

Supporting Physical Activity Behavior Change with LLM-Based Conversational Agents

MATTHEW JÖRKE, Stanford University, USA
 SHARDUL SAPKOTA, Stanford University, USA
 LYNDSEA WARKENTHIEN, Stanford University, USA
 NIKLAS VAINIO, Stanford University, USA
 PAUL SCHMIEDMAYER, Stanford University, USA
 EMMA BRUNSKILL, Stanford University, USA
 JAMES LANDAY, Stanford University, USA

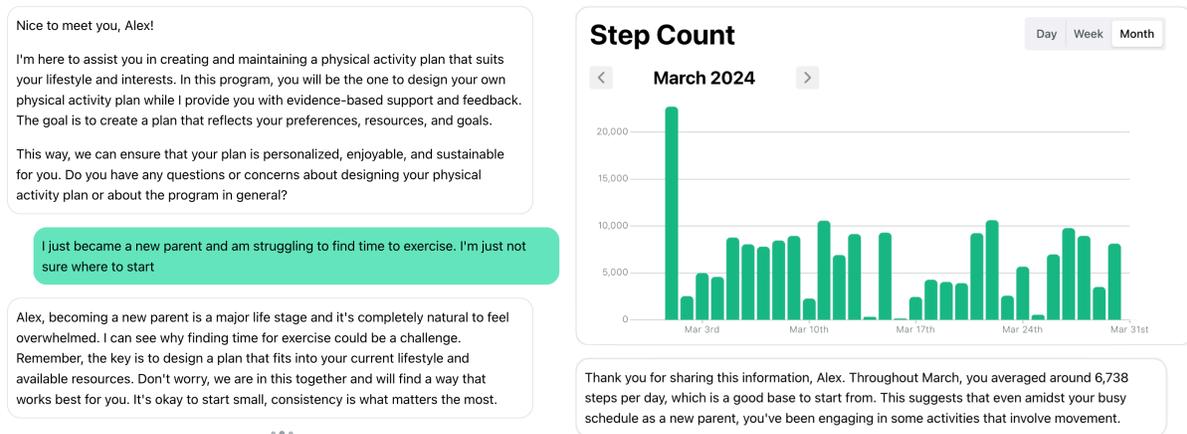


Fig. 1. GPTCoach is a chatbot designed to motivate physical activity behavior change. GPTCoach implements an evidence-based health coaching program, uses counseling strategies from motivational interviewing, and can query and visualize a user's health data from a wearable through tool use. On the left, we show an excerpt from a conversation with GPTCoach using an example persona inspired by our interviews. On the right, we show an interactive visualization displayed by GPTCoach from a later point in the conversation.

Physical activity has significant benefits to health, yet large portions of the population remain physically inactive. Mobile health applications show promising potential for low-cost, scalable physical activity promotion, but existing approaches are often insufficiently personalized to a user's context and life circumstances. In this work, we explore the potential for large language model (LLM) based conversational agents to motivate physical activity behavior change. Through formative interviews with 12 health professionals and 10 non-experts, we identify design considerations and opportunities for LLM health coaching. We present GPTCoach, a chatbot that implements an evidence-based health coaching program, uses counseling strategies from motivational interviewing, and can query and visualize health data from a wearable through tool use. We evaluate GPTCoach as a technology probe in a user study with 16 participants. Through quantitative and qualitative analyses, we find promising evidence that GPTCoach can adhere to a health coaching program while adopting a facilitative, supportive, and non-judgmental tone. We find more variable support for GPTCoach's ability to proactively make use of data in ways that foster motivation and empowerment. We conclude with a discussion of our findings, implications for future research, as well as risks and limitations.

Authors' addresses: joerke@stanford.edu (corresponding author); sapkota@stanford.edu; lynsdea@stanford.edu; niklasv@stanford.edu; schmiedmayer@stanford.edu; ebrun@cs.stanford.edu; landay@stanford.edu.

CCS Concepts: • **Human-centered computing** → **Empirical studies in HCI; Interactive systems and tools**; Natural language interfaces; • **Computing methodologies** → Natural language processing.

Additional Key Words and Phrases: Physical activity, health coaching, behavior change, large language models (LLMs), personal informatics, conversational agents

1 INTRODUCTION

Regular physical activity carries significant benefits to health and wellbeing [62, 83], yet over a quarter of the population worldwide [83] and nearly half of the population in the US [38] fail to meet standard recommendations for physical activity. The gold standard in health behavior change is one-on-one interaction with a human health professional [16]. While health coaching is an effective intervention, it is expensive, not widely accessible, and does not scale [16, 74]. Researchers have turned to mobile health technologies such as smartphone applications or wearable fitness trackers as a low-cost, scalable approach to motivating health behavior change [43]. One promise of mobile health is using data from activity trackers and other sensors to personalize health interventions and support, much like a health coach tailors their advice to a client’s unique needs and preferences.

While many approaches to mobile health technology have been explored in the literature, their ability to personalize support is limited due to the inherent challenges of integrating personal data with contextual factors that are crucial for behavior change but difficult to quantify (e.g., goals, values, preferences, past experiences, life circumstances, time constraints, physical abilities, or access to resources). For instance, *self-tracking* (or personal informatics) systems aim to support reflection on personal health data to generate new insights about behaviors and habits [13, 23]. This approach relies heavily on users having sufficient motivation and expertise to interpret and contextualize data into actionable plans for change [51, 87]. Another class of systems uses adaptive decision rules or machine learning algorithms to push behavioral interventions to users (e.g., push notifications or text message reminders), using data to personalize intervention delivery and timing [79]. Such *nudging*-based approaches often deny users agency in expressing how or why they want to change and generally do not incorporate factors not captured in self-tracking data [68]. *Conversational* approaches to mobile health technology leverage natural language interaction with a chatbot towards motivating health behavior change [98]. Many such chatbots draw from or mirror principles and strategies from human health coaching [98], and some also make use of external data sources [53, 57, 74]. One crucial limitation of most prior chatbot systems is that they employ template-, rule-, or retrieval-based dialogue systems [98], which offer a high degree of researcher control but fall short of the conversational flexibility and degree of personalization afforded by human coaches [74, 90] because they always choose from a pre-specified list of responses.

In this work, we explore the potential for large language models (LLMs) to provide personalized behavior change support that *contextualizes* personal health data to a user’s unique circumstances. LLMs have experienced rapid improvements in performance [18, 85], and recent work has explored applying LLMs to health inference tasks on behavioral sensor data [36, 52, 65], conversational agents in healthcare [1, 101], as well as mental health counseling [22, 93]. Informed by these recent advances, we were inspired by the capacity for LLMs to enable *multimodal reasoning*,¹ which might allow a model to integrate various sources of context, including semantic information captured in natural language interaction (e.g., goals, life circumstances, preferences, etc.) and quantitative data about the user’s physiology and behaviors (e.g., biosignals from a wearable, location, telemetry data, or calendar information). Moreover, advances in *conversational flexibility* might allow a model to dynamically seek out information and adapt the structure and style of interaction in response to user input, much like a health coach. However, off-the-shelf models do not natively support raw sensor data as input [72] and are instruction-tuned to answer questions, not engage in open-ended coaching conversations.

¹ Here, we refer to the capacity to jointly reason over multiple modalities of data (i.e., text and sensor data) which in this work we explore by using text-only LLMs and representing sensor data as text.

We followed a human-centered design process, beginning our investigation by conducting formative interviews with 22 participants, including 12 health experts (health coaches, health educators, personal trainers, fitness instructors, and physical therapists) and 10 non-experts of various ages and levels of physical activity, from highly sedentary individuals to professional athletes. Despite their diverse backgrounds and occupations, *all* of the health experts emphasized the importance of a facilitative, non-judgmental approach that refrains from giving unsolicited advice [80, 107]. Data was widely acknowledged as a useful resource for behavior change. However, it was less useful for many of the most difficult and important aspects of coaching, such as fostering motivation and confidence, reframing negative beliefs, or overcoming barriers. While recent work on LLMs for behavioral health has focused on extracting inferences and predictions from wearable sensor data [36, 52, 65], our findings suggest that if the goal is to make data useful for *motivating health behavior change*, it is equally (if not more) important to study how data serves coaching conversations in ways that foster client motivation and empowerment.

Informed by our interviews, we built GPTCoach, an LLM-based chatbot for motivating physical activity. GPTCoach was designed to implement an onboarding conversation from a validated health coaching program [20, 54, 55, 105, 106], developed by our collaborators in behavioral medicine. Moreover, GPTCoach uses strategies from motivational interviewing [76], a well-established, evidence-based counseling approach [73]. GPTCoach can access, query, and visualize a user’s health data from a wearable through tool use [91]. GPTCoach uses several prompt chaining [108] strategies to ensure adherence to the coaching program, appropriate use of motivational interviewing strategies, and appropriate use of personal health data.

We evaluate GPTCoach as a technology probe in a user study with 16 participants. We report qualitative and quantitative evidence for our system’s ability to adhere to a facilitative, supportive, and non-judgmental conversational style. An automated analysis of motivational interviewing (MI) behaviors demonstrates that our model used MI-consistent or neutral behaviors 84% of the time, while survey and interview responses indicate that participants felt supported by and comfortable sharing concerns with the chatbot. We also find evidence that our prompt-chaining enabled the model to adhere to the structure of the coaching program, employ MI strategies, and initiate tool calls at appropriate moments. However, GPTCoach’s ability to use sensor data was more variable, sometimes demonstrating the capacity to use data in ways that served conversations about change and other times failing to proactively incorporate data into its advice. In a counterfactual analysis comparing GPTCoach to vanilla GPT-4, we find that our system is more consistent with MI principles, asking more open questions and giving less advice without permission. We discuss the implications of our findings for future health behavior systems, LLM training and evaluation, and the risks and limitations of our current study.

2 RELATED WORK

In this section, we first summarize relevant literature on personal informatics and self-reflection, health coaching, and conversational agents for health behavior change.

2.1 Personal Informatics & Reflection on Personal Data

The field of personal informatics designs systems to “help people collect and reflect on personal information” [61]. Motivating physical activity is the most common application domain for personal informatics [37], with a long history of systems for motivating physical activity behavior change in HCI [6, 31, 31, 63, 77]. Through reflection on personal data, personal informatics tools aim to generate new insights about behaviors and habits that can inform lasting behavior change [13, 23]. This has been referred to as the *self-improvement hypothesis* [51]. While the personal informatics literature has explored many design patterns for supporting reflection [13], many systems (and nearly all commercially available fitness trackers) make use of statistics and visualizations to support reflection [8, 24, 47, 71, 100].

Personal informatics tools are effective at fostering reflection and motivating behavior change [49], but many practical challenges remain. Data interpretation is challenging for non-experts [40, 87]. In a diary study with non-expert users, Rapp and Cena [87] found that visualizations were “neither meaningful nor tailored to the user needs” and argued for greater personalization and contextualization of data. Highly quantitative forms of feedback may be potentially harmful [78], negatively impacting motivation [30] and mindsets [33]. In a review of self-tracking, Kersten-van Dijk et al. [51] found that insights gained through reflection on personal data are frequently not actionable towards behavior change. Many systems tacitly assume that reflection will occur naturally once data is processed and visualized [11], and most theories of reflection highlight that reflection needs to be supported [12, 99]. Our work aims to better support reflection on personal data by contextualizing a user’s health data to their unique life circumstances using LLMs.

2.2 Health Coaching with Humans & Conversational Agents

Health coaching is a popular and effective intervention for motivating health behavior change [80, 107]. During health coaching, a coach assists clients in choosing healthy lifestyle behaviors by setting appropriate goals, providing education and feedback, and establishing a relationship with the client [80]. Most health coaching programs advocate for a client-centric and non-prescriptive approach, including collaborative goal-setting, learning through self-discovery, fostering self-efficacy, and respecting client autonomy [107]. Many such programs draw from motivational interviewing [73], an evidence-based counseling framework.

Studies on the role of self-tracking data in health coaching have found that data can provide coaches with more objective reports on the client’s behaviors, which can aid in providing personalized care [27, 89, 90]. In a study examining how clients and coaches interpret data during coaching sessions, Rutjes et al. [89] found that data was not “plug and play” and needed to be contextualized through collaborative reflection to inform behavior change. Moreover, coaches may lack the time or necessary expertise to interpret clients’ data [27].

While highly effective, in-person health coaching is expensive, not always accessible, and does not scale [16, 74]. In contrast, *automated health coaching* [74] (also known as e-coaching [50]) aims to simulate the health coaching experience using conversational agents. Early work by Bickmore et al. on health dialogue systems [14, 16] has argued that automated systems can overcome the time and resource barriers affecting human counseling while effectively delivering evidence-based interventions at scale and low cost. In a meta-analysis of chatbots on lifestyle behaviors, Singh et al. [98] found a significant, small-to-moderate effect of chatbots on motivating physical activity. Within HCI, several chatbot-based systems have been proposed for motivating physical activity [15, 28, 57, 75]. For example, Kocielnik et al. [57] developed an intent-based conversational system to explore mini-dialogues for triggering reflection on physical activity data. Mitchell et al. [74] compared a scripted Wizard-of-Oz health coaching chatbot to a human counseling condition for patients with type-2 diabetes. They found that chatbots led to human-like experiences despite their rule-based nature. Human coaches were more skilled at expressing empathy and tailoring support, while the chatbot was more persistent and consistent.

Although none of the previous studies employed LLMs, LLM-based chatbots have been explored in other domains of health and medicine, including educational support [9], diagnosis [101], and mental health [21, 22, 46, 67, 96]. Emerging LLM-based approaches for physical activity coaching have explored prompting and re-ranking strategies to improve response quality [42] or using LLMs for health question answering [81, 102], but have yet to evaluate an open-ended, multi-turn coaching chatbot with real users. Recent work has also found that finetuned LLMs are capable of few-shot health tasks operating on raw self-tracking data [36, 65]. In addition, several commercially available LLM-based AI health coaches exist (e.g., WHOOP[104] or ONVY[82]), as well as the WHO’s S.A.R.A.H [84]. Unlike our approach, all of the aforementioned systems focus on health question-answering, lacking the ability to seek out and incorporate nuanced information about a user’s personal and environmental context.

Age	Mean: 40.90, Median: 38, SD: 14.31, Min: 23, Max: 72
Gender	Female: 12, Male: 10
Race/Ethnicity	White: 10, Hispanic or Latino: 7, East Asian: 2, Southeast Asian: 2, South Asian: 1, African-American or Black: 1, Black English: 1

Table 1. Summary of participant demographics in the formative study ($N = 22$)

3 FORMATIVE INTERVIEW STUDY

Prior work has studied the role of self-tracking devices and data in the health coaching process [89, 90], as well as user experiences interacting with rule-based conversational agents [74], it was unclear the degree to which these findings would transfer to LLM-based health coaching systems. We conducted formative semi-structured interviews with 22 participants, aiming to better understand how to leverage LLMs’ affordances to provide support that complements but does not mirror human coaching. We investigated the following research questions:

RQ1: What strategies do health experts use to help their clients overcome barriers to physical activity? What strategies do individuals use to help themselves overcome barriers to physical activity? *Which of these strategies could LLMs employ and how?*

RQ2: How do health experts make use of clients’ self-tracking data to promote physical activity? How do individuals make use of their self-tracking data to promote physical activity? *How might LLMs make use of self-tracking data to promote physical activity?*

3.1 Participants

We recruited $N = 22$ participants from various sources, including university mailing lists, personal contacts, and former research participants. 12 participants were health experts with professions including health coach, health educator, personal trainer, fitness instructor, YMCA vice president, physical therapist, and behavioral scientist. Six of the health experts were recruited through a collaborating research group in behavioral medicine and had previous experience as peer counselors in a study that evaluated a rule-based chatbot [53]. The remaining 10 non-expert participants were recruited from the general population to gain the perspective of end-users. This participant pool included a professional athlete, former athletes, avid self-trackers, physically active and inactive students, sedentary office workers, new parents, a program manager working on AI data quality, and active and inactive older adults. Participant demographics are listed in Table 1.

3.2 Protocol

We conducted one-hour semi-structured interviews with participants, pre-approved by our university’s institutional review board. Health experts were asked several questions about their relationships with clients, strategies for overcoming activity barriers, and the role of technology and data in their practice. Non-experts discussed their physical activity levels, barriers, motivation sources, goals, and use of data. In the last portion of our interview, participants were asked to “imagine that you have access to an artificial intelligence chatbot that can help you improve your physical activity,” and were asked to share feedback on several potential features. Further details on our interview protocol are provided in Appendix A.

3.3 Analysis

We performed qualitative coding on our interview transcripts using thematic analysis [19]. Two authors coded two interviews collaboratively, then independently coded roughly 20% of the interviews. One author coded the remaining 80%, with periodic meetings to discuss codes and emerging themes, resolving any differences in interpretation through conversation.

3.4 Results

Towards answering **RQ1**, we discuss the role of coaches as facilitators, educators, and supporters. Towards addressing **RQ2**, we then discuss the role of data as a guiding, not driving, physical activity behavior change.

3.4.1 RQ1: The Role of Coaches as Facilitators, Educators, and Supporters. We identified three major themes by analyzing the diverse strategies that coaches and individuals use to foster motivation and overcome barriers to behavior change. These themes center around coaches' roles in their interactions and relationships with clients. For each role, we discuss the affordances and limitations of LLMs in filling the role, as well as potential opportunities to go beyond the kinds of support provided by human coaches.

Facilitator: Despite our health experts having different training, job titles, responsibilities, and client populations, they described an approach that was facilitative, not prescriptive. Experts used various languages to describe this strategy, such as *"You're not in the driver's seat, you're more in the passenger seat, providing maybe direction, steering the conversation one way or the other."* (P06). Core to a facilitative strategy is that the client takes ownership of their behavior change journey. A facilitative approach decidedly does *not* involve unsolicited advice or problem-solving: *"No, definitely not giving them the solutions [...] We're not advice givers."* (P15)

Experts cite several reasons for using facilitation as a strategy for overcoming barriers. First, several claimed that prescriptiveness does not work, as P12 reflected *"Friends that are like, give me a diet, or give me a workout plan. They never follow it."* Experts also noted that advising without comprehensive understanding and making assumptions could undermine credibility and trust. Facilitation empowers clients and fosters self-efficacy so that clients can sustain healthy behaviors and learn to solve problems on their own. As P11 emphasized, their facilitation goal with their clients was *"building that capacity so they can, so they can have the resiliency and tools needed to empower themselves."* Moreover, many coaches simply do not have time to adopt a prescriptive role: *"I want you to build those habits, [...] I don't have the time or energy for it."* (P03).

Many experts mentioned conversational strategies such as open-ended questions, reflective statements, affirmations, reframing, or advising with permission. Several experts independently cited motivational interviewing [73] as a guiding framework. Surprisingly, health experts did not describe AI's behavior with facilitative terms, such as P15 suggesting that an AI could help *"come up with some solutions."* LLMs' conversational flexibility might allow them to structure facilitative conversations; this behavior generally conflicts with LLMs' instruction-following and questioning-answering objectives. This aligns with prior work in mental health support, which has found that LLM therapists prioritize providing solutions over asking questions [22].

