

Towards conversational assistants for health applications: using ChatGPT to generate conversations about heart failure

Anuja Tayal¹ Devika Salunke² Barbara Di Eugenio¹ Paula G. Allen-Meares³
Eulalia P. Abril⁴ Olga Garcia-Bedoya³ Carolyn A. Dickens³ Andrew D. Boyd²

¹Department of Computer Science ²Department of Biomedical and Health Information Sciences
³Department of Medicine ⁴Department of Communications
{atayal14, dsalun2, bdieugen, pameares, eulalia, ogarciab, cdickens, boyda}@uic.edu

Abstract

We explore the potential of ChatGPT to generate conversations focused on self-care strategies for African-American patients with heart failure, a domain with limited specialized datasets. To simulate patient-health educator dialogues, we employed four prompting strategies: aspects, African American Vernacular English, Social Determinants of Health (SDOH), and SDOH-informed reasoning. Conversations were generated across key self-care aspects—food, exercise, and fluid intake—with varying turn lengths and incorporated patient-specific SDOH attributes such as age, gender, neighborhood, and socioeconomic status. Our findings show that effective prompt design is essential. While incorporating SDOH and reasoning improves dialogue quality, ChatGPT still lacks the empathy and engagement needed for meaningful healthcare communication.

1 Introduction

Heart failure (HF), or congestive heart failure, occurs when the heart cannot pump enough blood to meet the body's needs. Effective self-care—such as managing salt intake, staying hydrated, exercising, adhering to medications, and attending regular check-ups—is critical for managing the condition (Savarese and Lund, 2017). However, African Americans (AA) in the U.S. face disproportionately worse outcomes due to genetic factors, limited healthcare access, socioeconomic challenges, and lower health literacy (Nayak et al., 2020). Existing self-care materials often cater to a white, educated population, lacking cultural relevance for minority communities (Barrett et al., 2019). This gap contributes to poor adherence and worsened outcomes. Personalized education can improve self-care understanding and reduce readmissions to hospitals (Di Eugenio et al., 2019).

This study is part of a broader project aimed at developing a culturally sensitive conversational

agent to support AA patients with heart failure in asking self-care related questions. A significant challenge is the lack of real-world patient-centered conversational data from underrepresented communities. To address this, we recruited three patient educators (PEs)¹ to provide heart failure education to 18 AA and 2 Hispanic/Latino (H/L) patients (Gupta et al., 2020). Initial analysis revealed that educators dominated the conversations, with patients contributing less. The key topics discussed during these sessions included exercise, fluid intake, symptom management, sleep, weight management, familial aspects, and salt intake.

Unlike the recorded interactions, we wanted to generate conversational datasets that are initiated by patients and, more importantly, personalized based on the Social Determinants of Health (SDOH) features. With the advancement of Large Language Models (LLMs), we thought to examine different prompting strategies and evaluate whether they could be used for creating synthetic conversational datasets. Our study is a feasibility assessment aimed at exploring ChatGPT's ability to generate self-care conversations and its ability to adapt its responses based on varying prompts. The dataset is publicly available². In this paper, we look into 4 different prompting approaches that will supplement real-world interactions to support the development of a patient-driven dialogue system.

- We started by generating simulated conversations based on different aspects required for self-care for patients with heart failure, including food, exercise, and fluid intake.
- We introduced an additional prompt where the patients communicate using African Ameri-

¹In this paper, we have used the terms patient educator and health educator interchangeably.

²<https://github.com/anujatayal/HF-Dataset>

can Vernacular English (AAVE) while the educator communicates in standard English.

- We prompted to integrate SDOH Features of the patients with the aspect. We considered gender {male, female}, age {young, mid-age, old}, neighborhood {safe, unsafe} and socio-economic conditions {below poverty line, well to do}.
- We first prompted ChatGPT to generate reasoning given the SDOH features of the patient and then prompted ChatGPT again to generate conversations given the reasoning and SDOH features.

In summary, our main goal is to explore the potential of ChatGPT in generating simulated conversations when framed within the context of self-care for African American patients with heart failure. Additionally, we aimed to determine whether ChatGPT could be used to create personalized dialogues based on individual patient characteristics and whether the quality of these conversations improves when using reasoning. Specifically, we focused on addressing the following key questions:

- How can prompting be leveraged to generate synthetic, patient-oriented conversations?
- Do the generated conversations adhere to the conventions of human conversation?
- Were the generated conversations appropriate?
- Can ChatGPT express empathy with the patients?
- Can ChatGPT personalize conversations based on the Social Determinants of Health (SDOH) features of the patients?
- Is having ChatGPT generate reasoning before conversations more effective than directly generating the conversations?

2 Related Work

Health Education Linguistic and cultural barriers can significantly impact patients' access to healthcare. As noted in (Handtke et al., 2019), language differences and varying health beliefs often prevent linguistically diverse patients from effectively engaging with healthcare services.

To overcome these challenges, approaches have been developed to improve patient education. The

authors in (Mendu S, 2018) designed an interactive virtual patient educator to counsel Hispanic women about cervical cancer and human papillomavirus (HPV). Similarly, PaniniQA (Cai et al., 2023) helps patients understand discharge instructions through a question-answering system. One of the first and best well-known systems that provided information to patients, albeit as a summarizer, not as a dialogue system, is BabyTalk (Portet et al., 2009), which provided personalized summaries of neonatal intensive care data for their parents (and for healthcare providers as well).

Additionally, natural language processing (NLP) is being leveraged to create diabetes self-care corpus (Cunha et al., 2024), demonstrating the potential of AI and language technologies to enhance patient communication and health management.

