

Generative AI for thematic analysis in a maternal health study: coding semistructured interviews using large language models

Shan Qiao^{1,6}  | Xingyu Fang^{2,3,4} | Junbo Wang⁵ |
Ran Zhang¹ | Xiaoming Li^{1,6,7} | Yuhao Kang^{2,3} 

¹Department of Health Promotion, Education, and Behavior, Arnold School of Public Health, University of South Carolina, Columbia, South Carolina, USA

²GISense Lab, Department of Geography and the Environment, The University of Texas at Austin, Austin, Texas, USA

³Department of Geography, University of South Carolina, Columbia, South Carolina, USA

⁴School of Civil, Environmental and Geomatic Engineering, University College London, London, UK

⁵Department of Geography and Sustainability, University of Tennessee, Knoxville, Tennessee, USA

⁶South Carolina SmartState Center for Healthcare Quality (CHQ), University of South Carolina, Columbia, South Carolina, USA

⁷Big Data Health Science Center (BDHSC), University of South Carolina, Columbia, South Carolina, USA

Correspondence

Yuhao Kang, GISense Lab, Department of Geography and the Environment, The University of Texas at Austin, Austin, Texas, USA.

Email: yuhao.kang@austin.utexas.edu

Abstract

Study Objectives: The coding of semistructured interview transcripts is a critical step for thematic analysis of qualitative data. However, the coding process is often labor-intensive and time-consuming. The emergence of generative artificial intelligence (GenAI) presents new opportunities to enhance the efficiency of qualitative coding. This study proposed a computational pipeline using GenAI to automatically extract themes from interview transcripts. **Methods:** Using transcripts from interviews conducted with maternity care providers in South Carolina, we leveraged ChatGPT for inductive coding to generate codes from interview transcripts without a predetermined coding scheme. Structured prompts were designed to instruct ChatGPT to generate and summarize codes. The performance of GenAI was evaluated by comparing the AI-generated codes with those generated manually. **Results:** GenAI demonstrated promise in detecting and summarizing codes from interview transcripts. ChatGPT exhibited an overall accuracy exceeding 80% in inductive coding. More impressively, GenAI reduced the time required for coding by 81%. **Discussion:** GenAI models are capable of efficiently processing

This is an open access article under the terms of the [Creative Commons Attribution-NonCommercial](https://creativecommons.org/licenses/by-nc/4.0/) License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited and is not used for commercial purposes.

© 2025 The Author(s). Applied Psychology: Health and Well-Being published by John Wiley & Sons Ltd on behalf of International Association of Applied Psychology.

Funding information

Population Research Center Seed Grant;
NIH, Grant/Award Numbers:
R01AI127203-05S2, R01AI174892

language datasets and performing multi-level semantic identification. However, challenges such as inaccuracy, systematic biases, and privacy concerns must be acknowledged and addressed. Future research should focus on refining these models to enhance reliability and address inherent limitations associated with their application in qualitative research.

KEYWORDS

Generative AI, Coding, Maternal health, Inductive coding, Thematic analysis

INTRODUCTION

Qualitative research commonly involves collecting and analyzing nonnumerical data (e.g. text, video, and audio) to understand concepts, perceptions, beliefs, attitudes, opinions, and lived experiences within specific contexts (Kuckartz, 2014; Mey, 2023). Qualitative studies have been widely used in social sciences and public health research to generate in-depth insights, discover latent patterns, and inform the development of new research directions (Lichtman, 2013; Mohajan, 2018; Tracy, 2024). In mixed-methods research, qualitative methods serve both exploratory and explanatory functions, uncovering salient constructs for quantitative measurement and providing detailed interpretations of quantitative results (Guest et al., 2015; Moser & Korstjens, 2018; Nardi, 2018). The choice of a qualitative research approach by researchers depends on the theoretical tradition and research paradigm that they adhere to.

As a form of inquiry, qualitative research is diverse and consists of many research paradigms that shape the assumptions a researcher adheres to when answering their research question(s) of interest (Ponterotto, 2005). For example, postpositivism as a research paradigm adheres to the worldview that even though there is a true reality, the ability to know everything about it is impossible, as there is always the possibility of new knowledge arising that disproves existing knowledge (Lincoln et al., 2011; Ponterotto, 2005). Epistemologically, postpositivists believe that knowledge is valid when it is generated by other researchers and scientific experts. Conversely, constructivism as a research paradigm assumes that multiple realities exist (Lincoln et al., 2011; Ponterotto, 2005). Rooted in hermeneutics, constructivism assumes that these realities are constructed within each individual's experience and can be uncovered through rich discussion between researchers and study participants (Domenici, 2008). Unlike postpositivism, constructivism epistemologically assumes that knowledge generation is a coconstructive process between researchers and study participants (Domenici, 2008). In a separate vein, the critical theory research paradigm stems from social movements and the Frankfurt School (Agger, 1991; Weaver & Olson, 2006). Ontologically, the critical theory paradigm assumes that an individual's reality is shaped by sociocontextual factors (i.e. social, political gender, cultural, and race/ethnicity) (Lincoln et al., 2011; Weaver & Olson, 2006). Epistemologically, the critical theory paradigm assumes that knowledge is cooperatively produced between researchers and the populations featured within their studies (Agger, 1991; Weaver & Olson, 2006). Regardless of which paradigm is selected by a researcher, adhering to their

underlying assumptions is important, as they provide the underlying rationale to assess the reasonableness of a study's methods, analyses, and conclusions (Domenici, 2008; Lincoln et al., 2011; Ponterotto, 2005; Weaver & Olson, 2006).

Thematic analysis can be used to analyze patterns across the respondents within a study sample, as outlined in a six-step process by Braun and Clarke (2006). It includes familiarizing, generating initial codes, searching for themes, reviewing themes, defining and naming themes, and producing the report. Within Braun and Clarke's approach to thematic analysis, interviews are one type of data collection strategy in which trained interviewers ask participants questions based on interview guides that cover key topics (Braun & Clarke, 2006; DiCicco-Bloom & Crabtree, 2006). Depending on the research team's onto-epistemology and qualitative research tradition, interviews can be structured, semistructured, or unstructured (Renjith et al., 2021). The process of analyzing interview data within this analytical approach typically requires researchers to transcribe audio recordings verbatim, review the transcripts, code the data, and group these codes into latent level themes. Within thematic analysis, coding is not only a critical step in the data analysis process but it also essential for developing familiarity with the data (Braun & Clarke, 2006).

In thematic analysis more generally, coding is a process of assigning labels to data (e.g. words, phrases, sentences, or paragraphs) that in relation to a research question of interest (Kawulich, 2016; Ruona, 2005). There are many different approaches to coding qualitative data; two approaches include an inductive approach and a deductive approach (Azungah, 2018). Deductive approaches to coding are often theory-driven, consisting of a predefined coding scheme (e.g. based on theory or literature review) that sorts data into preset categories (Bingham & Witkowsky, 2021). It relies on existing theories, concepts, or frameworks to guide the coding process (Azungah, 2018; Bingham & Witkowsky, 2021). Inductive approaches to coding do not use a predefined coding scheme but instead involve organically generating codes in relation to a research question (Azungah, 2018; Naeem et al., 2023).

