

CALLM: Understanding Cancer Survivors’ Emotions and Intervention Opportunities via Mobile Diaries and Context-Aware Language Models

ZHIYUAN WANG, Department of Systems and Information Engineering, University of Virginia, United States

KATHARINE E. DANIEL, Center for Behavioral Health and Technology, University of Virginia, United States

LAURA E. BARNES, Department of Systems and Information Engineering, University of Virginia, United States

PHILIP I. CHOW, Center for Behavioral Health and Technology, University of Virginia, United States

Cancer survivors face unique emotional challenges that impact their quality of life. Mobile diary entries—short text entries people record through their personal phone—provide a promising method for tracking emotional states, improving self-awareness, and promoting well-being outcome. This paper aims to, through mobile diaries, understand cancer survivors’ emotional states and key variables related to just-in-time intervention opportunities, including the desire to regulate emotions and the availability to engage in interventions. Although emotion analysis tools show potential for recognizing emotions from text, current methods lack the contextual understanding necessary to interpret brief mobile diary narratives. Our analysis of diary entries from cancer survivors (N=407) reveals systematic relationships between described contexts and emotional states—with administrative and health-related contexts associated with negative affect and regulation needs, while leisure activities promote positive emotions. We propose CALLM, a Context-Aware framework leveraging Large Language Models (LLMs) with Retrieval-Augmented Generation (RAG) to analyze these brief entries by integrating retrieved peer experiences and personal diary history. CALLM demonstrates strong performance with balanced accuracies reaching 72.96% for positive affect, 73.29% for negative affect, 73.72% for emotion regulation desire, and 60.09% for intervention availability, outperforming language model baselines. Post-hoc analysis reveals that model confidence strongly predicts accuracy, with longer diary entries generally enhancing performance, and brief personalization periods yielding meaningful improvements. Our findings demonstrate how contextual information in mobile diaries can be effectively leveraged to understand emotional experiences, predict key states, and identify optimal intervention moments for personalized just-in-time support.

CCS Concepts: • **Human-centered computing** → **Ubiquitous and mobile computing systems and tools**; • **Applied computing** → **Health care information systems**; • **Computing methodologies** → **Natural language processing**.

Additional Key Words and Phrases: Mobile diaries, cancer survivors, context awareness, language models, emotions, interventions

1 Introduction

In the context of cancer survivorship, understanding day-to-day emotional and behavioral fluctuations is critical to deliver tailored support and interventions. The National Cancer Institute (NCI) defines a cancer survivor broadly, encompassing individuals from diagnosis through the remainder of life [46]. This large and growing population often faces significant and unique emotional challenges [20, 93], including heightened psychological distress [4, 58], particularly in the initial years following diagnosis [92]. Addressing these emotional needs is a key component of comprehensive survivorship care, linked to long-term well-being and recovery [31, 63]. Despite the recognized importance of emotional support, many cancer survivors face barriers to accessing traditional mental health services, creating a need for accessible, low-burden approaches to emotional monitoring and intervention delivery [13, 29].

Authors’ Contact Information: Zhiyuan Wang, Department of Systems and Information Engineering, University of Virginia, Charlottesville, Virginia, United States, vmf9pr@virginia.edu; Katharine E. Daniel, Center for Behavioral Health and Technology, University of Virginia, Charlottesville, Virginia, United States, ked4fd@virginia.edu; Laura E. Barnes, Department of Systems and Information Engineering, University of Virginia, Charlottesville, Virginia, United States, lb3dp@virginia.edu; Philip I. Chow, Center for Behavioral Health and Technology, University of Virginia, Charlottesville, Virginia, United States, pic2u@virginia.edu.

Mobile diaries (e.g., brief text entries like “Pain in my legs” or “Visit from my daughter” in response to a prompt question “what impacted your mood?”), allow users to submit brief open-ended journal entries through smartphone interfaces, representing a growing modality in mobile health and ubiquitous computing [35]. Mobile diaries can be collected through ecological momentary assessment (EMA), which captures users’ experiences and states longitudinally, in-situ, as they occur in daily life via mobile phones [16, 70]. Crucially, these tools capture not just the *what* (emotional state) but potentially the *why* (the described context or ‘emotion driver’), providing rich contextual cues that can inform personalized interventions. If the qualitative data within these diaries can be effectively interpreted, it opens possibilities for numerous applications, including providing enhanced self-awareness [49], personalizing coaching feedback [38], identifying moments of risk [36, 81], and tailoring adaptive interventions [9].

For digital mental health interventions (DMHIs) to be effective, particularly within frameworks like just-in-time adaptive interventions (JITAs) [51], identifying and characterizing optimal “intervention opportunities” is crucial. Emotional states—like elevated negative affect or reduced positive affect—are well-established indicators of when psychological support may be most needed. In addition to emotional states, two other factors can further inform the timing and relevance of intervention delivery: (1) the subjective **desire** to regulate emotions—when individuals actively want to change their current emotional state, and (2) objective intervention **availability**—when individuals are practically able to engage with a mobile intervention. This dual requirement exists because even when individuals are available (not driving, in meeting, etc.) [51], without the desire to change their emotions, motivation to engage will be minimal—a challenge common in conditions like depression where emotional inertia may limit regulation attempts [77], despite distress. Conversely, strong motivation to regulate emotions is ineffective if practical constraints prevent intervention engagement. Understanding when either or both of these conditions are present can support more responsive and contextually appropriate intervention delivery [26]. Analyzing the emotion drivers captured in mobile diaries potentially offers a low-burden way to characterize and detect these intervention opportunity cues without requiring multiple separate assessments.

However, unlocking the potential of mobile diaries for triggering and tailoring JITAs hinges on accurately interpreting the states and their surrounding context reflected within the brief, often unstructured text entries. While computational techniques for analyzing health-related text show promise [66, 78], mobile diary data presents unique challenges that limit the direct application of methods developed for other text types (e.g., social media [19, 24], clinical notes [94]). Diary entries are characterized by their extreme brevity (in our dataset presented later, averaging only 6.83 words with a median word count of 4.0, as shown in Figure 1a), are subjective nature, possess an introspective focus [6], and rely on implicit context. Furthermore, entries often lack explicit emotional language (66.2% exhibited neutral sentiment polarity in our data, as shown in Figure 1b), making state inference difficult. Traditional text analysis methods, whether based on feature engineering or standard deep learning [80], often falter with such sparse data. Even recent LLMs, while powerful, typically require explicit prompting or fine-tuning for specific tasks and may struggle to capture the necessary personal and temporal context when processing entries in isolation [55, 75, 87]. There is a need for methods that can effectively leverage the limited text alongside personal contexts and trajectories to understand the emotional and behavioral states reflected in mobile diaries.

In this paper, our investigation is guided by the following research questions:

RQ1: What contextual information (e.g., described activities, social interactions, health factors) is conveyed in brief mobile diary entries from cancer survivors, and how does this described context relate to their concurrent emotional states, desire to regulate emotions, and availability to engage in interventions?

RQ2: Can individuals’ emotional states, desire to regulate emotions, and availability to engage in interventions be inferred from brief mobile diary text using a context-aware LLM framework?

RQ3: How do factors such as model confidence, diary entry length, the inclusion of personal temporal context, and personalization using initial user data influence the performance of LLM-based prediction of emotional states, desire to regulate emotions, and availability to engage in interventions from mobile diaries?

To address these research questions, we conducted a study involving the collection and analysis of a large-scale, longitudinal mobile diary dataset. First, to enable the contextual analysis required for RQ1, we collected and characterized 24,183 diary entries over five weeks from $N = 407$ cancer survivors, capturing their emotional experiences alongside self-reported states, as well as interpreting them with contextual cues captured from the diaries. Second, to investigate RQ2—predicting emotional states, desire to regulate emotions, and availability to do interventions from these brief entries—we developed and evaluated CALLM (Context-Aware language model framework for mobile diaries). This framework leverages LLMs augmented with context dynamically retrieved from peer experiences (via retrieval augmentation; RAG [95]) and the individual’s own temporal trajectory from diary history. Third, to explore the factors influencing performance (RQ3), we compared CALLM against various baselines and conducted extensive post-hoc analyses examining the impact of LLMs’ confidence, diary entry length, diaries’ temporal dependency, and personalization potential.

This work offers the following contributions:

1. We collected and analyzed a large-scale mobile diary dataset from cancer survivors, revealing systematic relationships between described contexts and emotional states. This characterization provides insights into how different daily activities and contexts influence cancer survivors’ emotional experiences and identifies patterns relevant for contextually-aware intervention delivery.
2. We designed and evaluated CALLM, a novel context-aware framework that enhances LLM-based analysis of brief mobile diary entries by integrating peer experiences and individual temporal trajectories. Our approach demonstrates how contextual enrichment can significantly improve the interpretation of sparse text data for understanding intervention opportunities by detecting cancer survivors’ self-reported emotional states and desire to regulate emotions, though availability detection still leaves room for improvement.
3. We identified key factors that influence prediction performance through comprehensive post-hoc analyses, exploring how model confidence, entry length, temporal context, and personalization affect accuracy. These findings provide practical guidance for implementing mobile diary analysis systems in real-world settings.