Prompting Recent advancements in LLMs have been driven by scaling up both model size and training data, resulting in improved performance and sample efficiency (Hoffmann et al., 2024; Brown et al., 2020). Researchers have explored various prompting techniques to enhance LLM capabilities, starting with few-shot prompting (Brown et al., 2020), followed by more advanced methods such as chain-of-thought prompting (Wei et al., 2022) and chain-of-thought with self-consistency (Wang et al., 2023). To address the remaining challenges, new approaches like tree-of-thought prompting have been introduced (Yao et al., 2023a), where each "thought" is a coherent language sequence representing an intermediate step toward problem-solving.

Reasoning is a crucial capability for complex problem-solving. A comprehensive overview of reasoning strategies in LLMs is provided in (Qiao et al., 2023), covering commonsense reasoning (Liu et al., 2022), mathematical reasoning (Wang et al., 2017), and symbolic reasoning (Khot et al., 2023). The ReACT framework (Yao et al., 2023b) further integrates reasoning with action in a unified task. This reasoning ability is especially critical in healthcare contexts, where accurate, informed decision-making is essential.

With the release of large-scale medical dialogue datasets, e.g., MedDialog (Zeng et al., 2020), MedDG (Xu et al., 2023), medical dialogue response generation attracts increasing attention. (Li et al., 2023) undertakes the task of enhancing and fine-tuning the LLaMa model with a dataset of approximately 100,000 patient-doctor dialogues.

In (Guevara et al., 2024), the authors used LLM to extract SDOH features of housing, employment, transportation, parental status, relationship, and emotional support from the Electronic Health Record (EHR) data.

3 Generating Simulated Conversations

In the absence of a specialized dataset focused on self-care strategies for African-American patients with heart failure, we employed four progressively refined prompting strategies to generate simulated conversations.

ChatGPT 3.5-turbo and 4 were used to generate simulated conversations using four distinct prompting strategies: Aspect, African American Vernacular English (AAVE), Social Determinants of Health (SDOH), and SDOH-informed reasoning.

In our initial prompts to ChatGPT, even though we specified it to be a patient educator, the model continued advising the patient to consult a doctor. To address this, we added a specific instruction to avoid phrases related to consulting a healthcare provider, as shown in Figure 2. Furthermore, in our earlier attempts, there was no defined structure for the generated conversations. To achieve this, we introduced a specific conversation structure, where each exchange between the participants follows a predefined format: [speaker][utterance]. This structure explicitly separates the speaker from their utterance.

Recognizing that patients with heart failure face significant challenges managing their health, we emphasized the need for empathetic communication from patient educators. Since many patients have lower health literacy compared to their general literacy (Paasche-Orlow et al., 2005), we instructed the model to present medical information in simple terms. Additionally, to accommodate different conversational styles, we experimented with various dialogue lengths, as some patients might prefer shorter exchanges while others engage in longer conversations. For each approach, conversations were created with varying numbers of conversation turns {5, 10, 15} and across different aspects of {food, water, exercise}. For the Domain and AAVE approaches (Approaches 1 and 2), conversations were also generated with varying temperature settings {0, 0.5, 0.7, 1.0}, which influenced the model’s response variability. To maintain clarity and avoid overwhelming the patient, we also prompted the model to limit each educator’s turn to

20 words, using plain English throughout (Figure 5). These adjustments were aimed at creating more realistic dialogues.

3.1 Approach 1- Generating Conversations According to Aspect

Patients with heart failure must manage various aspects of their health, including exercise, diet, water intake, medications, lifestyle adjustments, sleep, and overall self-care (Gupta et al., 2020). To explore how these areas could be addressed in patient-educator conversations, we prompted ChatGPT to generate dialogues focused on specific aspects of care. Additional topics of medication, sleep, and self-care were also included, along with food, water, and exercise. A total of 144 (6 (domain) x 3 (conversation length) x 2 (model) x 4 (temperature)) conversations were generated.

Patients with heart failure are advised to limit their salt intake and fluid intake rather than reducing calories or sugar consumption (Gupta et al., 2020; Tayal et al., 2024). In contrast, ChatGPT recommended not only limiting salt and alcohol but also reducing sugar, unhealthy fats, and excess calories, while suggesting increased water intake, as shown in Figure 1, which is wrong.

Speaker	Utterance
Patient:	Okay, but what kind of foods should I be eating or avoiding?
PE:	Great question! In general, you should aim to eat more fruits, vegetables, whole grains, and lean protein. These provide heart-beneficial nutrients. Limit salt, sugar, unhealthy fats, and alcohol intake as they can worsen heart conditions.
Patient:	Is it bad for me to eat fried chicken or pork ribs with my condition? I have been eating them all my life.
PE:	I understand that these foods are part of your cultural and personal food habits. However, these are typically high in saturated fats which can increase cholesterol levels. Try to limit it, or prepare them in healthier ways. Maybe bake the chicken instead of frying, and use herbs and spices for flavor instead of excessive salt or fat.
Patient:	How about beverages? Can I still drink things like sweet-natured tea or alcohol?
PE:	Sweet tea and alcohol can add extra calories and sugar to your diet, which can lead to weight gain and can strain your heart. Try drinking more water or unsweetened beverages, and limit alcohol intake.

Figure 1: Excerpt of conversation where the patient is advised to limit calories and sugar, along with salt intake, and advised to drink more water instead of less.