Coding interviews is a time-intensive process that involves an in-depth reading and understanding of each transcript. Depending on the decided unit of analysis (e.g. single lines of text, paragraphs, a complete document like a transcript) the amount of time needed to code an interview transcript can vary greatly (Elliott, 2018). In recent years, with the development of computer software, many qualitative data analysis software have begun to provide coding functions, helping researchers efficiently manage and organize their qualitative datasets (Paulus, 2023). However, the cost of expensive software and the opaque data processing workflow may be challenges for individual researchers or small research teams with limited budgets.

The emergence of generative artificial intelligence (GenAI), with the development of computer science tools and methods, offers promising opportunities for qualitative research (Stokel-Walker & Van Noorden, 2023). The field of Natural Language Processing (NLP) has made significant advancements in enabling computers to understand human languages and manipulate texts (Chowdhary, 2020). A notable example of this advancement refers to ChatGPT, a cutting-edge NLP model developed by OpenAI et al. (2024). ChatGPT is based on the generative pre-trained transformer (GPT) model, which utilizes the transformer architecture to generate coherent and contextually appropriate responses in humanlike conversational settings (Yenduri et al., 2024). By training on vast amounts of text data, ChatGPT can generate coherent, contextually relevant responses to a wide range of human language inputs. More importantly, the superior ability of ChatGPT to understand language allows it to engage in dynamic and natural conversations. This deep understanding enables ChatGPT to respond in ways that are

contextually appropriate, nuanced, and reflective of humanlike reasoning, creating opportunities across a wide range of domains.

Researchers are increasingly exploring the integration of ChatGPT into various fields, recognizing its potential to revolutionize the handling of qualitative data. For instance, researchers in psychology, education, and public health have begun investigating how ChatGPT can be leveraged to facilitate their tasks such as deductive coding, inductive coding, and thematic analysis (Biswas, 2023; Bryda & Sadowski, 2024; Demszky et al., 2023; Jang et al., 2024; Mathis et al., 2024; Tai et al., 2024). These pioneering efforts highlight the promise of large language models (LLMs) in enhancing the efficiency and scalability of qualitative research. Despite their success, prior research has also noted that LLMs, including ChatGPT, often lack domain-specific knowledge (Szymanski et al., 2025). There has been limited exploration of the applications of LLMs to support maternal health studies, highlighting the need for interdisciplinary collaboration to incorporate advanced technology with real-world practices.

To this end, our study examines the potential of ChatGPT, in particular, ChatGPT 4 model, to assist with the inductive coding process for thematic analysis in qualitative studies, specifically following Braun and Clarke's widely adopted framework (Braun & Clarke, 2006). Using semistructured interview data collected from maternity care providers in South Carolina, this study (1) proposes a novel computational workflow to use ChatGPT in inductive coding of thematic analysis, as well as (2) assesses the coding performance of ChatGPT in terms of its credibility and dependability, by comparing its coding results with those coded manually. Notably, OpenAI offers the access to ChatGPT through an application programming interface (API), which supports the automated processing of multiple input documents. Thus, we will leverage the ChatGPT API into our workflow to support thematic analysis and generate codes.

This study is situated within a postpositivist research paradigm. Therefore, it was conducted with the ontological assumption that while a single reality exists, it may not be possible to understand it in its entirety since there are factors that we cannot account for (Lincoln et al., 2011; Ponterotto, 2005). Likewise, this study epistemologically assumes that elements of quantitative forms of inquiry (e.g. statistics) can be used if needed to answer the study's research question (Lincoln et al., 2011; Ponterotto, 2005). Moreover, this study epistemologically believes that other researchers are viable assessors of the validity of our work (Lincoln et al., 2011; Ponterotto, 2005). We chose to situate this study within a postpositivist research paradigm because it allows us to flexibly conceptualize how to apply quantitative technology developed using positivist rationale (i.e. GenAI) within a qualitative research context. This study could contribute to this emerging area by illustrating how LLMs can be leveraged to enhance qualitative analysis in maternal health research, with potential implications for broader applications across public health domains.

METHODS

Interview raw data

The qualitative data used in this study were derived from a previous qualitative study that investigated racial disparities in maternal healthcare services and outcomes during the COVID-19 pandemic in South Carolina. In that study, we conducted semi-structured interviews with 39 women who gave birth between March 2020 and July 2021 and nine maternity care providers (i.e. physicians, nurses, and case managers) from clinics serving communities with a high

TABLE 1 Questions of the in-depth interview guide from the maternity care providers.

| # | Questions |
|----|---|
| 1 | What are the challenges in providing maternal care? |
| 2 | How have you dealt with those challenges? |
| 3 | Are there any differences in labor and delivery in your facilities during COVID-19 compared with prepandemic? |
| 4 | How did COVID-19 change the policies and practices in the department/unit where you work? |
| 5 | Do you think the changes have affected the quality of care provided to women? If so, how? |
| 6 | Do you think the changes have affected maternal health outcomes? If so, how? |
| 7 | Are there any changes in practice or procedure that you think should be continued beyond COVID-19? If so, which? |
| 8 | Based on your observations, what are the specific challenges that your patients have experienced in the pandemic? |
| 9 | What has made them stressed and anxious? Would you like to share any examples? |
| 10 | How would you describe your clients' psychological conditions during the pregnancy in general? |
| 11 | Did they complain about any problems related to mental health? |
| 12 | What about their psychological conditions after giving birth? |
| 13 | What do you think about the health disparity in maternal health outcomes in South Carolina? |
| 14 | In your opinion, what are the main factors that contribute to the disparities? |
| 15 | How did the COVID-19 pandemic affect these factors? |
| 16 | Do you have any suggestions to reduce health disparities in maternal health in South Carolina? |
| 17 | If there was an intervention, or a program, dedicated to addressing the maternal health of Black/African American or Latino women throughout the COVID-19 pandemic, what would you like to see from that program? |
| 18 | Do you have any suggestions for how that program could best support Black/African American or Latino women? |
| 19 | What questions do you have for me? Is there anything else you would like to add? |

proportion of Black and/or Hispanic populations. In the current study, only the data from the maternity care providers was used. More details about the dataset and study protocol can be referenced in Zhang et al. (2024). The audio recordings of the interviews were transcribed using Otter.ai (2024), and the transcripts were manually reviewed and refined by two research assistants to ensure accuracy. The interview guide questions are listed in Table 1.

