

Effects of Collaboration on the Performance of Interactive Theme Discovery Systems

Anonymous ACL submission

Abstract

NLP-assisted solutions have gained considerable traction to support qualitative data analysis. However, no unified evaluation framework exists which can account for the many different settings in which qualitative researchers may employ them. In this paper, we propose an evaluation framework to study the way collaboration settings may produce different outcomes across a variety of interactive systems. Specifically, we study the impact of synchronous vs. asynchronous collaboration using three different NLP-assisted qualitative research tools and present a comprehensive analysis of significant differences in the consistency, cohesiveness, and correctness of their outputs.

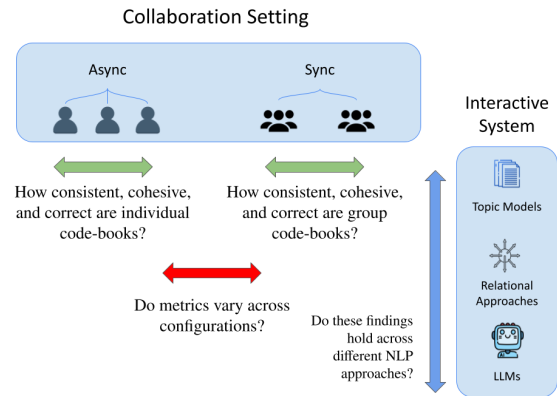


Figure 1: In this study, we measure the quality of coded themes using different interactive systems under different coding configurations.

1 Introduction

Making sense of large textual datasets is a common challenge across academic disciplines and is traditionally addressed through qualitative methods such as Thematic Analysis (Braun and Clarke, 2006) and Grounded Theory (Glaser et al., 1968). These approaches rely on manual *inductive coding*, in which researchers identify abstract themes by closely reading the data. However, as datasets grow in size, manual coding becomes impractical, motivating the use of Natural Language Processing (NLP) techniques to automate parts of the analysis process (Brady, 2019; Hilbert et al., 2019).

In recent years, a range of NLP-based systems have been developed to support qualitative research. These systems assist researchers by uncovering latent semantic structures through topic modeling (Smith et al., 2018; Fang et al., 2023), clustering documents and propagating limited human annotations across datasets (Pacheco et al., 2023; Chew et al., 2023), or offering real-time coding recommendations (Dai et al., 2023; Gao et al., 2024). To maintain researcher agency, such systems typically adopt human-in-the-loop (HitL) strategies that balance automation with manual interpretation.

Previous work typically evaluates HitL qualitative analysis tools in isolation, focusing on specific technical strengths and weaknesses. For example, by comparing topic coherence with and without human input (Fang et al., 2023) or contrasting machine-assisted and manual code-book generation (Dai et al., 2023). However, qualitative analysis in practice is often collaborative, with teams of researchers jointly coding and interpreting data (Flick, 2014), and supported by tools built on diverse methodologies that are applied to datasets with vastly differing characteristics (Baden et al., 2022). By abstracting away collaboration settings, methodological variation, and dataset diversity, existing evaluations risk misrepresenting how HitL systems operate in real-world settings and how broadly their findings apply. In this work, we seek to answer the following questions: (1) Does the collaboration setting measurably affect the quality of resulting code-books? (2) Do these findings hold across different NLP approaches? (3) How do dataset characteristics influence the outcomes of NLP-assisted inductive coding tools?

We focus on two common but contrasting collab-

065 oration settings: asynchronous coding, where in- 115
066 dividuals code independently before consolidating 116
067 results, and synchronous coding, where teams iden-
068 tify themes through live discussion. These settings
069 offer complementary strengths, with asynchronous
070 coding supporting flexibility across time and loca-
071 tion, and synchronous coding facilitating shared
072 understanding and efficient coordination. To com-
073 pare outcomes across these settings, we introduce
074 an evaluation framework that measures consistency
075 between synchronous and asynchronous coding, as
076 well as the cohesiveness and correctness of themes
077 produced within each setting.

078 We evaluate three NLP-assisted inductive coding 117
079 tools built on distinct methodological foundations: 118
080 a human-in-the-loop topic modeling system (Fang 119
081 et al., 2023), a concept-driven thematic modeling 120
082 approach (Pacheco et al., 2023), and an LLM-based 121
083 system for evaluating and propagating code defini- 122
084 tions (Chew et al., 2023). To assess generalizabil- 123
085 ity across data characteristics, we test all tools and 124
086 collaboration settings on two markedly different 125
087 datasets: a corpus of short social media posts and a 126
088 collection of advertising texts. We further comple- 127
089 ment our quantitative analysis with a small-scale 128
090 user study to capture coder experiences and derive 129
091 design recommendations. 130

092 Our contributions are twofold: (1) we demon- 131
093 strate that the collaboration setting affects the re- 132
094 sults of NLP-assisted inductive coding tools across 133
095 diverse methodologies and datasets, and (2) we 134
096 provide an evaluation strategy that captures multi- 135
097 ple dimensions of analysis quality. Together, these 136
098 findings aim to inform the design and evaluation of 137
099 language technologies that better align with real- 138
100 world qualitative research workflows. 139

101 2 Related Work 140

102 The overarching goal of the systems we investi- 141
103 gate is to partially automate the qualitative coding 142
104 process either by inducing topics in an interactive, 143
105 semi-supervised manner (Fang et al., 2023; Smith 144
106 et al., 2018), by learning user-defined themes inter- 145
107 actively (Pacheco et al., 2023; Gao et al., 2023), or 146
108 by prompting LLMs with natural language defini- 147
109 tions of the observed themes (Chew et al., 2023; 148
110 Dai et al., 2023). A separate but related line of work 149
111 exemplified by Gao et al. (2024) uses LLMs to gener- 150
112 ate label recommendations as users perform the 151
113 coding process. While this system is explicitly de- 152
114 signed for asynchronous collaboration, the systems

we study differ in their ability to annotate large por-
tions of the dataset without extensive supervision.

Our research addresses a real-world use case
for qualitative researchers using HiTL systems and
is informed by the Human-Computer Interaction
(HCI) literature (Jiang et al., 2021; Feuston and
Brubaker, 2021; Chen et al., 2018). Prior HCI work
shows that coders place particular emphasis on
identifying and resolving ambiguity. In traditional
settings, this is supported by an independent close
reading of the data. However, in large-scale, NLP-
assisted analysis, coders have limited visibility into
where such ambiguities arise. Solutions have been
proposed to either visualize codes (Drouhard et al.,
2017) or rank document disagreement (Zade et al.,
2018) regardless of dataset size. In contrast, we
look at different output qualities of the resulting
themes and their assignments, combining signal
from group overlaps and relationships in the seman-
tic embedding space, and perform manual post-hoc
evaluations. This method of evaluation highlights
areas where the coders diverge both with each other
and with the model, providing another perspective
on the ambiguity question.