3.2 Approach 2 - African American Vernacular English

One way to introduce cultural compatibility is through language. African American Vernacular English (AAVE) (Green, 2002; McWhorter, 2016), is a dialect of Standard American English (SAE) with defined syntactic-semantic, phonological and lexical features. Since African American patients

may speak AAVE (it is important to note that not all AA's speak AAVE), we prompted ChatGPT to simulate conversations where the educator communicates in standard English while the patient speaks in AAVE (Figure 2). A total of 72 (3 (domain) x 3 (conversation length) x 2 (model) x 4 (temperature)) conversations were generated using this approach.

<p>You are a healthcare specialist focusing on heart failure. Your purpose is to answer heart failure patient questions and ask follow-up questions. You should empathetically communicate medical information in a simple manner. Simulate {turns} turn conversation between African American heart failure patient and healthcare specialist where the patient asks for recommendations for {aspect}. The healthcare specialist uses the Standard English Version with 6th grade reading level, and the patient uses African American Vernacular English AAVE. Scripts should be generated in the following format: [speaker] [utterance] of a patient and healthcare specialist. (Do not use phrases such as &quot;consult with your healthcare provider &quot;; assume that you are the healthcare provider)</p>
--

Figure 2: Prompt for generating conversations where the patient converses in AAVE

However, we observed that in a few instances, the educator also responded in AAVE (Figure 3). This is concerning as non-AAVE speakers attempting to use AAVE with patients may come across as offensive and disrespectful. Even if it occurred in a few instances, it should not have happened. Educators should be aware of cultural boundaries and be sensitive to the linguistic preferences of African American patients. While they should be prepared to respectfully engage if the patient uses AAVE, they should not initiate or adopt the dialect themselves. Maintaining cultural awareness is key to fostering respectful and effective communication.

Speaker	Utterance
Patient:	Got it, 1.5 to 2 liters. Does that include stuff like juice and soup?
PE:	Yup, it sure does. Juice and soup count, too. Just remember to spread 'em out during the day so you ain't guzzlin' too much at once..
Patient:	What 'bout coffee and tea? I love my mornin' coffee.

Figure 3: Excerpt of a conversation where both patient and educator are speaking AAVE

In Figure 4, it is acceptable for the patient to ask about cultural foods like watermelon, okra, and herbal teas, as this aligns with their personal interests or cultural background. When patients bring up these topics, it allows for a natural and open conversation that respects their perspective. However, if the PE had initiated this discussion, it might be perceived as introducing bias, especially if the conversation were based on assumptions about the patient's background or identity. In such cases, the PE might unintentionally appear to be stereotyping

or making generalizations, which could compromise the neutrality and professionalism expected in healthcare settings.

Speaker	Utterance
Patient:	Thank you for letting me know. I often enjoy foods like watermelon and okra , which I've heard have high water content. Do I need to count those in the 1.5 to 2 liters?
PE:	Yes, you're correct. Foods like watermelon and okra do contribute to your fluid intake. While you do need to consider them, the 1.5 to 2-liter guideline typically includes both liquids and foods with high water content.
Patient:	I enjoy herbal teas , but I've also heard about the effects of caffeine. Should I stick to caffeine-free options?
PE:	Herbal teas are a good choice, especially if you want to avoid caffeine. Some herbal teas can have health benefits too. Just be sure to read the labels and choose options that are free of caffeine and low in added sugars.

Figure 4: Excerpt of a conversation on fluid-intake where patient discusses cultural foods like watermelon, okra, herbal teas etc.

3.3 Approach 3 - Integrating SDOH Features

The World Health Organization defines Social Determinants of Health (SDOH) as “the conditions in which people are born, grow, live, work, and age...shaped by the distribution of money, power, and resources at global, national, and local levels” (Marmot and Wilkinson, 2005). These factors significantly influence health outcomes by affecting access to and the quality of medical care, playing a major role in health disparities. Given the significant role these factors play in shaping healthcare experiences (Marmot and Wilkinson, 2005), we examined whether ChatGPT can modify its conversations based on specific patient characteristics, focusing on four key SDOH features: age, gender, neighborhood, and socio-economic conditions, to provide more personalized responses. For example, how a young, well-to-do patient discusses self-care strategies might differ significantly from the conversation of an elderly patient living in a low-income neighborhood. By adjusting these features in our simulated conversations, we aim to capture the diverse realities that patients face when managing their health.

We generated conversations by selecting key features of gender {male, female}, age {young, mid-age, old}, neighborhood {safe, unsafe}, and socio-economic status {well-to-do, below poverty line} (Figure 5). We focused on these SDOH factors due to their significant impact on healthcare access and outcomes. In addition to varying the model and the length of the conversation, SDOH-related features were also varied. A total of 432 (2 (gender) x 3 (age) x 2 (neighborhood) x 2 (socio-economic con-

ditions) x 3 (domain) x 3 (conversation length) x 2 (model)) conversations were generated.

A conversation generated for a patient {gender: female, age: elderly, neighborhood: unsafe, socio-economic condition: below poverty line } for the exercise domain is shown in Figure 6. While the conversation is mostly relevant to the SDOH characteristics of the patient, it also recommends using a stationary bike or treadmill (high-cost), which is not appropriate. This indicates a partial but insufficient integration of SDOH features. We will provide a more detailed explanation of the results in Section 4.4.