Research ethics of reusing data

All personal and identifiable information was removed from the transcripts to protect participants' privacy and confidentiality before conducting data analysis, following our standardized qualitative study protocol (Zhang et al., 2024). In addition, we applied for and obtained Institutional Review Board (IRB) approval for an amendment allowing the reuse of transcripts collected for a traditional qualitative study for use in the present context. Considering potential ethical concerns and privacy protection, we chose to exclude the transcripts or

data from ethnic minority women in this exploratory study. The reuse of the transcript data was approved as a study amendment by the IRB committee at the University of South Carolina (#Pro00115169).

Interview transcript and structure

Given the structure of an interview, the raw interview transcript typically follows a question-and-answer structure that allows researchers to consider a variety of responses from different participants to the same question, thereby identifying diverse opinions and lived experiences. To facilitate GenAI's coding, we first reformatted the text into structured, analyzable data by matching each question with its corresponding answer. Each question–answer pair (hereafter referred to as a *dialogue*) was assigned a unique identifier (ID#), making the data computer-readable and facilitating efficient data management and query. Figure 1 illustrates the reformatting process from raw text into structured data, where interviewers' questions were aligned with participants' responses. These structured data were saved in CSV files, with each participant's responses organized accordingly.

Coding of semistructured interview with GenAI

In this study, we aimed to investigate the potential of utilizing ChatGPT for thematic analysis in inductive coding for semistructured interviews. Specifically, through carefully designed instructions and question prompts, ChatGPT generated a set of codes by identifying and summarizing key points mentioned by each participant following thematic analysis approach. Braun and Clarke's six-phase coding framework widely adopted for thematic analysis of inductive coding was employed in the present study (Braun & Clarke, 2006). The framework involves the following steps:

1. Familiarizing: Read the data thoroughly, transcribe if necessary, and jot down initial ideas.
2. Generating initial codes: Systematically code features across the dataset, collating relevant data.

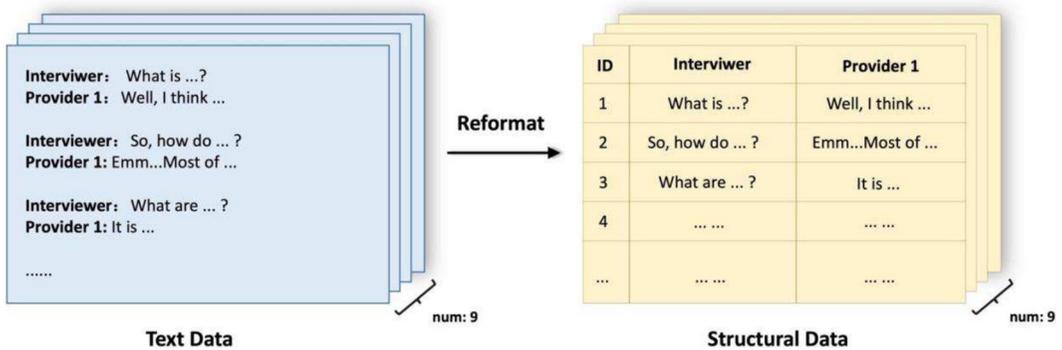


FIGURE 1 Reformatting raw text data into structured computer-readable data. Each interviewers' question is matched with its corresponding participants' answer as a dialogue.

3. Searching for themes: Organize codes into potential themes, gathering relevant data for each theme.
4. Reviewing themes: Check the application of themes, in relation to the coded extracts and the entire dataset, generating a thematic “map” of the analysis.
5. Defining and naming themes: Refine themes and the overall narrative, generating clear definitions and names.
6. Producing the report: Conduct the final analysis, select compelling examples, relate the analysis to the research question and literature, and produce a scholarly report.

In this study, we used ChatGPT 4 to simulate the thematic analysis procedure from the second step to the fourth step. The specific methods included *dialogue filtering*, *code generation*, and *code aggregation*, which are illustrated in Figure 2. During the development and testing of the proposed workflow, we first interacted with the ChatGPT 4 model through its web-based chatbot interface. This allowed us to iteratively refine prompts utilized for analysis and roughly evaluate model performance in generating appropriate codes. Once the prompt structure was established, we leveraged the ChatGPT API to automate the inductive coding process at scale.

Dialogue filtering

In-depth interviews often feature a deep, conversational nature where participants might address supporting points for one question while responding to another, either unintentionally or to reinforce their viewpoints. As a result, the same perspective might be repeated across different responses and intersect with multiple interview questions. Given this complexity, we performed *dialogue filtering* to locate participants' responses to specific questions. We first instructed ChatGPT to determine whether a specific question was asked during the interview. When the question was identified, ChatGPT provided the dialogue ID, allowing researchers to observe participants' responses to the question. This step ensured that only relevant sections of the dialogue are further coded, enhancing the accuracy and efficiency of the analysis. The prompt template used for dialogue filtering was illustrated in Figure 3. The prompt instructed ChatGPT to identify if a specific interview question is present within the dialogue and, if so, to return the ID of the relevant dialogue. If no relevant content was found, it returned a 0. By doing so, for each question, we identified all relevant dialogue from each participant for each question.

Code generation

After identifying the dialogues for each question, we instructed ChatGPT to perform *code generation* by summarizing and extracting key topics from the dialogues. Figure 4 illustrates the prompt template used for generating codes. In this phase, responses to specific questions from

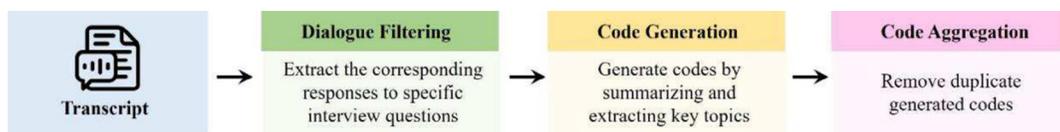


FIGURE 2 The workflow of using ChatGPT for inductive coding.



Identify if there is a dialogue asking the question:
'What are the challenges in providing maternal care?'
 (if yes tell me the number the dialogue, no then say 0, just give me the number):

- 1. "I'm 61. "
- 2. "Uh, two months. "
- ...
- {All dialogues in an interview transcript}



Yes, the question "What are the challenges in providing maternal care?" is asked in the dialogue. The number of the dialogue is 14.

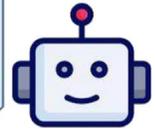


FIGURE 3 The prompt template used for dialogue filtering. If the given question is asked in the interview transcript, returns the row number of the dialogue, else then returns 0.



Can you extract the topics from the answers about the question **'What are the challenges in providing maternal care?'**:
'The challenges are getting people to get vaccinated and to get them to not come into our clinic sick and, and make other patients and family people sick. So it's, and it's very hard to provide care when they want to bring their family members into their visits, and they cannot. So it's been a very challenging time in the last two years because of that.'



The topics extracted from the response regarding the challenges in providing maternal care are:

- Vaccination: Ensuring that individuals are vaccinated.
- Clinic Safety: Preventing sick individuals from coming to the clinic and spreading illness.
- Visitor Restrictions: Challenges associated with limiting family members from attending visits.
- Pandemic Impact: The overall increased difficulty of providing care during the last two years, likely due to COVID-19 related challenges.
- ...