Previous evaluation methods introduced with
HiTL systems for qualitative coding have gener-
ally been ad hoc, with experiments conducted in
various group settings (Choo et al., 2013; Hoque
and Carenini, 2016; Smith et al., 2018), on individ-
ual participants (Rietz et al., 2020), and through
platforms such as MTurk (Zade et al., 2018). Our
contribution provides a standardized framework for
performing experiments regardless of the collabo-
ration modality, using a suite of metrics for evalu-
ating consistency, cohesiveness, and correctness in
experimental results.

151 3 Interactive Systems 151

152 We identify three categories of NLP techniques 152
153 used in interactive systems for qualitative coding 153
154 with large datasets: topic models, relational ap- 154
155 proaches, and LLMs. These techniques may be 155
156 applied in a variety of ways in interactive systems, 156
157 but we put special focus on their ability to help 157
158 code large datasets. To maximize coverage across 158
159 systems, we select a representative system from 159
160 each category to use in our experiments. In this 160
161 section, we briefly describe the unique aspects of 161
162 each category and introduce the selected system. 162

3.1 Topic Models

This category includes systems that use some variation of topic modeling to find emerging themes and facilitate document assignment. These systems benefit from the relative speed of the topic model, which allow users to quickly visualize and explore the dataset. Early exploration incorporated visualizations to help users adjust parameters (Chuang and McFarland, 2013), while later works implemented refinement operations that allow users to directly edit topic words and remove documents (Smith et al., 2018). However, topic modeling systems are limited by their lack of malleability and predictability. Refinement operations mostly edit topic words, which can have limited impact in the final results.

We select the HitL query-driven topic model (QDTM) system introduced by Fang et al. (2023). The topic model is initialized by providing input queries (i.e., words that represent concepts of interest for the user) which the model uses to generate the initial topics. In our experiments, users begin by iterating through each topic and naming them based on identified themes. Users then use a set of *refinement operations* to edit the topic model. They can merge and split topics based on topic words, add, remove, or reorder topic words, and remove documents from topics. The next iteration of the model is only produced when the users choose to apply refinements and the prior model is saved, allowing users to return to prior iterations to test different operations. Once satisfied with the state of the topic model, the user downloads the document distribution for that iteration.

To ensure comparable results, we use the same starting distribution of 13 topics for all our experiments using the same hyperparameters as Fang et al. (2021), which are $\alpha = 1.0$, $\beta = 0.5$, $\gamma = 1.5$. The QDTM also allows queries to be input prior to topic model initialization to produce partitions that follow prior knowledge, but we do not take advantage of this capability. The same initial topic model is provided for all experiments.

3.2 Relational Approaches

Relational approaches combine vector semantics and structured inference to model relationships between high-level concepts. Instead of treating themes as distributions over words (as topic models do), these frameworks define themes as distributions over generalized concepts. This reflects the

inductive coding process, where researchers identify patterns and concepts that are then synthesized into more abstract themes. However, their computational complexity grows with the number of dependencies considered, which hinders their ability to quickly adapt during coding sessions. Further, they rely on users to define informative concepts, making them less suited for inexperienced researchers.

We select the system introduced by Pacheco et al. (2023), which uses a two-stage relational framework. In the first stage, the system automatically partitions the dataset based on semantic similarity. The users explore each partition to identify themes, assign "good" and "bad" example documents for each theme, and input or correct supporting concepts for each example. In the second stage, the system uses the provided examples and concept relations to map the remaining dataset, only leaving documents unmapped if no theme is a sufficiently good match. The assignment procedure follows as a structured inference approach, where dependencies between concepts and themes are explicitly modeled. The unmapped documents are repartitioned as in the first stage and users are prompted to review unmapped partitions again. The process is iterated until all documents are mapped. More details about the framework and our experimental configuration can be found in App. C

3.3 Large Language Models

LLMs are ideal candidates for interactive systems, especially for tasks such as qualitative coding where the model can be prompted to produce themes or explanations without ad-hoc training (Kojima et al., 2024). They have been used for theme recommendation (Gao et al., 2023), for code conflict resolution (Gao et al., 2024), and for automated document assignment (Xiao et al., 2023). However, the flexibility of LLM outputs also leads to hallucinations, which are only partially addressed by prompt engineering. Models further suffer from biases in training which are difficult to identify and impact their ability to produce quality labels or recommendations (Chen et al., 2018). Additionally, their massive size is prohibitive when working with large datasets due to the high cost of inference.

We select the framework introduced by Chew et al. (2023). In their protocol, the human coder first manually codes a representative subset of the data and drafts definitions for each code. The

LLM is then prompted to label the data sample with the provided definitions. Agreement is calculated between human and model annotations using GWET’s AC_1 (Gwet, 2008). The prompt is then tweaked iteratively to achieve a satisfactory level of agreement, and the best-performing version is used to prompt the model to code the rest of the dataset. For our experiment, we select Llama 3.2 3B-Instruct as the base LLM for automated labeling. Additional details about our experimental settings can be found in App. E.

4 Study Design

To study the effects of different collaboration settings on the performance of the three selected systems, we design a protocol that can be used for both synchronous and asynchronous settings. For each system, we conduct three asynchronous experiments with one coder each and two synchronous experiments with three coders each for a total of 15 experiments. Evaluation metrics are calculated by comparing the resulting code-books within each experimental setting (e.g. the two code-books independently created by the two synchronous groups using the topic model). The rest of this section lays out the dataset, participant demographics and experimental protocol.

4.1 Datasets

We perform our experiments using two distinct datasets. The first consists of roughly 85,000 tweets about COVID-19 vaccines posted by users located in the US, uniformly distributed between Jan.-Oct. 2021 (Pacheco et al., 2022). The corpus also contains labels for vaccination stance (e.g. pro-vax, anti-vax) and morality frames (e.g. fairness/cheating and their actor/targets) (Roy et al., 2021), which are used as auxiliary concepts for the relational model. The second dataset consists of 5,471 climate related English language ads from the Facebook Ad Library (Islam et al., 2023). The ads focus on the US, were shown between Jan. 2021-Jan. 2022, and contain labels for climate change stance. The two datasets differ on a number of dimensions such as domain and purpose; a detailed analysis can be found in Appendix A.

4.2 Participants

We recruited a group of 33 researchers in NLP and Computational Social Science, 9 female and 24 male, between the ages of 20 and 45. This group

included professors at different levels of seniority, postdoctoral researchers, and graduate and undergraduate students from two different universities. This group covers the range of researchers likely to use interactive coding systems. All participants were either well-versed in qualitative data analysis, or were explicitly trained by senior researchers to perform the task. Due to the large number of experiments in our study, some participants took part in multiple experiments. These participants always performed the asynchronous experiment first to prevent external influence and took part in one synchronous experiment at most.