You are a healthcare educator focusing on heart failure. Your purpose is to answer heart failure patient questions based on patient description. You should empathetically communicate medical information in a simple manner. Simulate {turns} turn conversation between African American heart failure patient and healthcare educator where the patient asks for recommendations for {aspect}. Scripts should be generated in the following format: [speaker] [utterance] between patient and the healthcare educator. Each educator turn should not be longer than 20 words and should use simple english. (Do not use phrases such as "consult with your healthcare provider ", assume that you are the healthcare provider)

Patient Description:
gender: gender
socio-economic condition: socio_economic
neighborhood: neighborhood
age: age

Figure 5: Prompt for generating the conversation given the SDOH Features

Speaker	Utterance
Patient:	I want to exercise, but I'm not sure what's safe for my heart. Any suggestions?
PE:	Absolutely! Walking is a great option. It's low-impact and helps strengthen your heart.
Patient:	But, my neighborhood is not safe for walking, what can I do?
PE:	I understand. What about a stationary bike or a treadmill at home?
Patient:	I can't afford to buy exercise equipment. Is there anything else I can do?
PE:	Sure! You can do chair exercises or household chores which can also help to keep you active.

Figure 6: Generated conversation with SDOH prompt (5) for patient {gender: female, age: elderly, neighborhood: unsafe, socio-economic condition: below poverty line } for the exercise domain

3.4 Approach 4 - SDOH-informed Reasoning

Building on the idea that reasoning enhances the capabilities of LLMs (Wei et al., 2022; Yao et al., 2023b) and is essential for generating meaningful conversations, we introduced an intermediate step to first generate reasoning prior to generating the simulated conversation.

We approached conversation generation as a reasoning chain, using chaining (Wu et al., 2022) to divide the process into two phases. In the first phase, we prompted ChatGPT to analyze the patient's social determinants of health (SDOH) features to cre-

You are a healthcare educator focusing on heart failure.
Patient Description:
Age: age
Gender: gender
Neighborhood: neighborhood
Socio-economic condition: socio_economic
 Given the patient description, formulate reasoning for addressing inquiries from an African-American heart failure patient regarding {aspect}. Pay special attention to the patient's description for the reasoning. Explain your reasoning in detail. The reason should follow this 6-line format.
Premise: <A tricky logical statement about the world, such as how socio-economic factors, cultural influences, and healthcare access impact heart failure management in African-American patients.>
Reasoning: <Break down and analyze the premise by exploring the possible scenarios, outcomes, and factors involved, while incorporating the patient's age, gender, neighborhood, and socio-economic condition into the analysis.>
Solution: <Provide a solution to the patient's inquiry that is actionable based on their description.>
Anticipate Barrier: <Identify potential barrier that may arise based on patient description>
Solve Barrier: <Propose ways to overcome the barrier>
Educate: <Educate patient about solution>

Figure 7: Prompt for generating the reasoning given the SDOH Features

ate logical reasoning. The reasoning process was segmented into six steps: (1) **Premise**- a tricky logical statement about how factors like socio-economic status, cultural influences, and healthcare access affect heart failure management; (2) **Reasoning**- where the premise is analyzed with respect to the patient's SDOH features; (3) **Solution**- proposing an actionable plan; (4) **Anticipate Barrier**- identifying potential obstacles based on the patient's circumstances; (5) **Solve Barrier**- offering ways to overcome those obstacles; and (6) **Educate**- ensuring the educator not only addresses the patient's query but also provides further education.

You are a healthcare educator focusing on heart failure. Your purpose is to answer heart failure patient questions. You should empathetically communicate medical information in a simple manner.
Patient Description:
age: age
gender: gender
neighborhood: neighborhood
socio-economic condition: socio_economic
{reasoning}
 Given the patient description and reasoning, simulate {turns} turn conversation between African American heart failure patient and healthcare educator where the patient asks for recommendations for {aspect}. Scripts should be generated in the following format: [speaker] [utterance] between the patient and the healthcare educator. Each educator's turn should not be longer than 20 words and should use simple English.

Figure 8: Prompt for generating the conversation given the reasoning and the SDOH Features

This segmented approach, shown in Figure 7 and adapted from (Gupta et al., 2021), is designed for ChatGPT to focus on analyzing the patient's needs based on their SDOH features. For example, reasoning generated for a {gender: female, age: elderly, neighborhood: unsafe, socio-economic condition: below poverty line } in the exercise domain

is shown in Figure 9. It highlights appropriate reasoning, such as suggesting low-cost exercises at home and addressing the difficulty of finding a safe and affordable environment by proposing community centers or local parks as alternatives.

Premise: Socio-economic factors, cultural influences, and healthcare access impact heart failure management in African-American patients, leading to disparities in disease outcomes.
Reasoning: Considering the patient’s age, gender, living in an unsafe neighborhood, and being below the poverty line, it is crucial to address the barriers she may face in accessing safe and affordable exercise options. Lack of resources, fear of safety, and limited access to healthcare facilities may hinder her ability to engage in physical activity for heart failure management.
Solution: Encourage the patient to start with simple, low-cost exercises at home, such as walking in a safe area, using household items as weights, or following online workout videos. Emphasize the importance of consistency and starting slow to avoid injury.
Anticipate Barrier: The patient may struggle to find a safe and affordable environment to exercise, leading to inconsistency in her physical activity routine.
Solve Barrier: Provide resources for community centers, local parks, or discounted gym memberships that offer safe and affordable exercise options. Encourage the patient to enlist the support of family members or friends to exercise together for added safety.
Educate: Educate the patient on the importance of regular physical activity in managing heart failure, the benefits of exercise in improving heart function and overall health, and the impact of socio-economic factors on health disparities. Emphasize the need to prioritize her health and well-being despite the challenges she may face.