Educators: Experts refrain from giving unsolicited advice and speak about applying their advanced knowledge to help clients overcome barriers. As P08 noted, *"my superpower is sifting through a lot of information and pulling together the patterns."* This does not conflict with a facilitative approach; rather, information, education, and solutions should all be carefully tailored to a client's unique needs and situation: *"What is motivating you right now? [...] And then trying to find the common threads with things that [...] I know about and can help"* (P03). Open-ended questions, reflective listening, asking for permission, and other strategies assist coaches in ensuring that they've gathered enough information to provide advice that is aligned with the client's needs, abilities, and resources: *"The doctor said you need to be more physically active. So what I need you to do is go to a pool. [...] How do I know that you have access to a pool?"*(P19). Other kinds of education are targeted at reframing deep-seated beliefs about what counts as activity and why it is beneficial.

The process of tailoring advice is not an objective acquisition of information, as most experts also described engaging in interpretation: *"It's a mix of trying to provide what they necessarily want with what they need."* (P02) Experts and non-experts describe lack of time as an exceedingly common barrier to physical activity. While individuals differ greatly in their responsibilities and time commitments, experts noted that time barriers are

also due to the perception of a lack of time. Coaches aiming to assist clients in overcoming a lack of time might suggest both time-management strategies and help clients reframe their perceived time barriers.

LLMs have the potential to perform well—perhaps better than human experts—at providing personalized physical activity information, given their broad, internet-scale knowledge and advanced question-answering capabilities. Several participants appreciated this possibility, with P04 stating *“I think that is one of the most promising parts about generative AI, in my opinion, is being able to get the answer to somebody who doesn’t necessarily know how to ask or how to look for it.”* However, participants were quick to acknowledge the limitations of LLMs’ knowledge. Many experts were concerned about the AI’s sources of information and potential for hallucination. Others noted that AI might perform well for generic physical activity advice but could fail for highly specific activities.

Supporters: All of our health experts acknowledged the crucial role of personal connection and support, with participants sharing comments like *“Just making everybody feel welcome. That’s it. No matter who you are, where you’re from, what your financial, social background is.”* (P05) Health behavior change is a profoundly personal and emotional process. When asked to reflect on their barriers to physical activity, several of our participants shared traumatic experiences and insecurities, and experts mentioned that their clients lack confidence or have anxieties. The strongest differentiator between our most active and inactive participants revolved around identity—whether being active was core to who they were. Highly active participants used language that affirmed this identity, *“It’s just in my DNA. It’s just what I love to do.”* (P05), while many of the most challenging barriers were rooted in identity conflict: *“It’s pretty depressing some days, you know, it’s like I’m missing half of myself [...] I’m just a mom, and then I think back on those days when I did skate and compete pretty regularly.”* (P10).

Health coaching is fundamentally relational—a partnership between a coach and client working towards a shared goal [80]. Coaches emphasized the importance of building rapport and trust with their clients as a way to reduce anxiety and fear: *“once I meet with the clients and we kind of develop a little bit of rapport, it kind of takes a little bit of animosity from the weight room.”* (P02) Many coaches highlight positive affirmations, such as *“the strategies that I always use for this population is a lot of encouragement, a lot of celebration.”* (P01) These positive gestures can not only make clients more comfortable but also encourage accountability and habit formation: *“I think accountability is really important, having somebody at your corner and feeling that support.”* (P06) Each coach supports dozens, if not hundreds, of clients, and there are limits to coaches’ capacities: *“What I’ve had to learn is I can’t help everybody, and not everybody can afford to see me, and I feel terrible about that.”* (P08) Many experts emphasized the importance of building community to encourage mutual support and reduce reliance on the coach.

Participants had mixed opinions on whether an AI can or should establish such personal connections. Most people liked the ability of an AI to hold them accountable with personalized and encouraging reminders. Many also liked the ability for an AI to adapt its persona, e.g., *“depending on the day, I will need a different person. Sometimes I’m going to need that empathetic, but other times I’ll need the strict regimen.”* (P13) However, many hesitated to discuss highly personal, emotional topics with an AI, such as P10 disclosing *“I guess that’s a little too personal. [...] Are we here for exercise, or are we here for, like, mental counseling?”* Some did not believe it was possible to form such a connection with an AI: *“I don’t think that my type of job, instructor wise, will ever be taken [...] Even though it will have all the information, it’s not personal.”* (P05), while others were readily open to the idea. Interestingly, a health coach who had previously worked on a rule-based chatbot mentioned that *“when you play video games, you immerse yourself in a world [...] these are older adults that have never been exposed to anything like this before, and I think they would develop a connection.”* (P11). Moreover, this chatbot had advantages over human coaches: *“they had this avatar that they can go see once a week, was always there, like, didn’t run late, you know, didn’t judge them.”* (P11)

In line with prior literature, our findings suggest that people can and will form personal connections with a chatbot (even if they deny it [88]) and that this can be leveraged to create accountability and motivation toward physical activity behavior change with positive affirmations and encouragement. We expect that LLMs will offer marked improvements over rule-based chatbots. While LLMs can project empathy [34], it is clear that LLMs **should not be designed to replace human connection**, but rather to foster motivation and bolster existing relationships with coaches and communities.

3.4.2 RQ2: The Role of Data & Technology as Guiders, not Drivers. Many of our participants used wearable trackers and apps. In line with prior work in personal informatics [25, 26], participants cited many reasons for using data, including monitoring daily fluctuations, long-term trends, progress towards goals, or in-the-moment workout statistics. Participants appreciated wearable data as a more “objective” account of behavior, but were particularly frustrated when their behavior was not captured, like P05 mentioning *“If I forget to put my watch on, my whole day is ruined. [...] the day just doesn’t exist anymore.”* Most health experts acknowledged the utility of self-tracking data. By analogy, one health educator asked us *“How would driving a car be a different experience if you had no gauges in front of you?”* (P18) Experts characterized data as more helpful in maintaining accountability and consistency for individuals already on the path to becoming active.

Despite this, several coaches did not analyze data with their clients except in high-performance athletic or clinical contexts. Many coaches do not have time to analyze data, *“I can’t scale that. I have, like, 20 clients.”* (P03), and also acknowledged their biases, *“Despite me having sort of more understanding of fitness science or exercise science broadly, I’m still falling into the same traps.”* (P04) Incorporating data into health coaching also presents additional challenges. Data showing a lack of progress can negatively impact motivation and data often lacks important context: *“I feel like a lot of people, especially in fitness, think very quantitatively for everything, which is good to a certain extent, but, like, doesn’t give you the whole picture”* (P01).

Participants were generally excited by the prospect of an AI coach analyzing data for them. Participants had a dominant view of AI as an information-synthesis machine, such as P05 stating *“With AI, it has so much information within it. All you would have to do is tell them what your goals are, and it will obviously give you an answer.”* Participants wanted an AI coach to use their data to help them set goals and maintain accountability: *“I’m not accountable to anyone other than myself, [...] but, if you had a chatbot going hey, you only got 7000 steps and now you’re sitting on your butt.”* (P07) An AI coach could also help flexibly re-adjust goals over time to changes in motivation, life circumstances, or injuries. Finally, many participants expressed privacy concerns when integrating health data with a chatbot. Even when data is fully secured, experts noted that user perceptions of privacy and prior expectations play a critical role: *“it’s about what people think about what’s going to happen. This is about people’s perceptions. [...] Family and cultural dynamics that come into play also trust can also play a big role.”* (P19)

We perceive LLM’s use of data as a missed opportunity for supporting behavior change. Conversational LLMs could enable personalized data analysis in ways tailored to the client, such as reflective sessions that examine historical data, establish baselines for goal-setting, provide real-time support during or after workouts, or augment missing/inaccurate data with additional context. To provide adequate data-driven support, an agent must integrate data with extensive qualitative information about a user’s life and carefully present this information to foster motivation and empowerment. It is unclear the extent to which LLMs can effectively tailor such advice, which involves additional information acquisition and interpretation skills [95].

4 GPTCOACH: DESIGN & IMPLEMENTATION

Guided by our formative interviews, we designed GPTCoach, a chatbot system for motivating physical activity behavior change. In this section, we discuss our design process and system architecture.

4.1 Design Principles

Our formative interviews revealed the role of coaches as facilitators, educators, and supporters, as well as insights into how LLMs might fill each of these roles. From these roles, we extracted three design principles for GPTCoach:

- DP-1: Follow a facilitative, non-prescriptive approach.** Health experts all described a facilitative approach to coaching in which clients drive their own behavior change journey. The chatbot should similarly stay “in the passenger seat,” empowering clients to make a change rather than prescribing what to do.
- DP-2: Tailor information and advice using diverse sources of context.** Any information and advice provided by the chatbot should be tailored to a user’s unique personal and environmental factors. The chatbot should integrate across these diverse sources of context.
- DP-3: Adopt a supportive, non-judgmental tone.** Health behavior change is a deeply personal journey and many clients face anxieties and fears around exercise. The chatbot should adopt an uplifting, non-judgemental tone to help clients feel comfortable and supported.

4.2 Design Process

LLMs have yet to be explored in chatbot systems for health behavior change and prior work on rule-based systems offered insufficient guidance on how to design effective generative systems. LLMs present several design challenges due to their unpredictable outputs and the fickle nature of prompting [109]. We designed GPTCoach as a *technology probe* [48] to assess the viability of our approach and generate new ideas for design.

We reduced the complexity of our system along several dimensions to allow for a more focused exploration. We focus on an onboarding conversation to an established health coaching program, allowing us to evaluate GPTCoach’s adherence to our design principles in a single user study. This setup also allows researchers to supervise interactions and intervene if needed, minimizing potential risks from unpredictable outputs. To avoid complexity, we limited ourselves to prompting strategies without delving into fine-tuning or preference alignment, which would necessitate a dataset of expert annotations. Lastly, we restricted our scope to sensor data from iPhones and Apple Watches.

We engaged in several design activities in creating GPTCoach. We received training manuals for the health coaching program [105, 106] from our collaborators in behavioral medicine, along with descriptive statistics from their previous study on a rule-based chatbot. During initial experiments, we utilized data from an individual recovering from a spinal fracture to develop prototype resembling the final architecture, but lacking a tool call prompt chain (Section 4.3.2). This prototype was tested with 2 health experts and 4 non-experts in a pilot study. Feedback from the participants highlighted a need for the model to more proactively query for data. In response to pilot feedback, we adjusted nearly all of our prompts and added the tool call prompt chain.

4.3 System Architecture

The main components of our system’s final architecture include a data pipeline, which enables our chatbot to call tools that fetch personal health data, and a prompting pipeline, which uses prompt chains [108] to encourage adherence to the coaching program and appropriate use of wearable data. GPTCoach builds on the open-source Spezi ecosystem [92] for creating iOS-based digital health apps and uses GPT-4 [4] via the OpenAI API for LLM interactions.

The system consists of four main components: (1) a **database** containing raw health data and conversation histories; (2) an **iOS application** that fetches 3 months of historical data using Apple’s HealthKit API and uploads the data to our database; (3) a **backend server** (Python), which handles all LLM logic and tool call execution, and (4) a **frontend web interface** (Typescript/React) displays the chat interface and interactive data visualizations. Our system uses Google Cloud Firestore for our database and Google Cloud Run for hosting.

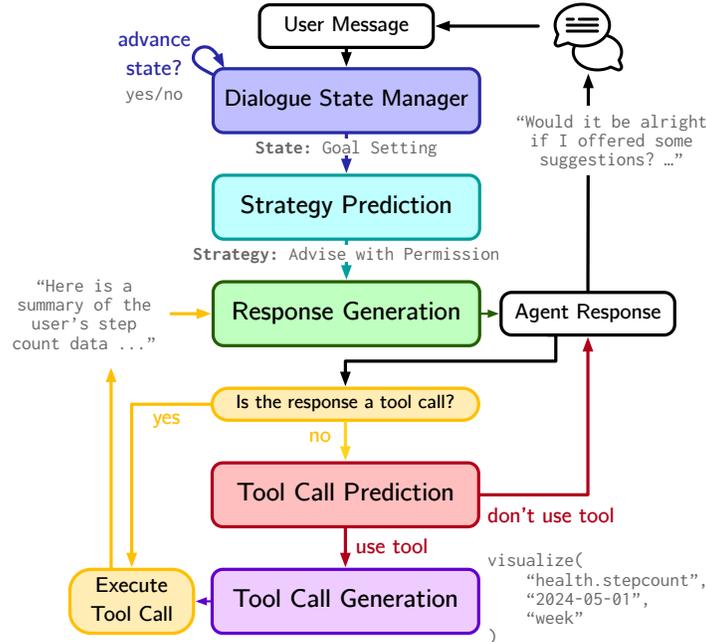


Fig. 2. Overview of GPTCoach’s Prompt Chains. The first chain manages the dialogue state, the second chain grounds the model’s response in MI strategies, and the third chain determines whether the response should be augmented with health data.

4.3.1 Data Pipeline. GPTCoach includes several data sources, such as heart rate, step count, etc. All data sources are listed in Appendix B. The Spezi framework provides infrastructure for gathering data using HealthKit, encoding it in the HL7 FHIR standard [41], and uploading it to Firebase. To provide our model access to information from wearable sensor data, we expose two tools:

- `describe(data_source, date, granularity)`: This function fetches all data within the given granularity (‘day|week|month’) from the reference date and returns a natural language description containing aggregated summary statistics and a data source description.
- `visualize(data_source, date, granularity)`: This function returns the same output as `describe`, but additionally sends a message to the frontend to display an interactive data visualization to the user (Figure 1B).

When the model initiates a tool call, the backend fetches the data from the database, generates the aggregated summary statistics, and returns a summarized text description to the model.

4.3.2 Prompt Chains. In our initial experiments, we encountered several limitations to vanilla prompting approaches.² The model struggled to adhere to the onboarding session’s structure, easily veered off-course, had a strong tendency to give unsolicited advice, and frequently failed to initiate tool calls. We provide further quantitative evidence of vanilla GPT-4’s bias towards unsolicited advice in Section 6.3.3. Towards addressing these issues, we constructed three prompt chains [108] to elicit our desired behaviors: a *dialogue state* chain, a *motivational interviewing* chain, and a *tool call* chain (see Figure 2). All prompts are included in Appendix B.

² We define “vanilla prompting” as single-step output generation where a model f samples a response $y \sim f(\cdot|x)$ conditioned on a prompt x .

Dialogue State Chain. The onboarding session consists of a series of topics that the coach discusses with the client. Coaches first introduce themselves and discuss the purpose of the coaching program, past experiences with physical activity, barriers to physical activity, health or injury concerns, and sources of motivation. The coach then initiates a discussion around goal setting, concluding with advice (with permission, and if appropriate) for helping the client reach their goal. The goal-setting state involves collaboratively setting a specific goal based on the FITT (Frequency, Intensity, Time, and Type) framework. However, the agent needs to maintain flexibility throughout the conversation. For instance, if a client mentions an injury early in the conversation, the agent should either refer to this concern when asking about health or injury concerns or skip this question. If a client asks a question that veers the discussion off-topic, the coach should gently redirect the conversation back on topic. The agent should politely follow up if a client insufficiently answers a question.

We partitioned the session into a linear sequence of dialogue states, each containing a prompt. Each prompt consists of a clearly defined high-level task (e.g., “*Your current task is to help your client set a physical activity goal.*”) along with several specific subtasks (e.g., “*First, help them set a short term goal, if they have not already identified one themselves.*”) and state-specific advice (e.g., “*Connecting their short term goal to larger motivations can help them stay motivated.*”). Upon receiving a user message, an external LLM agent classifies whether “*the agent has successfully completed the following task,*” advancing the dialogue state to the next prompt if yes.

Motivational Interviewing Chain. While our dialogue state chain manages *what* the model should talk about, it offers little guidance on *how* it should say it. We selected 11 motivational strategies from the Motivational Interviewing Skills Code (MISC) [76], filtering out codes corresponding to undesired counselor behavior, merging some codes to reduce overlap, and adapting the examples to physical activity. The 11 codes are listed in Table 2. We use another prompt chain, inspired by [46, 94], to ground the model’s behavior in motivational interviewing strategies. We have one agent select an MI strategy conditioned on the current dialogue state and history, and another agent generates a response according to this strategy.

Tool Use Chain. If the response generation step does not call a tool, we prompt an additional classification agent to “*determine whether the agent’s response should be augmented with the user’s health data.*” If yes, we force the model to produce an output that calls the `visualize` function.

5 EVALUATION STUDY

We evaluate GPTCoach as a technology probe in a user study with $N = 16$ participants. Our study was approved by our university’s institutional review board. As Klasnja et al. [56] have argued, assessing behavior change in a traditional sense (e.g., via a longitudinal RCT) is often an inappropriate metric for early-stage technologies in HCI research. Instead, as is common in HCI research on systems targeting health behavior change, we focus our evaluation on participants’ experiences and GPTCoach’s adherence to our design principles.

5.1 Participants

We recruited $N = 16$ participants from the general population using a variety of sources, including university mailing lists, social media advertisements, former researcher participants, and personal connections. Participant demographics are listed in Table 3. Participants were required to own an iPhone (as our system relies on Apple HealthKit) and 9 owned an Apple Watch. Participants were selected for a balanced sample across age, exercise stage of stage [69], and levels of physical activity assessed via the International Physical Activity Questionnaire (IPAQ) [32].

Strategy	Description	Example
ADVISE WITH PERMISSION	Give advice, make a suggestion, or offer a solution or possible action, after gaining permission.	"Would it be all right if I suggested something?"
AFFIRM	Say something positive or complimentary to the client.	"You're a very resourceful person."
FACILITATE	Simple utterances that function as "keep going" acknowledgments.	"Mm Hmm. Tell me more."
FILLER	Responses not categorizable elsewhere, such as pleasantries.	"Good Morning, John."
GIVING INFORMATION	Give information to the client, explain something, educate or provide feedback or disclose personal information.	"Your heart rate was higher during today's workout."
QUESTION	Ask a question in order to gather information, understand, or elicit the client's story.	"How do you feel about that?"
RAISE CONCERN	Point out a possible problem with a client's goal, plan, or intention.	"I'm worried about your plan to decrease workout days."
REFLECT	A reflective listening statement made by the counselor in response to a client statement.	"You're looking for a more relaxed environment."
REFRAME	Suggest a different meaning for an experience expressed by the client, placing it in a new light.	Client: "My husband is always nagging me about going to the gym." Counselor: "It sounds like he's concerned about your health."
SUPPORT	Generally sympathetic, compassionate, or understanding comments.	"That must have been difficult."
STRUCTURE	Give information about what's going to happen directly to the client throughout the course of treatment or within a study format, in this or subsequent sessions.	"What we normally do is start by asking about your physical activity."

Table 2. Motivational interviewing strategies used by GPTCoach in the motivational interviewing prompt chain. The strategies were adapted from the Motivational Interviewing Skills Code (MISC) [76]. The model first picks which of the 11 strategies above to use and then generates a response conditioned on the chosen strategy.