4.3 Coding Protocol

At the start of each experiment, participants were provided with a demonstration of all the operations in their respective systems. Every system starts with an initial partition of the data, so participants were instructed to read the first 25 samples in each partition, and manually create/name any themes they identified before freely exploring the rest of the dataset and start performing operations to find more themes.

In the topic model experiments, we suggested that participants merge and split topics based on their identified themes before making fine-grained refinements. They were then asked to refine the topic model based on their identified themes such that every topic corresponds to a unique theme. They kept re-running the model and making refinements until they were satisfied with the results, or until they failed to effect any meaningful changes.

In the relational system experiments, participants were tasked with selecting example documents for each identified theme, as well as determining concept relations for them. Following Pacheco et al. (2023), the supporting concepts considered were vaccination stances and morality frames (e.g., the identified theme “natural immunity” has an “anti-vax” stance, and is tied to the “purity” frame). Once participants were satisfied with their themes and selections, the system automatically coded the rest of the dataset. Unmapped examples were repartitioned and returned to the participants for a second (and last) round of coding.

In the experiments for the LLM-based system, participants produced natural-language definitions for each identified theme and selected a set of good examples for them. We then prompted the LLM with different task-prompt templates to find the

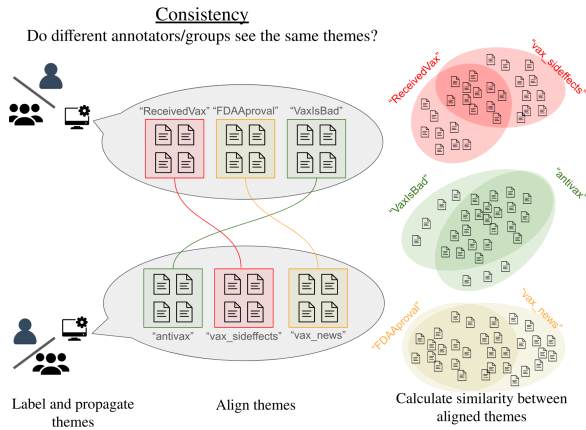


Figure 2: Two sets of coders use a particular HiTL system to find themes. Since the same theme can be named differently by different coders, we find the best match. In this example, the coder 1’s theme “VaxIsBad” has been matched with coder 2’s theme “antivax”. After aligning, we calculate the similarity between these two themes using methods like Jaccard Similarity or Centroid Distance.

best prompt for each set of participant-generated definitions, which was then used to code the rest of the dataset. Details of the templates, as well as the human-model agreement for the best template can be found in Appendix E.

5 Evaluation

We use both descriptive metrics and a user study to provide a comprehensive analysis on the differences when coding in synchronous and asynchronous settings. Our evaluation framework is comprised of three dimensions; **consistency**, **cohesiveness & distinctiveness**, and **correctness**, each of which uses metrics that are well-established in the literature (Ben-David and Ackerman, 2008; Hoyle et al., 2021; Pacheco et al., 2023).

5.1 Consistency

Coders risk overgeneralizing or overlooking key themes, leading to unsystematic results (Cornish et al., 2014). We address this by measuring consistency, defined as the extent to which different coders elicit the same themes from the same texts (Fig. 2). In semi-automated coding, consistency is nontrivial to assess: similarly named themes may cover different documents, while differently named themes may overlap substantially. We therefore measure consistency based on document overlap between themes. Specifically, we compute the maximum Jaccard similarity between each theme and

all themes produced by another coder, treating this maximum as the theme’s best alignment. Consistency is then calculated as the average similarity across all aligned theme pairs.

To account for semantically similar themes with differing document assignments, we also measure semantic consistency using S-BERT document embeddings (Reimers and Gurevych, 2019). We compute (1) centroid similarity, based on the cosine similarity between theme centroids, and (2) group average similarity, based on the average pairwise similarity across documents in two themes. As with Jaccard, we report maximum similarities per theme and average them for comparison across settings.

Table 1 reports average maximum Jaccard and embedding similarities across experiments. The high variance in Jaccard similarity across datasets, systems, and collaboration settings underscores the importance of semantics-based metrics. For the Covid dataset, synchronous groups produced more consistent themes for the topic modeling and relational systems, while the LLM-based system showed no statistically significant difference across collaboration settings. Results for the Climate dataset were less conclusive, likely due to greater semantic diversity and longer documents, which may encourage synchronous coders to surface more varied themes through discussion; as shown in the next section, these themes are also more cohesive.

Notably, the LLM-based system afforded the fewest opportunities for user intervention during coding. Unlike the other systems, which supported operations such as splitting and merging topics or defining relations between concepts, the LLM system only allowed users to edit theme definitions. We hypothesize that richer intervention mechanisms better enable coders to leverage the deliberation inherent to synchronous collaboration.

5.2 Cohesiveness and Distinctiveness

Another dimension for determining the systematicity and clarity of coding outcomes is by evaluating the similarities and differences between themes within the same code-book. We propose two metrics to measure this: cohesiveness and distinctiveness. A theme is said to be cohesive if its documents are similar to each other (measured by *intra*-theme similarity) and distinctive if it is dissimilar from documents in other themes within the same code-book (measured by *inter*-theme similarity). Intuitively, the purpose of grouping docu-

		Jaccard	Topic Model		Jaccard	Relational		Jaccard	LLM-Based	
			Centroid	Group Avg.		Centroid	Group Avg.		Centroid	Group Avg.
Covid	Sync	0.56(0.23)	0.98(0.05)**	0.52(0.10)	0.36(0.19)	0.98(0.01)*	0.52(0.07)*	0.14(0.08)	0.98(0.03)	0.44(0.03)
	Async	0.30(0.17)	0.96(0.05)**	0.51(0.09)	0.30(0.22)	0.94(0.07)*	0.44(0.10)*	0.17(0.11)	0.98(0.02)	0.45(0.03)
Climate	Sync	0.37(0.31)*	0.89(0.14)	0.43(0.14)	0.27(0.17)	0.95(0.05)**	0.43(0.08)	0.09(0.07)*	0.95(0.04)	0.37(0.07)
	Async	0.58(0.29)*	0.93(0.12)	0.46(0.16)	0.33(0.22)	0.94(0.09)**	0.42(0.08)	0.13(0.07)*	0.94(0.06)	0.36(0.06)

Table 1: **Avg. Consistency** between Best Theme Matches *across* Coder Groups. *Statistically significant using a two-sample unpaired t-test with $p < 0.05$. ** Near statistically significant with $p \approx 0.05$.