We envision this work informing the development of context-aware JITAI systems that leverage mobile data to deliver highly personalized and precisely timed healthcare interventions. By interpreting user-generated entries—and potentially integrating passive sensing data from smartphones and wearables in future studies—this approach supports interventions attuned to individuals’ emotional states, desire, and availability, aligning closely with their real-world motivations and contexts. Ultimately, this direction advances adaptive, user-centered digital health systems capable of supporting individuals through complex emotional and behavioral journeys.

2 Related Work

2.1 Mobile Diaries for In-Situ Health Monitoring

Mobile diaries, often implemented through smartphone applications as part of ecological momentary assessment (EMA) protocols, are increasingly utilized in mobile health research. Their primary strength lies in capturing data

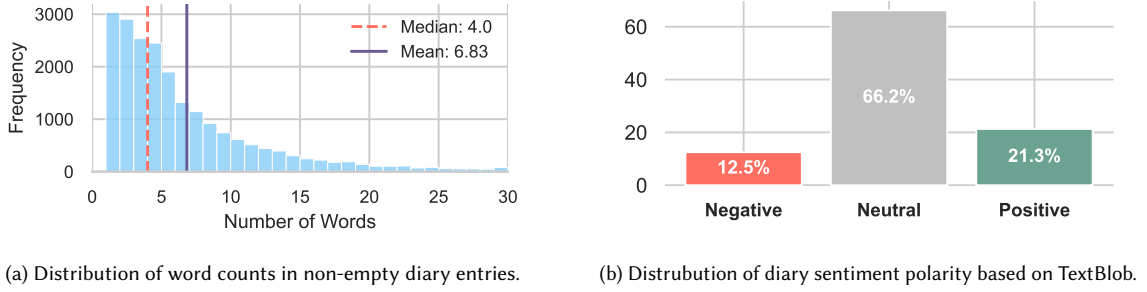


Fig. 1. Distributions of word count and sentiment polarity across collected mobile diaries from N=407 cancer survivors.

longitudinally and *in-situ*, reflecting experiences and states as they naturally occur in daily life [70]. This approach offers greater ecological validity and minimizes recall bias compared to retrospective self-reports [73]. Mobile diaries, particularly free-text formats, can achieve this with relatively lower subjective burden than frequent, lengthy structured questionnaires [50]. Common applications include monitoring symptoms in chronic illnesses [37], tracking mood and behavior patterns in mental health contexts [15], and understanding daily experiences in specific populations like cancer survivors [22, 47]. However, realizing the full potential of this rich qualitative data is often hindered by challenges including participant compliance over time, the inherent subjectivity of entries, and, most relevantly for computational analysis, the extreme brevity and sparsity of typical free-text diary entries [41, 71, 74]. This necessitates advanced methods capable of interpreting meaning from limited text input. These challenges, particularly the need to extract meaning from sparse text while considering individual context, motivate the exploration of context-aware systems (Section 2.2) and specialized text analysis techniques (Section 2.3).

2.2 Context-Aware Systems in Mobile Health

Understanding user context is a central theme in ubiquitous computing and mobile health [1, 25]. Context-aware systems aim to adapt their behavior based on situational information—such as location, activity, social setting, and human behavioral or emotional states—to provide more relevant and effective support [23, 33, 69]. In mobile health, context plays a critical role in understanding health behaviors and enabling timely interventions, such as JITAIs, where algorithms are designed to adapt dynamically to changing contextual factors [3, 51]. Previous work has also explored leveraging diverse signals for understanding user’s general contexts or contexts relevant to their disorders, including integrating passive sensor data (e.g., GPS, accelerometer, phone usage logs, heart rate, and audio signals) [48, 59, 76, 82] and employing multi-modal approaches that combine text with audio or visual cues [57]. Complementary to these approaches, our work focuses on the challenge of extracting rich contextual understanding directly *from the user’s own words* within brief, free-text mobile diary entries, potentially offering a low-friction method for capturing subjective context often missed by sensors alone. In particular, we characterize how cancer survivors as a unique population experience emotional deviations across different contextual situations (e.g., administrative, health-related, leisure activities) and express intervention opportunities through their diary narratives, enabling the identification of context-specific moments for targeted support without requiring additional sensing infrastructure.

2.3 Analyzing Mobile Diary Text with Language Models

Traditional emotion analysis has primarily targeted content-rich texts with explicit emotional cues, such as clinical interviews and social media posts (averaging 420 tokens¹ per post in Reddit datasets [79]), where contextual richness supports mental health assessments [66, 78]. These approaches leverage deep learning models [80], domain-specific frameworks [22], and increasingly, LLMs [87], thriving in contexts where users curate rich, expressive content [56]. However, mobile diary entries present fundamentally different challenges with their extreme brevity—averaging merely 7.57 tokens (6.83 words) per entry in this study’s dataset (Figure 1a)—and introspective nature. Unlike public-facing texts, these private reflections [6] often lack explicit emotional language (66.2% of entries demonstrate neutral sentiment polarity in our dataset, Figure 1b) and contextual redundancy. Existing machine learning methods struggle with this textual sparsity due to: (1) dependency on predetermined task designs unsuited for minimal text [19, 24], (2) inflexibility in adapting to personalized emotional trajectories [55], and (3) insufficient contextual awareness for interpreting ultra-brief entries [75].

Recent advances in LLMs have shown remarkable capabilities across various domains [89], particularly in mental health applications. These include detecting cognitive distortions [11, 87], analyzing depression through clinical interview data [62], supporting cognitive reframing of negative thoughts [68], and facilitating empathic peer-to-peer conversations in therapeutic settings [67, 83, 84]. RAG techniques further enhance LLMs by incorporating external knowledge, as demonstrated in personalized abusive language detection [91] and other knowledge-intensive tasks [39], while interpretable mental health analysis approaches offer explainable predictions [88] and well-being support [45]. In the specific domain of digital journaling, emerging work has begun exploring LLM integration: Jung et al.’s *MyListener* combines smartphone and Fitbit data to alleviate depression and loneliness [32], Kim et al.’s *MindfulDiary* improves journaling consistency through psychiatric-informed conversational interfaces [34], and Nepal et al. [52, 53] generates contextual journaling prompts from behavioral sensing data. While these approaches demonstrate LLMs’ potential for mental health applications, they primarily focus on enhancing the journaling experience rather than deriving rich insights from minimal text entries.

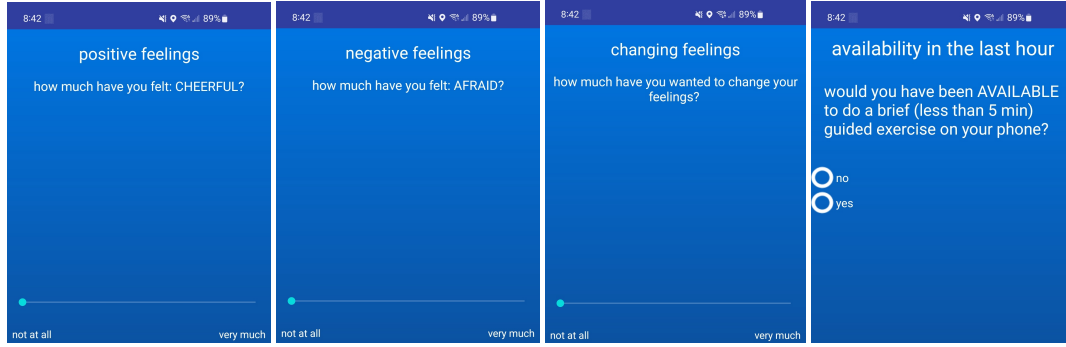
Building on existing work and addressing key challenges, the CALLM framework uniquely integrates peer experiences and personal temporal trajectories—both derived directly from diary entries—to extract meaningful insights from sparse text. It shows how even brief entries can illuminate cancer survivors’ life contexts and emotional needs, thereby enabling new applications of mobile diaries.

3 Methods

3.1 Participants

The present study involved data gathered from $N = 407$ US adults within five years of a cancer diagnosis (stages 0-4, referred to as ‘cancer survivors’) and who owned a smartphone. Participants were eligible regardless of cancer type and treatment status. Of the 426 who were enrolled, $N = 19$ were excluded ($n = 10$ did not initiate the EMA, $n = 9$ requested to withdraw), resulting in a final sample of $N = 407$. Participants were, on average, 48.73 years old ($SD = 12.23$); 9.09% male and 90.17% female; 86.67% White; and 92.52% non-Hispanic. Most participants (57.28%) had a primary diagnosis of breast cancer and 21.88% were actively receiving cancer treatment during the study. The study protocol has been fully approved by the Institutional Review Board at the University of Virginia IRB #HSR230080.

¹In language models, a token is the smallest unit of text produced by segmentation and is not necessarily a complete word. Typically, one word equals about 1.33 tokens, though this ratio can vary by tokenizer and model.



(a) Positive emotions. (b) Negative emotions. (c) Desire to regulate emotions. (d) Intervention availability.

Fig. 2. EMA interface for collecting emotional measures, desire to regulate emotions, and availability to do interventions.

3.2 Procedure

3.2.1 Recruitment Strategy. Participants were recruited via targeted online advertising. Recruitment was promoted by BuildClinical², which specializes in supporting robust and diverse recruitment for academic research trials. BuildClinical deploys advertisements according to data-driven strategy to target and engage specific patient populations through their large digital network (e.g., patient communities, Google, Facebook, Instagram, health websites, medical apps). Their advertising campaigns engage potentially eligible participants by deploying custom advertisements across relevant platforms using applied machine-learning and data mining to identify digital footprints of specific patient groups. During the study, algorithms are refined to engage specific groups and patient populations of interest.