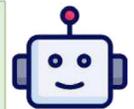


FIGURE 4 The prompt template used for code generation.

all interview transcripts were collected. ChatGPT was then tasked to extract and generate codes from the collected responses. This process enabled the identification of themes and topics across participants' responses.

Code aggregation

After GenAI generated a set of codes reflecting the key topics in response to each question, we further aggregated these codes into categories and developed themes. The first task was to

merge codes with similar semantic meanings. While ChatGPT generated the most fitting code for each participant's response, slight variations in expressions resulted in many repetitive codes that were redundant for analysis and required manual aggregation. The second task was to synthesize codes based on their conceptual meanings. The codes were combined and grouped into categories based on their definitions and scope in the conceptual system. The process of *code aggregation* ensured that the extracted codes were concise and useful for identifying themes in subsequent qualitative data analysis.

Evaluation

Reviewing the generated and aggregated codes was crucial to ensure the accuracy and applicability of the generated themes. Although ChatGPT is often efficient at identifying general themes present in dialogue, it may lack domain-specific knowledge. Thus, two evaluation approaches were performed to assess the validity of the generated codes: a human-centered approach and a machine-based approach. It is important to emphasize that the purpose of presenting these evaluation methods is not to position one as superior to the other. Rather, our goal is to assess and characterize the performance of ChatGPT through both lenses, recognizing that each approach has its own advantages and limitations.

Human-centered approach

We engaged domain experts in maternal health, with relevant background knowledge in manual coding of the interview transcripts. One faculty member and one doctoral student in public health, both with extensive experiences in qualitative studies, independently developed the codes after reviewing all the transcripts. The disagreements were sufficiently discussed and resolved before finalizing the codes assigned to each participant's response. The domain experts reviewed and compared the ChatGPT-generated codes with the human-generated codes. If the ChatGPT-generated codes were the same as and/or aligned with the human-generated codes, they were marked as "accepted." Two domain experts independently marked each code as accepted and discussed any disagreements. The accuracy rate for each response to a question was the number of "accepted" codes divided by the total number of codes generated by ChatGPT.

Machine-based metric approach

We also compared the human-generated codes and ChatGPT-generated codes and quantified the differences by assessing their semantic similarity. The underlying hypothesis is that if ChatGPT performs effectively, the semantic content of its generated codes should closely resemble that of human-generated codes. In particular, we utilized a cutting-edge text embedding method, Bidirectional Encoder Representations from Transformers (BERT), to convert the codes into high-dimensional vectors (Devlin et al., 2019). No preprocessing was conducted since the codes (whether generated by humans or ChatGPT) consisted of only a few words or phrases. These high-dimensional feature vectors captured the semantic meanings of the input

codes. We calculated the cosine similarity between the two high-dimensional vectors using the following formula.

$$S(V_G, V_H) = \frac{V_G \cdot V_H}{\|V_G\| \times \|V_H\|}$$

V_G represents the high-dimensional vector of a GenAI-generated code, V_H represents the high-dimensional vector of a human-defined code, $S(V_G, V_H)$ represents the cosine similarity of V_G and V_H . $S(V_G, V_H)$ ranges from 0 to 1, and the higher the $S(V_G, V_H)$, the more similar the two codes.

RESULTS

Inductive coding performance

We first show several examples of the codes generated by ChatGPT. Table 2 lists the ChatGPT-generated codes for three example questions. Each question was coded based on the participants' direct responses to the questions.

Overall, the codes created from the text are relevant and demonstrate alignment, indicating that ChatGPT has provided relatively meaningful results. For instance, the codes generated for the question about challenges in providing maternal care (Example A) cover a broad range of

TABLE 2 Three example questions and their corresponding codes created by ChatGPT.

| Question | Codes |
|--|--|
| Example A | |
| <i>What are the challenges in providing maternal care?</i> | (1) visitor restrictions, (2) telehealth limitations, (3) inconsistent care standards, (4) partner attendance issues, (5) missed appointments, (6) vaccination hesitancy, (7) difficulty managing COVID patients, (8) socio-economic disparities, (9) technology accessibility, (10) pandemic safety concerns |
| Example B | |
| <i>How did COVID-19 change the policies and practices in the department/unit where you work?</i> | (1) knee jerk reactions to changes, (2) communication issues, (3) visitation policy alterations, (4) COVID testing for patients, (5) push on vaccinations, (6) labor and delivery support changes, (7) personalized care challenges, (8) clothing policy changes, (9) donning and doffing PPE, (10) masking rules enforcement |
| Example C | |
| <i>Do you think the changes have affected the quality of care provided to women? If so, how?</i> | (1) reduced support for women, (2) less tolerance for birth plans, (3) nursing staff burnout, (4) early postpartum discharge, (5) fear of COVID in hospitals, (6) increased depression and anxiety, (7) missed appointments, (8) halted preventive screenings, (9) reduced physical touch in care, (10) implicit bias and fear |

relevant issues, including objective (e.g., visitor restrictions), personal (e.g., vaccination hesitancy), and social factors (e.g., socio-economic disparities). These results indicate that ChatGPT has successfully identified and categorized key themes from the participants' responses. On average, ChatGPT generated approximately 10 codes per interview question, with each code containing an average of five words. Notably, when inputting the same question into ChatGPT multiple times, the generated outcomes might be inconsistent. Thus, we performed an assessment to see if the generated results showed similar patterns. We leveraged the BERT and computed the similarity among generated codes. The results revealed that the codes generated by ChatGPT had a high similarity score of over 0.95, indicating that ChatGPT's outputs are relatively stable across repeated runs.

ChatGPT's performance in inductive coding was quantitatively evaluated by accuracy rate of codes for responses to each question. Figure 5 presents the number of codes generated by GenAI for each question, the number of "accepted" codes after manual review, and the corresponding accuracy rates. The overall accuracy rate observed was 85.15% across all questions. For individual questions, the accuracy rates ranged from a maximum of 100% to a minimum of 66.67%. Notably, 11 questions had coding accuracy rates equal to or greater than 90%, representing more than half of the total number of questions, with eight questions achieving a coding accuracy rate of 100%. Only three questions had coding accuracy rates less than 80%. After reviewing these three questions, we found that the responses are highly descriptive, and ChatGPT tended to extract more generalized information from these detailed responses. For example, in coding Question 6, "Do you think the changes have affected maternal health outcomes? How?," ChatGPT provided the broad summary code of "Changes affecting maternal health." For Question 10, "How would you like to describe your clients' psychological conditions during the pregnancy in general?," ChatGPT provided summarized codes like "Patients' health and baby's health concerns," while human-generated codes were more specific, such as "Psych," "Anxiety," and "Anger." This finding indicates that a small portion of questions did not achieve high accuracy in coding. Overall, the inductive coding results are reliable, and ChatGPT demonstrated promise in providing support for qualitative coding.

