxAvailabilityInfo. One new one, getting vaccines is difficult. Not related to local policy. <b>Decided against merging with previous theme</b>	VaxDistributionIssues
K-Means 4	A lot of talk about skepticism regarding the vaccine. Some good matches to previous MoreHarmThanGood, assign them. Mentions to profiting from the vaccine. <b>Look for similar instances to mentions of profits</b> <b>Text look up for "vaccine getting rich"</b> Mentions to redlining, implications of inequality <b>Text look up for "vaccine inequality"</b> Lots of mentions to racial and monetary inequalities in access to vaccine	VaxCapitalism VaxInequality
K-Means 5	Both PostVaxSymptoms and IGotMyVax examples, assign them.	-
K-Means 6	Mentions to vaccine safety. Weighting the safety/risks of the vaccine	VaxSafety
K-Means 7	A lot of discussion about the pandemic not being over Discussion on whether to open back up or not	CovidNotOver
K-Means 8	Repetitions, IGotMyVax. Assign them.	-
K-Means 9	Mentions to mandates. The vaccine should be a personal choice, mandates should not be there. Different reasons: personal choice, no proof of whether it works. For no proof, assign to previous MoreHarmThanGood	VaxPersonalChoice

Table 7: **Second Group's Second Iteration:** Patterns Identified in Subsequent Clusters and Resulting Themes

Cluster	Experts Rationale	New Named Themes
K-Means 0	Headlines, coverage. Some have an agenda (pro) Others are very academic and research-oriented Opinion pieces.	AcademicDiscussions
K-Means 1	Talking about apprehending immigrants at the border Some report about the border but no stance. Deportation. Leaning negative towards immigrants.	JustifiedDetainmentEnforce
K-Means 2	Less US-centric, more general. Talking about immigration as a global issue Humanitarian issues, mentions to refugees, forced migration Situation in country of origin that motivates immigration Mentions to how the west is responsible The role of the target countries in destabilizing countries Mentions to economic migrants. <b>Look up for "economic work migrants", "asylum seekers"</b>	EconomicMigrantsNotAsylumSeekers SituationCountryOfOrigin RoleOfWesternCountries
K-Means 3	About Trump. Trump immigration policy. Politicizing immigration.	TrumpImmiPolicy
K-Means 4	Attacking democrats. A lot of mentions to democrats wanting votes Common threads is democrats are bad	DemocratImmiPolicyBad
K-Means 5	Lacks context, lots of usernames. Not a cohesive theme. Both pro and con, and vague. Some mentions to invasion. <b>Look for "illegal immigrants invade"</b> Mentions to caravan, massive exodus of people. Mentions to crime. <b>Look for immigrants murder, immigrants dangerous.</b> A lot of tweets linking immigrants to crime	ImmigrantInvasion ImmigrantCrime
K-Means 6	Looks very varied. Not cohesive.	-
K-Means 7	Very cohesive. Mentions to detaining children, families.	DetainingChildren
K-Means 8	All tweets are about the UK and Britain. Both pro and anti immigration. Only common theme is the UK. Almost exclusively policy/politics	UKProImmiPolicy UKAntiImmiPolicy
K-Means 9	Economic cost of immigration. Immigration is bad for the US economy Some about crime, and democrats. Assign to existing themes.	FinacialCostOfImmigration

Table 8: **First Iteration Immigration:** Patterns Identified in Initial Clusters and Resulting Themes

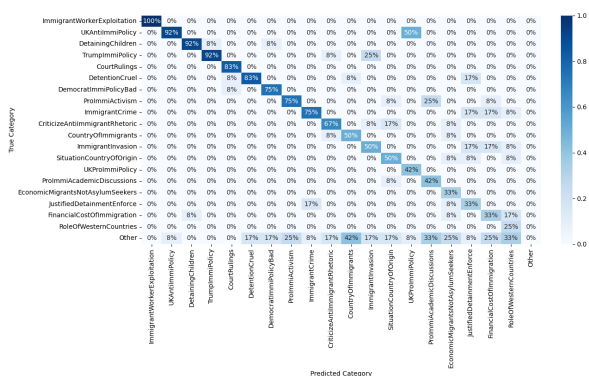


Figure 21: **Confusion matrix of Immigration themes after second iteration.** Values are normalized over the predicted themes (columns), and sorted from most accurate to least accurate.