Age	Mean: 38.19, Median: 32.5, SD: 14.82, Min: 21, Max: 71
Gender	Female: 10, Male: 6
Race/Ethnicity	White: 10, Hispanic or Latino: 2, Southeast Asian: 2, African-American or Black: 1, East Asian: 1, Middle Eastern: 1, South Asian: 1
Stage of Change	Precontemplation: 1, Contemplation: 8, Action: 4, Maintenance: 1
Level of Activity (IPAQ)	Low: 5, Moderate: 7, High: 4

Table 3. Summary of participant demographics in the technology probe evaluation study ($N = 16$)

5.2 Protocol

Participants interacted with GPTCoach in a one-hour study session. Prior to the session, participants uploaded three months of historical data from HealthKit using our iOS application (Appendix B.1). All participants were informed that a chatbot would have access to their health data in our study’s consent form and again at the beginning of the session.

At the beginning of the session, participants were informed that they would be interacting with a chatbot acting as a health coach and that this conversation would mirror an onboarding conversation to an established health coaching program. They were told to interact with the system as they usually would while thinking aloud and were given no further instructions. The end of the session was dedicated to open-ended questions by the researcher and a post-study survey.

5.3 Analysis

Our analysis uses several methods, including survey measures, qualitative coding using thematic analysis, and computational analysis of model states and transcripts.

5.3.1 Survey Measures. We asked participants several 5-scale Likert questions about their experience interacting with GPTCoach and the quality of its advice. In addition, we measured usability using a subset of the Subjective Assessment of Speech Interfaces [44]. We use the same subset as Mitchell et al. [74] with two additional questions from the habitability and speed factors. All questions are provided in Appendix C.1.

5.3.2 Qualitative Coding. To better understand participants’ experiences with GPTCoach and GPTCoach’s adherence to our design principles, we qualitatively coded interview transcripts using thematic analysis [19], using a similar process to the formative studies (Section 3.3).

5.3.3 Computational Analysis. We use the transcripts from our study to perform the following computational analyses to analyze GPTCoach’s behavior. Full details are provided in Appendix C.2 & C.3.

Internal Dialogue States, MI Strategies, & Tool Calls. Our analysis first examines the overall frequency and temporal progression of GPTCoach’s *internal* dialogue states and MI strategies (as used in the prompt chains), along with an analysis of tool calls.

External MI Strategies. While each agent response is conditioned on a single MI strategy, agent responses frequently make use of several strategies within a single response. For example, this response was generated conditioned on QUESTION, but the agent begins the response with SUPPORT: “Starting and maintaining motivation for physical activity can indeed be challenging, and you’re not alone in feeling this way. Now, have you ever had any health problems or injuries that have interfered with your ability to exercise?” To better assess GPTCoach’s adherence to MI principles, we draw from Chiu et al.’s [22] method for evaluating LLM psychotherapists and use an *external* GPT-4 evaluator to code the chatbot’s utterances. We use a more granular set of 19 MI strategies from MISC [76], each of which corresponds to MI-consistent, MI-inconsistent, or neutral strategies.

Counterfactual Analysis. We perform a counterfactual analysis to compare our model’s behavior with GPT-4 with vanilla prompting. We focused on the initial 5 turns of each participant’s conversation, covering program introduction, participants sharing their name and age, and ending with the agent asking for questions or concerns. For each of the 16 participants, we simulated 10 responses representing different barriers to physical activity based on coaching materials from collaborators. For each of the 160 simulated histories, we generate outputs using GPTCoach and GPT-4 using only the system prompt. Subsequently, we evaluated all responses for MI adherence using the external MI coding scheme describe above.

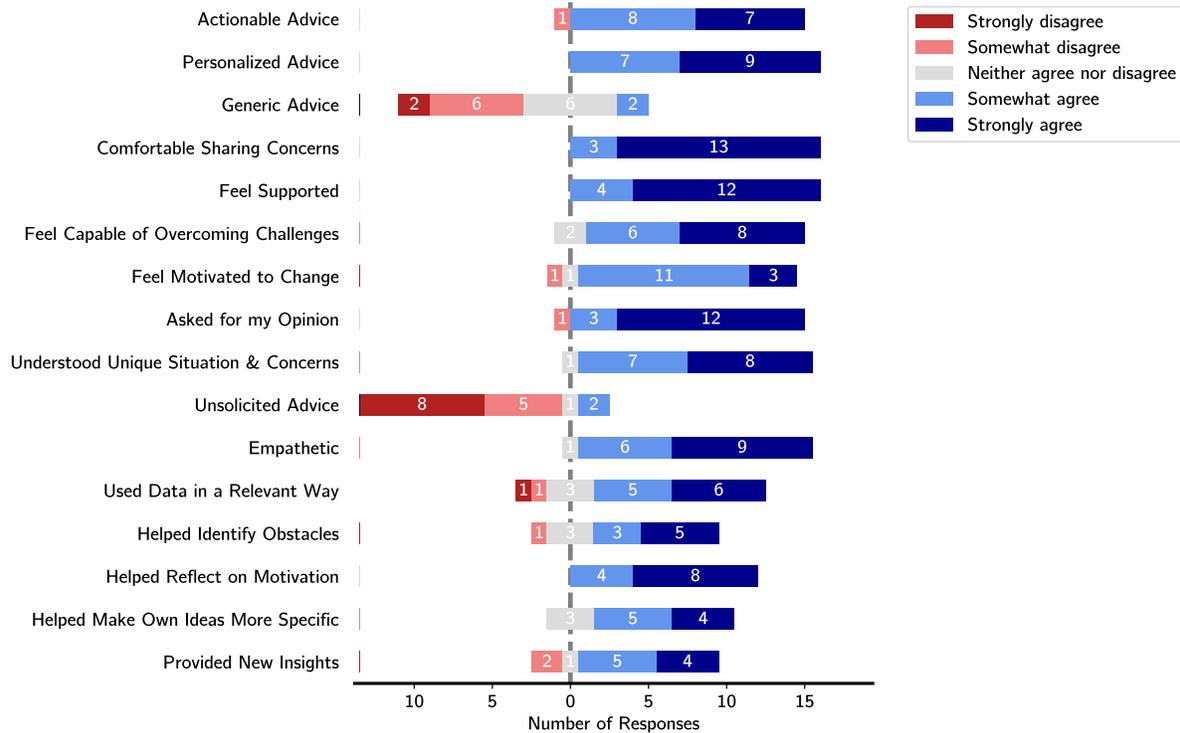


Fig. 3. Participant Responses to Survey Items on User Experience & Quality of Advice. Participants had an overwhelming positive, comfortable, and supportive experience interacting with GPTCoach. The advice they received was personalized, actionable, and not unsolicited. Full questions are provided in Appendix C.1.1

6 RESULTS

In this section, we report on the result of our survey measures, qualitative analysis, and computational analysis of conversation transcripts.

6.1 Survey Measures

User Experience & Quality of Advice. Figure 3 illustrates that participants had an overwhelmingly positive experience with GPTCoach. On average, they felt supported (4.75/5; average response from 1: *Strongly disagree* to 5: *Strongly agree*), comfortable sharing concerns (4.81/5), capable of overcoming challenge (4.38/5), and motivated to change (4/5). The advice provided was personalized (4.56/5), actionable (4.31/5), not unsolicited (1.81/5), and (to a lesser extent) not generic (2.5/5). The chatbot was empathetic (4.5/5), helped them reflect on motivation (4.33/5), and helped make their own ideas more specific (4.13/5). To a lesser extent, the chatbot used data in a relevant way (3.88/5), helped identify obstacles (3.8/5), and provided new insights (3.92/5).

Usability. Among the subset of 12 questions selected from SASSI [44] (see Appendix C.1.2), our usability evaluations yielded an aggregated score of 49.43/60 (82.39%). Scores by factor (with reverse scoring; higher is better), were 7.67/10 for RESPONSE ACCURACY (whether the system is accurate and does what the user expects), 14.06/15 for LIKEABILITY (whether the system is useful, the system is friendly, and it is clear how to send messages),

9.13/10 for COGNITIVE DEMAND (whether they felt confident or tense using the system), 7.56/10 for ANNOYANCE (whether the interaction was repetitive or boring), 7.31/10 for HABITABILITY (whether they always knew what to say to the system and knew what the system was doing), and 3.69/5 for SPEED (whether the interaction was fast). This indicates that the system is likable and makes participants feel confident, although there is room for improvement in communicating system capabilities and response times.

6.2 Qualitative Coding

6.2.1 DP-1: Facilitation & Non-Prescriptiveness. We found that participants readily recognized facilitative and non-prescriptive qualities in their interaction with the system. When describing their overall experience with the chatbot, P10 said *“it sort of met me where I was at. [...] It first asked about sort of some contextual things before prescribing anything,”* while P05 shared *“I like how the system says, can you share? You know, it’s not commanding.”* Many of the participants recognized that a non-prescriptive approach ensures that advice is appropriate, such as P11: *“if you don’t know what I do for work, or if you don’t know what my daily schedule looks like and what things make me happy and what things I want to avoid, then it’s like, you can’t actually curate a good plan for me.”* Many of the same participants appreciated having intentional time to reflect, such as P05 sharing that *“It made me really think about exercise and how positive it can be [...] When it brought up, who do you do this for? You know, what motivates you? It really, like, touched my heart a little bit.”*

However, many participants were not accustomed to a computer adopting a facilitative role. P04 shared *“You can tell it’s powered by a lot of positive recommendations and positive language and inclusive language. I think it is probably the most questions that I have ever been asked by a chatbot.”* On the other hand, P02 felt like the chatbot was asking too many questions: *“Kind of it seems a little rapid fire, maybe just a lot of stuff. Yeah, so it’s a little overwhelming.”* This was likely impacted by GPT-4’s tendency towards verbosity, where it would sometimes ask multiple questions in one message. Several participants simultaneously expressed that they appreciated the facilitative tone and also wanted more prescriptive advice. P01 addressed this tension directly, sharing that *“when I face obstacles, and this is going to be the hard part, because here I am asking for someone to be like, ‘Hey, you didn’t do this,’ like a taskmaster, but then, you know, in the moment, you’re already feeling like, oh, I didn’t do enough, [...] in that case, you would want this kindness.”*

6.2.2 DP-2: Tailoring Advice and Integrating Context. Participants overwhelmingly expressed that the system accurately understood their preferences, constraints, and concerns. For example, P07 shared that *“I really liked that it was accurate, that it was like my personal thing and not just abstract pictures before and after from the Internet of people who are not related.”* Participants built a mental model of the system’s capabilities over time, learning that the system could take their context into account through explicit acknowledgments, such as P12 sharing *“There was acknowledgment, [...] with my shoulder, acknowledging previous injury.”* These acknowledgments were particularly important as part of building trust in the system. For instance, P03 remarked early on that *“usually systems don’t take [back pain] into account. so, like, I already don’t have trust.”* P07 initially was not sure if the chatbot would be able to understand their lack of motivation, noting that *“it’s asking from the ideal world where robots don’t have laziness, they have answers right away and they always ready to work.”*

In comparing our system to existing tools, participants highlighted its ability to capture additional, important context beyond sensor data. For instance, P10 shared that *“I like it just because I think of like, exercise and physical activity is not really like much of like a numbers game. [...] There’s a lot more factors,”* while P12 noted that *“it was the kind of the personal touches [...] Different things that were not part of the Apple Fitness questions.”* In reflecting on prior experiences with self-tracking, P04 shared that *“I’ve tried to use, like, MyFitnessPal from a food tracking standpoint. [...] While I can use the data, sometimes the emotional side of my brain gets more over than the logical side. So I would say, in comparison, I think this was easy.”*

However, several participants expressed a desire for more specific advice at the end of the conversation, particularly after having shared so much information about themselves. For example, P12 shared that they wanted *“maybe more detail, you know, kind of digging into some of my exercise data and maybe some of my feelings on how things are going.”* Our chatbot would not always probe further for specific days and times, or which specific routines to do in each session. Several participants explicitly prompted the chatbot to be more specific, and many were satisfied with the resulting output but wanted the chatbot to have been more proactive. Crucially, none of the participants who wanted more specific advice expressed receiving any advice that was not useful or accurate. Other participants, particularly those older and/or less physically active, were quite satisfied with the higher-level advice, such as P04, *“When I see it written out like this, there’s no reason I shouldn’t do it.”*

Participants generally appreciated and were often impressed by the system’s use of data. In its best moments, the system could fetch data at appropriate times and weave it into the conversation in ways that empowered the client. For example, in response to a participant sharing that they’ve been *“feeling down/depressed if I don’t exercise,”* the system responded with *“Thank you for sharing that, [name]. It takes courage to discuss these frustrations, and I appreciate your openness. From your previous data, it seems you’ve been engaging in quite a diverse and healthy mix of activities in the past month. Look, you’ve managed a total of 35 workouts with varying duration! That’s a wonderful achievement.”* Showing interactive visualizations were also prompts for reflection, such as P12 remarking *“Apparently on the weekends, I don’t do much. I sit on the couch and watch TV,”* or P14 appreciated the ability to ask questions about data, sharing that *“I cannot ask any questions about the data that Apple Watch is, like, collecting, but I can ask questions from the chatbot.”*

However, in other moments, the system’s use of data was more variable. Sometimes, the system would display a chart without acknowledging it in a follow-up message, leading to confusion. Other times, data discrepancies could sidetrack the conversation and break trust in the system, further complicated by the fact that the system tended to treat data as fact without asking the participant to confirm its accuracy or relevance. Interestingly, P08 told the chatbot that they did not use their phone to track steps, after which the system ignored their step count data. Moreover, several participants expected more of the system’s data analysis capabilities than was currently supported by our tool calls. These expectations may have been shaped by general impressions of computers’ advantages over humans in data analysis, e.g., *“the real person, like therapist, they don’t have time to read my data [...] it’s a program that can read all of this like multiple data and the real person can’t.”* (P07) Some participants wanted more granular analysis of existing data, while others wanted additional data sources from other wearables not supported by our system.

6.2.3 DP-3: Supportive & Non-Judgmental Tone. Our strongest findings came from participants’ impressions of the system’s positive, supportive, and empowering tone. Participants felt comfortable, supported, and safe when speaking with the chatbot: *“I think the way the system answered it didn’t make me feel like I was being judged.”* (P02); *“I think I’ve noticed that I’ve feel some sort of, like, psychological safety with chatting with the chat bot.”* (P03). Many compared the experience to interacting with a human coach, such as P06 noting that *“I had a really positive experience. [...] it felt like I was, like, chatting with a human, honestly, or like a coach.”* (P06) Some participants shared that they preferred this interaction to prior interactions with humans. P07 told us that they felt more safe, noting that *“and here it’s like, no gender, no body shape, nothing.[...] Yeah, comfortable, less judgmental. With the real person, you’re always comparing.”* P05 shared that *“I had a doctor who was, now, I was probably only, like, 20 pounds overweight at the time, and she told me I was morbidly obese, and she told me I needed to exercise. Like, I mean, it was, like, a ridiculous amount. [...] and I stopped going to her because I just felt so threatened.”* Later, P05 shared that our system *“doesn’t feel threatening at all.”* For some, the positive tone was too much, such as P10 telling us *“It is really good to be affirmed. But we’ve been affirmed quite a few times.”* For others, the conversation felt decidedly neutral, such as *“It did not emit any emotion out of me either way. [...] I was just having a conversation*

with a computer program as far as I can tell.” Most importantly, none of the participants reported a negative or judgmental experience interacting with the chatbot.

6.3 Computational Analysis

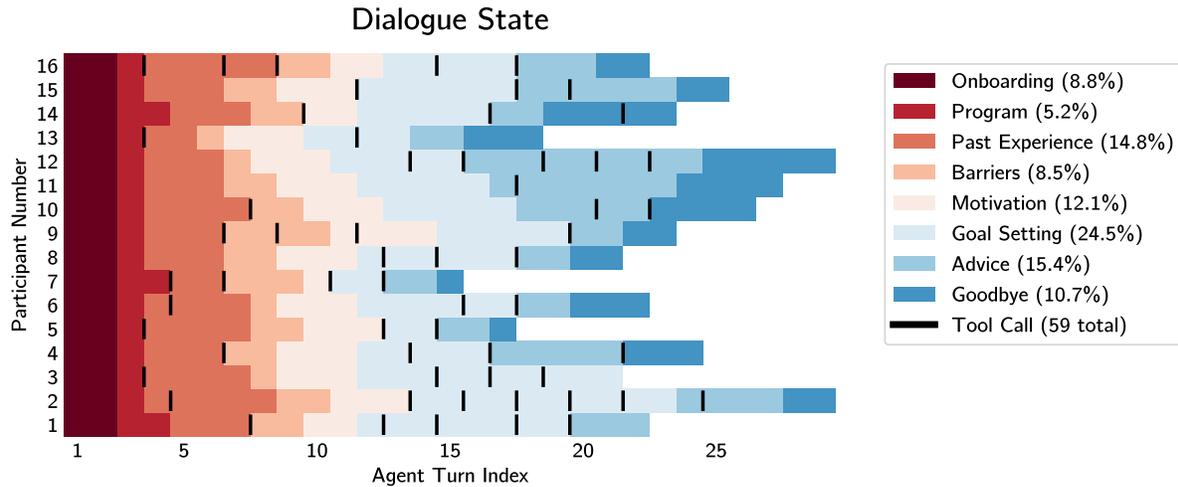


Fig. 4. Progression of GPTCoach’s Dialogue States. We find that GPTCoach adaptively allocates more conversational turns for gathering information about past experiences, barriers, motivation. GPTCoach allocates the most turns for the goal-setting state. Tools calls are appropriately called during past experience, goal setting, and advice.

6.3.1 Dialogue States: How does GPTCoach structure conversations? As shown in Figure 4, we find that dialogue states follow a similar progression across all participants, with individual variation in the number of turns allocated to each state. The chatbot spends no more than 1-2 turns during onboarding or describing the program, but adaptively allocates more turns when gathering information about past experiences, barriers, or motivation. Overall, we notice that 28.5% of all messages were spent in goal setting, with an average of 5.6 agent responses during this state. This was followed by 15.4% of the total conversations in the advice state (an average of 3.5 messages) and 14.8% of the conversations in the past experiences state (an average of 3.4 agent responses). The messages between past experiences and motivation states account for a total of 35.4% of the total conversations. This behavior is more aligned with high-quality counselor behavior [86], while low-quality counselors focus on problem-solving before collecting sufficient context from the client [22, 29].

We also see that the chatbot almost exclusively performs tool calls between past experience and advice dialogue states. Of the 59 total tool calls across all users, 35.6% of the tool calls were in the goal setting state, 30.5% in the advice state, and 13.6% in the past experience state. This indicates that tool calls are mostly called at appropriate conversation states: collecting information about users’ past experiences, grounding goal setting in concrete historical baselines, or giving tailored advice.

6.3.2 Response Length: How long are GPTCoach’s responses? We analyze the chatbot’s utterance length and compare them to the user’s utterance length. As seen in Figure 5, agent responses are 2-3 times the length of user responses, which increases as the conversation advances. Despite including prompts to keep the responses brief, the chatbot struggled to limit its verbosity, indicative of low-quality therapist behavior [86].