3.2.2 Enrollment Process. Interested individuals responded to an online advertisement and completed an online pre-screen survey to verify eligibility (i.e., diagnosed with cancer in the last 5 years and owns a smartphone). Pre-eligible individuals had to pass a background check to verify their identity and the consent process was conducted on the phone. Prior to initiating the EMA phase, participants completed a battery of self-report questionnaires to assess demographics, cancer diagnosis and treatment history, and psychological functioning.

3.2.3 EMA Protocol. During the five-week EMA phase, participants received three surveys per day via the Effortless Assessment Research System (EARS) mobile application. Surveys were delivered randomly within three 2-hour time windows: morning (8-10am), afternoon (1-3pm), and evening (7-9pm). This schedule was designed to capture how psychological, behavioral, and emotional processes unfold throughout the day. Participants were notified of new surveys via push notifications, and surveys expired if not answered within 2 hours.

3.2.4 Compensation Structure. Participants received up to \$100 in gift cards for completing the study, with prorated compensation:

- \$20 for completing the baseline questionnaire battery
- Additional \$50 for completing between 50%-75% of the EMA surveys
- Or, additional \$80 for completing at least 75% of the EMA surveys

²<https://www.buildclinical.com/>

3.3 EMA Measures

3.3.1 Primary Outcome Measures.

- **Emotion driver diary entry:** Participants responded to the prompt regarding the last hour, “What has had the biggest impact on your mood?” using an open text box on their smartphone.
- **State affect:** Participants rated the degree to which they felt specific emotions in the past hour using a 0 (*not at all*) to 10 (*very much*) scale:
 - **Positive Affect:** Sum of ratings for happy, cheerful, and pleased (range: 0-30)
 - **Negative Affect:** Sum of ratings for sad, afraid, and miserable (range: 0-30)
- **Desire to regulate emotions:** Participants were asked regarding the last hour, “How much have you wanted to change your feelings?” using a 0 (*not at all*) to 10 (*very much*) scale.
- **Intervention availability:** Participants were asked regarding the last hour, “Would you have been AVAILABLE to do a brief (less than 5 minute) guided exercises on your phone?” using a yes/no response format.

3.3.2 *Secondary Outcome Measures.* Additional variables were measured using a 0 (*not at all*) to 10 (*very much*) Likert scale: social interaction quality, pain, worried, lonely, and grateful.

3.3.3 *Binary State Indicators.* To standardize analysis across participants and account for individual differences in baseline levels and response biases, we created binary state indicators for each emotional and behavioral measure. These indicators reflect whether a participant’s momentary rating exceeds their own mean level (True) or not (False) for that measure, allowing us to identify periods of elevated states relative to each participant’s typical experience.

4 Characterizing Contextual Information in Mobile Diaries

To answer **RQ1** (concerning contextual information in mobile diary entries and its relationship to emotional states), we systematically analyzed 24,183 diary entries collected from $N = 407$ cancer survivors. Our analysis approach was multi-faceted, examining (1) the relationship between diary topic categories and emotional/behavioral states, (2) how sentiment expressed in diary entries relates to these states, and (3) lexical patterns characterizing high-intensity emotional experiences. These analyses collectively provide insights into what contextual cues are present in brief diary entries and how they might inform the triggering and personalization of mobile interventions.

4.1 Topic-Based Analysis of Emotional States

To analyze how different activity contexts relate to emotional and behavioral measures, we categorized diary entries into context groups through keyword matching. The groups include: administrative, cancer & health, family events, home activities, leisure, outdoor activity, pet activities, routine, sleep & rest, social/holiday, transportation, work & study, and other. Each entry could be assigned to multiple groups based on its content (e.g., “walked dog in park” would be categorized under both “outdoor activity” and “pet activities”). The ‘Others’ category includes diaries that are too brief or lack the necessary keywords to provide sufficient information for topic matching.

Analysis of emotional and behavioral measures across different context categories reveals distinct patterns in how various activities and contexts relate to cancer survivors’ emotional experiences (Figure 3). For positive affect (Figure 3a), leisure activities (mean=19.90), social/holiday events (mean=19.11), and family events (mean=18.73) demonstrate notably higher scores compared to administrative tasks (mean=16.18) or cancer/health-related activities (mean=15.87). Mixed effects models, accounting for participant-level clustering (which explained 61.6% of variance), confirmed these

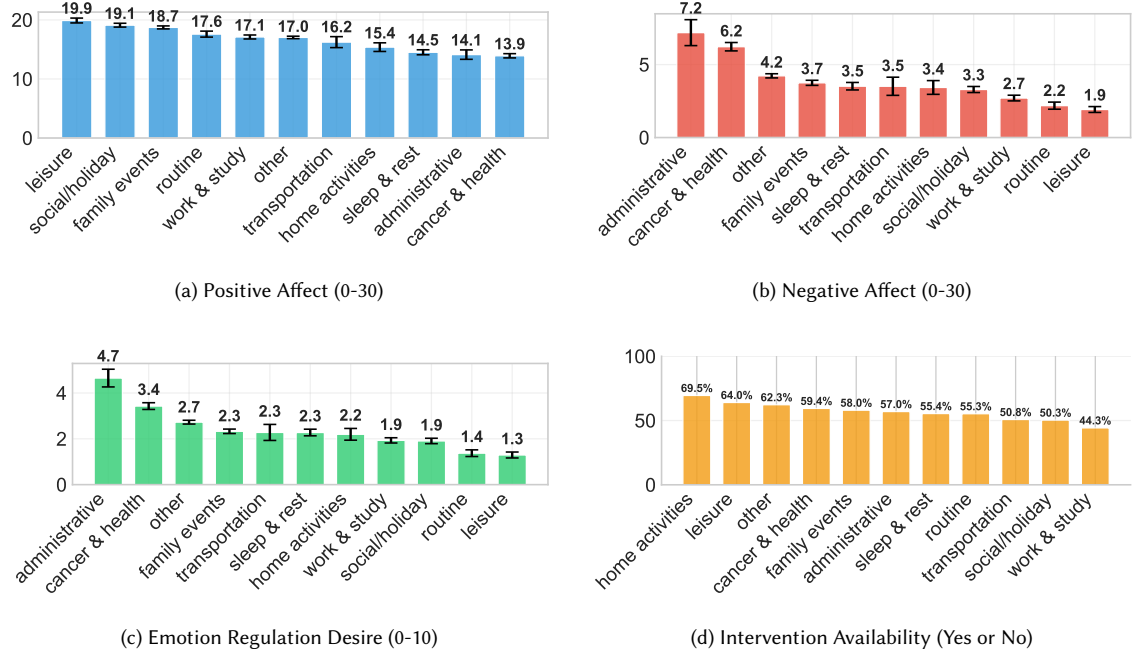


Fig. 3. Emotional and Behavioral Measures Across Different Emotion Driver Categories.

differences were statistically significant. Compared to administrative tasks (reference category), significantly higher scores were observed for contexts related to leisure activities (coef=3.97, adjusted $p < 0.001$), social/holiday events (coef=4.16, adjusted $p < 0.001$), and family events (coef=4.39, adjusted $p < 0.001$).

Negative affect distributions (Figure 3b) show an inverse pattern, with cancer/health-related categories (mean=6.16) and administrative tasks (mean=5.72) exhibiting higher mean scores and greater variability, particularly evident in the numerous outliers. Work and study activities show moderate negative affect levels (mean=2.71), highlighting the emotional challenges cancer survivors face in managing daily responsibilities. After controlling for individual differences (which explained 59.0% of variance), significantly lower negative affect scores were found for leisure (coef=-3.23, adjusted $p < 0.001$), cancer/health-related categories (coef=-3.63, adjusted $p < 0.001$), social/holiday events (coef=-2.94, adjusted $p < 0.001$), and work & study (coef=-2.90, adjusted $p < 0.001$) compared to administrative tasks.

Emotion regulation desire (Figure 3c) peaks during administrative tasks (mean=4.44) and cancer/health-related activities (mean=3.50), while being notably lower during leisure (mean=1.29), outdoor activities (mean=1.34), and routine activities (mean=1.37). Statistical analysis, controlling for individual differences (which explained 48.9% of variance), confirmed that all categories showed significantly lower scores compared to administrative tasks: cancer & health (coef=-0.94, adjusted $p < 0.001$), leisure (coef=-2.65, adjusted $p < 0.001$), outdoor activities (coef=-2.53, adjusted $p < 0.001$), social/holiday events (coef=-2.44, adjusted $p < 0.001$), and work & study (coef=-2.02, adjusted $p < 0.001$).