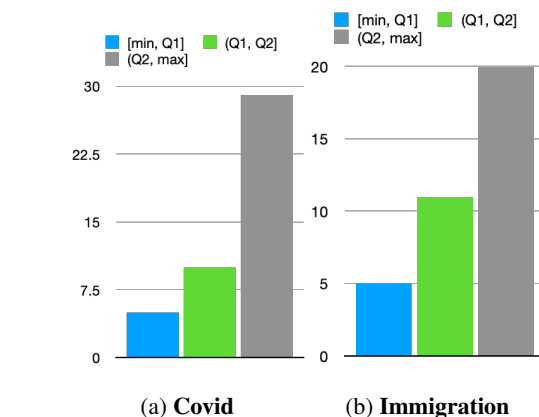


Figure 22: Tweets that Do Not Match Current Set of Themes (True Category is "Other") at Different Intervals

### A.7 Shifting Predictions between Iterations

Heatmaps of shifting predictions for Covid can be seen in Fig. 23. The distribution of the unmatched predictions for both Covid and Immigration, according to their similarity intervals can be seen in



Cluster	Experts Rationale	New Named Themes
K-Means 0	Legal decisions and rulings. Both pro and anti immigration rulings Not a single event, but cohesively talking about rulings	CourtRulings
K-Means 1	The same tweet reworded and tweeted at different people Talks about worker exploitation, and Cesar Chavez. <b>Look up for "exploitation"</b> . Mentions to workers and wages <b>Look up for "cheap labor"</b>	ImmigrantWorkerExploitation
K-Means 2	Blaming Trump for being irresponsible Criticizing his rhetoric. Mentions to hateful speech About the rhetoric rather than policy. Mentions to racist language Others about policy, added to previous TrumpImmiPolicy theme	CriticizeAntiImmigrantRhetoric
K-Means 3	Nation of immigrants. Identity, we are all immigrants	CountryOfImmigrants
K-Means 4	Organizing. Call to action. Skews pro. language of rights and liberties. We are here, we demand, sign here. <b>Look up "ACLU", "rights for immigrants"</b>	ProImmiActivism
K-Means 5	A lot of mentions to numbers and stats. Short URLs. Headlines.	-
K-Means 6	A lot of usernames. Bad policies, criticizing policies on both sides. Send them to either DemocratImmiPolicyBad or TrumpImmiPolicy	-
K-Means 7	Very messy. Links.	-
K-Means 8	European headlines and news. Some about the UK. Send the ones that are relevant to UK policy themes	-
K-Means 9	Detention, detention centers, solitary confinement as cruel.	DetainmentCruel

Table 9: **First Iteration Immigration**: Patterns Identified in Initial Clusters and Resulting Themes

Fig. 24. Additionally, some examples of shifting predictions for the two themes with the most movement for the Immigration case can be seen in Tabs. 10 and 11.

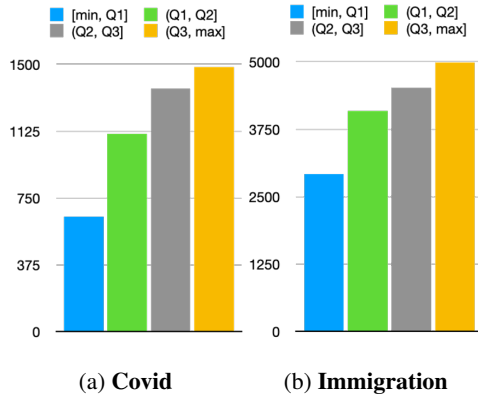


Figure 24: Unmatched Predictions (Shifting from Named Theme to Unknown) at Different Intervals