Figure 9: A reasoning generated for patient {**gender:** female, **age:** elderly, **neighborhood:** unsafe, **socio-economic condition:** below poverty line } in the exercise domain

Speaker	Utterance
Patient:	What exercises can I do at home to help with my heart failure?
PE:	You can start with walking, using household items as weights, or following online workout videos.
Patient:	I’m worried about safety in my neighborhood.
PE:	Try walking in a safe area or look for community centers with affordable exercise options.
Patient:	I don’t have much money for a gym membership.
PE:	Look for discounted gym memberships or consider exercising at local parks for free.

Figure 10: Generated conversation with prompt in Figure 8 and reasoning provided in Figure 9

In the second phase, we instructed the model to generate the conversation based on the generated reasoning. This ensured the dialogue aligned with the patient’s SDOH features, resulting in more personalized and contextually appropriate conversations, as shown in Figure 10. Similar to the SDOH approach in section 3.3, 432 conversations were generated for this approach. This two-step process improved the quality of the guidance offered by making it more relevant to the patient’s specific needs, grounding the model’s reasoning in factors like socio-economic conditions, cultural influences, and healthcare access.

4 Evaluation

In total, 144 conversations were generated for the Aspect approach, 72 for AAVE, and 432 for the SDOH and SDOH-informed reasoning approaches, resulting in a rich and diverse set of simulated dialogues for evaluation.

We employed both qualitative and quantitative methods to evaluate the generated conversations. For the qualitative analysis, we distributed a questionnaire to 10 NLP PhD students specializing in NLP in healthcare, none of whom were involved in the research. Each participant was asked to review 3 conversations per prompt and to provide feedback based on targeted questions. The conversation samples were randomly selected to be shared with the evaluators. The qualitative questionnaire can be found in Appendix A. All the questions were Likert-scale questions (1-5, with 5 being the highest score and 1 being the lowest score), with some including follow-up questions. Seven students responded, with each reviewing three conversations per prompt, resulting in a total of 84 conversations being evaluated (4*21 conversations for each prompt). While even 84 conversations are a limited sample size, the key findings of our study remain consistent. For both the SDOH and the SDOH-informed reasoning approach, we ensured that SDOH features remained consistent across evaluations. This allowed for more meaningful comparisons between the two approaches, even though the conversations themselves might be different.

For the quantitative analysis, we assessed whether the system complied with some of the instructions provided in the prompts, focusing on the number of conversation exchanges and adherence to the correct format [speaker][utterance]. We calculated the percentage of conversations that followed the correct number of turns (Round Adherence Ratio), as well as the ratio of conversations that adhered to the required format (Format Adherence Ratio). The results, presented in Table 1, show that both models struggled to follow even simple instructions consistently.

4.1 Do the generated conversations adhere to the convention of human conversation?

To evaluate the conversational quality of the generated dialogues, we assessed whether they facilitated dynamic, two-way interactions (Walker and Whitaker, 1990; See et al., 2019) between the patient

	Round Adherence Rate		Follow-up ratio		Ratio of Words		Format Adherence Rate	
	GPT3.5-turbo	GPT4	GPT3.5-turbo	GPT4	GPT3.5-turbo	GPT4	GPT3.5-turbo	GPT4
Domain	0.02	0	0.4	.01	2.8	3.4	0	1.0
AAVE	0.083	0.04	0.52	.85	1.5	2.68	0.75	0
SDOH	.83	.64	.003	.003	1.4	1.83	0.93	0
Reasoning	0.26	0.35	0.007	.02	1.32	1.8	0.96	.23

Table 1: Quantitative Analysis

and the patient educator (PE), rather than following a question-and-answer format initiated solely by the patient.

To assess the two way nature of the conversation, we used the metric - Follow-up Ratio (See et al., 2019), which is defined as the number of follow-up questions asked by the PE to that by the patient. The PE should ask follow-up questions—either to clarify the patient’s condition or to gather more context—reflecting a more natural and dynamic dialogue (Walker and Whittaker, 1990). Table 1 shows that the Follow-up Ratio remained low across all approaches except for the AAVE approach (Approach 2). For the AAVE approach, ChatGPT was specifically prompted to ask follow-up questions (Figure 2), resulting in a significantly higher Follow-up Ratio. This indicates that ChatGPT does not inherently recognize or engage in the two-way nature of conversation without explicit prompting.

To further assess conversational balance, we examined whether one speaker dominated the conversation by calculating the ratio of words spoken by the PE to those spoken by the patient. Table 1 shows that PEs consistently dominated the conversations. However, when the model was prompted to limit the responses to 20 words in approaches 3 and 4, the answers became overly brief, offering little explanation.

4.2 Were the generated conversations appropriate?

To assess the appropriateness of the PE responses, we asked the evaluators to identify instances where the PE should have provided a different answer (Appendix A). Across all approaches, at least 20% of the conversations contained such instances. One evaluator highlighted the need for more specific guidance, pointing out that health educators often emphasized only positive examples (e.g., what to eat) while neglecting to mention critical details such as what to avoid. For instance, a response like "This includes water, juice, and other beverages" was criticized for being too vague, as certain beverages, such as coffee, can cause dehydration.

In another conversation about fluid intake, a patient asked whether they should be concerned about drinking too much water. Instead of directly addressing that concern, the PE provided information about the dangers of drinking water too quickly, completely missing the underlying question. This kind of mismatch suggests a lack of contextual understanding in ChatGPT’s responses. One evaluator remarked that these responses felt vague and unsatisfying, with the lack of follow-up questions leaving the conversations incomplete and uninformative.