			Topic Model		Relational		LLM-based	
			Intra-Theme	Inter-Theme	Intra-Theme	Inter-Theme	Intra-Theme	Inter-Theme
Covid	All	Sync	0.52(0.10)	0.40(0.04)	0.51(0.08)*	0.42(0.05)*	0.44(0.06)	0.40(0.04)
		Async	0.52(0.10)	0.40(0.04)	0.45(0.10)*	0.34(0.11)*	0.43(0.05)	0.39(0.04)
	Top 25%	Sync	0.56(0.11)	0.39(0.05)	0.70(0.09)*	0.52(0.07)*	0.63(0.07)	0.55(0.05)
		Async	0.56(0.11)	0.39(0.05)	0.64(0.09)*	0.46(0.13)*	0.63(0.05)	0.54(0.05)
Climate	All	Sync	0.57(0.23)*	0.25(0.08)*	0.44(0.10)*	0.30(0.07)*	0.39(0.11)*	0.29(0.05)*
		Async	0.50(0.19)*	0.24(0.07)*	0.43(0.09)*	0.29(0.07)*	0.38(0.08)*	0.28(0.05)*
	Top 25%	Sync	0.72(0.21)*	0.26(0.09)*	0.66(0.11)*	0.39(0.10)*	0.63(0.11)*	0.40(0.09)*
		Async	0.65(0.20)*	0.26(0.08)*	0.65(0.10)*	0.39(0.11)*	0.62(0.12)*	0.40(0.09)*

Table 2: Group Avg. Similarity *within* Coder Groups. Themes are considered to be more **cohesive** if intra-theme similarity is high and more **distinctive** if inter-theme similarity is low. *Statistically significant using a two-sample unpaired t-test with $p < 0.05$.

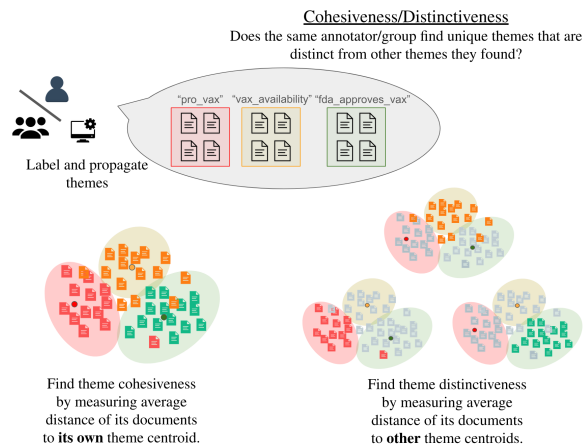


Figure 3: Once a coder has identified themes and they have been propagated the full dataset, we calculate *intra-theme similarity* by measuring the avg. of the pairwise distances between each document within a theme (left). We calculate *inter-theme similarity* by measuring the avg. of pairwise distances between each document in a theme and documents assigned to all other themes (right)

ments by theme is to create abstract representations of a dataset, where each theme represents a distinct facet of the data. If themes are not cohesive and distinctive, then it becomes hard to tell which theme a given document should belong to and the code-book falls apart.

Figure 3 shows how to evaluate these metrics for a single coder (or coder group). We calculate

both the intra-theme similarity and the inter-theme similarity for all the themes in the code-book. *Intra-theme similarity* is calculated by taking the average of pair-wise similarity between all documents of the same theme. *Inter-theme similarity* for a given theme is calculated by taking the average pair-wise similarity of documents in that theme with documents in all other themes.

A confounding factor in these measures is that all systems provide broad coverage of documents such that even distantly related documents may be assigned to a theme. To more accurately represent the cohesiveness and distinctiveness of the themes in each experiment, we perform the same calculations on a subset comprised of only the top 25% of documents most closely related to each theme. For the relational and LLM-based systems, this top quartile is selected using the distance from the centroid. For the interactive topic model, we use the weights assigned by the model.

Table 2 shows results for both the whole dataset as well as the subset of the documents closest to each theme. Overall, we find that the intra-theme similarities are always higher than inter-theme similarities, which means that themes are at least moderately cohesive and distinctive across the board. For the Covid dataset, we find that themes may be more cohesive but not more distinctive in the synchronous setting especially for the relational system. This may be due to the homogeneity of the

dataset, which further explains the uniformity of results in the topic model and LLM experiments.

The results from the Climate dataset present a more compelling finding, where the synchronous experiments uniformly produce much more cohesive themes. When only the top 25% of documents are considered, the increase in cohesiveness is not correlated with a decrease in distinctiveness, unlike in the Covid case. Our results strongly suggest that synchronous collaboration facilitates deeper analysis to uncover clearer themes in the data. This is further supported by the greater distinctiveness in the Climate data experiments: we expect coders to find more distinct themes in a more complex and heterogenous dataset.

5.3 Correctness

Interactive systems allow users to automate large portions of the coding process at the risk of producing inaccurate theme assignments. To estimate how correct the outputs of each system are, we conduct a post-hoc analysis by manually checking a randomly selected sample of 1,200 document-theme pairs (200 per experimental setting). To ensure that our sample is representative of the overall dataset for each experimental setting, we split the data into quartiles based on document relatedness to each theme and select a uniform sample of themes and relatedness scores. As before, relatedness is calculated using the theme weight distribution for the interactive topic model and distance from the centroid for the other two systems. To assess reliability, each assignment is evaluated by two annotators, with a third weighing in for tie breaks. We also find that human evaluators have moderate-to-high agreement when assessing system outputs (with an overall Krippendorff’s α of 0.632), suggesting that we can trust these estimations.

Figure 4 shows the correctness results for each quartile sample per experiment in the Covid dataset. First, we observe that the relational system is not only the most accurate, but it shows negligible correctness differences in synchronous vs. asynchronous configurations. This an encouraging result, given that this system took the most advantage of synchronous deliberation based on our other metrics of quality. The other two systems showed marked differences, with the topic modeling approach producing more accurate assignments in asynchronous operation, and the LLM producing more accurate assignment in a synchronous

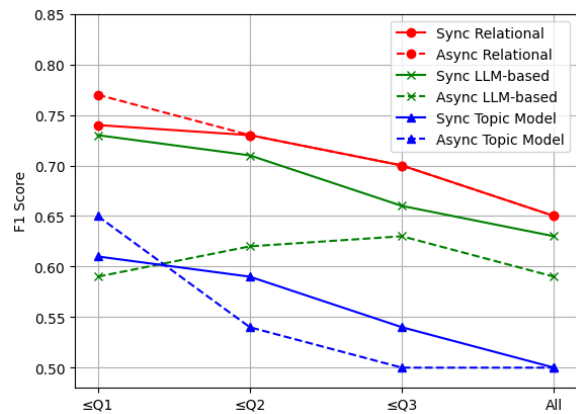


Figure 4: Correctness w.r.t. distance from theme in Covid data experiments.

paradigm. This observation leads us to believe that interactive systems that introduce less and weaker inductive biases from interaction are more sensitive to choices of configuration – and researchers should be aware of these possibilities when designing experiments. Second, we find that the correctness of different approaches tapers differently, based on the distance of the example from the theme. In other words, as the data samples resemble less like the human labeled ones, some models handle it more robustly than others. In our observation, relational approaches and LLM methods outperform topic models when it comes to assigning themes to distant examples.