The availability to engage with digital interventions (Figure 3d) varies considerably across contextual categories, with home activities showing the highest availability (69.5%), while work and study contexts demonstrate the lowest (44.3%). Chi-square tests confirmed that these differences in intervention availability across topic categories were statistically significant ($\chi^2(12) = 313.5$, adjusted $p < 0.001$, Cramer's $V = 0.12$), indicating a small to medium association between

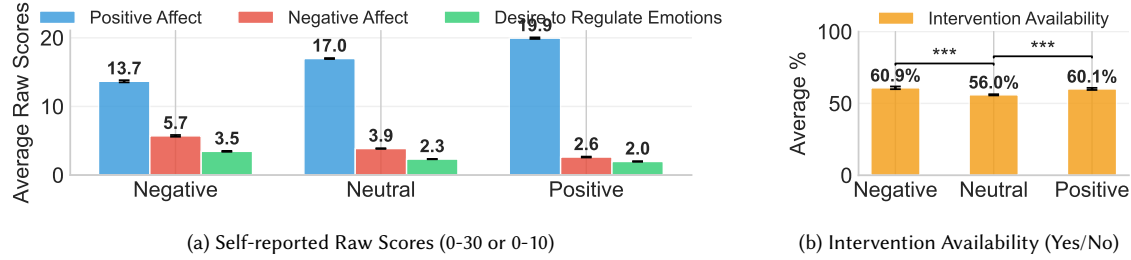


Fig. 4. Emotional States and Intervention Availability Across Different Sentiment Categories

activity context and receptiveness to digital intervention. These findings have direct implications for the targeting of just-in-time interventions, suggesting that both emotional need and practical availability should be considered when determining optimal intervention moments.

4.2 Sentiment Analysis and Emotional States

While topic analysis provides insights into how activity contexts relate to emotional states, the sentiment expressed within diary entries offers another dimension of contextual understanding. We performed sentiment analysis on diary entries using TextBlob toolkit [44], resulting in continuous sentiment scores ranging from -1.0 to 1.0. Entries were categorized as negative (score < -0.1, 12.5%), neutral (score = -0.1 to 0.1, 66.2%), or positive (score > 0.1, 21.3%). The predominance of neutral sentiment entries (66.2%) highlights the challenge of emotion inference from brief diary texts, as participants often describe situations factually without explicit emotional polarity.

Analysis of emotional self-reports by sentiment category confirms significant differences in experienced emotions that align with the detected text sentiment (Figure 4a). For affective states, positive sentiment entries showed the highest positive affect scores (mean=19.95), followed by neutral (mean=16.98) and negative (mean=13.66) entries, while negative sentiment entries showed the highest negative affect (mean=5.71), followed by neutral (mean=3.86) and positive (mean=2.62) entries. Mann-Whitney U tests confirmed these differences were statistically significant for both positive affect (negative vs. neutral: $U = 13.8M$, $p < 0.001$, $r = 0.24$; positive vs. neutral: $U = 37.6M$, $p < 0.001$, $r = -0.21$; negative vs. positive: $U = 3.3M$, $p < 0.001$, $r = 0.44$) and negative affect (negative vs. neutral: $U = 21.3M$, $p < 0.001$, $r = -0.20$; positive vs. neutral: $U = 26.2M$, $p < 0.001$, $r = 0.13$; negative vs. positive: $U = 7.6M$, $p < 0.001$, $r = -0.33$).

The desire to regulate emotions varied similarly across sentiment categories, with the highest scores in negative sentiment entries (mean=3.45), followed by neutral (mean=2.32) and positive (mean=1.98) entries, with Mann-Whitney U tests confirming statistical significance (negative vs. neutral: $U = 22.0M$, $p < 0.001$, $r = -0.21$; negative vs. positive: $U = 7.6M$, $p < 0.001$, $r = -0.29$). Interestingly, intervention availability (Figure 4b) showed a U-shaped pattern, with higher availability during moments where diary entries exhibit any polarity, either negative (60.9%) or positive (60.1%), compared to neutral moments (56.0%). Chi-square tests confirmed these differences were statistically significant ($\chi^2(2) = 38.2$, $p < 0.001$), though the effect size was small (Cramer's $V = 0.04$). This finding suggests that participants may be more receptive to interventions during emotionally salient moments (either positive or negative) compared to neutral states.



Fig. 5. Word clouds with high (above individual mean) affective scores reported.

4.3 Lexical Patterns in Emotional Experiences

Beyond topics and sentiment, examining the specific vocabulary used in diary entries provides deeper insights into the contextual factors driving cancer survivors' emotional experiences. To visualize these patterns, we generated word clouds based on entries associated with high emotional intensity (defined as above individual mean plus one standard deviation for both positive and negative affect). Common English stop words (e.g., "the", "is", "at") and emotion-specific terms (e.g., "feel") were removed to focus on content-bearing words.

Analysis of the word clouds (Figure 5) reveals distinct vocabulary patterns between positive and negative affect contexts. In positive affect entries (Figure 5a), daily activities and social connections dominate the discourse: work, day, and time appear as central terms, while social terms like family, daughter, and friend feature prominently. Activity-related words such as dinner, vacation, and walking suggest the importance of both routine and recreational activities in positive experiences.

The negative affect word cloud (Figure 5b) presents a markedly different emotional landscape. Health-related terms are most prominent, with cancer appearing as one of the largest words, accompanied by terms like pain, tired, and health. Relationship terms also feature differently here - while family remains present, husband appears with notable prominence (notably, 89.93% of participants are reported as female, $n=375$). The presence of terms like worried, anxiety, and stress directly reflects emotional challenges, while words related to medical experiences (hospital, treatment, doctor) indicate the ongoing impact of health concerns on daily life.

5 Predicting Emotional States and Intervention Opportunities from Mobile Diary Text

To address **RQ2** (whether nuanced emotional states, desire and availability relevant to intervention opportunities can be accurately predicted from brief mobile diary text using a context-aware LLM framework), we developed and evaluated CALLM (Context-Aware Language Model). This section details both the design of our framework and its evaluation results, demonstrating how brief text entries can be leveraged to predict emotional states, desire to regulate emotions, and availability to engage in interventions for personalized intervention delivery.

5.1 CALLM: A Context-Aware Framework for Mobile Diary Analysis

Building upon our contextual analysis findings from **RQ1**, we designed CALLM to leverage three key insights: (1) the significant impact of contextual factors on emotional states, (2) the substantial role of individual differences in emotional experiences, and (3) LLMs' pre-trained knowledge regarding cancer survivors' emotional experiences and potential

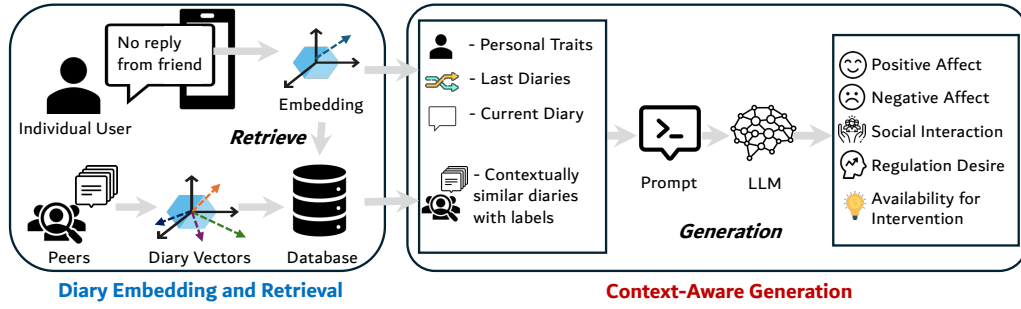


Fig. 6. Illustration of the CALLM Framework: The left side represents the Diary Embedding and Retrieval Module, while the right side illustrates the Context-Aware Generation Module.

needs. As illustrated in Figure 6, the framework consists of two main modules: the Diary Embedding and Retrieval Module and the Context-Aware Generation Module. By embedding and retrieving contextually relevant emotion diary data from peers, along with integrating individual traits and historical trajectories, CALLM incorporates these contextual cues into the LLM generation process to identify current emotional and behavioral states. This integration improves both the interpretability and predictive accuracy for real-world emotional states, even from extremely brief diary entries.

5.1.1 Diary Embedding and Retrieval Module. The Diary Embedding and Retrieval Module begins by converting each diary entry into a high-dimensional text embedding using OpenAI’s text-embedding-3-small model³. These embeddings capture semantic similarities among entries and are stored in a vector database. Each vector is accompanied by metadata, including self-reported outcome scores and participant traits, ensuring rich contextual alignment.

To facilitate efficient retrieval, a FAISS (Facebook AI Similarity Search⁴) index is constructed over these embeddings. This index allows for rapid similarity searches, leveraging L2 distance computation to identify the most relevant entries, where the L2 distance $d(\mathbf{a}, \mathbf{b})$ between two vectors \mathbf{a} and \mathbf{b} is defined as: $d(\mathbf{a}, \mathbf{b}) = \sqrt{\sum_{i=1}^n (a_i - b_i)^2}$

When a new diary entry is embedded, the system queries the FAISS index to retrieve the top- K semantically similar entries from the database. These entries are selected based on their proximity in the embedding space, indicating contextual relevance.

The retrieved entries, along with their associated ground truth labels, serve as contextual examples for the subsequent generation process. This retrieval mechanism ensures that predictions are grounded in contextually relevant historical data from peers, enhancing the model’s ability to understand and predict emotional states in context despite entries containing limited content.