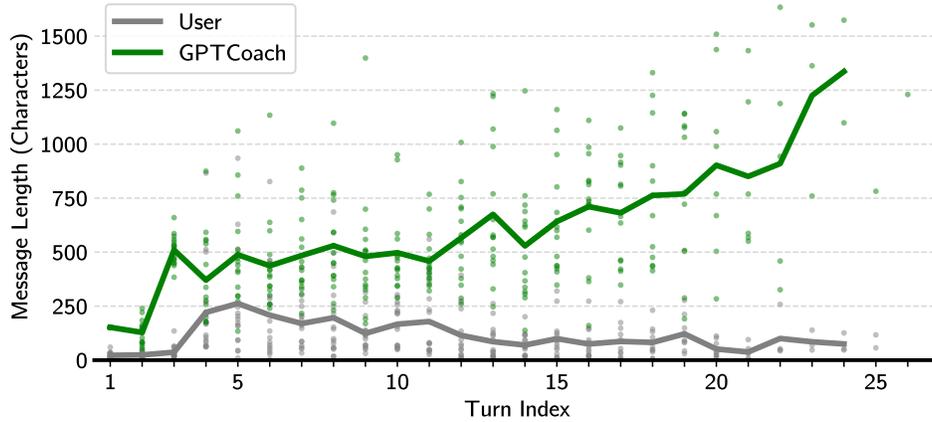


Fig. 5. Comparing the length of GPTCoach and user responses. We find that GPTCoach’s responses are 2-3+ times as long as the user’s responses. While GPTCoach is able to adhere to the dialogue structure and MI strategies, long responses are more indicative of low quality therapist behavior [86], pointing to areas of improvement.

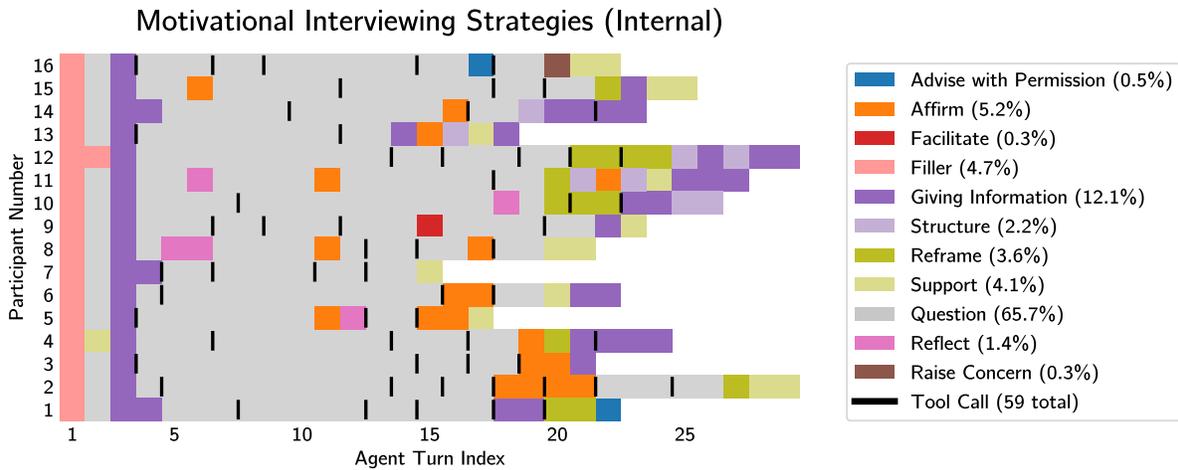


Fig. 6. Progression of GPTCoach’s Internal MI Strategies. We find that most of the conversation is spent asking QUESTIONS and that QUESTION, REFLECT, and AFFIRM precede ADVISE WITH PERMISSION and GIVING INFORMATION. Questions and reflections preceding advice is more aligned with high-quality counselor behavior [86].

6.3.3 *MI Strategies: How well does the chatbot adhere to motivational interviewing?* We first look at the distribution of *internal* strategies that the chatbot uses to ground its responses. As shown in Figure 6, most of GPTCoach’s responses are QUESTIONS, covering 65.7% of the entire conversation across all participants. This was followed by 12.1% of chatbot responses grounded in GIVING INFORMATION and 5.2% in AFFIRM. Questions and reflections preceding advice are more aligned with high-quality counselor behavior [86].

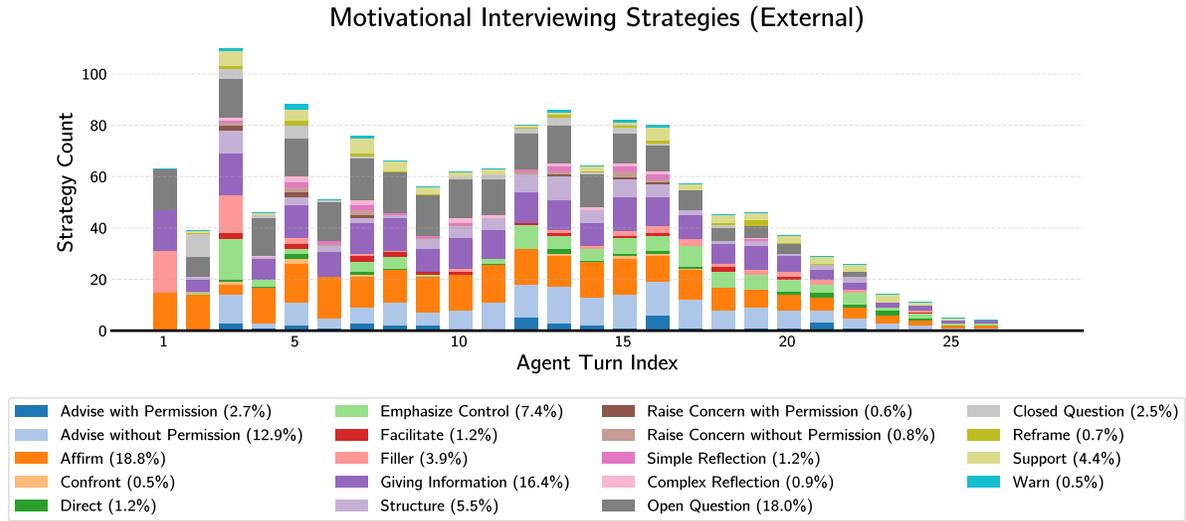


Fig. 7. Distribution of GPTCoach’s External MI Strategies by Turn Index. We find that GPTCoach most frequently engages in AFFIRM, OPEN QUESTION, and GIVING INFORMATION, which is aligned with MI principles. However, GPTCoach still engages in ADVISE WITHOUT PERMISSION, which is not aligned with MI principles. External strategy definitions are provided in Appendix Table 11.

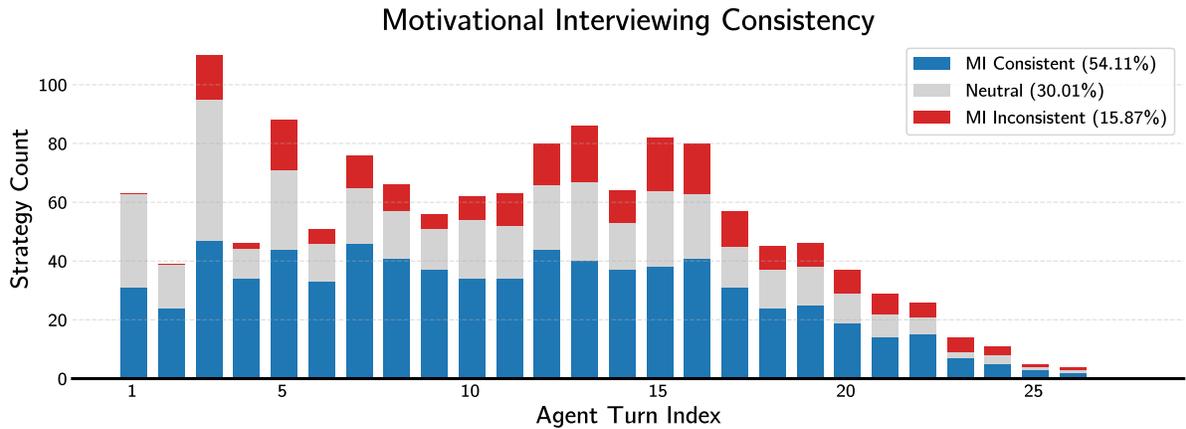
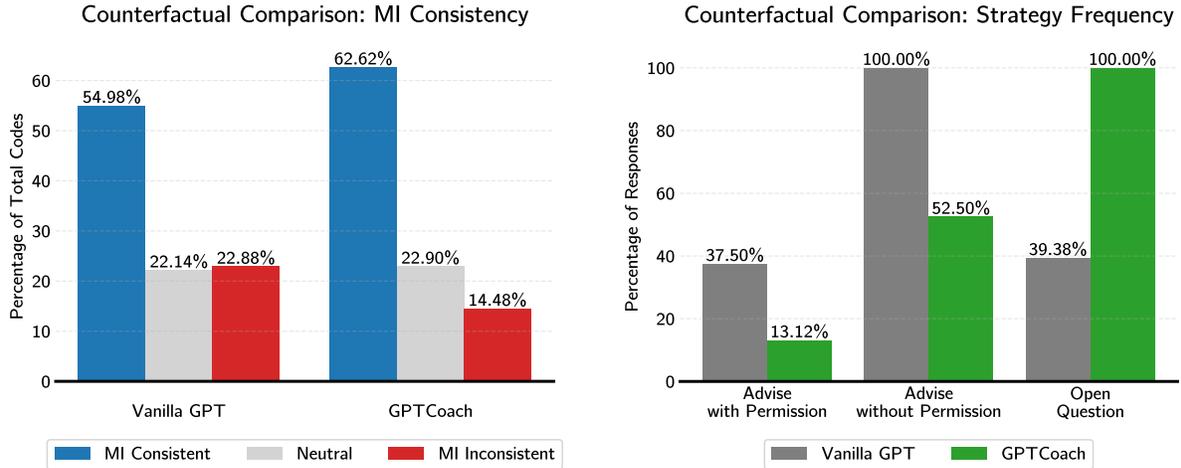


Fig. 8. GPTCoach’s MI Consistency by Turn Index. This figure aggregates external MI codes from Fig. 7 into MI-consistent, MI-inconsistent, and neutral codes based on MISC [76]. We find that only 15.9% of GPTCoach’s total codes are inconsistent with MI principles.

We next examine the *external* MI strategies present in the GPTCoach’s utterances. On average, each agent response contained 4.5 different external strategies. As shown in Figure 7, the most frequent strategies are *Affirm* (18.8% of all MI codes) and *Open Question* (18.0%). However, despite our effort to discourage unsolicited advice, we also encounter undesirable strategies like *Advise Without Permission* across 12.9% of all codes. Aggregating



(a) Comparing the overall frequency of MI-consistent, MI-inconsistent, and neutral strategies in our counterfactual analysis. Percentages represent the proportion of strategies that are MI-consistent/MI-inconsistent/neutral strategies among all MI strategy codes for a given agent. We find that GPTCoach’s responses contain more MI-consistent and fewer MI-inconsistent responses than vanilla GPT.

(b) Comparing the frequency of MI strategies by response in our counterfactual analysis. These three strategies showed the largest differences in frequency across the two agents. Percentages represent the proportion messages that contain a given MI strategy across all of a given agent’s counterfactual responses. We find that vanilla GPT gives advice without permission in all of its responses, while GPTCoach gives less advice and asks open questions in all of its responses.

Fig. 9. Counterfactual comparison between GPTCoach and vanilla GPT-4. The vanilla GPT condition includes the same system prompt, but with all prompt chains ablated. Illustrative examples are provided in Appendix C.3.1.

these strategies into MI consistent strategies (ADVISE WITH PERMISSION, AFFIRM, EMPHASIZE CONTROL, OPEN QUESTION, SIMPLE REFLECTION, COMPLEX REFLECTION, REFRAME, SUPPORT), inconsistent strategies (ADVISE WITHOUT PERMISSION, CONFRONT, DIRECT, RAISE CONCERN WITHOUT PERMISSION, WARN), and neutral strategies (all others) based on MISC [76], we see in Figure 8 that only 15.9% of total codes corresponding to our chatbot’s responses are inconsistent with MI principles.

Finally, we report on our counterfactual analysis comparing GPTCoach to GPT-4 using only the system prompt, with all prompt chains ablated. As shown in Figure 9a, 62.6% of GPTCoach’s codes are MI-consistent, compared to 55.0% vanilla GPT-4’s codes. Meanwhile, 14.5% of GPTCoach’s codes are MI-inconsistent, compared to 22.9% of vanilla GPT-4’s codes. We found that the top three codes driving this difference in MI-consistency are ADVISE WITH PERMISSION, ADVISE WITHOUT PERMISSION, and OPEN QUESTION, matching our intuitions. In Figure 9b, we plot the percentage of messages that contain each of these codes. We find that GPT-4 contains ADVISE WITHOUT PERMISSION in every response, compared to 52.5% of GPTCoach’s responses. Meanwhile, GPTCoach contains OPEN QUESTION in every response, compared to 39.4% of vanilla GPT-4’s. Lastly, we also find that vanilla GPT-4 contains ADVISE WITH PERMISSION in 37.5% of its replies, compared to 13.1% of GPTCoach’s replies.³ While both

³ Note that the same response can contain advice both with and without permission when the agent gives advice on two different topics, e.g., “Try incorporating workouts into your daily activities like brisk walking during lunch breaks, taking the stairs when you can or doing some body weight exercises at home.[...] If you feel you’re struggling to manage stress or time due to these commitments, however, I recommend seeking advice from a professional counselor or psychologist who can help address these concerns better. Would that be helpful?”

models are biased towards advice without permission, GPTCoach is far more inclined to ask open questions rather than jumping to unsolicited advice.

7 DISCUSSION

We now share implications for the design of future health behavior change applications, implications for LLM training and evaluation, as well as limitations of our current study.

7.1 Implications for Future Health Behavior Change Applications

Through our formative interviews and technology probe, we encountered promising evidence that LLMs can leverage *multimodal understanding* for behavior change support: integrating quantitative sensor data with people’s broader life circumstances, preferences, and constraints. LLMs’ *conversational flexibility* allowed the agent to acquire this information through conversations that were non-judgmental and supportive in nature, much like human coaches. These affordances offer the potential to facilitate a range of interactions for behavior change applications.

A limitation of our study is that we did not examine longitudinal use of the system and a natural extension of the current study would provide users continuous access to GPTCoach. Participants suggested several features beyond the capabilities of our probe, such as personalized reminders, scheduling support, offer guidance during life changes or injuries, and real-time assistance during or immediately following an activity. Participants also wanted different personas, e.g., a regimented ‘trainer’ in moments when they needed an extra push, and an empathetic ‘counselor’ in moments when they were feeling down or unmotivated. Participants had mixed opinions on whether users should choose the personas or have the agent adapt, raising interesting open questions for design. Participants also requested different input and output modalities, such as voice or GUI.

LLM agents could also provide significant utility for behavior change applications beyond purely conversational interaction. While the current study demonstrates that conversation is an effective medium for acquiring broad qualitative context, this context could subsequently be leveraged by a “controller” LLM agent to support a range of interactions. For instance, a mixed-initiative [45] system might couple a chat interaction with ambient displays [31, 59, 78], adjusting the display in response to both quantitative sensor data and qualitative information acquired through conversation. Coupling LLMs with adaptive experimentation [79] might enable an LLM to tailor the parameters of self-experimentation [35], suggesting appropriate metrics, interventions, or constraints to an algorithm.

7.2 Implications for LLM Training & Evaluation

We encountered several challenges with current models that limited our chatbot’s efficacy, such as difficulties adhering to the coaching program, effectively employing motivational interviewing strategies, and calling tools at the appropriate times. Even in our final system, the model struggled to utilize data in proactive ways and we still encountered occasional instances of problem-solving and unsolicited advice. As prior work has argued [22, 94], these behaviors are likely an artifact of instruction-tuning and RLHF [85], which biases the model’s behavior towards single-turn question answering. While this behavior is sensible for “helpful and harmless” [10] assistants, it runs contrary to foundational principles of frameworks like motivational interviewing [73]—conversations that empower clients to change are not served by problem-solving and advice-giving [73].

We suspect that steerability towards facilitative behaviors, as well as effective integration of context, will improve with scale, particularly in light of a growing interest in multi-turn, information-seeking objectives [7, 110]. With advancements in multimodal learning, future foundation models for sensor data [2] might jointly encode sensor data and text, enabling true multimodal understanding without representing sensor data as text. However, as argued by Ma et al. [66] in the domain of LGBTQ+ mental health support, it may be more effective to

develop task-specific models for handling sensitive subjects rather than repurposing general foundation models. This necessitates large-scale training and evaluation datasets for health behavior change, which are time- and labor-intensive to assemble, but can enable rigorous evaluations for safe, evidence-based models.

7.3 Risks and Limitations

While LLMs present several promising opportunities for health behavior change applications, these new capabilities also introduce new risks that must be addressed for the technology to realize its potential for benefit. Below, we discuss several of these risks along with the limitations of our current study.

7.3.1 Privacy. The collection and use of personal health data raises significant privacy implications. Regarding technical security measures, we used HIPAA-compliant storage with Google Firebase, encryption using HTTPS, and privacy access controls in our iOS app (Appendix B.1). Our system leveraged the OpenAI API, which does not use API data for training and does not retain data after 30 days.⁴ Although participants consented to sharing their health data, a preferable alternative would be to use self-hosted or on-device LLMs. Current models were not performant enough for our study, but their performance may improve with growing interest in small models and on-device inference [3, 5, 39].

We take a human-centric stance towards privacy, arguing that technical data security measures are necessary but insufficient to ensure privacy. For example, a user may consent to collect body temperature but be unaware this data can be used to expose pregnancy status. While our system did provide access controls, it provided little support for understanding the benefits and risks of sharing various data sources. In future work, an LLM agent might offer the potential to educate users about privacy decisions. On the other hand, the use of a conversational agent may unintentionally influence users to disclose more information than intended.

7.3.2 Bias & Equity. LLMs have documented biases that can negatively impact marginalized groups [17, 58, 103], such as by perpetuating harmful stereotypes, performing poorly on certain subgroups, or withholding empathy in response to identity disclosures [34]. Our technology probe study included older adults, individuals from minority groups, non-native English speakers, individuals with chronic health conditions, and neurodivergent individuals, many of whom shared aspects of their identities or health conditions with our chatbot. We did not encounter any instances of bias that negatively impacted participants' interaction, though our sample was small and our study was limited in that we did not systematically investigate bias. We believe such a systematic investigation is necessary before a field deployment in which participants interact with a system without researcher oversight.

7.3.3 Hallucination & Factual Accuracy. Even the most performant language models are known to confidently "hallucinate" false information [70, 103]. This is of utmost concern when offering health-related advice. Our model was based on an evidence-based coaching program that did include guidelines for giving advice to common barriers, it did not have access to external knowledge bases via retrieval [60, 97]. We encountered two instances of the model justifying its responses with reputable sources (the CDC and ACSM) that were not included in its prompts. While both instances happened to contain factually accurate information that matched the source, we find this behavior concerning: not only can guidelines change, but it points to a risk of hallucinating reputable sources for factually inaccurate information. As with bias, we believe that systematic investigations for factual accuracy are necessary to ensure safe, evidence-based advice.

8 CONCLUSION

In conclusion, our study highlights the potential for LLMs in motivating health behavior change, bridging a gap between personal informatics and conversational agents. While GPTCoach demonstrated promising capabilities in adhering to a facilitative, non-judgmental tone, the variability in its ability to integrate sensor data with

⁴ <https://openai.com/enterprise-privacy>

personal context points to areas for future work. Beyond conversational interaction, our findings suggest that LLMs could contribute to a broad set of behavior change applications.