5.4 Dataset Differences

To what extent do structural and semantic characteristics of a dataset affect the coding process? We carried out a comprehensive analysis across multiple dimensions, including readability, templatic or formulaic patterns, semantic theme ambiguity, and semantic similarity distribution (details in App. A). The higher formulaic content in climate ads (20.1%) may facilitate easy pattern recognition, but could also induce coder fatigue through repetitive content. The greater linguistic complexity (lower Flesch scores, longer sentences) and semantic diversity (lower pairwise similarity) of climate ads suggest higher cognitive load per coding decision.

Conversely, the semantic homogeneity of COVID tweets, evidenced by higher pairwise similarity and centroid proximity, may enable more efficient real-time discussion in synchronous coding settings, as coders share common semantic reference points. The less formulaic nature of the COVID discourse may also sustain coder engage-

ment through content variety. This may explain why the average consistency scores tends to favor synchronous groups for COVID discourse, while for Climate data, asynchronous coders tend to either slightly outperform or are at par with synchronous groups.

These findings suggest that coding protocol selection (synchronous vs. asynchronous) should account for corpus-specific characteristics: semantically diverse, linguistically complex corpora may benefit from asynchronous approaches allowing extended reflection, while homogeneous corpora may be efficiently processed through synchronous collaborative coding.

6 User Study and Recommendations

We conducted semi-structured interviews to understand the participant experiences with the task and tools, with a focus on synchronous versus asynchronous coding (see App. D for interview script). Several themes emerged.

Synchronous teamwork eased coding and improved outcomes. Participants reported that working synchronously helped them contextualize data, resolve disagreements quickly, and “break ties” through discussion. These experiences align with our quantitative findings showing higher consistency and cohesiveness in synchronous settings, suggesting the value of systems that explicitly support real-time deliberation.

Asynchronous coders were more sensitive to tool limitations. Because they worked largely in isolation, asynchronous participants focused more on usability issues and tool constraints, highlighting the need to improve support for independent coding workflows.

Limited control reduced user trust, particularly in topic modeling. Participants reported a loss of agency when using the topic modeling system, citing insufficient control over operations and difficulty tracking theme evolution – *‘the merge process did not offer the ideal amount of control and made it difficult to keep track of the theme groups’*. Initial topics also sometimes conflated opposing themes due to lexical similarity, frustrating users’ attempts to refine results. One participant commented – *“Many Anti-Vax and Pro-Vax standpoints use the same words/phrases in their tweets, which the Fang et al. (2023) model groups together despite the stark difference in message between the two.”* While some appreciated the model’s initial

theme induction, participants desired greater control and clearer explanations of cluster structure. These findings highlight the need for systems to maximize the degree of control afforded to users.

LLM-based coding was costly and unreliable at scale. Despite strong reasoning capabilities, LLMs proved inefficient for large-scale annotation due to high computational costs and inconsistent classification performance. This opens an opportunity for NLP researchers to make LLMs more reliable inductive reasoners and to come up with prompting strategies that can allow LLMs to reliably classify documents in bulk, especially when working at scale.

Overall, our findings suggest that interactive systems should balance automation with user control and provide tailored workflows that support both real-time collaborative deliberation and independent coding at scale.

7 Conclusion and Future Work

We examined three categories of NLP-assisted qualitative research tools in different collaboration modalities, and conducted inductive and deductive coding on two very different datasets of English texts. We designed an evaluation framework that describes the quality of the induced themes and their resulting document assignments under synchronous and asynchronous collaboration. We find that the collaboration modality is a significant factor in determining if the quality of a system’s output. This is particularly true for systems where users have a wider range of interaction that can benefit from group consensus. We also observe that solutions based on topic modeling, although popular in the data analysis literature, can struggle with inducing cohesive themes and accurate code assignments. Finally, we show that while LLM-based solutions show promise, they pose significant challenges when it comes to coding at scale.

While this study focuses on collaboration modalities, there are numerous other variables that can affect a tool’s efficacy for qualitative coding. We believe that our proposed evaluation framework can be repurposed and expanded to evaluate a wide range of interventions, such as the underlying NLP technology, the interactive interface, the expertise of the coders, and the type of data being annotated. We hope to inform the development of more robust evaluations of NLP tools for qualitative research in realistic settings.

662 Limitations

663 The study presented in this paper has two main
664 limitations.

665 (1) While we selected three distinct, representa-
666 tive tools to perform our analysis of synchronous
667 vs. asynchronous settings, as well as two datasets
668 with distinct characteristics, the list is of course non-
669 exhaustive. A larger study incorporating more tools
670 and datasets could yield additional insights.

671 (2) While we look at an important variable in
672 qualitative research settings (collaboration modal-
673 ity), there are several other variables that can influ-
674 ence the outcome of NLP-assisted solutions (e.g.,
675 choice of tool, expertise and live experience of an-
676 notators, type of data being annotated, etc.). In
677 addition to this, we did not explore the many dif-
678 ferent consolidation strategies that are often used
679 to bring together the perspectives of asynchronous
680 coders. We leave the explorations of these ques-
681 tions for future work.

682 Ethical Considerations

683 To the best of our knowledge, no code of ethics was
684 violated during the development of this project. We
685 used publicly available tools and datasets according
686 to their licensing agreements. For our annotation
687 experiments, we followed IRB protocol and did not
688 retain any personally identifiable information.

689 All information needed to replicate our experi-
690 ments is presented in the paper. We reported all ex-
691 perimental settings, as well as any pre-processing
692 steps, learning configurations, hyper-parameters,
693 and additional technical details. Due to space con-
694 straints, some of this information was relegated to
695 the Appendix. In addition to this, we will make the
696 results of the annotation experiment available to
697 the community, as well as the code to produce all
698 of our reported results..

699 References

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

705 Shai Ben-David and Margareta Ackerman. 2008. [Mea-](#)
706 [sures of clustering quality: A working set of axioms](#)
707 [for clustering.](#) In *Advances in Neural Information*
708 *Processing Systems*, volume 21. Curran Associates,
709 Inc.