5.1.2 Context-Aware Generation Module. The Context-Aware Generation Module integrates the retrieved data with additional contextual elements to improve prediction accuracy. As illustrated by Figure 7, this module constructs a comprehensive prompt by incorporating:

- **Individual Traits:** Participant-specific demographic and clinical information (e.g., “60-year-old male, stage II Kidney cancer”) to enable personalized analysis. These traits help contextualize the emotional experience within

³<https://platform.openai.com/docs/guides/embeddings> OpenAI - Vector Embeddings

⁴<https://ai.meta.com/tools/faiss/> Faiss: A library for efficient similarity search

the individual’s specific cancer journey and life circumstances, allowing for more nuanced interpretations of emotional states.

- **Temporal Trajectory Context:** Historical diary entries from the participant’s own timeline (e.g., since “Current Day” or “Last Day”, if available), constructing an emotional trajectory from their past diary entries. This temporal information helps capture emotional patterns and transitions, providing insight into how current emotions relate to recent experiences and potentially indicating emerging trends in emotional experiences.
- **Retrieved Peer Experiences:** The RAG module retrieves the k most similar cases (where k ranges from 1 to 20 in our experiments) using FAISS vector similarity search with L2 distance computation. Each retrieved case consists of an emotion diary text and its associated emotional outcomes, serving as concrete examples to ground the model’s predictions.

This structured approach ensures that each prediction is grounded in both individual-specific context and broader patterns observed across the participant population. For each analysis, the new diary entry serves as the primary input, while the LLM is instructed to act as an emotion analysis assistant specifically trained for cancer survivors’ emotional experiences. The prompt is formatted to generate predictions in a consistent JSON structure, enclosed within <PREDICTIONS> tags for reliable parsing. This standardized output format facilitates downstream processing and integration with other components of the emotion analysis pipeline.

5.1.3 Inference and Postprocessing. Using the constructed prompt, the LLM predicts both continuous affective scales and binary emotional states. The model outputs probabilities for each emotional state, ranging from 0.0 to 1.0, where 0.0 indicates the LLM assesses the specific state as unlikely, and 1.0 indicates it as likely. These probabilistic estimates are then converted into binary classifications using a 0.5 threshold. This step ensures compatibility with evaluation metrics and facilitates practical deployment of the framework. In our experiments, the framework utilizes GPT-4o-mini as the base LLM with 0.3 temperature and 1000 max tokens.

5.2 Evaluation Methodology

To assess whether CALLM can accurately predict emotional and behavioral states from brief mobile diary text (RQ2), we established a comprehensive evaluation protocol. We evaluated the framework using metrics including balanced accuracy (i.e., the mean of sensitivity and specificity) and area under the receiver operating characteristic curve (ROC-AUC) across binary emotional and behavioral states.

Notably, both metrics maintain a naive baseline of guessing the majority as 50% (referred as ‘majority baseline’ later), reflecting the expected performance of majority guessing in binary classification, regardless of class imbalance.

The evaluation was conducted using 5-fold grouped cross-validation, with each fold representing a distinct group of participants (which means participants were stratified into five non-overlapping groups, ensuring that all entries from the same participant remained within a single fold to prevent data leakage and maintain the independence of training

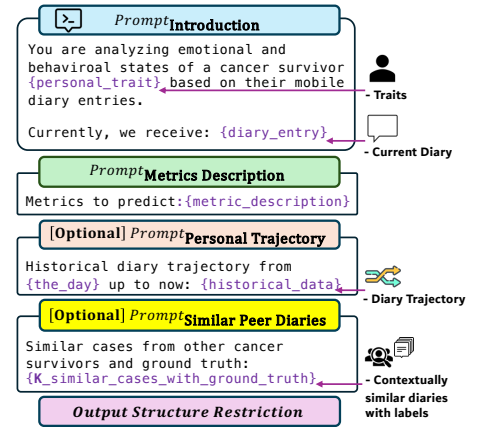


Fig. 7. Illustration of simplified prompt design, with additional details in Figure 12 in Appendix A.1.

and testing sets). This validation approach ensures that the model’s performance metrics reflect its true generalizability to new, unseen participants rather than merely its ability to recognize patterns from familiar individuals.

We compared CALLM to several baseline approaches:

- 1) **Traditional text classification methods:** We implemented Bag-of-Words [27], TF-IDF [64], and LIWC [7] feature extraction, each paired with logistic regression classifiers to enable simultaneous prediction across all target variables.
- 2) **General-purpose transformer models:** We fine-tuned state-of-the-art transformer architectures including BERT [18], SentenceBERT [61], XLNet [90], and RoBERTa [42]. Each model was configured with shared representations and multiple classification heads to simultaneously predict all emotional/behavioral states, leveraging cross-task learning for improved performance.
- 3) **Emotion-specialized transformer models:** We evaluated pre-trained emotion-specific models including EmoBERT [17] (trained on the GoEmotions dataset), RoBERTa-Emotion [42], DeBERTa-Emotion [28], and Emotion-Transformer based on DistilBERT [65]. These models were specifically pre-trained on emotion recognition tasks, providing a strong benchmark for affective text analysis.
- 4) **LLM-based classification:** We tested zero-shot and few-shot approaches using ‘GPT-4o-mini’ (same as the base LLM for CALLM evaluation reported), with standard prompting and random few-shot retrieval from peers’ diaries paired with ground-truth labels (versus semantic similarity-based retrieval used in CALLM). This allowed us to evaluate the effectiveness of CALLM’s design of peer-experience-based retrieval augmentation and individual trajectory incorporation.

All models were evaluated under identical conditions using our 5-fold grouped cross-validation approach as the outer validation approach within nested cross validations for hyper-parameter tuning detailed in Appendix B, ensuring a fair comparison across fundamentally different architectures and learning paradigms.

CALLM differs from the plain LLM baselines along two orthogonal axes:

- (1) **Retrieval-augmented few-shot context:** Rather than sampling examples at random, CALLM retrieves the $k \in \{0, 5, 20\}$ most semantically similar peer diaries paired with its ground-truth labels using FAISS similarity search. Setting $k = 0$ yields a pure zero-shot prompt, while $k > 0$ conducts a similarity-based few-shot scheme, shown as ‘CALLM 0/5/20-shot’ in Table 1.
- (2) **Personal temporal history:** We optionally prepend the participant’s own recent diary trajectory so the model can exploit short-term temporal dependencies. We test three scopes: *none*, *entries since today’s morning*, and *entries since yesterday’s morning*, shown as ‘None Diary’, ‘Diaries Since Current Day (Morning)’, and ‘Diary Since Last Day (Morning)’ in Table 1, respectively.

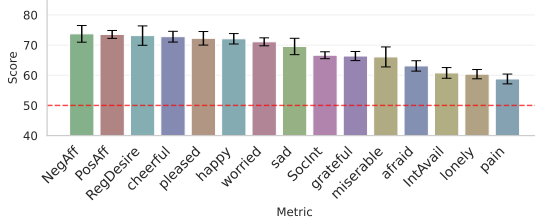
Different combinations of these two factors yields the six CALLM variants in Table 1. Crucially, every other setting—base model (GPT-4o-mini), temperature = 0.3, oken limit = 1000, and JSON output format—remains identical to the LLM baselines, allowing us to isolate the incremental benefit of (i) similarity-based retrieval and (ii) personal history conditioning.

5.3 Prediction Performance Results

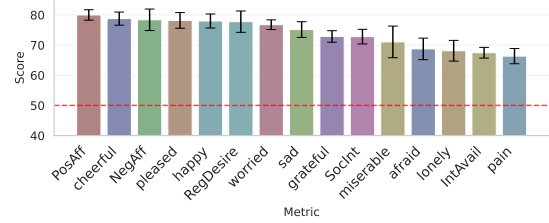
Table 1 reports balanced-accuracy scores for all models on the four target constructs. Across the conventional baselines—including sparse text features paired with logistic regression and fine-tuned or emotion-specialised transformers—the highest performance plateaus at roughly 68–69% for positive affect, 67–68% for negative affect and regulation-desire,

Table 1. Balanced Accuracy (% , mean \pm SD) on four prediction targets. Note, LLM rows show a vanilla GPT-4o-mini prompted with random 0/5/20-shot diaries paired with ground truth measures. CALLM rows add two enhancements: (i) similarity-based retrieval of k peer diaries paired with ground truth measures (k -shot) and (ii) optional personal diary history (*Diary History* = None / Since Current Day / Since Last Day). Bold text indicates the best performance.