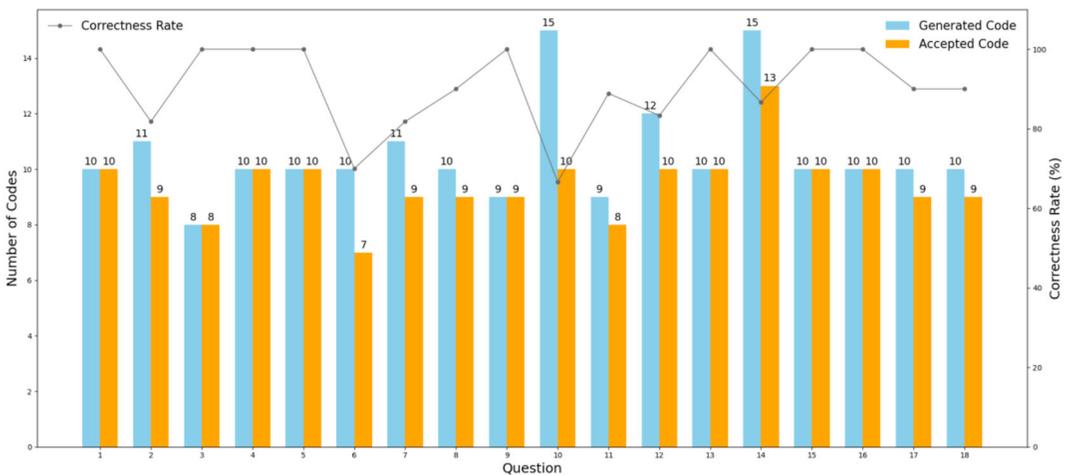


FIGURE 5 The number of generated codes, the number of "accepted" codes after manual review, and the accuracy rate for each question.

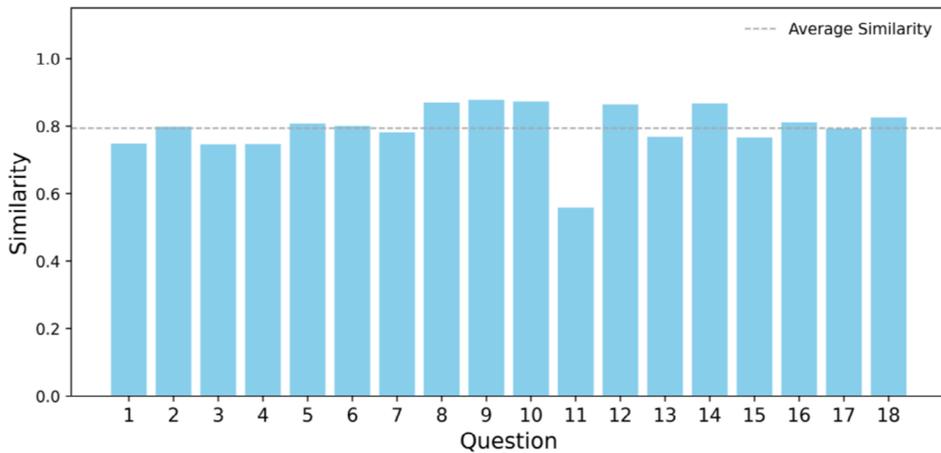


FIGURE 6 Text similarity between the ChatGPT-generated codes and human-generated codes for each question. The x-axis represents text similarity, and the y-axis represents the question number.

Moreover, we compared the ChatGPT-generated codes and human-generated codes for each question using the text semantic similarity score. Figure 6 shows the cosine similarity between the high-dimensional vectors of the two sets of texts. The highest similarity score was 0.879, with an overall similarity of 0.795. It should be noted that there is no established threshold for determining semantic similarity between two groups of codes. Thus, to provide greater interpretability of these results, we present three sample codes, as illustrated in Table 3 with high, medium, and low semantic similarity scores. These findings indicate that the semantic meaning of the codes generated by ChatGPT is very similar to those by human coders, illustrating the proficiency of using ChatGPT for the inductive coding of interview texts.

Time efficiency

In addition to demonstrating accuracy, the GenAI-based coding approach has also shown significant time efficiency compared to traditional methods. The total estimated time spent by human coders for inductive coding was approximately 30.5 hours, of which 4.5 hours were spent getting familiar with the nine transcripts and 11.25 hours were spent in generating and assigning codes to the dialogues.

The usage of ChatGPT was found to significantly reduce the time required to complete inductive coding. The time expenditure for GenAI-based coding primarily reflects two activities: *setup* and *coding*. The first is the time spent setting up the environment, including data preparation for ChatGPT coding, designing appropriate prompts, and reformatting the transcript files. The second is the time required for the actual coding, by ChatGPT. The first part, setting up the environment, requires most of the time in the work pipeline, around 5 hours. In particular, a human coder takes 3 hours to reformat the original transcript into Excel to prepare the materials for use in the GenAI coding. Then, it takes around 2 hours to test different prompts and reformat the transcript files to be fed into GenAI. Once the setup is complete, the actual coding process by GenAI is much faster, taking only around 20 minutes to generate and assign codes following the inductive approach.

TABLE 3 Similarity of inductive codes generated by human coders and ChatGPT.

| Inductive code generated by human | Inductive code generated by ChatGPT | Similarity |
|--|---|------------|
| Example A | | |
| COVID vaccine, job security/work, domestic situation, risk for COVID, mask wearing, childcare, hard jobs, vaccine hesitancy coupled with a high risk work environment, childcare/work/balance, pregnancy stress coupled with COVID stress, vaccine and mask discord, financial, transportation issues, unsafe work environment, having to homeschool or manage virtual learning for children, anger and anxiety, concerns about pregnancy during the pandemic among better off patients, financial concerns, isolation, fear of COVID itself, vaccine decisions, domestic violence maybe but screening was not optimal, scared to come in to appointment, family members dying from COVID, and vaccine hesitancy | Jobs and employment, family safety concerns, difficulty breathing in masks, pregnancy mask concerns, deciding on vaccination, isolation from loved ones, financial worries, fear of the virus, and vaccine anxiety | 0.879 |
| Example B | | |
| Visitation restrictions, differences in care for COVID positive patients, less triages, policy changes affected patient experience, PPE and infection control measures, telemedicine, center shut down, and nothing stopped at their clinic | Visitor restrictions, barriers to care, reduced one-on-one support, PPE requirements, changes in room availability, staffing and nursing shortages, patient experiences affected, and supply chain issues | 0.746 |
| Example C | | |
| Depression, general life stressors, COVID fear, anxiety, mood, and sleep issues | No more mental health complaints; pregnancy-related insomnia; increased domestic violence; emotional, physical, sexual abuse; overwhelmed mothers; small children at home; challenges of work-from-home; unemployment effects; and pandemic-related fatigue | 0.559 |

It should be noted that additional human-centered tasks are still necessary with current GenAI techniques. For example, these methods require domain experts to review the coding results and assess the quality of how each transcript was coded, which results in additional time costs. However, with the development of more robust and advanced GenAI models in the future, these review steps may be simplified and the process further streamlined. Consequently, the total time required for GenAI-based coding might be reduced to just 19% of the time required for traditional coding methods.