## A.8 LDA vs. our Themes

An overlap coefficient heatmap between LDA topics with Variational Bayes sampling and our themes for the first iteration of Covid can be seen in Fig. 25. Similarly, they can be seen for the second iterations of both Covid and Immigration in Fig. 26. We also include these heatmaps for LDA with Gibbs sampling in Figs. 27, 28 and 29

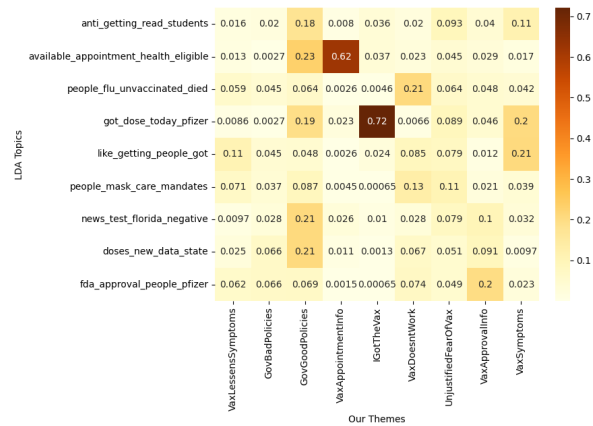


Figure 25: Overlap Coefficients between LDA Var. Bayes and our Themes (First Iteration for Covid).

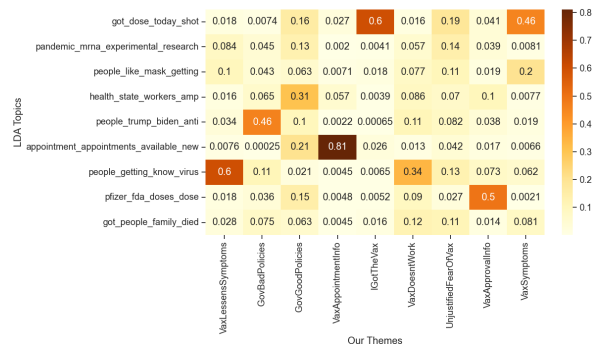


Figure 27: Overlap Coefficients between LDA Gibbs Sampling and our Themes (First Iteration for Covid).

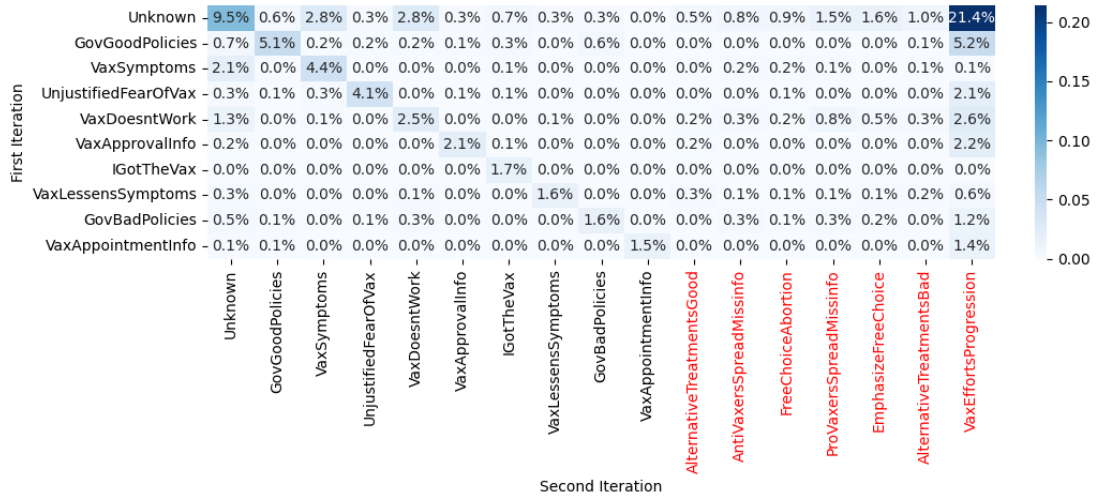


Figure 23: **Shifting predictions for Covid.** Themes added during second iteration are shown in red, and values are normalized over the full population.