Model	PosAff (%)	NegAff (%)	RegDesire (%)	IntAvail (%)
Majority Baseline	50.00 \pm 0.00	50.00 \pm 0.00	50.00 \pm 0.00	50.00 \pm 0.00
BoW	66.20 \pm 1.09	62.05 \pm 1.47	62.17 \pm 1.05	50.12 \pm 1.27
TF-IDF	67.34 \pm 1.67	62.75 \pm 1.16	62.52 \pm 0.83	50.69 \pm 1.29
LIWC	62.67 \pm 0.81	60.29 \pm 1.11	60.69 \pm 1.32	50.00 \pm 0.00
XLNet	61.55 \pm 0.91	57.98 \pm 1.14	59.75 \pm 1.38	52.20 \pm 0.82
BERT	67.99 \pm 1.38	64.39 \pm 2.15	64.81 \pm 1.29	51.42 \pm 1.14
SentenceBERT	68.76 \pm 1.33	67.14 \pm 1.46	66.92 \pm 2.33	52.52 \pm 0.91
EmoBERT	65.81 \pm 1.18	64.18 \pm 1.38	65.54 \pm 2.03	52.41 \pm 0.92
RoBERTa	68.32 \pm 1.42	66.27 \pm 1.50	67.23 \pm 1.80	53.73 \pm 1.20
RoBERTa-Emotion	68.06 \pm 1.42	67.20 \pm 1.25	66.61 \pm 2.08	52.13 \pm 1.10
DeBERTa-Emotion	68.60 \pm 1.90	64.30 \pm 2.40	65.90 \pm 1.52	51.90 \pm 0.70
Emotion-Transformer	68.25 \pm 1.62	66.13 \pm 1.32	66.20 \pm 1.92	55.33 \pm 1.89
LLM Zero-shot	70.14 \pm 1.34	69.98 \pm 2.89	70.22 \pm 2.97	53.16 \pm 1.20
LLM 5-Shot	69.95 \pm 1.42	69.87 \pm 2.76	70.10 \pm 3.05	53.05 \pm 1.18
LLM 20-Shot	69.99 \pm 1.50	69.75 \pm 2.90	69.95 \pm 2.95	52.98 \pm 1.22
CALLM 5-shot + None Diary	72.35 \pm 1.56	72.56 \pm 2.80	71.26 \pm 3.64	53.08 \pm 1.17
CALLM 0-shot + Diaries Since Current Day	70.25 \pm 1.37	72.18 \pm 2.69	73.15 \pm 3.22	55.55 \pm 1.22
CALLM 5-shot + Diaries Since Current Day	72.31 \pm 1.37	72.75 \pm 2.85	71.77 \pm 3.40	56.04 \pm 1.15
CALLM 0-shot + Diaries Since Last Day	70.42 \pm 1.47	72.30 \pm 2.68	73.72 \pm 3.40	58.85 \pm 1.77
CALLM 5-shot + Diaries Since Last Day	72.26 \pm 1.46	72.65 \pm 2.65	71.45 \pm 3.60	60.09 \pm 1.15
CALLM 20-shot + Diaries Since Last Day	72.96 \pm 1.54	73.29 \pm 2.89	71.86 \pm 3.12	56.31 \pm 1.23



(a) Balanced Accuracy.



(b) ROC-AUC.

Fig. 8. Model performance across different emotional and behavioral states. The red dashed line indicates the naive majority baseline of 50%.

and 55% for intervention availability. These figures are only a few points above the 50% majority baseline of balanced accuracy, underscoring a clear ceiling when contextual information is absent.

For LLM baselines with no or random context awareness, the zero-shot setting lifts the upper bound for the three subjective constructs to $\approx 70\%$, yet adding five or twenty randomly selected peer exemplars provides no additional benefit and in some folds yields a small decline. The lack of improvement suggests that exemplars unaligned with the index entry’s situational context introduce distracting cues that offset any nominal few-shot advantage.

In contrast, every CALLM variant surpasses baseline configurations. The best balanced-accuracy scores reach 72.96% for positive affect, 73.29% for negative affect, 73.72% for regulation-desire, and 60.09% for intervention availability,

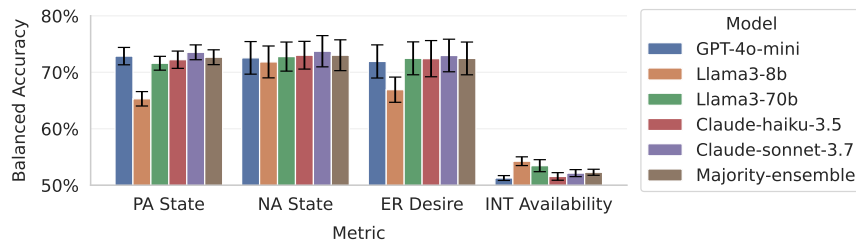


Fig. 9. Model Performance Comparison across Metrics. Error bars represent standard deviation across 5-fold grouped cross-validation.

improving the ceiling compared to the baseline methods. For affective states, the largest boost arises when similarity-matched peer exemplars are combined with the participant’s diary trajectory since the prior day, implying that survivors share sufficiently convergent affective contexts for cross-participant retrieval to be informative. Regulation desire, however, improves significantly only when personal trajectory is present; adding peer examples alone offers no benefit and occasionally degrades performance, pointing to the idiosyncratic nature of moment-to-moment motivation to alter one’s desire to regulate emotions. Availability detection remains the most challenging task, while CALLM raises accuracy to about 60%, indicating that this objective behavioral constraint may require additional multimodal signals beyond diary text.

Taken together, the results confirm that large language models can already decode brief diary entries more accurately than traditional text-classification pipelines, but that carefully curated contextual prompts—rather than randomly chosen exemplars—are essential for pushing performance beyond the 70 % ceiling and for extracting the nuanced temporal cues embedded in participants’ recent histories.

5.4 Generalizability Discussion

To assess whether CALLM generalizes to the full spectrum of nuanced emotional and behavioral states captured in our dataset, we extended the evaluation to a wider set of constructs beyond the four core targets reported in Table 1.

As shown in Figure 8, the model achieves balanced accuracies ranging from 58.8% to 73.2% and ROC-AUC scores from 66.4% to 79.5% across different states, consistently outperforming the majority baseline (50%). Beyond the core constructs reported above, the model shows robust performance in detecting specific emotional states such as cheerfulness (72.8%), pleasure (71.6%), and worry (71.1%). Even for more challenging states like fear (63.1%), loneliness (60.4%), and physical pain (58.8%), the model maintains above-chance performance, highlighting its broad applicability across diverse emotional and behavioral dimensions in cancer survivors.

We also evaluated CALLM’s generalizability across different base LLMs. See Figure 9, comparing GPT-4o-mini (OpenAI), Llama-3 models (8b and 70b), and Claude models (3.5-haiku and 3.7-sonnet), all models demonstrated consistent effectiveness in analyzing cancer survivors’ diary entries. Ensembling the models (aggregating predictions from all models to conduct majority voting) does not contribute to better performance. Claude 3.7-sonnet slightly outperformed others across most constructs. Llama3-8b exhibited the lowest performance, perhaps due to the smaller model size which may not have enough power to capture the nuances of the data. Besides Llama3-8b, models’ performance differences were minimal for most constructs (within 2-3%), with high prediction result correlations between models (97%), demonstrating CALLM’s robustness across different LLM implementations.

6 Analysis of Performance Determinants

To address **RQ3** concerning the factors that influence model performance, we conducted several post-hoc analyses exploring how model confidence, diary entry length, temporal context, and personalization affect prediction accuracy. These analyses provide valuable insights for real-world deployment of LLM-based mobile diary analysis systems.

6.1 Model Confidence and Prediction Accuracy

Our first analysis examined the relationship between the LLM’s confidence (probability estimates) and prediction accuracy. By varying confidence thresholds (defined as $|\text{probability} - 50\%|$), we found that all states except intervention availability showed improved prediction accuracy as confidence increased. For positive affect, accuracy improved dramatically from 73% to 88% at the highest confidence levels (Figure 10). This finding suggests that in practical applications, the system could use confidence scores to selectively trigger interventions only when predictions reach a certain reliability threshold.

6.2 Effect of Diary Entry Length

We also investigated how diary entry length affects prediction performance by analyzing entries containing 3–15 words. Spearman rank correlations revealed that for most emotional states (9 out of 15 analyzed), longer entries significantly improved prediction accuracy. Positive affect showed a strong positive correlation ($\rho = 0.874$, adjusted $p < 0.01$) with entry length, while intervention availability demonstrated a strong negative correlation ($\rho = -0.996$, adjusted $p < 0.001$), suggesting that shorter entries may be more effective for this particular prediction. Negative affect showed no significant relationship with text length ($\rho = -0.276$, adjusted $p = 0.472$). These findings provide practical guidance for optimizing prompt length in mobile diary applications, potentially with construct-specific instructions.

6.3 Temporal Prediction Capabilities

To explore the potential for forecasting emotional states—valuable for proactive intervention—we tested whether current-day diary entries could predict next-day emotional states. The analysis showed modest but consistent predictive power, with balanced accuracies between 54–56% across constructs, significantly above chance level (Wilcoxon signed-rank test, adjusted $p < 0.001$). As shown in Figure 11, prediction accuracy varied by time of day, with morning states being marginally more predictable than afternoon or evening periods. This temporal pattern suggests that emotional

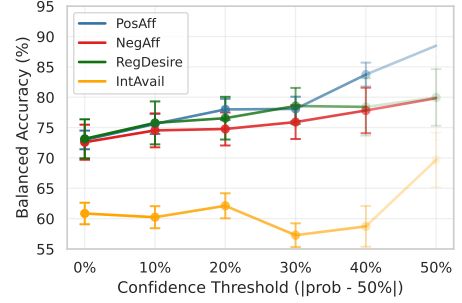


Fig. 10. Balanced Accuracy performance for predicted states at different LLM confidence thresholds ($|\text{probability} - 50\%|$; ranging from 0 to 50%, higher confidence thresholds yield better accuracy but fewer samples). Transparency indicates sample retention rate; error bars show standard deviations across 5-fold grouped cross-validation.