Henry E. Brady. 2019. [The challenge of big data and](#)
[data science.](#) *Annual Review of Political Science*,
22(1):297–323. 710
711
712

Virginia Braun and Victoria Clarke. 2006. [Using the-](#)
[matic analysis in psychology.](#) *Qualitative Research*
in Psychology, 3:77–101. 713
714
715

Nan-Chen Chen, Margaret Drouhard, Rafal Kocielnik,
Jina Suh, and Cecilia R. Aragon. 2018. [Using ma-](#)
[chine learning to support qualitative coding in social](#)
[science: Shifting the focus to ambiguity.](#) *ACM Trans.*
Interact. Intell. Syst., 8(2). 716
717
718
719
720

Robert Chew, John Bollenbacher, Michael Wenger, Jes-
sica Speer, and Annice Kim. 2023. [Llm-assisted con-](#)
[tent analysis: Using large language models to support](#)
[deductive coding.](#) *Preprint*, arXiv:2306.14924. 721
722
723
724

Jaegul Choo, Changhyun Lee, Chandan K. Reddy, and
Haesun Park. 2013. [Utopian: User-driven topic mo-](#)
[deling based on interactive nonnegative matrix fac-](#)
[torization.](#) *IEEE Transactions on Visualization and*
Computer Graphics, 19(12):1992–2001. 725
726
727
728
729

Jason Chuang and Daniel A. McFarland. 2013. [Docu-](#)
[ment exploration with topic modeling : Designing](#)
[interactive visualizations to support effective analysis](#)
[workflows.](#) 730
731
732
733

Flora Cornish, Alex Gillespie, and Tania Zittoun. 2014.
The SAGE Handbook of Qualitative Data Analysis.
SAGE Publications Ltd, London; London. 734
735
736

Shih-Chieh Dai, Aiping Xiong, and Lun-Wei Ku. 2023.
[LLM-in-the-loop: Leveraging large language model](#)
[for thematic analysis.](#) In *Findings of the Association*
for Computational Linguistics: EMNLP 2023, pages
9993–10001, Singapore. Association for Computa-
tional Linguistics. 737
738
739
740
741
742

Margaret Drouhard, Nan-Chen Chen, Jina Suh, Rafal
Kocielnik, Vanessa Peña-Araya, Keting Cen, Xiangyi
Zheng, and Cecilia R. Aragon. 2017. [Aeonium: Vi-](#)
[sual analytics to support collaborative qualitative cod-](#)
[ing.](#) In *2017 IEEE Pacific Visualization Symposium*
(PacificVis), pages 220–229. 743
744
745
746
747
748

Zheng Fang, Lama Alqazlan, Du Liu, Yulan He, and
Rob Procter. 2023. [A user-centered, interactive,](#)
[human-in-the-loop topic modelling system.](#) In *Pro-*
ceedings of the 17th Conference of the European
Chapter of the Association for Computational Lin-
guistics, pages 505–522, Dubrovnik, Croatia. Associ-
ation for Computational Linguistics. 749
750
751
752
753
754
755

Zheng Fang, Yulan He, and Rob Procter. 2021. [A query-](#)
[driven topic model.](#) In *Findings of the Association*
for Computational Linguistics: ACL-IJCNLP 2021,
pages 1764–1777, Online. Association for Computa-
tional Linguistics. 756
757
758
759
760

Jessica L. Feuston and Jed R. Brubaker. 2021. [Putting](#)
[tools in their place: The role of time and perspective](#)
[in human-ai collaboration for qualitative analysis.](#)
Proc. ACM Hum.-Comput. Interact., 5(CSCW2). 761
762
763
764

765	Uwe Flick. 2014. The sage handbook of qualitative data analysis .	
766		
767	Jie Gao, Kenny Tsu Wei Choo, Junming Cao, Roy Ka-Wei Lee, and Simon Perrault. 2023. Coaicoder: Examining the effectiveness of ai-assisted human-to-human collaboration in qualitative analysis . <i>ACM Trans. Comput.-Hum. Interact.</i> , 31(1).	
768		
769		
770		
771		
772	Jie Gao, Yuchen Guo, Gionnieve Lim, Tianqin Zhang, Zheng Zhang, Toby Jia-Jun Li, and Simon Tangi Perrault. 2024. Collabcoder: A lower-barrier, rigorous workflow for inductive collaborative qualitative analysis with large language models . <i>Preprint</i> , arXiv:2304.07366.	
773		
774		
775		
776		
777		
778	Barney G Glaser, Anselm L Strauss, and Elizabeth Strutzel. 1968. The discovery of grounded theory; strategies for qualitative research. <i>Nursing research</i> , 17(4):364.	
779		
780		
781		
782	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, and 542 others. 2024. The llama 3 herd of models . <i>Preprint</i> , arXiv:2407.21783.	
783		
784		
785		
786		
787		
788		
789		
790	Kilem Li Gwet. 2008. Computing inter-rater reliability and its variance in the presence of high agreement. <i>British Journal of Mathematical and Statistical Psychology</i> , 61(1):29–48.	
791		
792		
793		
794	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: A methodological catalyzer for a maturing discipline . <i>International Journal of Communication</i> , 13(0).	
795		
796		
797		
798		
799		
800		
801		
802		
803		
804	Enamul Hoque and Giuseppe Carenini. 2016. Interactive topic modeling for exploring asynchronous online conversations: Design and evaluation of convisit . <i>ACM Trans. Interact. Intell. Syst.</i> , 6(1).	
805		
806		
807		
808	Alexander Hoyle, Pranav Goel, Denis Peskov, Andrew Hian-Cheong, Jordan Boyd-Graber, and Philip Resnik. 2021. Is automated topic model evaluation broken? the incoherence of coherence. In <i>Proceedings of the 35th International Conference on Neural Information Processing Systems, NIPS '21</i> , Red Hook, NY, USA. Curran Associates Inc.	
809		
810		
811		
812		
813		
814		
815	Tunazzina Islam, Ruqi Zhang, and Dan Goldwasser. 2023. Analysis of climate campaigns on social media using bayesian model averaging . In <i>Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society, AIES '23</i> , page 15–25, New York, NY, USA. Association for Computing Machinery.	
816		
817		
818		
819		
820		
	Jialun Aaron Jiang, Kandrea Wade, Casey Fiesler, and Jed R. Brubaker. 2021. Supporting serendipity: Opportunities and challenges for human-ai collaboration in qualitative analysis . <i>Proc. ACM Hum.-Comput. Interact.</i> , 5(CSCW1).	821
		822
		823
		824
		825
	Xin Jin and Jiawei Han. 2010. <i>K-Means Clustering</i> , pages 563–564. Springer US, Boston, MA.	826
		827
	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2024. Large language models are zero-shot reasoners. In <i>Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS '22</i> , Red Hook, NY, USA. Curran Associates Inc.	828
		829
		830
		831
		832
		833
	Julia Mendelsohn, Ceren Budak, and David Jurgens. 2021. Modeling framing in immigration discourse on social media . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 2219–2263, Online. Association for Computational Linguistics.	834
		835
		836
		837
		838
		839
		840
	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 <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 5821–5839, Seattle, United States. Association for Computational Linguistics.	841
		842
		843
		844
		845
		846
		847
		848
		849
	Maria Leonor Pacheco, Tunazzina Islam, Lyle Ungar, Ming Yin, and Dan Goldwasser. 2023. Interactive concept learning for uncovering latent themes in large text collections . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 5059–5080, Toronto, Canada. Association for Computational Linguistics.	850
		851
		852
		853
		854
		855
		856
	Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERT-networks . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.	857
		858
		859
		860
		861
		862
		863
		864
	Tim Rietz, Peyman Toreini, and Alexander Maedche. 2020. Cody: An interactive machine learning system for qualitative coding . In <i>Adjunct Proceedings of the 33rd Annual ACM Symposium on User Interface Software and Technology, UIST '20 Adjunct</i> , page 90–92, New York, NY, USA. Association for Computing Machinery.	865
		866
		867
		868
		869
		870
		871
	Shamik Roy, Maria Leonor Pacheco, and Dan Goldwasser. 2021. Identifying morality frames in political tweets using relational learning . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 9939–9958, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	872
		873
		874
		875
		876
		877
		878

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 *Proceedings of the 23rd International Conference on Intelligent User Interfaces, IUI '18*, page 293–304, New York, NY, USA. Association for Computing Machinery.