Distance to Centroid	Example Tweets Kept on <i>Role of Western Countries</i>	Example Tweets Shifted to <i>Unknown</i>
0.27	The U.S. Helped Destabilize Honduras. Now Honduran Migrants Are Fleeing Political and Economic Crisis	Interesting that your problem is with "migrants", where the U.S. has issues with illegal aliens, that even our legal migrants wish to be rid of.
0.29	These people are fleeing their countries DIRECTLY because of U.S. ForeignPolicy. If you don't like refugees. Don't create 'em.	The root causes of migration aren't being addressed ASAP, as they must be. The governments are all busy talking about stopping the consequences without concrete plans to solve the causes.
0.30	Don't want migrants? Stop blowing their countries to pieces	What's missing in the US corporate news on migrants is the way American "aid" is used to overturn democracies, prop up strongmen and terrify the opposition.

Table 10: *Role of Western Countries*: Examples of tweets kept on theme (Left) and shifted to unknown (Right) between the first and second iteration. On Right are the tweets closest to the theme centroid that shifted to *Unknown*. On Left are tweets that did *not* shift, but have the same distance.

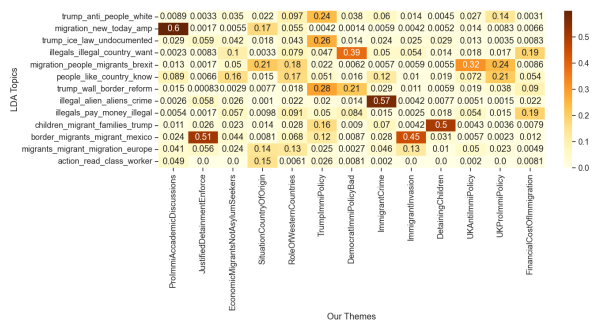


Figure 28: Overlap Coefficients between LDA Gibbs Sampling and our Themes (First Iteration for **Immigration**).

Distance to Centroid	Example Tweets Kept on <i>Trump Immigration Policy</i>	Example Tweets Shifted to <i>Unknown</i>
0.24	Racist realDonaldTrump wastes our tax money on locking up little kids in #TrumpConcentrationCamps and steals from our military to waste money on his #ReElectionHate-Wall and spends little on anything else.	The anti-migrant cruelty of the Trump Admin knows no bounds. This targeting of migrant families is meant to induce fear and doesnt address our broken immigration system. We should be working to make our immigration system more humane, not dangerous and cruel.
0.25	Trump promises immigration crackdown ahead of U.S. election	This is unlawful and is directed at mothers with their children! He had no remorse for separating immigrants earlier, now he's threatening their lives! It's heart wrenching, but Trumpf has no heart! He's void of feeling empathy! Read they are in prison camps? WH ignoring cries
0.26	Trump to end asylum protections for most Central American migrants at US-Mexico border	BBC News - Daca Dreamers: Trump vents anger on immigrant programme

Table 11: *Trump Immigration Policy*: Examples of tweets kept on theme (Left) and shifted to unknown (Right) between the first and second iteration. On Right are the tweets closest to the theme centroid that shifted to *Unknown*. On Left are tweets that did *not* shift, but have the same distance.