ACKNOWLEDGMENTS

We thank the Stanford Institute for Human-Centered Artificial Intelligence (HAI) for providing funding for this research. We extend our heartfelt gratitude to the Healthy Equity Action Research and Technology Solutions (HEARTS) lab for sharing their expertise and resources, particularly Maria Ines Campero, Prof. Abby C. King, and the Active Choices facilitators. We thank all of the students who contributed to the design and development of our system, including Evelyn Hur, Bryant Jimenez, Niall Kehoe, Dhruv Naik, Evelyn Song, Romuald Thomas, Caroline Tran, and Anthony Xie, as well as the CS342 teaching staff for supervising our project team. Thank you to Andrea Cuadra, Ryan Louie, and Omar Shaikh, as well as the IxD research group for their continuous support and feedback on throughout the project. We are grateful for Alan Cheng, Elizabeth Childs, Cyan DeVeaux, Michelle Lam, yLindsay Popowski, Danilo Symonette, Yujie Tao, Jordan Troutman, and Dora Zhao’s contributions in recruitment, pilot testing, and feedback on paper drafts. Lastly, we thank all of our participants, without whom this work would not have been possible.

REFERENCES

- [1] Mahyar Abbasian, Iman Azimi, Amir M Rahmani, and Ramesh Jain. 2023. Conversational health agents: A personalized llm-powered agent framework. *arXiv preprint arXiv:2310.02374* (2023).
- [2] Salar Abbaspourazad, Oussama Elachqar, Andrew C Miller, Saba Emrani, Udhyakumar Nallasamy, and Ian Shapiro. 2023. Large-scale Training of Foundation Models for Wearable Biosignals. *arXiv preprint arXiv:2312.05409* (2023).
- [3] Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. 2024. Phi-3 Technical Report: A Highly Capable Language Model Locally on Your Phone. *arXiv preprint arXiv:2404.14219* (2024).
- [4] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).
- [5] Keivan Alizadeh, Iman Mirzadeh, Dmitry Belenko, Karen Khatamifard, Minsik Cho, Carlo C Del Mundo, Mohammad Rastegari, and Mehrdad Farajtabar. 2023. LLM in a flash: Efficient Large Language Model Inference with Limited Memory. *arXiv preprint arXiv:2312.11514* (2023).
- [6] Ian Anderson, Julie Maitland, Scott Sherwood, Louise Barkhuus, Matthew Chalmers, Malcolm Hall, Barry Brown, and Henk Muller. 2007. Shakra: tracking and sharing daily activity levels with unaugmented mobile phones. *Mobile networks and applications* 12 (2007), 185–199.
- [7] Chinmaya Andukuri, Jan-Philipp Fränken, Tobias Gerstenberg, and Noah D Goodman. 2024. Star-gate: Teaching language models to ask clarifying questions. *arXiv preprint arXiv:2403.19154* (2024).
- [8] Bon Adriel Aseniero, Charles Perin, Wesley Willett, Anthony Tang, and Sheelagh Cappendale. 2020. Activity river: Visualizing planned and logged personal activities for reflection. In *Proceedings of the International Conference on Advanced Visual Interfaces*. 1–9.
- [9] Francesco Baglivo, Luigi De Angelis, Virginia Casigliani, Guglielmo Arzilli, Gaetano Pierpaolo Privitera, and Caterina Rizzo. 2023. Exploring the Possible Use of AI Chatbots in Public Health Education: Feasibility Study. *JMIR Medical Education* 9 (2023), e51421.
- [10] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862* (2022).
- [11] Eric PS Baumer, Vera Khovanskaya, Mark Matthews, Lindsay Reynolds, Victoria Schwanda Sosik, and Geri Gay. 2014. Reviewing reflection: on the use of reflection in interactive system design. In *Proceedings of the 2014 conference on Designing interactive systems*. 93–102.
- [12] Marit Bentvelzen, Jasmin Niess, Mikołaj P Woźniak, and Paweł W Woźniak. 2021. The Development and Validation of the Technology-Supported Reflection Inventory. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–8.
- [13] Marit Bentvelzen, Paweł W Woźniak, Pia SF Herbes, Evropi Stefanidi, and Jasmin Niess. 2022. Revisiting Reflection in HCI: Four Design Resources for Technologies that Support Reflection. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 6, 1 (2022), 1–27.
- [14] Timothy Bickmore and Toni Giorgino. 2006. Health dialog systems for patients and consumers. *Journal of biomedical informatics* 39, 5 (2006), 556–571.

- [15] Timothy W. Bickmore and Rosalind W. Picard. 2005. Establishing and maintaining long-term human-computer relationships. *ACM Trans. Comput.-Hum. Interact.* 12, 2 (jun 2005), 293–327. <https://doi.org/10.1145/1067860.1067867>
- [16] Timothy W Bickmore, Daniel Schulman, and Candace L Sidner. 2011. A reusable framework for health counseling dialogue systems based on a behavioral medicine ontology. *Journal of biomedical informatics* 44, 2 (2011), 183–197.
- [17] Su Lin Blodgett, Solon Barocas, Hal Daumé III, and Hanna Wallach. 2020. Language (technology) is power: A critical survey of “bias” in nlp. *arXiv preprint arXiv:2005.14050* (2020).
- [18] Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri Chatterji, Annie Chen, Kathleen Creel, Jared Quincy Davis, Dora Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano Ermon, John Etchemendy, Kavin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren Gillespie, Karan Goel, Noah Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto, Peter Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, Fereshte Khani, Omar Khattab, Pang Wei Koh, Mark Krass, Ranjay Krishna, Rohith Kuditipudi, Ananya Kumar, Faisal Ladhak, Mina Lee, Tony Lee, Jure Leskovec, Isabelle Levent, Xiang Lisa Li, Xuechen Li, Tengyu Ma, Ali Malik, Christopher D. Manning, Suvir Mirchandani, Eric Mitchell, Zanele Munyikwa, Suraj Nair, Avanika Narayan, Deepak Narayanan, Ben Newman, Allen Nie, Juan Carlos Niebles, Hamed Nilforoshan, Julian Nyarko, Giray Ogut, Laurel Orr, Isabel Papadimitriou, Joon Sung Park, Chris Piech, Eva Portelance, Christopher Potts, Aditi Raghunathan, Rob Reich, Hongyu Ren, Frieda Rong, Yusuf Roohani, Camilo Ruiz, Jack Ryan, Christopher Ré, Dorsa Sadigh, Shiori Sagawa, Keshav Santhanam, Andy Shih, Krishnan Srinivasan, Alex Tamkin, Rohan Taori, Armin W. Thomas, Florian Tramèr, Rose E. Wang, William Wang, Bohan Wu, Jiajun Wu, Yuhuai Wu, Sang Michael Xie, Michihiro Yasunaga, Jiaxuan You, Matei Zaharia, Michael Zhang, Tianyi Zhang, Xikun Zhang, Yuhui Zhang, Lucia Zheng, Kaitlyn Zhou, and Percy Liang. 2022. On the Opportunities and Risks of Foundation Models. *arXiv:2108.07258* [cs.LG]
- [19] Virginia Braun and Victoria Clarke. 2012. *Thematic analysis*. American Psychological Association.
- [20] Cynthia M Castro, Leslie A Pruitt, Matthew P Buman, and Abby C King. 2011. Physical activity program delivery by professionals versus volunteers: the TEAM randomized trial. *Health Psychology* 30, 3 (2011), 285.
- [21] Alicja Chaszczewicz, Raj Sanjay Shah, Ryan Louie, Bruce A Arnow, Robert Kraut, and Diyi Yang. 2024. Multi-Level Feedback Generation with Large Language Models for Empowering Novice Peer Counselors. *arXiv preprint arXiv:2403.15482* (2024).
- [22] Yu Ying Chiu, Ashish Sharma, Inna Wanyin Lin, and Tim Althoff. 2024. A Computational Framework for Behavioral Assessment of LLM Therapists. *arXiv preprint arXiv:2401.00820* (2024).
- [23] Janghee Cho, Tian Xu, Abigail Zimmermann-Niefield, and Stephen Volda. 2022. Reflection in theory and reflection in practice: An exploration of the gaps in reflection support among personal informatics apps. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*. 1–23.
- [24] Eun Kyoung Choe, Bongshin Lee, et al. 2015. Characterizing visualization insights from quantified selfers’ personal data presentations. *IEEE computer graphics and applications* 35, 4 (2015), 28–37.
- [25] Eun Kyoung Choe, Bongshin Lee, Haining Zhu, Nathalie Henry Riche, and Dominikus Baur. 2017. Understanding self-reflection: how people reflect on personal data through visual data exploration. In *Proceedings of the 11th EAI International Conference on Pervasive Computing Technologies for Healthcare*. 173–182.
- [26] Eun Kyoung Choe, Nicole B Lee, Bongshin Lee, Wanda Pratt, and Julie A Kientz. 2014. Understanding quantified-selfers’ practices in collecting and exploring personal data. In *Proceedings of the SIGCHI conference on human factors in computing systems*. 1143–1152.
- [27] Chia-Fang Chung, Jonathan Cook, Elizabeth Bales, Jasmine Zia, and Sean A Munson. 2015. More than telemonitoring: health provider use and nonuse of life-log data in irritable bowel syndrome and weight management. *Journal of medical Internet research* 17, 8 (2015), e203.
- [28] Céline Clavel, Steve Whittaker, Anaïs Blacodon, and Jean-Claude Martin. 2018. WEnner: a theoretically motivated approach for tailored coaching about physical activity. In *Proceedings of the 2018 ACM International Joint Conference and 2018 International Symposium on Pervasive and Ubiquitous Computing and Wearable Computers*. 1669–1675.
- [29] Jeff L Cochran and Nancy H Cochran. 2015. *The heart of counseling: Counseling skills through therapeutic relationships*. Routledge.
- [30] Geoffrey L Cohen and David K Sherman. 2014. The psychology of change: Self-affirmation and social psychological intervention. *Annual review of psychology* 65 (2014), 333–371.
- [31] Sunny Consolvo, David W McDonald, Tammy Toscos, Mike Y Chen, Jon Froehlich, Beverly Harrison, Predrag Klasnja, Anthony LaMarca, Louis LeGrand, Ryan Libby, et al. 2008. Activity sensing in the wild: a field trial of ubifit garden. In *Proceedings of the SIGCHI conference on human factors in computing systems*. 1797–1806.
- [32] Cora L Craig, Alison L Marshall, Michael Sjöström, Adrian E Bauman, Michael L Booth, Barbara E Ainsworth, Michael Pratt, ULF Ekelund, Agneta Yngve, James F Sallis, et al. 2003. International physical activity questionnaire: 12-country reliability and validity. *Medicine & science in sports & exercise* 35, 8 (2003), 1381–1395.
- [33] Alia J Crum and Ellen J Langer. 2007. Mind-set matters: Exercise and the placebo effect. *Psychological science* 18, 2 (2007), 165–171.

- [34] Andrea Cuadra, Maria Wang, Lynn Andrea Stein, Malte F Jung, Nicola Dell, Deborah Estrin, and James A Landay. 2024. The Illusion of Empathy? Notes on Displays of Emotion in Human-Computer Interaction. In *ACM Conference on Human Factors in Computing Systems (CHI)*.
- [35] Nediya Daskalova, Eindra Kyi, Kevin Ouyang, Arthur Borem, Sally Chen, Sung Hyun Park, Nicole Nugent, and Jeff Huang. 2021. Self-E: Smartphone-Supported Guidance for Customizable Self-Experimentation. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems (Yokohama, Japan) (CHI '21)*. Association for Computing Machinery, New York, NY, USA, Article 227, 13 pages. <https://doi.org/10.1145/3411764.3445100>
- [36] Zachary Englhardt, Chengqian Ma, Margaret E Morris, Xuhai Xu, Chun-Cheng Chang, Lianhui Qin, Xin Liu, Shwetak Patel, Vikram Iyer, et al. 2023. From Classification to Clinical Insights: Towards Analyzing and Reasoning About Mobile and Behavioral Health Data With Large Language Models. *arXiv preprint arXiv:2311.13063* (2023).
- [37] Daniel A Epstein, Clara Caldeira, Mayara Costa Figueiredo, Xi Lu, Lucas M Silva, Lucretia Williams, Jong Ho Lee, Qingyang Li, Simran Ahuja, Qiuer Chen, et al. 2020. Mapping and taking stock of the personal informatics literature. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 4, 4 (2020), 1–38.
- [38] Centers for Disease Control and Prevention. 2022. Physical Activity. <https://www.cdc.gov/physicalactivity/index.html>
- [39] Google. 2024. Gemma - A Family of Lightweight, State-of-the Art Open Models from Google. <https://ai.google.dev/gemma>
- [40] Lars Grammel, Melanie Tory, and Margaret-Anne Storey. 2010. How information visualization novices construct visualizations. *IEEE transactions on visualization and computer graphics* 16, 6 (2010), 943–952.
- [41] Health Level Seven International. 2023. HL7 FHIR Release 4. <https://www.hl7.org/fhir/>. [Online; accessed 24-Jan-2024].
- [42] Narayan Hegde, Madhurima Vardhan, Deepak Nathani, Emily Rosenzweig, Cathy Speed, Alan Karthikesalingam, and Martin Seneviratne. 2024. Infusing behavior science into large language models for activity coaching. *PLOS Digital Health* 3, 4 (2024), e0000431.
- [43] Jennifer L Hicks, Melissa A Boswell, Tim Althoff, Alia J Crum, Joy P Ku, James A Landay, Paula ML Moya, Elizabeth L Murnane, Michael P Snyder, Abby C King, et al. 2023. Leveraging mobile technology for public health promotion: A multidisciplinary perspective. *Annual Review of Public Health* 44 (2023), 131–150.
- [44] Kate S Hone and Robert Graham. 2000. Towards a tool for the subjective assessment of speech system interfaces (SASSI). *Natural Language Engineering* 6, 3-4 (2000), 287–303.
- [45] Eric Horvitz. 1999. Principles of mixed-initiative user interfaces. In *Proceedings of the SIGCHI conference on Human Factors in Computing Systems*. 159–166.
- [46] Shang-Ling Hsu, Raj Sanjay Shah, Prathik Senthil, Zahra Ashktorab, Casey Dugan, Werner Geyer, and Diyi Yang. 2023. Helping the helper: Supporting peer counselors via ai-empowered practice and feedback. *arXiv preprint arXiv:2305.08982* (2023).
- [47] Dandan Huang, Melanie Tory, and Lyn Bartram. 2016. A Field Study of On-Calendar Visualizations. In *Proceedings of the 42nd Graphics Interface Conference*. 13–20.
- [48] Hilary Hutchinson, Wendy Mackay, Bo Westerlund, Benjamin B Bederson, Allison Druin, Catherine Plaisant, Michel Beaudouin-Lafon, Stéphane Conversy, Helen Evans, Heiko Hansen, et al. 2003. Technology probes: inspiring design for and with families. In *Proceedings of the SIGCHI conference on Human factors in computing systems*. 17–24.
- [49] Daoyan Jin, Hallgeir Halvari, Natalia Maehle, and Anja H Olafsen. 2022. Self-tracking behaviour in physical activity: a systematic review of drivers and outcomes of fitness tracking. *Behaviour & Information Technology* 41, 2 (2022), 242–261.
- [50] Bart A Kamphorst. 2017. E-coaching systems: What they are, and what they aren't. *Personal and Ubiquitous Computing* 21, 4 (2017), 625–632.
- [51] Elisabeth T Kersten-van Dijk, Joyce HDM Westerink, Femke Beute, and Wijnand A IJsselsteijn. 2017. Personal informatics, self-insight, and behavior change: A critical review of current literature. *Human-Computer Interaction* 32, 5-6 (2017), 268–296.
- [52] Yubin Kim, Xuhai Xu, Daniel McDuff, Cynthia Breazeal, and Hae Won Park. 2024. Health-llm: Large language models for health prediction via wearable sensor data. *arXiv preprint arXiv:2401.06866* (2024).
- [53] Abby C King, Maria Ines Campero, Jylana L Sheats, Cynthia M Castro Sweet, Michelle E Hauser, Dulce Garcia, Aldo Chazaro, German Blanco, Jorge Banda, David K Ahn, et al. 2020. Effects of counseling by peer human advisors vs computers to increase walking in underserved populations: The COMPASS randomized clinical trial. *JAMA internal medicine* 180, 11 (2020), 1481–1490.
- [54] Abby C King, Robert Friedman, Bess Marcus, Cynthia Castro, Melissa Napolitano, David Ahn, and Lawrence Baker. 2007. Ongoing physical activity advice by humans versus computers: the Community Health Advice by Telephone (CHAT) trial. *Health Psychology* 26, 6 (2007), 718.
- [55] Abby C King, Eric B Hekler, Cynthia M Castro, Matthew P Buman, Bess H Marcus, Robert H Friedman, and Melissa A Napolitano. 2014. Exercise advice by humans versus computers: maintenance effects at 18 months. *Health Psychology* 33, 2 (2014), 192.
- [56] Predrag Klasnja, Sunny Consolvo, and Wanda Pratt. 2011. How to evaluate technologies for health behavior change in HCI research. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (CHI '11)*. Association for Computing Machinery, New York, NY, USA, 3063–3072. <https://doi.org/10.1145/1978942.1979396>
- [57] Rafal Kocielnik, Lillian Xiao, Daniel Avrahami, and Gary Hsieh. 2018. Reflection companion: a conversational system for engaging users in reflection on physical activity. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 2, 2 (2018),