Laurens van der Maaten and Geoffrey Hinton. 2008. [Visualizing data using t-sne](#). *Journal of Machine Learning Research*, 9(86):2579–2605.

Ziang Xiao, Xingdi Yuan, Q. Vera Liao, Rania Abdelghani, and Pierre-Yves Oudeyer. 2023. [Supporting qualitative analysis with large language models: Combining codebook with gpt-3 for deductive coding](#). In *Companion Proceedings of the 28th International Conference on Intelligent User Interfaces, IUI '23 Companion*, page 75–78, New York, NY, USA. Association for Computing Machinery.

Himanshu Zade, Margaret Drouhard, Bonnie Chinh, Lu Gan, and Cecilia Aragon. 2018. [Conceptualizing disagreement in qualitative coding](#). In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems, CHI '18*, page 1–11, New York, NY, USA. Association for Computing Machinery.

A Dataset Analysis

The two corpora exhibit substantial structural differences as seen in Table 3. The climate dataset comprises 5,471 documents with a mean length of 53.1 words, while the COVID dataset contains 85,799 documents averaging 27.5 words. This approximately two-fold difference in document length is statistically significant (Mann-Whitney U, $p < 0.001$, effect size $r = -0.353$). Climate advertisements also demonstrate greater variability in length ($CV = 1.47$) compared to the more constrained COVID tweets ($CV = 0.50$). Lexical diversity metrics reveal nuanced differences. While type-token ratios are comparable (0.060 vs. 0.051), the Measure of Textual Lexical Diversity (MTLD) indicates higher lexical sophistication in COVID tweets (119.6 vs. 95.6).

Readability analyses indicate that climate advertisements present greater cognitive demands. The Flesch Reading Ease score for climate ads ($M = 47.4$) falls within the “difficult” range, whereas COVID tweets score higher ($M = 54.9$), indicating moderately easier comprehension ($p < 0.001$, $r = 0.211$). Correspondingly, climate ads require higher grade-level reading ability (Flesch-Kincaid Grade: 10.7 vs. 9.2; Gunning Fog Index: 13.2 vs. 11.2). See Table 4.

We also performed an analysis over a sample of 5000 documents from both datasets to identify

Metric	Covid	Climate
Number of Documents	85799	5471
Words per Document	27.50	53.10
Sentences per Document	2.50	3.40
Type-Token Ratio	0.051	0.060
Measure of Textual		
Lexical Diversity (threshold = 0.72)	119.65	95.58

Table 3: Corpus-level statistics for the Covid and Climate datasets.

Metric	Covid	Climate
Flesch Reading Ease Score	54.85	47.38
Flesch Kincaid Grade	9.18	10.68
Gunning Fog Index	11.20	13.22
Automated Readability Index	12.34	11.85

Table 4: Measures of Linguistic Complexity of the Covid and Climate datasets.

the extent of formulaic or template centric content. DBSCAN clustering on TF-IDF similarity matrices (threshold=0.8) identified 137 template clusters in climate ads, encompassing 20.1% of documents, compared to only 31 clusters (3.7%) in COVID tweets. The proportion of high-similarity document pairs (≥ 0.8 cosine similarity) is an order of magnitude greater in climate ads (0.059% vs. 0.006%). See Table 5.

Characteristic n-grams reflect domain-specific discourse patterns. Climate ads prominently feature phrases such as “oil natural gas”, “clean energy jobs”, and “fight climate change”, while COVID tweets center on vaccine-related language (“covid 19 vaccine”, “getting covid vaccine”). These patterns suggest that climate advertising employs stan-

Metric	Covid	Climate
Template Clusters	31	137
Percentage of Documents in Template Clusters	3.70	20.10
Percentage of High Similarity Document Pairs	0.006	0.059
Average Cluster Size	5.96	7.34
Largest Cluster Size	31	75

Table 5: Template patterns across Covid and Climate datasets.

Covid	Climate
'covid 19 vaccine', 'getting covid vaccine', 'covid vaccine https', 'got covid vaccine', 'covid 19 vaccines'	'https bit ly', 'oil natural gas', 'clean energy jobs', 'fight climate change', 'oil gas industry'

Table 6: Top n-gram phrases in Covid and Climate datasets.

Metric	Covid	Climate
Average Entropy across Documents	0.042	0.023
Average Confidence in Primary Cluster Assignment	0.989	0.994
Percentage of Documents Where Top 2 Clusters are within 0.2 Probability	0.660	0.280
Silhouette Score	0.022	0.018

Table 7: Semantic theme ambiguity for Covid and Climate datasets.

948 dardized messaging templates, whereas COVID
949 discourse largely centers around vaccines. See Ta-
950 ble 6.

951 Gaussian Mixture Model analysis with SBERT
952 embeddings (Reimers and Gurevych, 2019) reveals
953 that both corpora exhibit high thematic clarity,
954 with mean maximum cluster probabilities exceed-
955 ing 0.98. However, COVID tweets demonstrate
956 marginally higher semantic entropy ($M=0.042$ vs.
957 0.023), indicating slightly greater theme ambiguity.
958 The proportion of documents spanning multiple
959 clusters (probability difference < 0.2 between top-2
960 clusters) is higher for COVID tweets (0.66% vs.
961 0.28%). Silhouette scores are low for both corpora
962 (0.022 vs. 0.018), suggesting that while documents
963 cluster clearly into dominant themes, inter-cluster
964 boundaries are not sharply defined. See Table 7.

965 Mean pairwise cosine similarity in the SBERT
966 embedding space is substantially higher for COVID
967 tweets (0.429 vs. 0.261), indicating greater se-
968 mantic homogeneity. The inter-quartile range con-
969 firms this pattern: COVID tweets exhibit similarity
970 scores between $0.365 - 0.518$, while climate ads
971 range from $0.167 - 0.348$. Notably, 62.6% of cli-
972 mate ad pairs fall below 0.3 similarity, compared
973 to only 14.1% of COVID tweet pairs. Document-
974 to-centroid similarity further corroborates this find-
975 ing (0.654 vs. 0.510), demonstrating that COVID
976 tweets cluster more tightly around their corpus cen-
977 troid. See Table 8.