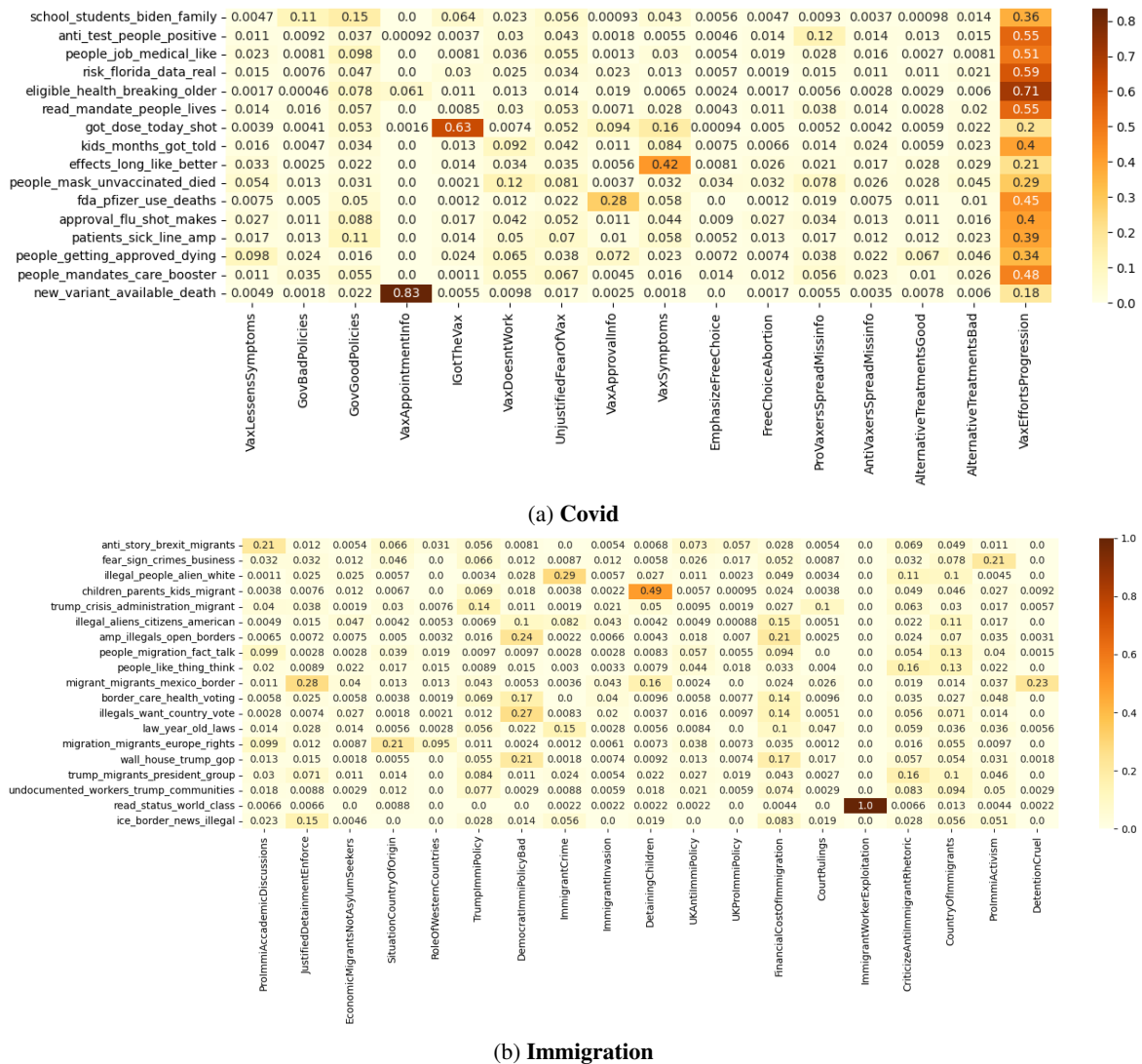
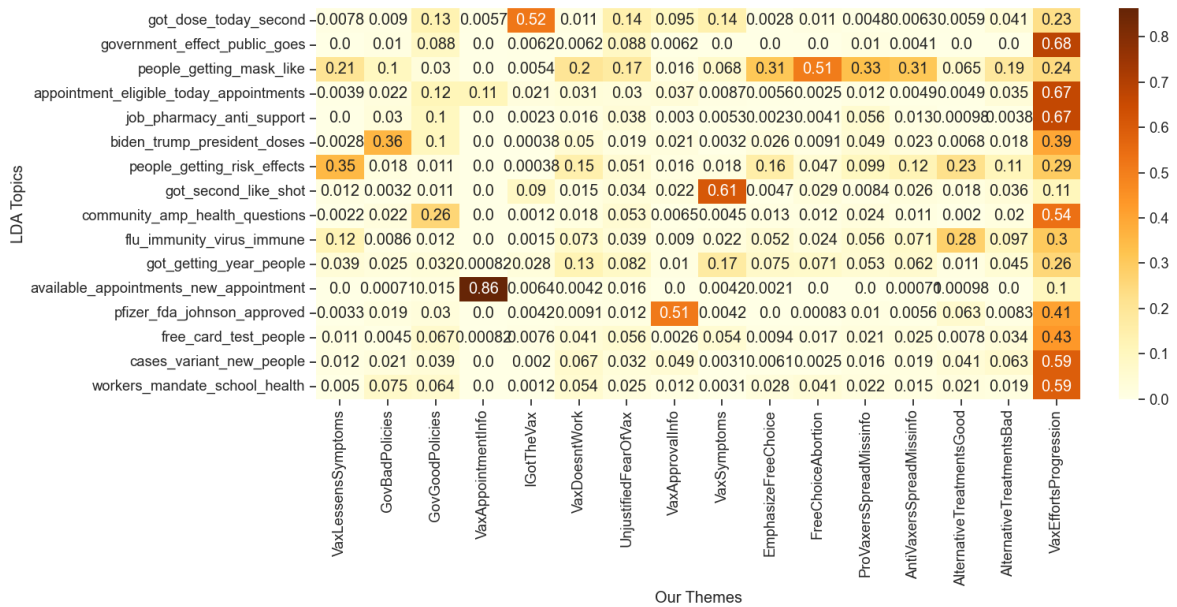