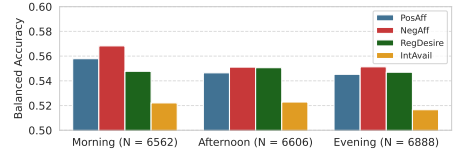


Fig. 11. Next-day prediction performance by time period. The figure shows balanced accuracy across different emotional and behavioral states when predicting morning, afternoon, and evening states of the following day.

states may be more stable and predictable in the morning, possibly before the accumulated effects of daily stressors. While the predictive capability is limited, it offers potential for anticipatory intervention planning, particularly for morning-focused interventions based on the previous day’s diary content.

6.4 Benefits of Personalization

Finally, we examined whether incorporating individual-level ground truth data could enhance prediction accuracy. By simulating a scenario where the system collects one week of user self-reports before operating independently, we found statistically significant improvements for key emotional states: positive affect (1.53% improvement, adjusted $p < 0.05$), negative affect (1.42% improvement, adjusted $p < 0.05$), and social interaction quality (3.62% improvement, adjusted $p < 0.01$). These results suggest that a short calibration period can meaningfully improve the personalization of emotional state predictions in a practical deployment scenario.

7 Discussion

Here we discuss the findings and takeaways from the present study with respect to the research questions, clinical implications we inferred from the findings, practical deployment considerations towards actionable intervention systems.

7.1 Findings and takeaways

Our research addressed three key questions about analyzing cancer survivors’ mobile diary entries. For **RQ1**, our analysis revealed systematic relationships between described contexts and emotional states—administrative and health-related contexts were associated with greater negative affect and regulation needs, while leisure and social activities promoted positive affect. Even without explicit emotional language, these brief diary entries reflected cancer survivors’ unique experiences and indicated their emotional states. Importantly, intervention availability varied significantly across contexts (higher at home, lower during work/study), highlighting opportunities for context-sensitive support. The significant proportion of variance explained by individual random effects underscores the importance of personalization, as participants often responded differently under similar contextual circumstances.

For **RQ2**, we demonstrated that CALLM successfully leverages both peer experiences and temporal diary context to analyze emotional states from brief diary entries, enabling mobile health applications to predict these states from the emotion drivers found in diary entries alone. This approach minimizes user burden by harnessing the information in one brief diary response rather than having users repeatedly report on their current emotional and behavioral states across many distinct questionnaire items. By combining LLMs’ pre-trained knowledge, personal diary history, with RAG-based peer exemplars, CALLM delivers flexible, context-aware predictions—even from sparse data—consistently outperformed baselines with balanced accuracies of 72.96% for positive affect, 73.29% for negative affect, and 73.72% for emotion regulation desire. Its consistent performance across various emotional constructs attests to both its effectiveness and generalizability. Moreover, the RAG-enabled peers’ contextually-similar exemplar retrieval is particularly valuable in privacy-sensitive healthcare settings, as it enables learning from an anonymized database of peer experiences and personal diary history maintained on local devices.

Addressing **RQ3**, our post-hoc analyses revealed four practical insights that influence prediction performance: (1) Model confidence strongly predicts accuracy, with improvements up to 15 percentage points when filtering for high-confidence predictions, suggesting confidence thresholds in deployed systems could significantly enhance reliability. (2) Longer diary entries generally enhance predictive performance for affective states, but overly long entries offer

diminishing returns, indicating a trade-off between information gain and user burden—especially when analytical efficiency is prioritized. This observation aligns with prior findings on optimizing assessment length [40]. (3) Diaries provide modest yet consistent predictive value for next-day emotional states, particularly for morning periods, suggesting value for short-term rather than long-range forecasting. and (4) Brief personalization periods—such as collecting one week of self-reported ground truth measures and incorporating them into the retrieval process—can yield meaningful improvements in prediction accuracy, particularly for emotional states.

Both our pre-hoc and post-hoc analyses demonstrate that context-aware language models can illuminate the complex, dynamic emotional experiences of cancer survivors through their open-ended diary entries. These findings collectively demonstrate the potential of extracting rich contextual understanding from brief mobile diary entries to inform timing, content, and personalization of mobile interventions while respecting individual differences in emotional responses to similar contexts—with important implications for developing more effective, less burdensome mobile health applications.

7.2 Clinical Implications

DMHIs can increase access to needed mental health care for cancer survivors by overcoming the significant barriers to care that they otherwise face (e.g., social stigma [30], financial costs [54], time burdens [8], and lack of available providers [10]). Despite cancer survivors finding digital interventions acceptable [14], sustained engagement with digital interventions is a widespread problem [21]. Factors that can impact user engagement relate to the timing of interventions, their dosage, and the burden placed on individuals to use them. Just-in-time adaptive intervention (JITAI) frameworks offer a way to improve intervention engagement and effectiveness by increasing the personal relevance of intervention timing, content, and dosage [51]. However, researchers have yet to identify a convenient and low-burden method for identifying when to intervene and what intervention to provide to patients. This is in part due to a tension between gathering enough information from participants to accurately inform when to provide an intervention, and a logical desire to minimize burden by not requiring them to report on their emotions multiple times per day using the same questionnaire items. Minimizing user burden is critical to mitigate the risk of individuals losing interest and dropping out of an intervention altogether.

The CALLM framework proposed here represents an important step towards reducing user burden without sacrificing the personalized, contextual information required to guide optimized intervention delivery decisions. For instance, we demonstrate that brief emotion-driver diary entries—each comprising just a few words—can predict whether an individual is feeling more negative than usual, more positive than usual, or desires to change their emotional state with approximately 73% accuracy, rising further when LLMs report higher confidence levels. Although digital diaries have been widely adopted for emotional self-regulation, symptom monitoring, and reflective journaling, CALLM extends these capabilities by providing a lightweight channel for real-time inference of user states and intervention opportunities. Moreover, passive sensing of behavioral and contextual cues—such as GPS-derived location (e.g., workspace, home, or outdoor), ambient sound patterns (indicating social activity versus solitude), and phone-usage logs (reflecting availability for interaction)—offers an objective complement to subjective desire measurements [25]. By integrating these signals, CALLM can potentially supplement the roughly 60% accuracy limitation of intervention-availability detection based solely on diaries, enabling more precisely timed, personalized interventions that enhance both relevance and engagement.

The ability to reasonably approximate these emotion factors and intervention opportunities in survivors through their brief diaries is a drastic improvement over the traditional method of asking individuals about each of these states in separate questionnaire items, to determine whether an intervention is warranted. Furthermore, our findings suggest

complementary integration with passive sensing data (such as location, activity, and device usage patterns) could create even more comprehensive context awareness, as we found systematic relationships between certain contexts (e.g., home activities, administrative tasks) and intervention opportunity components. Such multi-modal approaches could further reduce user burden while improving prediction accuracy, especially for constructs like intervention availability that showed more modest prediction performance from text alone. CALLM’s potential to improve JITAI delivery systems while minimizing user input is an exciting path forward to better address the mental health needs of cancer survivors.

7.3 Deployment Considerations

Healthcare applications of emotion analysis require privacy safeguards under HIPAA [2] and GDPR [60], including encrypted storage, clear consent processes, and data minimization practices. Our RAG-based approach offers an additional privacy advantage by using anonymized peer experiences rather than requiring extensive personal data collection, aligning with privacy-by-design principles.

Computational efficiency is another practical consideration. Taking the GPT-4o-mini as an example, processing costs (\$0.0005 per diary entry) remain feasible at scale, with potential for further reduction through smaller, specialized models. Recent advances with sub-billion parameter models like MobileLLM [43] suggest on-device deployment is viable, enhancing both privacy and accessibility.

Implementing our post-hoc insights would further optimize real-world deployment. For instance, adaptive confidence thresholds could dynamically balance prediction reliability against coverage, triggering interventions only when confidence exceeds certain thresholds. Similarly, contextual guidance about optimal diary entry length could improve user experience while maintaining prediction quality. For healthcare systems integration, the framework could provide confidence scores alongside predictions, giving clinicians transparency into reliability when reviewing automated assessments. These practical considerations help bridge the gap between research findings and sustainable implementation in clinical settings.

7.4 Future Work

This study represents a first step in using mobile health data and LLMs to identify cancer survivors’ emotional and behavioral states and inform personalized DMHIs, with limitations to acknowledge which open up future directions. It is worth noting that our study utilized GPT-4o-mini as the base LLM - while this demonstrates the potential of even smaller language models in healthcare applications, a comprehensive benchmark of different LLM architectures was not our focus and remains an open direction for future research.

Notably, future work should address demographic limitations of our study, with data skewed toward female (90.17%) and White (86.67%) participants. Extending validation to more diverse populations is critical, as recent research suggests language models may perform differently across demographic groups [5], potentially affecting equitable emotion analysis.

While this study focuses on proactive diary entries, future work could incorporate complementary “passive sensing” data such as phone usage, location, and physiological signals [25, 82]. Such passive data could supplement analysis when diary entries are unavailable, creating a more resilient system while minimizing user burden. Additionally, while the diary entries were collected in participants’ everyday lives, participants were compensated for completing the study. As such, future diary collection efforts without participant compensation may encounter different engagement patterns.