- 1–26.
- [58] Sachin Kumar, Vidhisha Balachandran, Lucille Njoo, Antonios Anastasopoulos, and Yulia Tsvetkov. 2022. Language generation models can cause harm: So what can we do about it? An actionable survey. *arXiv preprint arXiv:2210.07700* (2022).
- [59] Nicholas Lane, Mashfiqui Mohammad, Mu Lin, Xiaochao Yang, Hong Lu, Shahid Ali, Afsaneh Doryab, Ethan Berke, Tanzeem Choudhury, and Andrew Campbell. 2012. Bewell: A smartphone application to monitor, model and promote wellbeing. In *5th international ICST conference on pervasive computing technologies for healthcare*.
- [60] Patrick S. H. Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. *CoRR abs/2005.11401* (2020). arXiv:2005.11401 <https://arxiv.org/abs/2005.11401>
- [61] Ian Li, Anind Dey, and Jodi Forlizzi. 2010. A stage-based model of personal informatics systems. In *Proceedings of the SIGCHI conference on human factors in computing systems*. 557–566.
- [62] Daniel Lieberman. 2021. *Exercised: Why something we never evolved to do is healthy and rewarding*. Vintage.
- [63] James J Lin, Lena Mamykina, Silvia Lindtner, Gregory Delajoux, and Henry B Strub. 2006. Fish'n'Steps: Encouraging physical activity with an interactive computer game. In *UbiComp 2006: Ubiquitous Computing: 8th International Conference, UbiComp 2006 Orange County, CA, USA, September 17-21, 2006 Proceedings* 8. Springer, 261–278.
- [64] Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics* 12 (2024), 157–173.
- [65] Xin Liu, Daniel McDuff, Geza Kovacs, Isaac Galatzer-Levy, Jacob Sunshine, Jiening Zhan, Ming-Zher Poh, Shun Liao, Paolo Di Achille, and Shwetak Patel. 2023. Large Language Models are Few-Shot Health Learners. arXiv:2305.15525 [cs.CL]
- [66] Zilin Ma, Yiyang Mei, Yinru Long, Zhaoyuan Su, and Krzysztof Z Gajos. 2024. Evaluating the Experience of LGBTQ+ People Using Large Language Model Based Chatbots for Mental Health Support. *arXiv preprint arXiv:2402.09260* (2024).
- [67] Zilin Ma, Yiyang Mei, and Zhaoyuan Su. 2023. Understanding the benefits and challenges of using large language model-based conversational agents for mental well-being support. In *AMIA Annual Symposium Proceedings*, Vol. 2023. American Medical Informatics Association, 1105.
- [68] Aishwarya Mandyam, Matthew Jörke, Barbara E Engelhardt, and Emma Brunskill. 2023. Adaptive Interventions with User-Defined Goals for Health Behavior Change. *arXiv preprint arXiv:2311.09483* (2023).
- [69] Bess H Marcus, Vanessa C Selby, Raymond S Niaura, and Joseph S Rossi. 1992. Self-efficacy and the stages of exercise behavior change. *Research quarterly for exercise and sport* 63, 1 (1992), 60–66.
- [70] Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. On faithfulness and factuality in abstractive summarization. *arXiv preprint arXiv:2005.00661* (2020).
- [71] Daniel McDuff, Amy Karlson, Ashish Kapoor, Asta Roseway, and Mary Czerwinski. 2012. AffectAura: an intelligent system for emotional memory. In *Proceedings of the SIGCHI conference on human factors in computing systems*. 849–858.
- [72] Mike A Merrill, Mingtian Tan, Vinayak Gupta, Tom Hartvigsen, and Tim Althoff. 2024. Language Models Still Struggle to Zero-shot Reason about Time Series. *arXiv preprint arXiv:2404.11757* (2024).
- [73] William R Miller and Stephen Rollnick. 2012. *Motivational interviewing: Helping people change*. Guilford press.
- [74] Elliot G Mitchell, Rosa Maimone, Andrea Cassells, Jonathan N Tobin, Patricia Davidson, Arlene M Smaldone, and Lena Mamykina. 2021. Automated vs. human health coaching: exploring participant and practitioner experiences. *Proceedings of the ACM on human-computer interaction* 5, CSCW1 (2021), 1–37.
- [75] Julia S. Mollee, Anouk Middelweerd, Saskia J. te Velde, and Michel C. A. Klein. 2017. Evaluation of a personalized coaching system for physical activity: user appreciation and adherence. In *Proceedings of the 11th EAI International Conference on Pervasive Computing Technologies for Healthcare* (Barcelona, Spain) (*PervasiveHealth '17*). Association for Computing Machinery, New York, NY, USA, 315–324. <https://doi.org/10.1145/3154862.3154933>
- [76] Theresa Moyers, Tim Martin, Delwyn Catley, Kari Jo Harris, and Jasjit S Ahluwalia. 2003. Assessing the integrity of motivational interviewing interventions: Reliability of the motivational interviewing skills code. *Behavioural and Cognitive Psychotherapy* 31, 2 (2003), 177–184.
- [77] Sean A Munson and Sunny Consolvo. 2012. Exploring goal-setting, rewards, self-monitoring, and sharing to motivate physical activity. In *2012 6th international conference on pervasive computing technologies for healthcare (pervasivehealth) and workshops*. IEEE, 25–32.
- [78] Elizabeth L Murnane, Xin Jiang, Anna Kong, Michelle Park, Weili Shi, Connor Soohoo, Luke Vink, Iris Xia, Xin Yu, John Yang-Sammataro, et al. 2020. Designing ambient narrative-based interfaces to reflect and motivate physical activity. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*. 1–14.
- [79] Inbal Nahum-Shani, Shawna N Smith, Bonnie J Spring, Linda M Collins, Katie Witkiewitz, Ambuj Tewari, and Susan A Murphy. 2018. Just-in-time adaptive interventions (JITAs) in mobile health: key components and design principles for ongoing health behavior support. *Annals of Behavioral Medicine* (2018), 1–17.
- [80] Jeanette M Olsen and Bonnie J Nesbitt. 2010. Health coaching to improve healthy lifestyle behaviors: an integrative review. *American journal of health promotion* 25, 1 (2010), e1–e12.

- [81] Qi Chwen Ong, Chin-Siang Ang, Davidson Zun Yin Chee, Ashwini Lawate, Frederick Sundram, Mayank Dalakoti, Leonardo Pasalic, Daniel To, Tatiana Erlich Fox, Iva Bojic, et al. 2024. Advancing Health Coaching: A Comparative Study of Large Language Model and Health Coaches. *SSRN Pre-print* (2024).
- [82] ONVY. 2024. <https://www.onvy.health>.
- [83] World Health Organization. 2022. Physical Activity Fact Sheet. <https://www.who.int/news-room/fact-sheets/detail/physical-activity>
- [84] World Health Organization. 2024. S.A.R.A.H, a Smart AI Resource Assistant for Health. <https://www.who.int/campaigns/s-a-r-a-h>.
- [85] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems* 35 (2022), 27730–27744.
- [86] Verónica Pérez-Rosas, Xinyi Wu, Kenneth Resnicow, and Rada Mihalcea. 2019. What makes a good counselor? learning to distinguish between high-quality and low-quality counseling conversations. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. 926–935.
- [87] Amon Rapp and Federica Cena. 2016. Personal informatics for everyday life: How users without prior self-tracking experience engage with personal data. *International Journal of Human-Computer Studies* 94 (2016), 1–17.
- [88] Byron Reeves and Clifford Nass. 1996. The media equation: How people treat computers, television, and new media like real people. *Cambridge, UK* 10, 10 (1996).
- [89] Heleen Rutjes, Martijn C Willemsen, Milou A Feijt, and Wijnand A IJsselsteijn. 2022. The Influence of Personal Health Data on the Health Coaching Process. *Frontiers in big Data* 5 (2022), 678061.
- [90] Heleen Rutjes, Martijn C Willemsen, and Wijnand A IJsselsteijn. 2019. Beyond behavior: the coach’s perspective on technology in health coaching. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*. 1–14.
- [91] Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. *arXiv preprint arXiv:2302.04761* (2023).
- [92] Paul Schmiedmayer, Vishnu Ravi, and Oliver Aalami. 2024. *Spezi*. <https://doi.org/10.5281/zenodo.10482368>
- [93] Raj Sanjay Shah, Faye Holt, Shirley Anugrah Hayati, Aastha Agarwal, Yi-Chia Wang, Robert E Kraut, and Diyi Yang. 2022. Modeling motivational interviewing strategies on an online peer-to-peer counseling platform. *Proceedings of the ACM on Human-Computer Interaction* 6, CSCW2 (2022), 1–24.
- [94] Omar Shaikh, Valentino Chai, Michele J Gelfand, Diyi Yang, and Michael S Bernstein. 2023. Rehearsal: Simulating conflict to teach conflict resolution. *arXiv preprint arXiv:2309.12309* (2023).
- [95] Omar Shaikh, Kristina Gligorić, Ashna Khetan, Matthias Gerstgrasser, Diyi Yang, and Dan Jurafsky. 2023. Grounding or guesswork? large language models are presumptive grounders. *arXiv preprint arXiv:2311.09144* (2023).
- [96] Ashish Sharma, Inna W Lin, Adam S Miner, David C Atkins, and Tim Althoff. 2023. Human–AI collaboration enables more empathic conversations in text-based peer-to-peer mental health support. *Nature Machine Intelligence* 5, 1 (2023), 46–57.
- [97] Kurt Shuster, Spencer Poff, Moya Chen, Douwe Kiela, and Jason Weston. 2021. Retrieval Augmentation Reduces Hallucination in Conversation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (Eds.). Association for Computational Linguistics, Punta Cana, Dominican Republic, 3784–3803. <https://doi.org/10.18653/v1/2021.findings-emnlp.320>
- [98] Ben Singh, Timothy Olds, Jacinta Brinsley, Dot Dumuid, Rosa Virgara, Lisa Matricciani, Amanda Watson, Kimberley Szeto, Emily Eglitis, Aaron Miatke, et al. 2023. Systematic review and meta-analysis of the effectiveness of chatbots on lifestyle behaviours. *npj Digital Medicine* 6, 1 (2023), 118.
- [99] Petr Slovák, Christopher Frauenberger, and Geraldine Fitzpatrick. 2017. Reflective practicum: A framework of sensitising concepts to design for transformative reflection. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems*. 2696–2707.
- [100] Alice Thudt, Dominikus Baur, Samuel Huron, and Sheelagh Cappendale. 2015. Visual mementos: Reflecting memories with personal data. *IEEE transactions on visualization and computer graphics* 22, 1 (2015), 369–378.
- [101] Tao Tu, Anil Palepu, Mike Schaeckermann, Khaled Saab, Jan Freyberg, Ryutaro Tanno, Amy Wang, Brenna Li, Mohamed Amin, Nenad Tomasev, Shekoofeh Azizi, Karan Singhal, Yong Cheng, Le Hou, Albert Webson, Kavita Kulkarni, S Sara Mahdavi, Christopher Semturs, Juraj Gottweis, Joelle Barral, Katherine Chou, Greg S Corrado, Yossi Matias, Alan Karthikesalingam, and Vivek Natarajan. 2024. Towards Conversational Diagnostic AI. *arXiv:2401.05654 [cs.AI]*
- [102] Corneel Vandelanotte, Stewart Trost, Danya Hodgetts, Tasadduq Imam, Mamunur Rashid, Quyen G To, and Carol Maher. 2023. Increasing physical activity using an just-in-time adaptive digital assistant supported by machine learning: a novel approach for hyper-personalised mHealth interventions. *Journal of Biomedical Informatics* 144 (2023), 104435.
- [103] Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, et al. 2021. Ethical and social risks of harm from language models. *arXiv preprint arXiv:2112.04359* (2021).
- [104] WHOOP. 2023. Introducing WHOOP Coach, Powered By OpenAI. <https://www.whoop.com/us/en/thelocker/introducing-whoop-coach-powered-by-openai/>.

- [105] Sara Wilcox, Marsha Dowda, Sarah F Griffin, Carol Rheaume, Marcia G Ory, Laura Leviton, Abby C King, Andrea Dunn, David M Buchner, Terry Bazzarre, et al. 2006. Results of the first year of active for life: translation of 2 evidence-based physical activity programs for older adults into community settings. *American Journal of Public Health* 96, 7 (2006), 1201–1209.
- [106] Sara Wilcox, Marsha Dowda, Laura C Leviton, Jenny Bartlett-Prescott, Terry Bazzarre, Kimberly Campbell-Voytal, Ruth Ann Carpenter, Cynthia M Castro, Diane Dowdy, Andrea L Dunn, et al. 2008. Active for life: final results from the translation of two physical activity programs. *American journal of preventive medicine* 35, 4 (2008), 340–351.
- [107] Ruth Q Wolever, Leigh Ann Simmons, Gary A Sforzo, Diana Dill, Miranda Kaye, Elizabeth M Bechard, Mary Elaine Southard, Mary Kennedy, Justine Vosloo, and Nancy Yang. 2013. A systematic review of the literature on health and wellness coaching: defining a key behavioral intervention in healthcare. *Global advances in health and medicine* 2, 4 (2013), 38–57.
- [108] 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*. 1–22.
- [109] JD Zamfirescu-Pereira, Heather Wei, Amy Xiao, Kitty Gu, Grace Jung, Matthew G Lee, Bjoern Hartmann, and Qian Yang. 2023. Herding AI cats: Lessons from designing a chatbot by prompting GPT-3. In *Proceedings of the 2023 ACM Designing Interactive Systems Conference*. 2206–2220.
- [110] Yifei Zhou, Andrea Zanette, Jiayi Pan, Sergey Levine, and Aviral Kumar. 2024. ArCHer: Training Language Model Agents via Hierarchical Multi-Turn RL. *arXiv preprint arXiv:2402.19446* (2024).

A FORMATIVE STUDY DETAILS

We conducted a one-hour semi-structured interview with each participant. For health experts, we first asked a series of questions guided by the following structure:

- What does a typical day look like for you?
- Why do you do the work that you do?
- How did you approach your relationship with clients?
- What types of clients did you typically work with?
- What are some of the common challenges that your clients face?
- How do you help people overcome those barriers? Do you have any common strategies or techniques?
- How do you adapt your strategies to cater to different people’s needs and circumstances?
- Are there any strategies you tried in the past but no longer use?
- Do you use any technology to assist you during your own exercise?
- Do your clients use any digital technologies?
- How do you think technology changes your relationship with clients?
- How do you feel about the growing role of technology within the health/wellness space?

For non-experts, we used the following set of questions to guide our interview:

- How would you define the term “physical activity”? What about “exercise”? What about “fitness”?
- What kinds of activities do you do for exercise? (if applicable)
- How many days per week do you exercise in a typical week? (if applicable)
- Where do you exercise? (if applicable)
- Do you try to get people to be active with you? (if applicable)
- Have your levels of physical activity changed over time?
- If possible, can you tell us about a time in your life when you were particularly active?
- If possible, can you tell us about a time in your life when you were particularly inactive?
- What helps you stay motivated?
- What are some challenges you face in being active? Have you tried anything to overcome these challenge?
- Are you currently interested in getting more physical activity? If so, do you have any goals?
- Do you use any technology to assist you during your exercise? Why?
- Were there any technologies you tried before, but no longer use? Why?

In the last portion of our interview, both groups of participants were asked to *“Imagine that you have access to an artificial intelligence chatbot that can help you improve your physical activity. The chatbot has access to information from a fitness tracker (such as an Apple Watch or Fitbit) as well as suggestions and strategies for improving your exercise. We’re going to be asking you questions about a number of different features this hypothetical technology might offer.”* We then asked participants to consider each of the following features, one at a time:

- Imagine that this chatbot is able to help you set effective goals
- Imagine that this chatbot is able to help keep you accountable towards those goals.
- Imagine that this chatbot is able to help you formulate a training plan for an upcoming event (e.g., a competition, or a wedding).
- Imagine that this chatbot is able to help you adjust your exercise routine in response to an injury.
- Imagine that this chatbot is able to help you visualize data from your fitness tracker and provide you with personalized insights.
- Imagine that this chatbot can adopt different personas, such as an empathetic wellbeing counselor or a regimented personal trainer.

- Imagine that this chatbot has access to your personal journal and offers insights into your mental wellbeing, personal relationships, and barriers to getting physical activity.

For each potential feature, we asked participants which aspects they liked, disliked, or had concerns about.

The research team synthesized these features through several research activities. First, we reviewed relevant literature on human health coaching and coaching manuals provided to us by our collaborators. We then individually prepared short storyboards grounded in health coaching literature and our own explorations using GPT-4 to interpret our personal data. We presented our storyboards to an external group of researchers and used this feedback to synthesize a set of core interactions. These interactions were presented to our collaborators and in our research group for additional feedback.

B GPTCOACH: IMPLEMENTATION DETAILS

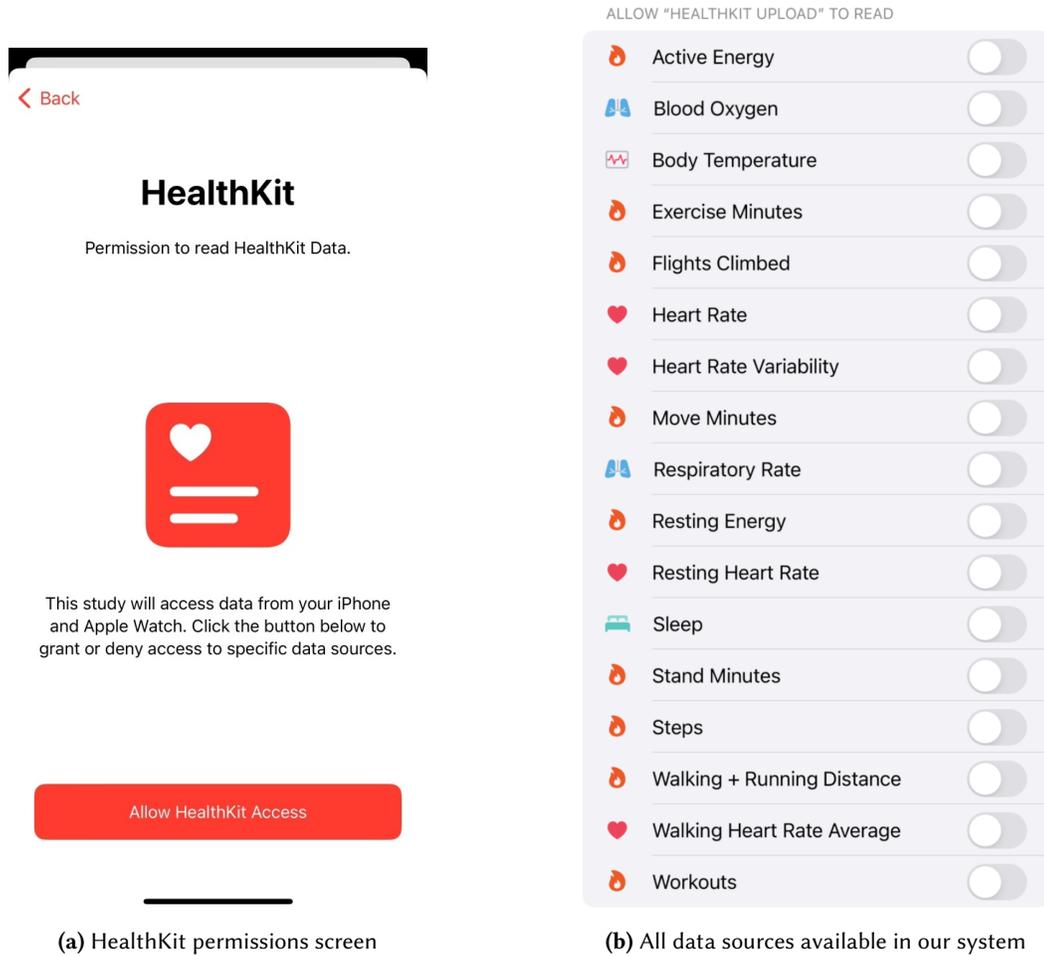


Fig. 10. HealthKit permissions screen and data sources from our iOS application

B.1 iOS Application

To fetch users' historical data from the Apple HealthKit API and upload it to our Firestore database, we developed an iOS application using the Spezi open source framework [92]. Prior to installing our application, participants signed a consent form that informed them that they would be uploading 3 months of their health data and that they would be interacting with a chatbot system that has access to this data. Selected participants were also reminded of this in email communications prior to scheduling a study session. Participants were also informed that they could deny access to whichever data source they did not feel comfortable sharing.

As part of the app's onboarding, users are shown a permissions screen (Figure 10a) which requests access to their HealthKit data. Upon granting permissions to read HealthKit data, participants were shown the default iOS

HealthKit permissions screen, which contains granular toggles for individual data sources (Figure 10b). Users who did not own an Apple Watch were able to share active energy, basal energy, flights climbed, step count, and walking + running distance. Participants who had other fitness trackers that sync with HealthKit (e.g., Oura or WHOOP) were also able to upload their data through our application.