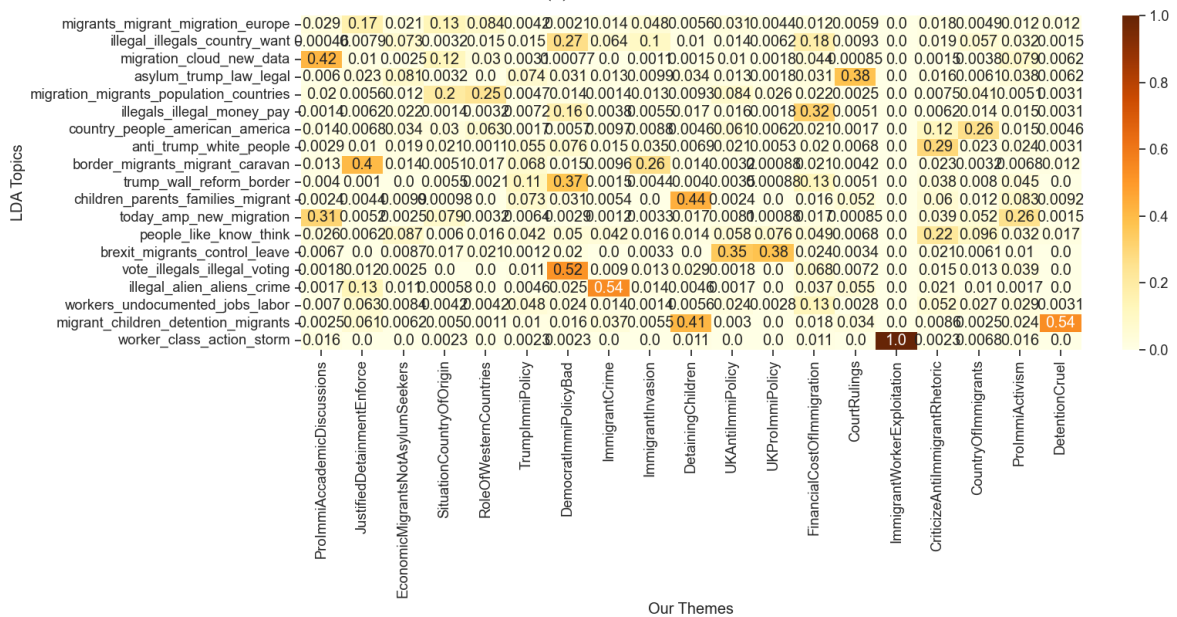