Metric	Covid	Climate
Average Pairwise Similarity between documents	0.43	0.26
25th Percentile Pairwise Similarity Score	0.37	0.17
75th Percentile Pairwise Similarity Score	0.52	0.35
Average Similarity of documents with the centroid	0.65	0.51
Percentage of Pairs with Similarity > 0.7	0.70	0.20
Percentage of Pairs with Similarity > 0.3	14.14	62.55

Table 8: Semantic Similarity Distribution for Covid and Climate datasets.

B Topic Model Experimental Settings

978

979 To ensure comparable results, we use the same
980 starting distribution of 13 topics for all our exper-
981 iments using the same hyperparameters as Fang
982 et al. (2021), which are $\alpha = 1.0$, $\beta = 0.5$, $\gamma = 1.5$.
983 The QDTM also allows queries to be input prior
984 to topic model initialization to produce partitions
985 that follow prior knowledge, but we do not take
986 advantage of this capability. The same initial topic
987 model is provided for all experiments.

C Experimental Settings and Interactive System Details

988
989

990 During the interactive coding process, researchers
991 are provided initial clusters to identify themes. We
992 use $K = 10$ means clustering to generate initial
993 partitions and the same partitions are provided for
994 all experiments. Users select positive and negative
995 examples for each theme, which are used to create a
996 distributional representations to calculate semantic
997 similarity with unlabeled documents. Simultane-
998 ously, supporting concepts are defined using logical
999 rules which researchers can use to label example
1000 documents. The relationship between themes and
1001 concepts are then used as a structured inference
1002 task to predict label assignments. By using mul-
1003 tiple sources of information, the model can much
1004 more efficiently extrapolate from a small set of
1005 manually labeled data.

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 (Mendelsohn et al., 2021) 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 9: Interactive Operations for the Pacheco et al. (2023) System

C.1 Relational Approach Operational Details

D Semi-Structured Interview

D.1 Interviewing

We administered interviews after annotation sessions. Asynchronous annotators were asked questions individually about their experience, whereas synchronous annotator groups were asked questions with their fellow annotators.

D.2 Script

1. Have you worked on annotation projects before? Did these annotation projects use qualitative coding strategies (ex: grounded theory)? How experienced are you as an annotator?
2. How was your experience on the COVID-19 vaccine annotation session we conducted on

Sunday? Particularly, we are interested in your thoughts and feelings over the session.

3. You annotated in a group, working together as a team. Did you find this setup to be beneficial? What were some of the limitations you faced, both individually and as a group, when working synchronously?
4. On a similar line, what would you consider to be the pros and cons if you were to annotate alone?

The last questions would be flipped based on if we are posing it to synchronous or asynchronous annotators.

E LLM-based Experimental Configuration and Prompt Details

For this study, we use the same starting partitions as used in the relational approach experiment in App. C. The original work uses existing codes from a theoretical framework whereas we use codes defined by the user during this step, but this does not affect the overall process.

code-book	Gwet’s AC_1	# Unlabeled Docs
Sync 1	0.42	5,548(6.5%)
Sync 2	0.49	9,766(11.4%)
Async 1	0.61	611(0.7%)
Async 2	0.62	2,506(2.9%)
Async 3	0.46	13,816(16.1%)

Table 10: Results for the selected prompt for each coding session using the LLM-based system. Gwet’s AC_1 is used to select the best prompt for running the full dataset. The number of unlabeled documents represent documents where the LLM produced a label not created by human annotators after running the full dataset (percentage of the dataset unlabeled).

Using the LLM, we generate 3 additional templates based on a prompt structure provided in the original work. We used the Llama 3.2 3B-Instruct (Grattafiori et al., 2024) model for all generation tasks:

E.1 LLM Hyperparameters

```

Batch size: 32
Model: Llama 3.2 3B-Instruct
GPU: A100 40GB VRAM
Average Compute Time: 24hrs per
job
Number of jobs: 5

```

E.2 LLM Prompts

```

To code this tweet, do the fol-
lowing:
- First, read the codebook and
the tweet.
- Next, decide which code is most
applicable and explain your rea-
soning for the coding decision.
- Finally, generate json with
your code and your reason for the
coding decision. The response
MUST be formatted as JSON.
Codes:
-
<codes>
-
Codebook:
-
<codebook>
-
Tweet:
-
<tweet>
-
JSON Output:
-
"code" : "",
"reason" : ""
-

```

```

To code this tweet, do the fol-
lowing:

First, read the codebook and the
tweet.
Next, decide which code is most
applicable based on the tweet's
content and explain your reason-
ing for the coding decision.
Finally, generate a JSON object
with the selected code and pro-
vide a brief explanation for your
coding decision.
The response MUST be formatted as
JSON.
Codebook: Themes: <"theme":
"definition">
Tweet: < "text": "<text>" >
JSON Output: < "code": "",
"reason": "" >

```

To generate code for this tweet, provide a step-by-step explanation of how to approach the task.

First, analyze the tweet's content and identify key concepts, such as the type of object or class being described, any specific behaviors or requirements, and relevant keywords.

Next, evaluate the codebook options and determine which one is most applicable. Explain your reasoning for your decision, including any similarities between the tweet and the code definitions, or any specific requirements mentioned in the tweet that align with a particular code.

Finally, generate a JSON object with the selected code and provide additional context, including:

A clear explanation of how you arrived at your chosen code

Any relevant notes or comments about the code's functionality and requirements

A brief comparison to other codes in the book, if applicable

The response MUST be formatted as JSON.

Codebook: <codebook>

Tweet: <tweet>

JSON Output: < "code": "", "reasoning": "", "context": "" >

To analyze this tweet and select a relevant theme, follow these steps:

First, read the tweet and identify key concepts, such as emotions, objects, or ideas mentioned in the text.

Next, evaluate the theme options and determine which one is most applicable. Explain your reasoning for your decision, including any connections you see between the tweet's content and the theme definitions.

Then, generate a JSON object with the selected theme and provide additional insight into your analysis. Include:

A clear explanation of how you arrived at your chosen theme

Any specific characteristics or keywords from the tweet that support your decision

A brief comparison to other themes, if applicable

The response MUST be formatted as JSON.

Themes: <Codebook>

Tweet: <tweet>

JSON Output: <"theme": "", "insight": "">

1053

1052

F Use of Pre-Existing Artifacts

1054

All pre-existing artifacts utilized in this study, including datasets, software libraries, models, and computational tools, are publicly available under open-source and open-access licenses. This academic work adheres to all intended use guidelines and terms of service for the respective resources.

1055

1056

1057

1058

1059

1060