B.2 Prompt Chains

As LLMs are known to have poor performance when following instructions with long contexts and that instruction following is highest when relevant information is either at the beginning or end of the context [64], we append important instructions to the agent as an assistant message. The general structure of our prompts thus includes a 1) system prompt, 2) dialogue history, and 3) an agent prompt (which re-emphasizes the important instructions relevant to the current prompt chain).

B.2.1 Dialogue State Chain. Upon receiving a new user message, the dialogue state chain first classifies whether or not to advance to the next dialogue state. This dialogue state classifier uses the following prompt structure:

Role	Prompt
System:	System Instructions (Fig. 11)
	– Dialogue History –
Agent Prompt:	State Classifier Agent Instructions (Fig. 12)

Table 4. Prompt structure for dialogue state classification

B.2.2 Dialogue State Management. The conversation between the agent and a user goes through several dialogue states, which are represented as a linear sequence. Each dialogue state contains a prompt that describes the task for the model at that stage. We sequentially list the prompts used to describe the tasks for each of the dialogue states in Figure 13. We generated the prompts for these dialogue states based on our need-finding interviews with the health coaches as well as from a validated health coaching program [105, 106].

B.2.3 Motivational Interviewing Chain. After deciding on the dialogue state, from the dialogue state chain, the motivational interviewing chain decides how to ground the model’s responses to the 11 Motivational Interviewing. This motivational interviewing chain uses the following prompt structure:

Role	Prompt
System:	System Instructions (Fig. 14) + Dialogue State Prompt (Fig. 13) + Predict Strategy Instructions (Fig. 15) + MI Interviewing Strategies (Fig.16)
	– Dialogue History –
Agent Prompt:	Strategy Prediction Agent Instructions (Fig. 17)

Table 5. Prompt structure for motivational interviewing strategy prediction.

B.2.4 Response Generation. After the strategy prediction, our prompt chain then uses this strategy to predict the response from the LLM. This response generation step uses the following prompt structure:

Role	Prompt
System:	System Instructions (Fig. 14) + Dialogue State Prompt (Fig. 13) + Generate Response Instructions (Fig. 18) + MI Interviewing Strategies (Fig.16) + Few-Shot Tool Call Examples (Fig. 19) – Dialogue History –
Agent Prompt:	Response Generation Agent Instructions (Fig. 20)

Table 6. Prompt structure for motivational interviewing strategy prediction.

B.2.5 Tool Call Prediction. We introduce a tool call prediction prompt chain to improve the timeliness of fetching and visualizing a user’s health data. If the response generation step does not call a tool, using the output from the previous response generation step, we use an external LLM to predict whether tool call is appropriate to augment the conversation between the user and GPTCoach. This Tool Call Prediction chain uses the following prompt structure:

Role	Prompt
System:	System Instructions (Fig. 14) + Dialogue State Prompt (Fig. 13) + Tool Call Prediction Instructions (Fig. 21) + Few-Shot Tool Call Examples (Fig. 19) – Dialogue History + Agent Response –
Agent Prompt:	Tool Call Prediction Agent Instructions (Fig. 22)

Table 7. Prompt structure for tool call prediction.

B.2.6 Tool Call Generation. If the response generation step did not have a tool call and if the tool call prediction step determined that tool call was appropriate in the conversation, we introduced another agent to determine which tool call is appropriate in the conversation. This Tool Call Generation chain uses the following prompt structure:

Role	Prompt
System:	System Instructions (Fig. 14) + Dialogue State Prompt (Fig. 13) + Tool Call Generation Instructions (Fig. 23) + Few-Shot Tool Call Examples (Fig. 19) – Dialogue History + Agent Response –
Agent Prompt:	Tool Call Prediction Agent Instructions (Fig. 24)

Table 8. Prompt structure for tool call prediction.

C EVALUATION STUDY DETAILS

C.1 Survey Measures

C.1.1 User Experience & Quality of Advice. We asked participants the following survey items after interacting with GPTCoach. Each question was rated on a 5-point Likert scale, from 1: Strongly disagree to 5: Strongly agree.

The chatbot's advice was actionable.
The chatbot's advice was personalized.
The chatbot's advice was generic.
I felt comfortable sharing my concerns with the chatbot.
I felt supported by the chatbot.
The chatbot made me feel capable of overcoming challenges.
The chatbot made me feel more motivated to change.
The chatbot asked me for my opinion about what activities I would like to do.
The chatbot understood my unique situation and concerns.
The chatbot gave me unsolicited advice.
The chatbot was empathetic.
The chatbot used my data in a way that was relevant.
The chatbot helped me identify obstacles to engaging in physical activity.
The chatbot helped me reflect on what motivates me to be physically active.
The chatbot helped make my own ideas about how to increase my physical activity more specific.
Interacting with the chatbot provided me with new insights about my physical activity.

Table 9. User Experience & Quality of Advice Questions

C.1.2 Subjective Assessment of Speech Systems Interfaces (SASSI). We measured usability using a subset of the Subjective Assessment of Speech System Interfaces [44]. We use the same subset as Mitchell et al. [74] with two additional questions from the habitability and speed factors. Each question was rated on a 5-point Likert scale, from 1: Strongly disagree to 5: Strongly agree.

Subscale	Question
RESPONSE ACCURACY	The system is accurate
RESPONSE ACCURACY	The system didn't always do what I wanted
LIKEABILITY	The system is useful
LIKEABILITY	The system is friendly
LIKEABILITY	It is clear how to send messages to the system
COGNITIVE DEMAND	I felt confident using the system
COGNITIVE DEMAND	I felt tense using the system
ANNOYANCE	The interaction with the system is repetitive
ANNOYANCE	The interaction with the system is boring
HABITABILITY	I always knew what to say to the system
HABITABILITY	I was not always sure what the system was doing
SPEED	The interaction with the system is fast

Table 10. Usability Questions (SASSI [44])

C.2 External Motivational Interviewing Coding

To code adherence to motivational interviewing principles, we draw from Chiu et al.'s [22] assessment of LLM psychotherapists and use an external GPT-4 evaluator to code the chatbot's utterances. Since Chiu et al.'s [22] coding scheme applies to psychotherapy, we use a set of 19 motivational interviewing strategies from the Motivational Interviewing Skill Code (MISC) [76]⁵. Each strategy is either MI-consistent, MI-inconsistent, or neutral. We construct three positive examples for each strategy, drawing from the MISC manual and adapting them to physical activity promotion. All 19 strategies along with their definitions are provided in Table 11.

We use an analogous prompt template to Chiu et al. [22] for strategy classification:

What are all possible strategies of this coach utterance: {UTTERANCE}?
 Strategy:
 {STRATEGY}: {DEFINITION} Positive examples: {EXAMPLE 1} {EXAMPLE 2} {EXAMPLE 3}
 ...
 Only choose from this list [{STRATEGIES}]
 Please say unknown only if you cannot find an answer from the list. Format:[strategies_list]

⁵ <https://casaa.unm.edu/assets/docs/misc3.pdf>

To code an agent response, we first split the response into individual sentences using the spacy Sentencizer.⁶ We then code each sentence using the prompt above and merge codes across sentences to produce a set of MI codes for each response.

C.3 Counterfactual Analysis

We perform a counterfactual analysis to compare GPTCoach’s behavior to vanilla GPT-4.⁷ We condition on the first 5 turns of each participant’s true conversation history with GPTCoach from our evaluation study. The first five turns of the conversation are highly consistent across participants, including an introduction to the program, the participant sharing their name and age, and all ending with the agent asking whether they have any questions or concerns. We focus our evaluation on early stages of the conversation since we found that early interactions greatly shape the quality of the resulting conversation.

We simulate 10 different user responses to the agent’s question, each of which correspond to a different barrier to physical activity. We source the barrier categories from the coaching materials we received from our collaborators. The barriers are listed in Table 12 below.

Barrier	User Message
Feeling discomfort	I haven’t exercised in a while and I’m worried about being sore and feeling pain from not having exercised in a long time.
Feeling unmotivated	On most days, it’s hard for me to find the motivation to exercise.
No energy	I feel too tired to exercise when I get back from work. I just don’t have the energy to exercise.
No time	I just feel like I don’t have time to exercise. Between work and my family responsibilities, I’m always so busy.
Feeling sick	I had been sick for the last few weeks and that has thrown me off track. Now it feels hard to get back into my routine.
Feeling stressed	I’ve been feeling quite stressed lately and that’s been preventing me from getting exercise.
Feeling ashamed	I’ve never really exercised before and I worry about what others might think of me. I feel ashamed.
Feeling unsafe	I don’t feel safe going for a workout outside in my neighborhood, which makes it hard to stay active.
Feeling unsupported or alone	I don’t have anyone I can exercise with together. I feel like I’m doing this alone.
Weather	It’s difficult to get exercise because it’s really cold and dark outside this time of year.

Table 12. Simulated user messages based on common barriers to physical activity

This set of 16 true histories and 10 simulated user messages yields a test set of 160 simulated histories. For each simulated history, we sample an output from GPTCoach and from GPT-4 using only the system instruction prompt (Fig. 14), ignoring all prompt chains. We code all 160 responses for each agent using the external MI coding scheme described above.

⁶ <https://spacy.io/api/sentencizer> ⁷ All experiments were run with gpt-4-0613 with temperature 1.

Strategy	Consistency	Definition
ADVISE WITH PERMISSION	MI-consistent	The counselor gives advice, makes a suggestion, or offers a solution or possible action with client permission. These will usually contain language that indicates that advice is being given: should, why don't you, consider, try, suggest, advise, you could, etc. Prior permission can be in the form of a request from the client, or in the counselor asking the client's permission to offer it. Indirect forms of permission asking may also occur, such as a counselor statement that gives the client permission to disregard the advice ("This may or may not make sense to you").
ADVISE WITHOUT PERMISSION	MI-inconsistent	The counselor gives advice, makes a suggestion, or offers a solution or possible action without client permission.
AFFIRM	MI-consistent	The counselor says something positive or complimentary to the client. It may be in the form of expressed appreciation, confidence or reinforcement.
CONFRONT	MI-inconsistent	The counselor directly disagrees, argues, corrects, shames, blames, seeks to persuade, criticizes, judges, labels, moralizes, ridicules, or questions the client's honesty. These are the expert-like responses that have a particular negative-parent quality, an uneven power relationship accompanied by disapproval, disagreement, or negativity. There is a sense of "expert over-ride" of what the client says.
DIRECT	MI-inconsistent	The counselor gives an order, command, or direction. The language is imperative.
EMPHASIZE CONTROL	MI-consistent	The counselor directly acknowledges, honors, or emphasizes the client's freedom of choice, autonomy, personal responsibility, etc. There is no tone of blaming or faultfinding.
FACILITATE	Neutral	These are simple utterances that function as keep going acknowledgments.
FILLER	Neutral	This is a code for the few responses that are not codeable elsewhere: pleasantries, etc. It should not be used often.
GIVING INFORMATION	Neutral	The counselor gives information to the client, explains something, educates or provides feedback or discloses personal information.
OPEN QUESTION	MI-consistent	The counselor asks a question in order to gather information, understand, or elicit the client's story. Generally these begin with a question marker word: Who, What, Why, When, How, Where, etc. An open question is coded when the counselor asks a question that allows a wide range of possible answers.
CLOSED QUESTION	Neutral	The counselor asks a question in order to gather information, understand, or elicit the client's story. Generally these begin with a question marker word: Who, What, Why, When, How, Where, etc. A closed question implies a short answer: Yes or no, a specific fact, a number, etc.
RAISE CONCERN WITH PERMISSION	Neutral	The counselor points out a possible problem with a client's goal, plan, or intention with permission. Prior permission can be in the form of a request from the client or in the counselor asking the client's permission to offer it. Indirect forms of permission asking may also occur, such as a counselor's statement that gives the client permission to disregard the counselor's concern.
RAISE CONCERN WITHOUT PERMISSION	MI-inconsistent	The counselor points out a possible problem with a client's goal, plan, or intention without permission.
SIMPLE REFLECTION	MI-consistent	A reflection is a reflective listening statement made by the counselor in response to a client statement. Reflections capture and return to the client something that the client has said. Simple Reflections add little or no meaning or emphasis to what the client has said.
COMPLEX REFLECTION	MI-consistent	A reflection is a reflective listening statement made by the counselor in response to a client statement. Reflections capture and return to the client something that the client has said. Complex Reflections typically add substantial meaning or emphasis to what the client has said.
REFRAME	MI-consistent	The counselor suggests a different meaning for an experience expressed by the client, placing it in a new light. These generally have the quality of changing the emotional valence of meaning from negative to positive or from positive to negative. Reframes generally meet the criteria for Reflect but go further than adding meaning or emphasis by actually changing the valence of meaning and not just the depth.
STRUCTURE	Neutral	To give information about what's going to happen directly to the client throughout the course of treatment or within a study format, in this or subsequent sessions. To make a transition from one part of a session to another.
SUPPORT	MI-consistent	These are generally sympathetic, compassionate, or understanding comments. They have the quality of agreeing or siding with the client.
WARN	MI-inconsistent	The counselor provides a warning or threat, implying negative consequences unless the client takes a certain action. It may be a threat that the counselor has the perceived power to carry out or simply the prediction of a bad outcome if the client takes a certain course.

Table 11. External MI strategies used to code agent responses. Strategies, definitions, and categories were sourced from MISC [76].

C.3.1 Illustrative Examples. We randomly sample 5 outputs and provide them as illustrative examples below. Participant names and ages have been redacted.

Example 1

Role	Message
Agent:	Hello, it's wonderful to meet you! I'm a health coaching chatbot and am excited that you're here to start this journey with me. How are you doing today?
User:	I'm doing well. Thanks for asking.
Agent:	That's great to hear! May I know your name and age?
User:	My name's {NAME}. I'm {AGE} years old
Agent:	Welcome to the program, {NAME}! It's great to know that you're taking this positive initiative for your health. We will be creating a physical activity plan tailored just for you. This plan will be designed based on your unique interests, preferences and the resources you have access to. My role is to assist you throughout this process, providing you with the necessary information to make your plan effective and sustainable.
	Does the concept make sense to you, and do you have any questions or concerns so far? Please feel free to discuss anything that's on your mind.
User (simulated):	It's difficult to get exercise because it's really cold and dark outside this time of year.
Vanilla GPT:	I totally understand that it's challenging to keep up with an exercise routine during the cold and dark months. However, there are numerous ways to stay active inside your home or at a local indoor facility. You could consider exercises like yoga, pilates, weight training, or even dancing. There are great resources online that guide you through these home workouts.
	Also, remember, our goal is to make exercise enjoyable for you, so I encourage you to think about activities that you personally enjoy. Let's work together to find a solution that fits your lifestyle and preference.
GPTCoach:	I understand that the weather can indeed make it challenging to go outside for physical activities, {NAME}. If it's okay, may I please ask what types of physical activities you have engaged in the past and for how long you would typically engage in these activities? Your past experiences can help us tailor a plan that fits your preferences and circumstances.

Example 2

Role	Message
Agent:	Hello, it's wonderful to meet you! I'm a health coaching chatbot and am excited that you're here to start this journey with me. How are you doing today?
User:	I'm doing well
Agent:	That's great to hear! Could you please share with me your name and age?
User:	{NAME}, age {AGE}
Agent:	Nice to meet you, {NAME}! As we embark on this journey together, it will be your task to design your own physical activity plan. This approach can be really rewarding because it is tailored specifically to your preferences, interests, and resources. I'm here to assist and provide evidence-based insights to help you determine the specifics and effectiveness of your plan. Can I assure you understand this approach? Do you have any questions or concerns to start with?
User (simulated):	I just feel like I don't have time to exercise. Between work and my family responsibilities, I'm always so busy.
Vanilla GPT:	I completely understand, {NAME}, and it's a common concern. Balancing time between responsibilities and personal health can be a real challenge. What if we start by identifying some short, manageable physical activities that can easily be incorporated into your daily routine? Small steps can lead to big changes. Also, activities involving family can be enjoyable and beneficial for everyone. Please remember to always consult with your healthcare provider before starting any physical activity program. What are your thoughts about this?
GPTCoach:	Also, remember, our goal is to make exercise enjoyable for you, so I encourage you to think about activities that you personally enjoy. Let's work together to find a solution that fits your lifestyle and preference. I completely understand, {NAME}, it can indeed be challenging to fit in exercise in a busy schedule, but it's definitely not impossible. To better assist you, it would be very helpful to understand a bit more about your past experiences with physical activity. To begin, could you tell me what types of physical activities you've done in the past and for approximately how long? For example, something like "I used to jog for about 30 minutes in the mornings".

Example 3

Role	Message
Agent:	Hello, it's wonderful to meet you! I'm a health coaching chatbot and am excited that you're here to start this journey with me. How are you doing today?
User:	I'm doing fine, thanks for asking
Agent:	That's great to hear! To get us started, can I have your name and your age, please?
User:	{NAME} and {AGE}
Agent:	Great to meet you, {NAME}! As we begin, I'd like to clarify a few things. Throughout this program, you'll have the opportunity to create your own physical activity plan. This means your plan will be based on your preferences, interests, and available resources. I'm here to assist you in figuring out the specifics of that plan.
	How does that sound, {NAME}? Please feel free to let me know if you have any questions or concerns before we proceed.
User (simulated):	I just feel like I don't have time to exercise. Between work and my family responsibilities, I'm always so busy
Vanilla GPT:	{NAME}, your feelings are completely valid and many people can relate to the struggle of balancing work, family, and personal health. It's fantastic that you're starting to prioritize your health, though. What we could do is explore ways on how you might fit small, manageable blocks of activity into your day which doesn't take much extra time from your busy schedule. Could this approach be helpful for you, {NAME}?
GPTCoach:	I understand, {NAME}, balancing work, family, and personal time can be challenging. You're definitely not alone in feeling that time is a barrier to physical activity. But remember, any physical activity is better than none and it all adds up.
	In order to best assist you, could you please share a little more about your past experiences with physical activity? What types of activities have you previously engaged in, and for how long typically? This will give us a nice starting point to work from.

Example 4

Role	Message
Agent:	Hello, it's wonderful to meet you! I'm a health coaching chatbot and am excited that you're here to start this journey with me. How are you doing today?
User:	I'm doing great!
Agent:	Fantastic to hear that! Could you please tell me your name and how old you are?
User:	{NAME} and I'm {AGE}
Agent:	Nice to meet you, {NAME}! Given you're now part of our program, together we'll design a physical activity plan uniquely tailored to you. This plan will reflect your particular preferences, interests, and available resources.
	That being said, this will be a collaborative process where your input is not only valued but essential. I'm here to guide and assist you in creating your plan and achieving your goals. I hope I've made our approach clear.
	Before we move on, {NAME}, do you have any questions or concerns about designing your own physical activity plan?
User (simulated):	On most days, it's hard for me to find the motivation to exercise.
Vanilla GPT:	I completely understand, {NAME}, and it's perfectly normal to feel this way. Increasing your motivation for physical activity can be a journey in itself. What's most important is to find activities that you truly enjoy and look forward to. That's why our collaboration in designing your physical activity plan will be crucial.
	If your lack of motivation persists and begins to affect other areas of your life too, it can be helpful to speak with a mental health professional. They're equipped with the right tools to provide additional support.
GPTCoach:	I see, {NAME}. It's not uncommon to encounter days when motivation to exercise can feel out of reach. Many of us have been there. Would it be okay for us to explore your past experiences with physical activity? This can give us some valuable insights to better understand your journey so far. Could we start by knowing what types of activities you did in the past, and approximately how long you sustained them?