DISCUSSION

These results of using ChatGPT for inductive coding highlight the feasibility of GenAI in qualitative studies that use a postpositivist approach to Braun and Clarke's thematic analysis

(Braun & Clarke, 2006). ChatGPT achieved relatively high performance in inductive coding. Traditionally, coding interview transcripts manually is a time-consuming process. Using GenAI in coding can greatly shorten this process and maintain relatively high accuracy. Thus, it could enhance efficiency in qualitative research. The findings from our study provide preliminary evidence for both opportunities and challenges of using GenAI in qualitative data coding.

Opportunities brought by GenAI for qualitative study

The emergence of GenAI technology provides unprecedented opportunities to enhance the coding of qualitative data to support thematic analysis.

User friendly and accessible

GenAI is easily accessible in terms of cost and ease of use, accommodating users with varying levels of technical expertise. This accessibility allows for broad applicability across various research domains without requiring extensive prior knowledge or domain-specific input. For example, ChatGPT offers both web-based chatbot user interface and APIs. Thus, users could interact with it as they would in human conversations. Also, ChatGPT APIs facilitate integration of GenAI into custom workflows and enable large-scale, automated processing of textual data. The advanced natural language understanding ability makes GenAI automatically screen interview transcripts and identify specific questions or topics related to a given question. When designing prompts in ChatGPT, there is no need to input a large training dataset. Even when using a small number of sample dialogues, ChatGPT promisingly achieved high levels of accuracy in identifying and categorizing topics in the data. In our study, although the main research topic is related to maternal health, ChatGPT successfully identified discussions related to cannabis, opioid, and alcohol use from the participants' responses, showing its flexibility across different public health areas. This suggests GenAI can be effectively applied across public health research.

Multi-level semantic structure identification

GenAI demonstrates an ability to identify multi-level semantic structures when coding, which is well represented in the code aggregation process. ChatGPT can generate fine-resolution codes and merge those with similar semantic meanings into broader categories. This granularity of analysis provides a more detailed understanding of qualitative data, enabling researchers to capture the complexity and diversity of participants' perspectives. Traditionally, qualitative researchers have to manually merge and aggregate codes across different levels. With GenAI, subtopics can be automatically grouped under their main topics, enhancing the overall coherence of the analysis.

Scalability and efficiency in large datasets

GenAI excels at handling large datasets, providing a scalable solution for analyzing large datasets. Analyzing large datasets is often a challenge due to time and human resource

limitations. Once the environment setup is completed, GenAI can process transcripts in parallel by leveraging its APIs, which can dramatically improve coding efficiency. Therefore, GenAI can serve as a tool, completing most of the initial coding tasks, thereby allowing human coders and experts to focus their efforts on the careful review and discussion of outputs as well as the applicability, interpretation, and dissemination of the findings. Thus, collaboration between GenAI and human coders can potentially reduce the time required for qualitative data analysis while preserving the quality of the data analysis.

Challenges and concerns when using GenAI

Despite the promising potential of GenAI for coding in qualitative studies, several challenges and concerns need careful consideration, including **inaccuracy, systematic biases, and privacy issues**. The inaccuracy issues in GenAI-based coding stem from two aspects. First, the lack of maternity knowledge, due to the sensitive nature of maternal health research, may result in inaccurate results. We observed that ChatGPT tends to produce more general codes than human coders and lacks the domain-specific terminology and professional knowledge unique to maternal experts. For example, for Question 10, “How would you like to describe your clients’ psychological conditions during the pregnancy in general?,” one of the codes generated by ChatGPT was “Post-traumatic stress from having COVID,” while the human-generated code was “Some form of PTSD symptoms in patients who were really sick with COVID,” which used more professional medical terminology and provided a more detailed description of the psychological issues described. This discrepancy may be attributed to ChatGPT being trained on a general corpus with limited exposure to public health knowledge, especially maternal health knowledge. Also, maternal health data are often highly sensitive, which hinders ChatGPT’s ability to learn maternity-related knowledge. This implies that qualitative studies seeking to use GenAI-based coding to assist within their data analysis will need to consider selecting a model that can be trained on their public health topic of interest. Second, the “hallucination” phenomena are commonly observed in current GenAI models. Since current GenAI models, such as ChatGPT, are trained on a general corpus rather than specific domain knowledge, they can produce inaccurate content and false information. Therefore, these hallucinations further demonstrate the necessity of human coders and experts to carefully review all outputs generated by GenAI during the coding process.

Additionally, systematic biases are rooted in representation issues when training datasets. Since the underlying generation processes of GenAI are inaccessible to prompt engineers, it is difficult to quantify errors and debug them. The datasets used to train GenAI are not readily available to users and may underrepresent certain minority population groups, thereby overlooking their lived experiences, sociolinguistic differences, and potentially leading to the amplification or replication of existing social biases (Gal, 2016). Being able to capture and understand these nuances are essential within qualitative research, as it is not possible to understand a phenomenon of interest without doing so. Since thematic analysis looks at patterns across a dataset, researchers conducting latent level analysis using GenAI-developed codes are particularly vulnerable to unknowingly drawing conclusions that are decontextualized and incorrect (Braun & Clarke, 2006). Anticipating and preventing these instances of error are vital as studies using GenAI-developed codes could inform clinical practices and influence public health policies.

Moreover, using GenAI systems, such as ChatGPT, poses privacy concerns, particularly regarding the handling of sensitive or personal interview data. Current GenAI systems may

retain and use input data for further model training, which could lead to misuse or sharing of private information. Researchers must sufficiently assess potential ethical issues before entering sensitive information into GenAI models and adhere strictly to IRB guidelines to ensure compliance with data protection protocols. Implementing robust data security measures, including data deidentification, and stringent access controls, can mitigate privacy risks. In addition, potential technical solutions may help address this issue in the future. The development of privacy-preserving machine learning techniques, such as differential privacy and secure multi-party computation, can enhance data protection without compromising the utility of GenAI models in qualitative research. Furthermore, techniques such as machine unlearning may protect sensitive data by preventing models from retaining or memorizing sensitive input data, thereby contributing to more secure and ethical use of LLMs. Ultimately, ensuring the responsible use of GenAI in public health research necessitates a multidisciplinary approach that integrates technological advancement, regulatory compliance and further discussions on ethics. This paper highlights such importance and advocates for more interdisciplinary collaborations to harness the benefits of GenAI and maintain the highest standards of data integrity, participant confidentiality, and rigor in qualitative methodology.