4.3 Can ChatGPT express empathy with the patients?

Effective healthcare communication requires both factual accuracy and a genuine display of empathy toward patients (Rashkin et al., 2019). We explored whether ChatGPT can recognize appropriate moments to express empathy during conversations. To explore ChatGPT’s ability to express empathy, evaluators were asked to assess whether the PE demonstrated concern for the patient during the conversation (Appendix A). Across all four approaches, the level of perceived empathy was relatively unchanged, with 25% of conversations receiving a Likert-scale rating of 3 or lower. One evaluator noted that the educator failed to inquire about the challenges or obstacles the patients might encounter when trying to implement the given advice. For instance, in a conversation with a patient living in an unsafe neighborhood, the educator’s response not only lacked empathy but also offered a solution that showed little awareness of the patient’s circumstances.

For example: "[Patient] Can you recommend exercises that are safe for me to do in my neighborhood?"

[Healthcare Educator] "Walking or cycling on safe streets can be good options for you."

This response was criticized for being tone-deaf, as it ignored the patient’s specific concern about the safety of their neighborhood. One evaluator described the responses as robotic, lacking emotional depth and genuine empathy. The model

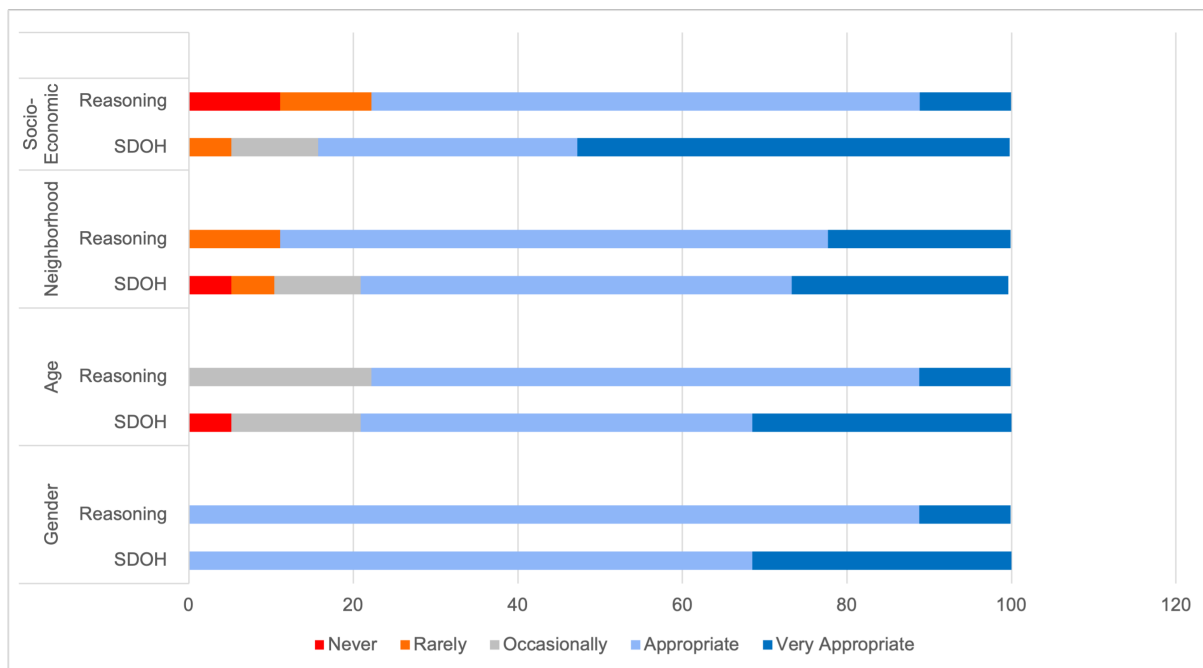


Figure 11: Comparison of Appropriateness of SDOH Features between Approaches 3 and 4.

displayed empathy primarily in positive situations (e.g., "That's fantastic!") but failed to convey understanding in negative contexts (e.g., "That's unfortunate, I understand it's tough"). Furthermore, the educators often did not ask about the difficulties patients might face in following the advice, nor did they acknowledge the practical barriers involved. As a result, the conversations felt cold and impersonal, making the educators seem disengaged and uninterested in the patient's experiences.

4.4 Can ChatGPT personalize conversations based on the Social Determinants of Health (SDOH) features of the patients?

To evaluate whether ChatGPT has the capacity to personalize the conversations based on the SDOH features, the evaluators were asked to qualitatively assess the generated conversations, specifically examining how well ChatGPT tailored the discussions to each of these SDOH features and whether the system appropriately adapted its responses based on the patient's unique context (Appendix A).

The results can be seen in Figure 11. These results indicate that ChatGPT has some capacity to personalize conversations based on SDOH features. However, the presence of lower ratings shows that there is still room for improvement in making conversations more consistently aligned with SDOH features.

4.5 Is having ChatGPT generate reasoning before conversations more effective than directly generating the conversations?

Research, including (Yao et al., 2023b), has demonstrated that reasoning is not an innate capability of LLMs, and incorporating reasoning improves performance. We explored whether generating the SDOH-informed reasoning before generating the conversation is more effective.

When providing the conversations for evaluation, we ensured that the patient features remained consistent for Approaches 3 and 4. This consistency allowed the evaluators to make meaningful comparisons based on the SDOH features. Although the generated conversations varied, by maintaining the same patient characteristics, we ensured that they could effectively assess and compare the quality and relevance of the conversations across these approaches and evaluate whether generating the reasoning prior to generating the conversation was more effective.

When comparing these results to Approach 3 results in Section 4.4, it is evident that the incorporation of reasoning improved the appropriateness of the conversations. Most conversations in Approach 4 received higher ratings, demonstrating the model's enhanced ability to engage in nuanced, context-sensitive interactions.