Figure 26: Overlap Coefficients between LDA Var. Bayes and our Themes (Second Iteration).



(a) Covid



(b) Immigration

Figure 29: Overlap Coefficients between LDA Gibbs Sampling and our Themes (Second Iteration).

## ACL 2023 Responsible NLP Checklist

---

### A For every submission:

- A1. Did you describe the limitations of your work?  
*Section 5: Limitations*
- A2. Did you discuss any potential risks of your work?  
*Not applicable. Left blank.*
- A3. Do the abstract and introduction summarize the paper's main claims?  
*Abstract, Section 1: Introduction*
- A4. Have you used AI writing assistants when working on this paper?  
*Left blank.*

### B Did you use or create scientific artifacts?

*Yes, Section 3: Framework*

- B1. Did you cite the creators of artifacts you used?  
*Yes, Section 3: Framework*
- B2. Did you discuss the license or terms for use and / or distribution of any artifacts?  
*Linked to gitlab, where this information is provided.*
- B3. Did you discuss if your use of existing artifact(s) was consistent with their intended use, provided that it was specified? For the artifacts you create, do you specify intended use and whether that is compatible with the original access conditions (in particular, derivatives of data accessed for research purposes should not be used outside of research contexts)?  
*Linked to gitlab, where this information is provided.*
- B4. Did you discuss the steps taken to check whether the data that was collected / used contains any information that names or uniquely identifies individual people or offensive content, and the steps taken to protect / anonymize it?  
*Not applicable. No new data introduced.*
- B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?  
*Section 4: Case Studies*
- B6. Did you report relevant statistics like the number of examples, details of train / test / dev splits, etc. for the data that you used / created? Even for commonly-used benchmark datasets, include the number of examples in train / validation / test splits, as these provide necessary context for a reader to understand experimental results. For example, small differences in accuracy on large test sets may be significant, while on small test sets they may not be.  
*Section 4: Case Studies*

### C Did you run computational experiments?

*Section 4: Case Studies*

- C1. Did you report the number of parameters in the models used, the total computational budget (e.g., GPU hours), and computing infrastructure used?  
*Section 4: Case Studies*

---

*The Responsible NLP Checklist used at ACL 2023 is adopted from NAACL 2022, with the addition of a question on AI writing assistance.*

- C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values?  
*Section 4: Case Studies*
  - C3. Did you report descriptive statistics about your results (e.g., error bars around results, summary statistics from sets of experiments), and is it transparent whether you are reporting the max, mean, etc. or just a single run?  
*Section 4: Case Studies*
  - C4. If you used existing packages (e.g., for preprocessing, for normalization, or for evaluation), did you report the implementation, model, and parameter settings used (e.g., NLTK, Spacy, ROUGE, etc.)?  
*Linked to gitlab, where this information is provided.*
- D**  **Did you use human annotators (e.g., crowdworkers) or research with human participants?**  
*Section 3: Framework, Section 4: Case Studies*
- D1. Did you report the full text of instructions given to participants, including e.g., screenshots, disclaimers of any risks to participants or annotators, etc.?  
*Section 3: Framework, Section 4: Case Studies*
  - D2. Did you report information about how you recruited (e.g., crowdsourcing platform, students) and paid participants, and discuss if such payment is adequate given the participants' demographic (e.g., country of residence)?  
*Section 3: Framework, Section 4: Case Studies*
  - D3. Did you discuss whether and how consent was obtained from people whose data you're using/curating? For example, if you collected data via crowdsourcing, did your instructions to crowdworkers explain how the data would be used?  
*Section 3: Framework, Section 4: Case Studies*
  - D4. Was the data collection protocol approved (or determined exempt) by an ethics review board?  
*Section 3: Framework, Section 4: Case Studies*
  - D5. Did you report the basic demographic and geographic characteristics of the annotator population that is the source of the data?  
*Section 3: Framework, Section 4: Case Studies*