Limitations and future directions

This study has several limitations that need to be considered. First, we analyzed transcripts from individual semi-structured interviews with maternity care providers conducted during the COVID-19 pandemic. The specificity of this dataset may limit the transferability of our findings to other interview formats (e.g., structured interviews, unstructured interviews, or focus groups) and different research contexts (e.g. different geographic regions, various healthcare settings, or alternative qualitative analytical approaches). Future studies should apply the computational workflow across a wider range of datasets and research contexts to showcase its adaptability and potential in qualitative research.

Second, in the current study, we tested only one GenAI model (i.e. ChatGPT). Variability in performance, accuracy, and interpretability across different GenAI models (e.g., Gemini, DeepSeek) remains an open question. Future studies may explore the potential to incorporate more advanced GenAI approaches, such as agent-based GenAI methods, or different GenAI models (e.g., commercial vs. open-sourced models). Comparing different models will provide insights into their relative strengths, weaknesses, and applicability to qualitative research.

Third, in our study, we used accuracy rate based on two human coders' judgment to assess the credibility of AI-generated codes against the human-generated codes within our dataset. Within qualitative research using thematic analysis, many tools exist to assess the accuracy of manually coded transcripts (e.g., reflexivity journals, triangulation, internal audits, prolonged engagement, member checks, and thick descriptions), with the selection of each tool depending on qualitative research paradigm in which a study is situated. In addition, the potential for anchoring bias among human coders must be considered during the evaluation of AI-generated codes. When reviewing outputs produced by LLMs, human evaluators may be more inclined to accept the suggested codes rather than critically proposing or independently generating alternative codes. Therefore, future research needs to explore the best practices for assessing the quality and validity of qualitative coding when incorporating LLMs.

Fourth, we did not train the LLMs on codes relevant to maternal health. This limitation may have resulted in the model overlooking critical domain-specific nuances in the coding

process. Future research is needed to assess the degree to which coding quality changes when the LLMs are trained versus untrained.

CONCLUSION

In summary, we proposed a computational framework for coding qualitative data for thematic analysis by leveraging advanced GenAI methods. Our findings demonstrated the potential of GenAI (i.e., ChatGPT 4) in automating the inductive coding processes, thereby enhancing efficiency within qualitative studies using Braun and Clarke (2006)'s approach to thematic analysis. In our case study on maternal health, ChatGPT could generate meaningful codes from the data to uncover key themes and topics from the semi-structured interview transcripts. Our evaluation showed a relatively high overall accuracy of GenAI-generated codes and a significant reduction in the time spent on coding processes. These findings suggest that GenAI can effectively support qualitative coding, provide reliable results, and offer substantial benefits for various public health studies. Despite its promise, we have also noted the challenges in applying GenAI in qualitative research, including inaccuracies, systematic biases, and privacy issues. Therefore, it is crucial to incorporate more domain knowledge into GenAI and ensure that domain experts carefully review the results of GenAI-based coding given the current technologies. With advancements in AI technologies, we believe qualitative research will continue to benefit from these cutting-edge methods.

ACKNOWLEDGMENTS

The authors would like to thank Camryn Garrett and Ariele N'Diaye at the University of South Carolina who helped us preprocess the maternal dataset and proofread the manuscript. Y.K. acknowledges the funding support provided by the Population Research Center Seed Grant, The University of Texas at Austin. M.L. acknowledges the funding support provided by NIH grant #R01AI127203-05S2. S.Q. acknowledges the funding support provided by NIH grant #R01AI174892. Any opinions, findings, conclusions, or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Institutes of Health (NIH).

CONFLICT OF INTEREST STATEMENT

The authors declare no conflicts of interest.

DATA AVAILABILITY STATEMENT

Author elects to not share data.

ORCID

Shan Qiao  <https://orcid.org/0000-0003-1834-1834>

Yuhao Kang  <https://orcid.org/0000-0003-3810-9450>

REFERENCES

- Agger, B. (1991). Critical theory, poststructuralism, postmodernism: Their sociological relevance. *Annual Review of Sociology*, 17(1), 105–131. <https://doi.org/10.1146/annurev.so.17.080191.000541>
- Azungah, T. (2018). Qualitative research: Deductive and inductive approaches to data analysis. *Qualitative Research Journal*, 18(4), 383–400. <https://doi.org/10.1108/QRJ-D-18-00035>