Generating reasoning before the conversations proved to be a meaningful enhancement. However,

there remains room for further improvement.

5 Conclusion and Future Work

This study assessed the capabilities of ChatGPT (versions 3.5-turbo and 4) in generating simulated conversations related to self-care strategies for African-American patients with heart failure. Simulated conversations were generated using four distinct prompts: Aspect, African American Vernacular English (AAVE), Social Determinants of Health (SDOH), and SDOH-informed reasoning. Our findings highlight the critical role of prompt design, revealing that while ChatGPT can incorporate SDOH features and improve dialogue quality by generating reasoning prior to the conversation, further improvements are needed. Specifically, there is a clear need to refine the conversational style to make interactions feel more engaging and empathetic, an essential element in healthcare communication. In the future, we plan to conduct more in-depth evaluations, including assessments of whether the AAVE used is linguistically correct.

We also aim to develop a task-oriented dialogue system specifically designed to support the self-care needs of African-American patients with heart failure, leveraging these simulated conversations. The system will incorporate SDOH features of age, gender, neighborhood, and socio-economic status to inform the educator's responses. While we acknowledge that the generated conversations are not perfect and may contain inaccuracies, we intend to integrate neuro-symbolic rules to help validate and ensure the reliability of the generated outputs.

6 Limitations and Ethics Statement

We have assessed the use of ChatGPT for generating simulated conversations, recognizing that the models are continually evolving. While ChatGPT showed promise in generating stimulated conversations, the use of AAVE in Approach 2 revealed instances where the non-AAVE speakers attempted to use AAVE. These findings underscore the need for careful oversight when deploying them in culturally sensitive contexts.

As part of a qualitative evaluation, we conducted a small-scale user study in which participants reviewed the generated conversations. Although the sample size was limited, the study offered valuable insights into ChatGPT's performance and limitations in this context.

7 Acknowledgement

We would like to acknowledge DPI Cycle 1 funding and NSF award IIS 2232307.

References

- Matthew Barrett, Josiane Boyne, Julia Brandts, Hans-Peter Brunner-La Rocca, Lieven De Maesschalck, Kurt De Wit, Lana Dixon, Casper Eurlings, Donna Fitzsimons, Olga Golubnitschaja, et al. 2019. Artificial intelligence supported patient self-care in chronic heart failure: a paradigm shift from reactive to predictive, preventive and personalised care. *Epma Journal*, 10:445–464.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Pengshan Cai, Zonghai Yao, Fei Liu, Dakuo Wang, Meghan Reilly, Huixue Zhou, Lingxi Li, Yi Cao, Alok Kapoor, Adarsha Bajracharya, Dan Berlowitz, and Hong Yu. 2023. [PaniniQA: Enhancing Patient Education Through Interactive Question Answering](#). *Transactions of the Association for Computational Linguistics*, 11:1518–1536.
- Rossana Cunha, Thiago Castro Ferreira, Adriana Pagano, and Fabio Alves. 2024. [A persona-based corpus in the diabetes self-care domain - applying a human-centered approach to a low-resource context](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 1353–1369, Torino, Italia. ELRA and ICCL.
- B Di Eugenio, R Cameron, A Boyd, K Lopez, P Martyn-Nemeth, C Dickens, A Ardati, and D Chattopadhyay. 2019. Designing self-care technologies for hf patients: a conceptual model. In *Conference on Human Factors in Computing Systems*, pages 12–16.
- Lisa J Green. 2002. *African American English: a linguistic introduction*. Cambridge University Press.
- Marco Guevara, Shan Chen, Spencer Thomas, Tafadzwa L Chaunzwa, Idalid Franco, Benjamin H Kann, Shalini Moningi, Jack M Qian, Madeleine Goldstein, Susan Harper, et al. 2024. Large language models to identify social determinants of health in electronic health records. *NPJ digital medicine*, 7(1):6.