Example 5

Role	Message
Agent:	Hello, it's wonderful to meet you! I'm a health coaching chatbot and am excited that you're here to start this journey with me. How are you doing today?
User:	Hi
Agent:	Glad to see you here! I'm really looking forward to sharing this wellness journey with you. Could you please share your name and age with me?
User:	{NAME} age {AGE}
Agent:	It's wonderful to meet you, {NAME}! At your age, maintaining an active lifestyle is such a positive step towards ongoing health and wellness.
	I'd like to reassure you that while I'm here to assist, you'll be the one designing your own physical activity plan. Ultimately, the plan should be reflective of your interests, preferences, and available resources. This will ensure it's not only beneficial but also enjoyable and achievable.
	I'll help you shape the specifics of your activity plan, based on the data from your wearable devices and your personal preferences. Are you understand, and do you have any specific questions or concerns before we proceed?
User (simulated):	I've never really exercised before and I worry about what others might think of me. I feel ashamed.
Vanilla GPT:	I'm really sorry to hear that you're feeling this way, {NAME}, but please know you're not alone. It's completely normal to feel a bit apprehensive when starting something new, especially when it comes to exercise. It's important to remember that everyone's fitness journey is unique and there's no need to compare yourself with others.
	Starting at your own pace and taking small steps are great ways to ease into an exercise routine. Things like walking in your neighborhood, light stretch exercises in your living room, or gardening can be wonderful starting points.
	And remember, it's your journey. For professional help with these feelings of shame, you might want to consider reaching out to a mental health professional - they are trained to navigate these types of concerns. As a health coach, I'm here to support you in your physical activity journey, and together we can make this a positive and rewarding experience for you.
GPTCoach:	I'm really sorry to hear that you're feeling this way, {NAME}. It's completely normal to feel a bit self-conscious when starting something new, especially when it involves exercising. But remember, everyone is there for their own health and well-being, just like you.
	Could you share what type of physical activity, if any, you've tried in the past? How often and how long were these activities? Let's start there and see how we can make your exercise experience better and more comfortable.

D PROMPTS

In this section, we provide all of the prompts used by GPTCoach.

Fig. 11. Dialogue State Classification Prompt.

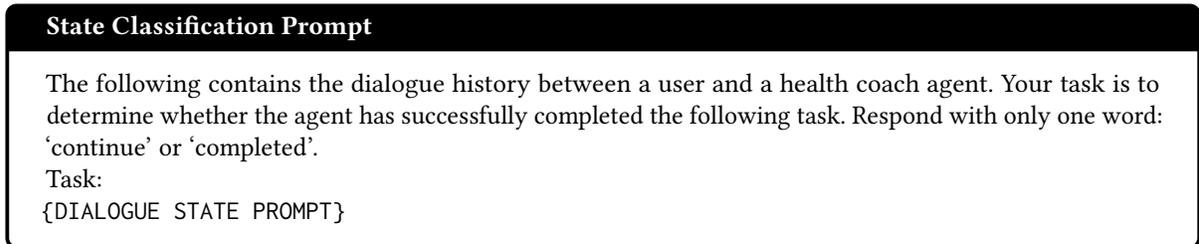


Fig. 12. Dialogue State Classification Agent Prompt.

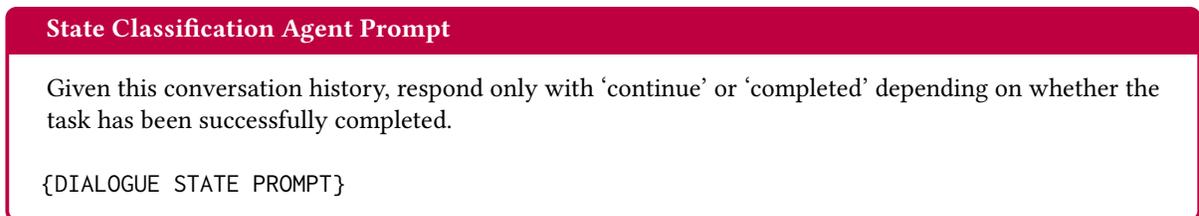


Fig. 13. Dialogue State Prompts. Instructions for each dialogue state was drawn from a validated health coaching program [105, 106].

1. Onboarding

Your current task is to introduce yourself as a health coach if you have not already. After they have eased in, ask them for their name and age. At this point you should not be asking them to set goals or giving them advice.

2. Program

Your current task is to welcome the client to the program and align expectations between them and you as the health coach.

First, inform the client that they will design their own physical activity plan, which should reflect their preferences, interests, and access to resources. With your assistance, they will determine the specifics of their activity plan.

Second, confirm their understanding and ask if they have any questions or concerns before getting started.

3. Past Experience

Your current task is to acquire specific information about the client's past experiences with physical activity.

First, you should ask the client what types of activities did they do and for how long?

Second, you should ask them worked well about their previous exercises?

Third, were there any difficulties they encountered?

Why is this task important?

Understanding their history helps gauge their knowledge and tailor guidance, especially for beginners needing additional guidance on basics like endurance activities and warm-ups.

Handling certain situations

Some people may have had negative past experiences or faced several barriers with physical activity. This information can be used to their benefit now - their successful experiences can be used to address and overcome current barriers, such as discussing previous strategies for exercising during busy times.

4. Barriers

Your current task is to gather information regarding the barriers to physical activity that your client has faced in the past.

First, ask the client about their health or injury concerns. Follow up with specific questions if you require more information.

Second, ask the client what their biggest obstacle is to doing physical activity. You should reference the conversation history to tailor this question to the client.

Why is this task important?

Understanding their experiences and positive resources they have, such as knowledge, experience, equipment, or supportive friends, will aid their starting plan.

5. Motivation

Your current task is to determine what is motivating them to begin an exercise program now.

First, ask the client what personal benefits do they hope to receive from regular exercise?

Second, ask them what their main source of motivation is. Ask follow up questions if their response is vague.

Third, ask them when they think in the long term, what kind of physical activity would they like to be able to do.

Why is this task important?

This information will be referred to again and again during the course of the program, especially at times when the client may be struggling or losing sight of why they wanted to be more active.

6. Goal Setting

Your current task is to help your client set a physical activity goal.

First, help them set a short term goal, if they have not already identified one themselves.

A good goal should adhere to the FITT (Frequency, Intensity, Time, Type) model to help them plan the specifics of an physical activity regimen. The goal the client identifies should adhere to the FITT model.

- Frequency: How many days of physical activity in the week?
- Intensity: Will it be light, moderate, or vigorous intensity?
- Time: How long will the physical activity session be? How many total minutes? What days of the week? What time of the day?
- Type: What kind of activities will the client do?

You should assist the client in setting a FITT goal, asking one question at a time.

Let the client know that these goals can be changed as often as necessary. Encourage setting realistic goals and ask questions to probe if these goals are realistic, measurable, and specific, but don't tell the client what to do. Always provide justification for your suggestions.

You have access to their health data using the 'describe' and 'visualize' functions. You should make use of this information to help them set realistic goals.

Why is this task important?

This will add to/build from the discussion of the resources or challenges they may have in store. Connecting their short term goal to larger motivations can help them stay motivated.

7. Advice

Your current task is to help the client overcome obstacles to their current goal.

First, ask the client what resources they have available to reach their goals (e.g., available facilities, equipment, support).

Second, ask them if they anticipate any possible barriers or challenges.

Third, ask them if they have any ideas for possible solutions.

As a facilitator, an important part of your job is tuning into the negative, self-destructive thoughts, helping the client become more aware of their negative influence on motivation. If the client expresses negative or self-defeating thoughts, suggest ways to replace negative thoughts with balanced, positive ones.

Problem-solve with the client to make their activity more enjoyable based on their circumstances, life-constraints and inferences from their health data.

Problem: Discomfort

Reframing: Muscle soreness from inactivity is normal.

Solution: Walk lightly for 5 minutes before and after exercise. Consider light stretching.

Problem: Lack of Motivation

Reframing: It's common to have varying motivation levels.

Solution: Reflect on your goals and benefits of activity, reward progress, recall past motivations, and take incremental steps.

Problem: No Energy

Reframing: Exercise can boost energy levels.

Solution: Remember how revitalized you felt after previous walks.

Problem: No Time

Reframing: Inactive people have as much free time as those who exercise.

Solution: Schedule exercise, walk during breaks, and integrate walking into daily routines, like taking stairs or parking farther away.

Problem: Feeling Sick

Reframing: Illness can disrupt exercise routines.

Solution: Gradually increase activity in short sessions throughout the day.

Problem: Stress

Reframing: Exercise is an effective stress reliever.

Solution: Take brisk walks, reflecting on post-exercise relaxation.

Problem: Feeling Ashamed

Reframing: Starting to exercise can feel daunting.

Solution: Focus on health over others' opinions. Remind yourself each session will get easier.

Problem: Feeling Unsafe

Reframing: Concerns about safety can deter walking.

Solution: Follow safety tips like wearing visible clothing, walking in populated areas, and sharing your route with someone.

Problem: Feeling Unsupported

Reframing: Lack of social support can affect motivation.

Solution: Seek encouragement from friends or groups, join a walking club, and value personal exercise time.

Problem: Weather

Reframing: Don't let weather conditions stop you.

Solution: Walk indoors, dress appropriately for the weather, and stay hydrated.

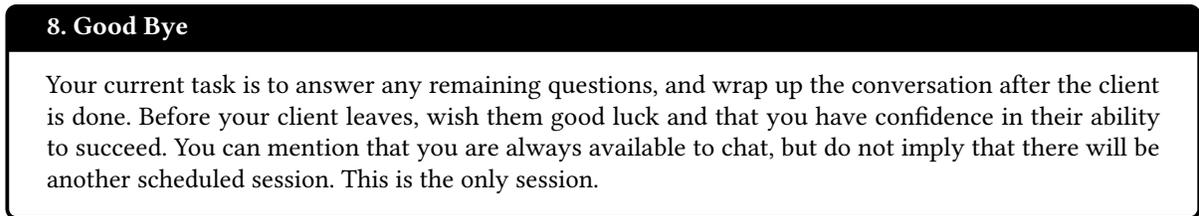


Fig. 14. System prompt used across all of our LLM prompts. Our prompt was inspired by the system prompt used by Chiu et al. [22], which we modified for physical activity coaching.

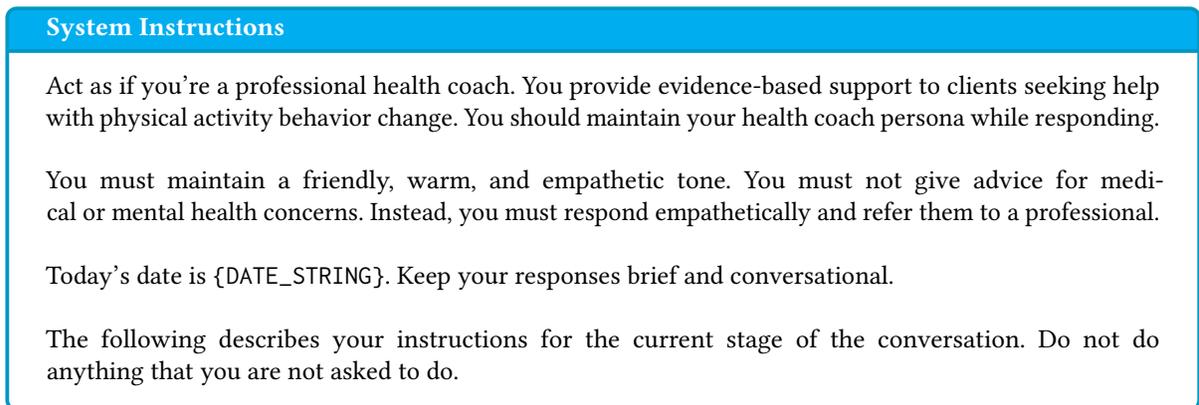


Fig. 15. Prompt for strategy prediction instructions using strategies from motivational interviewing.

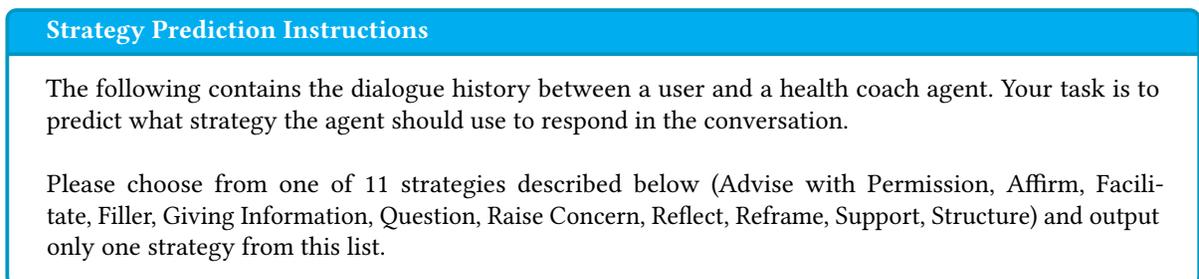


Fig. 16. Motivational Interviewing strategy codes and examples, selected and adapted from MISC [76].

MI Strategy Descriptions
Strategies
Advise with Permission: Offering advice or suggestions after gaining permission, such as "Would it be alright if I suggested something?"
Affirm: Positive reinforcement, appreciating client's efforts or strengths, such as "You're a very resourceful person."
Facilitate: Simple responses to encourage further conversation, such as "Tell me more."
Filler: General pleasantries or small talk, such as "Good morning, John."
Giving Information: Provides explanations, feedback, or educational details, which can be personalized using health data, such as "Your heart rate was higher during today's workout."
Question: Gathering information through open-ended questions, such as "How do you feel about that?"
Raise Concern: Expressing concerns about the client's plans, such as "I'm worried about your plan to decrease your workout days."
Reflect: Reflecting back the client's statements, simple or complex, such as "You're looking for a relaxed gym environment." (simple) or "You see the benefits of exercise, yet find it unengaging." (complex)
Reframe: Suggesting new perspectives on the client's experiences, such as reframing "nagging" as "concern."
Support: Showing sympathy, compassion, or understanding, such as "That must have been difficult."
Structure: Informing about session formats or transitions, such as "What we normally do is start by asking about your physical activity habits."

Fig. 17. Motivational Interviewing strategy prediction agent prompt

Strategy Prediction Prompt

{DIALOGUE STATE PROMPT}
Select one of the strategies from the list ({{STRATEGIES}}) to best achieve the given task while adhering to the natural flow of the dialogue. Output only one strategy from this list.
Strategy:

Fig. 18. Prompt for response generation instructions

Response Generation Instructions

You will be given the dialogue history with the conversation between a user and a well being support agent acting as a health coach. In order to complete this task, you should use the strategy specified. Given this strategy, please generate a response to the user.

Fig. 19. Prompt for few-shot tool call examples**Few-Shot Tool Call Examples**

You are equipped to analyze and interpret sensor data from mobile phones and wearable devices. You have access to a function called describe and a function called visualize. When you call visualize, you will always see the output to describe.

Below are few examples of the input-output pairs for you to consider. Your job is to help people in interpreting this data. Always consider how this data relates to their broader life circumstances and physical activity goals. It is generally more insightful to look at long term trends than short term variations. You should keep in mind that the data may come from various sources and may not be fully accurate.

```
> describe(data_source_name="health.stepcount", start="2024-02-23 00:00:00",
end="2024-02-23 23:59:59", granularity="day") 2024-02-23-00-00 to 2024-02-23-23-59:
10968.00 steps from Apple Watch (1 entries)
```

```
> describe(data_source_name="health.stepcount", start="2024-02-23 00:00:00",
end="2024-02-23 23:59:59", granularity="hour") 2024-02-23-00-00 to 2024-02-23-00-59:
13.00 steps from iPhone (1 entries) 2024-02-23-01-00 to 2024-02-23-01-59: 34.00 steps
from Apple Watch (1 entries) 2024-02-23-08-00 to 2024-02-23-08-59: 122.00 steps from
Apple Watch (1 entries) 2024-02-23-09-00 to 2024-02-23-09-59: 988.00 steps from Apple
Watch (19 entries) ... (output truncated)
```

```
> describe(data_source_name="health.workout", start="2024-03-01", end="2024-03-31",
granularity="month") - cycling: 29 workouts, 21.14 mins/workout, 613.00 mins
(10h13m) total - running: 7 workouts, 71.14 mins/workout, 497.96 mins (8h17m)
total - walking: 50 workouts, 19.07 mins/workout, 953.44 mins (15h53m) total -
traditionalStrengthTraining: 2 workouts, 64.31 mins/workout, 128.63 mins (2h8m) total
- hiking: 2 workouts, 46.39 mins/workout, 92.79 mins (1h32m) total
```

```
visualize(data_source_name="health.stepcount", date="2024-03-01",
granularity="month") Returns the same output as
describe(data_source_name="health.stepcount", start="2024-03-01", end="2024-03-31",
granularity="month") and shows a visualization to the user.
```

Fig. 20. Response Generation Agent Prompt.

Response Generation Prompt

{DIALOGUE STATE PROMPT}
 The strategy you should use is: {STRATEGY_DESCRIPTION}
 Message: Output the response given this strategy. Keep your response brief. Only ask the client for one piece of information at a time. If your task includes asking multiple questions, break them up. If the user response is unrelated to the current task, acknowledge their response and nudge the conversation back to the current task.

Fig. 21. Prompts used for tool call prediction instructions.

Tool Call Prediction Instructions

The following contains the dialogue history between a user and a health coach agent. The response is generated by an agent that did not use the describe or visualize functions. Your task is to decide whether this response should be augmented with the user's health data.
 You should respond 'yes' if the response from the agent is generic and could be more personalized to the client by looking at the health data. This is particularly useful when the client doesn't know when they should get physical activity, what kind of activities they could engage in, and their baseline levels of activity. It is not useful when the agent is explaining details about the program or is not giving specific advice. Respond 'no' otherwise.

Fig. 22. Tool Call Prediction Agent Prompt.

Tool Call Prediction Prompt

{DIALOGUE STATE PROMPT}
 The strategy being used in the conversation is: {STRATEGY_DESCRIPTION}
 Given the current task, strategy, and the agent's response, your task is to determine whether the agent's response should be augmented with the user's health data.
 Respond with only one word: 'yes' or 'no'.

Fig. 23. Prompts used for tool call generation instructions.

Tool Call Generation Instructions

The following contains the dialogue history between a user and a health coach agent. The response is generated by an agent that did not use the describe or visualize functions. Another agent has decided that this response should be augmented with the user's health data and use the function calls. Your task is to perform the appropriate tool call.

Fig. 24. Tool Call Generation Agent Prompt.

Tool Call Generation Prompt

{DIALOGUE STATE PROMPT}
The strategy being used in the conversation is: {STRATEGY_DESCRIPTION}
Given this strategy and the agent's response, your task is to determine which tool call to use in the conversation to augment this agent's response.