- Bingham, A. J., & Witkowsky, P. (2021). Deductive and inductive approaches to qualitative data analysis. In *Analyzing and interpreting qualitative research: After the interview* (pp. 133–146). SAGE.
- Biswas, S. S. (2023). Role of chat GPT in public health. *Annals of Biomedical Engineering*, *51*(5), 868–869. <https://doi.org/10.1007/s10439-023-03172-7>
- Braun, V., & Clarke, V. (2006). Using thematic analysis in psychology. *Qualitative Research in Psychology*, *3*(2), 2. <https://doi.org/10.1191/1478088706qp0630a>
- Bryda, G., & Sadowski, D. (2024). From words to themes: AI-powered qualitative data coding and analysis. In J. Ribeiro, C. Brandão, M. Ntsohi, J. Kasperuniene, & A. P. Costa (Eds.), *Computer supported qualitative research* (pp. 309–345). Springer Nature Switzerland. https://doi.org/10.1007/978-3-031-65735-1_19
- Chowdhary, K. R. (2020). Fundamentals of artificial intelligence. In *Natural language processing* (pp. 603–649). Springer India. https://doi.org/10.1007/978-81-322-3972-7_19
- Demszky, D., Yang, D., Yeager, D. S., Bryan, C. J., Clapper, M., Chandhok, S., Eichstaedt, J. C., Hecht, C., Jamieson, J., Johnson, M., Jones, M., Krettek-Cobb, D., Lai, L., JonesMitchell, N., Ong, D. C., Dweck, C. S., Gross, J. J., & Pennebaker, J. W. (2023). Using large language models in psychology. *Nature Reviews Psychology*. <https://doi.org/10.1038/s44159-023-00241-5>
- Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding,” *arXiv*: arXiv:1810.04805. <https://doi.org/10.48550/arXiv.1810.04805>
- DiCicco-Bloom, B., & Crabtree, B. F. (2006). The qualitative research interview. *Medical Education*, *40*(4), 314–321. <https://doi.org/10.1111/j.1365-2929.2006.02418.x>
- Domenici, D. J. (2008). Implications of hermeneutic constructivism for personal construct theory: Imaginally construing the nonhuman world. *Journal of Constructivist Psychology*, *21*(1), 25–42. <https://doi.org/10.1080/10720530701503868>
- Elliott, V. (2018). Thinking about the coding process in qualitative data analysis. *The Qualitative Report*, *23*(11). Retrieved September 05, 2024, from <https://ora.ox.ac.uk/objects/uuid:5304bf7f-6214-4939-9f1b-b64415d4fac1>
- Gal, S. (2016). Sociolinguistic differentiation. In *Sociolinguistics: Theoretical debates* (Vol. 113, pp. 113–124). Cambridge University Press.
- Guest, G., Namey, E. E., Guest, G., & Fleming, P. (2015). Mixed Methods Research. In *Public health research methods* (pp. 581–614). SAGE Publications, Inc. <https://doi.org/10.4135/9781483398839>
- Jang, K. M., Chen, J., Kang, Y., Kim, J., Lee, J., Duarte, F., & Ratti, C. (2024). Place identity: A generative AI’s perspective. *Humanities and Social Sciences Communications*, *11*(1), 1–16. <https://doi.org/10.1057/s41599-024-03645-7>
- Kawulich, B. B. (2016). *The BERA/SAGE handbook of educational research* (pp. 1–1170).
- Kuckartz, U. (2014). *Qualitative text analysis: A guide to methods, practice and using software*. SAGE.
- Lichtman, M. (2013). *Qualitative research for the social sciences*. SAGE Publications.
- Lincoln, Y. S., Lynham, S. A., & Guba, E. G. (2011). Paradigmatic controversies, contradictions, and emerging confluences, revisited. In *The SAGE handbook of qualitative research* (pp. 97–128). SAGE Publications. https://www.miguelangelmartinez.net/IMG/pdf/2018_denzin_lincoln_handbook_qualitative_research-213-263.pdf
- Mathis, W. S., Zhao, S., Pratt, N., Weleff, J., & De Paoli, S. (2024). Inductive thematic analysis of healthcare qualitative interviews using open-source large language models: How does it compare to traditional methods? *Computer Methods and Programs in Biomedicine*, *255*, 108356. <https://doi.org/10.1016/j.cmpb.2024.108356>
- Mey, G. (2023). Qualitative methodology. In J. Zumbach, D. A. Bernstein, S. Narciss, & G. Marsico (Eds.), *International handbook of psychology learning and teaching* (pp. 453–478). Springer International Publishing. https://doi.org/10.1007/978-3-030-28745-0_22
- Mohajan, H. K. (2018). Qualitative research methodology in social sciences and related subjects. *Journal of Economic Development, Environment and People*, *7*(1), 23–48. <https://doi.org/10.26458/jedep.v7i1.571>
- Moser, A., & Korstjens, I. (2018). Series: Practical guidance to qualitative research. Part 3: Sampling, data collection and analysis. *The European Journal of General Practice*, *24*(1), 9–18. <https://doi.org/10.1080/13814788.2017.1375091>
- Naem, M., Ozuem, W., Howell, K., & Ranfagni, S. (2023). A step-by-step process of thematic analysis to develop a conceptual model in qualitative research. *International Journal of Qualitative Methods*, *22*, 16094069231205789. <https://doi.org/10.1177/16094069231205789>

- Nardi, P. M. (2018). *Doing survey research: A guide to quantitative methods* (4th ed.). Routledge. <https://doi.org/10.4324/9781315172231>
- OpenAI, Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altschmidt, J., Altman, S., Anadkat, S., Avila, R., Babuschkin, I., Balaji, S., Balcom, V., Baltescu, P., Bao, H., Bavarian, M., Belgum, J. ..., Zoph, B., (2024). "GPT-4 Technical Report," *arXiv*: arXiv:2303.08774. Retrieved September 05, 2024, from <http://arxiv.org/abs/2303.08774>
- Otter.ai. (2024). "Voice Meeting Notes & Real-time Transcription." Retrieved October 15, 2024, from <https://otter.ai/>
- Paulus, T. M. (2023). Using qualitative data analysis software to support digital research workflows. *Human Resource Development Review*, 22(1), 139–148. <https://doi.org/10.1177/15344843221138381>
- Ponterotto, J. G. (2005). Qualitative research in counseling psychology: A primer on research paradigms and philosophy of science. *Journal of Counseling Psychology*, 52(2), 126–136. <https://doi.org/10.1037/0022-0167.52.2.126>
- Renjith, V., Yesodharan, R., Noronha, J. A., Ladd, E., & George, A. (2021). Qualitative methods in health care research. *International Journal of Preventive Medicine*, 12, 20. https://doi.org/10.4103/ijpvm.IJPVM_321_19
- Ruona, W. E. (2005). Analyzing qualitative data. In *Research in organizations: Foundations and methods of inquiry* (pp. 223–263). Berrett-Koehler.
- Stokel-Walker, C., & Van Noorden, R. (2023). What ChatGPT and generative AI mean for science. *Nature*, 614(7947), 214–216. <https://doi.org/10.1038/d41586-023-00340-6>
- Szymanski, A., Ziems, N., Eicher-Miller, H. A., Li, T. J.-J., Jiang, M., & Metoyer, R. A. (2025). Limitations of the LLM-as-a-Judge Approach for Evaluating LLM Outputs in Expert Knowledge Tasks. In *Proceedings of the 30th international conference on intelligent user interfaces, in IUI'25* (pp. 952–966). Association for Computing Machinery. <https://doi.org/10.1145/3708359.3712091>
- Tai, R. H., Bentley, L. R., Xia, X., Sitt, J. M., Fankhauser, S. C., Chicas-Mosier, A. M., & Monteith, B. G. (2024). An examination of the use of large language models to aid analysis of textual data. *International Journal of Qualitative Methods*, 23, 16094069241231168. <https://doi.org/10.1177/16094069241231168>
- Tracy, S. J. (2024). Qualitative research methods: Collecting evidence, crafting analysis. In *Communicating impact*. John Wiley & Sons.
- Weaver, K., & Olson, J. K. (2006). Understanding paradigms used for nursing research. *Journal of Advanced Nursing*, 53(4), 459–469. <https://doi.org/10.1111/j.1365-2648.2006.03740.x>
- Yenduri G., Srivastava G., Maddikunta P. K., Jhaveri R. H., Wang W. Vasilakos A. V., & Gadekallu T. R. (2024). GPT (Generative Pre-trained Transformer): A Comprehensive Review on Enabling Technologies, Potential Applications, Emerging Challenges, and Future Directions. IEEE Access. Retrieved September 05, 2024, from <https://ieeexplore.ieee.org/document/10500411>
- Zhang, T. Byrd, S. Qiao, M. E. Torres, X. Li, and J. Liu, (2024). "Maternal care utilization and provision during the COVID-19 pandemic: Voices from minoritized pregnant and postpartum women and maternal care providers in deep south". *PLoS ONE*. Retrieved September 05, 2024, from <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0300424>

How to cite this article: Qiao, S., Fang, X., Wang, J., Zhang, R., Li, X., & Kang, Y. (2025). Generative AI for thematic analysis in a maternal health study: coding semistructured interviews using large language models. *Applied Psychology: Health and Well-Being*, 17(3), e70038. <https://doi.org/10.1111/aphw.70038>