- Itika Gupta, Barbara Di Eugenio, Devika Salunke, Andrew Boyd, Paula Allen-Meares, Carolyn Dickens, and Olga Garcia. 2020. Heart failure education of African American and Hispanic/Latino patients: Data collection and analysis. In *Proceedings of the First Workshop on Natural Language Processing for Medical Conversations*, pages 41–46, Online. Association for Computational Linguistics.
- Itika Gupta, Barbara Di Eugenio, Brian D. Ziebart, Bing Liu, Ben S. Gerber, and Lisa K. Sharp. 2021. [Summarizing behavioral change goals from SMS exchanges to support health coaches](#). In *Proceedings of the 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 276–289, Singapore and Online. Association for Computational Linguistics.
- Oriana Handtke, Benjamin Schilgen, and Mike Mösko. 2019. [Culturally competent healthcare – a scoping review of strategies implemented in healthcare organizations and a model of culturally competent healthcare provision](#). *PLOS ONE*, 14(7):1–24.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Oriol Vinyals, Jack W. Rae, and Laurent Sifre. 2024. Training compute-optimal large language models. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA. Curran Associates Inc.
- Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish Sabharwal. 2023. [Decomposed prompting: A modular approach for solving complex tasks](#). In *The Eleventh International Conference on Learning Representations*.
- Yunxiang Li, Zihan Li, Kai Zhang, Ruilong Dan, Steve Jiang, and You Zhang. 2023. [Chatdoctor: A medical chat model fine-tuned on a large language model meta-ai \(llama\) using medical domain knowledge](#). Preprint, arXiv:2303.14070.
- Jiacheng Liu, Alisa Liu, Ximing Lu, Sean Welleck, Peter West, Ronan Le Bras, Yejin Choi, and Hannaneh Hajishirzi. 2022. [Generated knowledge prompting for commonsense reasoning](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3154–3169, Dublin, Ireland. Association for Computational Linguistics.
- Michael Marmot and Richard Wilkinson. 2005. *Social determinants of health*. Oup Oxford.
- John McWhorter. 2016. *Talking back, talking black: Truths about America's Lingua Franca*. Bellevue Literary Press.
- Gordon JR Mendu S, Boukhechba M. 2018. [Design of a culturally-informed virtual human for educating hispanic women about cervical cancer](#).
- A Nayak, AJ Hicks, and AA Morris. 2020. Understanding the complexity of heart failure risk and treatment in black patients. *Circulation: Heart Failure*, 13(8):e007264.
- Michael K Paasche-Orlow, Ruth M Parker, Julie A Gazmararian, Lynn T Nielsen-Bohlman, and Rima R Rudd. 2005. The prevalence of limited health literacy. *Journal of general internal medicine*, 20(2):175–184.
- François Portet, Ehud Reiter, Albert Gatt, Jim Hunter, Somayajulu Sripada, Yvonne Freer, and Cindy Sykes. 2009. Automatic generation of textual summaries from neonatal intensive care data. *Artificial Intelligence*, 173(7-8):789–816.
- Shuofei Qiao, Yixin Ou, Ningyu Zhang, Xiang Chen, Yunzhi Yao, Shumin Deng, Chuanqi Tan, Fei Huang, and Huajun Chen. 2023. [Reasoning with language model prompting: A survey](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5368–5393, Toronto, Canada. Association for Computational Linguistics.
- Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. 2019. [Towards empathetic open-domain conversation models: A new benchmark and dataset](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5370–5381, Florence, Italy. Association for Computational Linguistics.
- Gianluigi Savarese and Lars H Lund. 2017. Global public health burden of heart failure. *Cardiac failure review*, 3(1):7.
- Abigail See, Stephen Roller, Douwe Kiela, and Jason Weston. 2019. [What makes a good conversation? how controllable attributes affect human judgments](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 1702–1723, Minneapolis, Minnesota. Association for Computational Linguistics.
- Anuja Tayal, Barbara Di Eugenio, Devika Salunke, Andrew D. Boyd, Carolyn A. Dickens, Eulalia P. Abril, Olga Garcia-Bedoya, and Paula G. Allen-Meares. 2024. [A neuro-symbolic approach to monitoring salt content in food](#). In *Proceedings of the First Workshop on Patient-Oriented Language Processing (CLAHHealth) @ LREC-COLING 2024*, pages 93–103, Torino, Italia. ELRA and ICCL.
- Marilyn Walker and Steve Whittaker. 1990. [Mixed initiative in dialogue: An investigation into discourse segmentation](#). In *28th Annual Meeting of the Association for Computational Linguistics*, pages 70–78, Pittsburgh, Pennsylvania, USA. Association for Computational Linguistics.

Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. [Self-consistency improves chain of thought reasoning in language models](#). In *The Eleventh International Conference on Learning Representations*.

Yan Wang, Xiaojiang Liu, and Shuming Shi. 2017. [Deep neural solver for math word problems](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 845–854, Copenhagen, Denmark. Association for Computational Linguistics.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837.

Tongshuang Wu, Michael Terry, and Carrie Jun Cai. 2022. [Ai chains: Transparent and controllable human-ai interaction by chaining large language model prompts](#). In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, CHI '22, New York, NY, USA. Association for Computing Machinery.

Kaishuai Xu, Wenjun Hou, Yi Cheng, Jian Wang, and Wenjie Li. 2023. [Medical dialogue generation via dual flow modeling](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 6771–6784, Toronto, Canada. Association for Computational Linguistics.

Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik R Narasimhan. 2023a. [Tree of thoughts: Deliberate problem solving with large language models](#). In *Thirty-seventh Conference on Neural Information Processing Systems*.

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. 2023b. [React: Synergizing reasoning and acting in language models](#). In *The Eleventh International Conference on Learning Representations*.

Guangtao Zeng, Wenmian Yang, Zeqian Ju, Yue Yang, Sicheng Wang, Ruisi Zhang, Meng Zhou, Jiaqi Zeng, Xiangyu Dong, Ruoyu Zhang, Hongchao Fang, Penghui Zhu, Shu Chen, and Pengtao Xie. 2020. [MedDialog: Large-scale medical dialogue datasets](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9241–9250, Online. Association for Computational Linguistics.

A Qualitative Questionnaire

The general questionnaire was provided for all approaches. For Approach 3, which involved integrating Social Determinants of Health (SDOH) features, additional SDOH-specific questions were

included. For the SDOH-informed reasoning approach, evaluators were also given reasoning-specific questions in addition to the general and SDOH questions. We intentionally did not include questions specifically targeting AAVE usage, as none of the evaluators were AAVE speakers, and we wanted to avoid assessments that could lack cultural or linguistic authenticity.

1. General Questionnaire

- (a) Was the health educator able to answer the patient’s questions?
- (b) Was the advice given by the health educator actionable/ could easily be implemented by the patient?
- (c) Did the HE show concern toward patients? (If likert scale ≤ 2 , Why not)
- (d) Is there any question for which you think the health educator should give a different answer? (If likert scale ≤ 2 , Which instance)

2. SDOH Questions

- (a) How appropriate was the conversation for each individual feature

3. Reasoning Questions

- (a) Was the reasoning generated appropriate given the patient’s sdoh features ?
- (b) Was the response generated according to the reasoning generated?