Additionally, while LLMs demonstrate impressive performance in our emotion analysis tasks, their “black box” nature presents challenges for clinical applications where understanding prediction rationale is essential. Future work

should explore better interpretability through prediction explanations, user-oriented feedback loop, and alignment with established clinical scales. Meanwhile, our finding that model confidence correlates strongly with prediction accuracy offers a practical mechanism for communicating reliability to clinicians and patients.

Towards more personalized emotion management solutions, such as JITAIs [85] and embodied conversational chatbot [12], the next step might be longer-term context studies to understand how LLM systems adapt to changing emotional patterns and life circumstances of cancer survivors. Future work should explore mechanisms for continuous learning and adaptation of the model to individual users over extended periods, particularly in handling evolving emotional needs and identifying shifting patterns in intervention opportunities during different phases of survivorship. We envision this study as a feasibility test toward building LLM-empowered agentic workflows that incorporate memory of individual users' context trajectories, sensor- and diary-based context understanding, retrieval-augmented external knowledge integration, and theory-informed workflow automation, leveraging LLM systems' neural-symbolic and task execution potential [72, 86]. This could ultimately lead to more personalized and context-aware emotional management solutions for cancer survivors and potentially broader populations.

8 Conclusion

This paper introduced CALLM, a context-aware framework for analyzing brief mobile diary entries from cancer survivors. Using only diary entries collected via smartphones and validated with data from 407 cancer survivors, CALLM performed well in predicting emotional states and variables relevant to personalized intervention opportunities (desire to regulate emotions and availability to engage in an interventions). By uniquely integrating retrieved peer experiences with individual temporal trajectories, our framework demonstrates how contextual information in mobile diaries can be effectively leveraged to understand users' emotional experiences (RQ1), predict key states from brief text with reasonable accuracy (RQ2), and identify factors that optimize prediction performance like confidence thresholds and diary length (RQ3). These contributions highlight how context-aware computational approaches can transform lightweight user inputs into rich behavioral insights while addressing the critical challenge of identifying optimal moments for intervention delivery—potentially reducing assessment burden while enabling more responsive and personalized mobile health applications.

References

- [1] Gregory D Abowd, Anind K Dey, Peter J Brown, Nigel Davies, Mark Smith, and Pete Steggles. 1999. Towards a better understanding of context and context-awareness. In *Handheld and Ubiquitous Computing: First International Symposium, HUC'99 Karlsruhe, Germany, September 27–29, 1999 Proceedings 1*. Springer, 304–307.
- [2] Accountability Act. 1996. Health insurance portability and accountability act of 1996. *Public law* 104 (1996), 191.
- [3] Nabil Alshurafa, Jayalakshmi Jain, Rawan Alharbi, Gleb Iakovlev, Bonnie Spring, and Angela Pfammatter. 2018. Is More Always Better? Discovering Incentivized mHealth Intervention Engagement Related to Health Behavior Trends. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 2, 4, Article 153 (Dec. 2018), 26 pages. doi:10.1145/3287031
- [4] Michael A Andrykowski, Emily Lykins, and Andrea Floyd. 2008. Psychological health in cancer survivors. In *Seminars in oncology nursing*, Vol. 24. Elsevier, 193–201.
- [5] 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).
- [6] Niall Bolger, Angelina Davis, and Eshkol Rafaeli. 2003. Diary methods: Capturing life as it is lived. *Annual review of psychology* 54, 1 (2003), 579–616.
- [7] Ryan L Boyd, Ashwini Ashokkumar, Sarah Seraj, and James W Pennebaker. 2022. The development and psychometric properties of LIWC-22. *Austin, TX: University of Texas at Austin* 10 (2022).
- [8] Katherine Bradbury, Mary Steele, Teresa Corbett, Adam WA Geraghty, Adele Krusche, Elena Heber, Steph Easton, Tara Cheetham-Blake, Joanna Slodkowska-Barabasz, Andre Matthias Müller, et al. 2019. Developing a digital intervention for cancer survivors: an evidence-, theory- and person-based approach. *NPJ digital medicine* 2, 1 (2019), 85.

- [9] Michael S Businelle, Scott T Walters, Eun-Young Mun, Thomas R Kirchner, Emily T Hébert, and Xiaoyin Li. 2020. Reducing drinking among people experiencing homelessness: protocol for the development and testing of a just-in-time adaptive intervention. *JMIR Research Protocols* 9, 4 (2020), e15610.
- [10] Louise Camm-Crosbie, Louise Bradley, Rebecca Shaw, Simon Baron-Cohen, and Sarah Cassidy. 2019. ‘People like me don’t get support’: Autistic adults’ experiences of support and treatment for mental health difficulties, self-injury and suicidality. *Autism* 23, 6 (2019), 1431–1441.
- [11] Zhiyu Chen, Yujie Lu, and William Yang Wang. 2023. Empowering psychotherapy with large language models: Cognitive distortion detection through diagnosis of thought prompting. *arXiv preprint arXiv:2310.07146* (2023).
- [12] 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).
- [13] Briana K Clifford, David Mizrahi, Carolina X Sandler, Benjamin K Barry, David Simar, Claire E Wakefield, and David Goldstein. 2018. Barriers and facilitators of exercise experienced by cancer survivors: a mixed methods systematic review. *Supportive care in cancer* 26 (2018), 685–700.
- [14] Teresa Corbett, Karpaul Singh, Liz Payne, Katherine Bradbury, Claire Foster, Eila Watson, Alison Richardson, Paul Little, and Lucy Yardley. 2018. Understanding acceptability of and engagement with Web-based interventions aiming to improve quality of life in cancer survivors: a synthesis of current research. *Psycho-oncology* 27, 1 (2018), 22–33.
- [15] Andy Davies, Eiko Fried, Omar Costilla-Reyes, and Hane Aung. 2023. Individual Behavioral Insights in Schizophrenia: A Network Analysis and Mobile Sensing Approach. In *International Conference on Pervasive Computing Technologies for Healthcare*. Springer, 18–33.
- [16] Lianne P De Vries, Bart ML Baselmans, and Meike Bartels. 2021. Smartphone-based ecological momentary assessment of well-being: A systematic review and recommendations for future studies. *Journal of Happiness Studies* 22, 5 (2021), 2361–2408.
- [17] Dorottya Demszky, Dana Movshovitz-Attias, Jeongwoo Ko, Alan Cowen, Gaurav Nemade, and Sujith Ravi. 2020. GoEmotions: A Dataset of Fine-Grained Emotions. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 4040–4054.
- [18] Jacob Devlin. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [19] Johannes C Eichstaedt, Robert J Smith, Raina M Merchant, Lyle H Ungar, Patrick Crutchley, Daniel Preotiuc-Pietro, David A Asch, and H Andrew Schwartz. 2018. Facebook language predicts depression in medical records. *Proceedings of the National Academy of Sciences* 115, 44 (2018), 11203–11208.
- [20] Jenny Firkins, Lissi Hansen, Martha Driessnack, and Nathan Dieckmann. 2020. Quality of life in “chronic” cancer survivors: a meta-analysis. *Journal of Cancer Survivorship* 14 (2020), 504–517.
- [21] Theresa Fleming, Lynda Bavin, Mathijs Lucassen, Karolina Stasiak, Sarah Hopkins, and Sally Merry. 2018. Beyond the trial: systematic review of real-world uptake and engagement with digital self-help interventions for depression, low mood, or anxiety. *Journal of medical Internet research* 20, 6 (2018), e199.
- [22] Burkhardt Funk, Shiri Sadeh-Sharvit, Ellen E Fitzsimmons-Craft, Mickey Todd Trockel, Grace E Monterubio, Neha J Goel, Katherine N Balantekin, Dawn M Eichen, Rachael E Flatt, Marie-Laure Firebaugh, et al. 2020. A framework for applying natural language processing in digital health interventions. *Journal of medical Internet research* 22, 2 (2020), e13855.
- [23] Chenyang Gao, Ivan Marsic, Aleksandra Sarcevic, Waverly Gestrinch-Thompson, and Randall S. Burd. 2023. Real-time Context-Aware Multimodal Network for Activity and Activity-Stage Recognition from Team Communication in Dynamic Clinical Settings. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 7, 1, Article 12 (March 2023), 28 pages. doi:10.1145/3580798
- [24] Sharath Chandra Guntuku, Anneke Buffone, Kokil Jaidka, Johannes C Eichstaedt, and Lyle H Ungar. 2019. Understanding and measuring psychological stress using social media. In *Proceedings of the international AAAI conference on web and social media*, Vol. 13. 214–225.
- [25] Gabriella M Harari and Samuel D Gosling. 2023. Understanding behaviours in context using mobile sensing. *Nature Reviews Psychology* 2, 12 (2023), 767–779.
- [26] Wendy Hardeman, Julie Houghton, Kathleen Lane, Andy Jones, and Felix Naughton. 2019. A systematic review of just-in-time adaptive interventions (JITAs) to promote physical activity. *International Journal of Behavioral Nutrition and Physical Activity* 16 (2019), 1–21.
- [27] ZS Harris. 1954. Distributional structure.
- [28] Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. DEBERTA: DECODING-ENHANCED BERT WITH DISENTANGLED ATTENTION. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=XPZiaotutsD>
- [29] Maria Hewitt and Julia H Rowland. 2002. Mental health service use among adult cancer survivors: analyses of the National Health Interview Survey. *Journal of Clinical Oncology* 20, 23 (2002), 4581–4590.
- [30] Jimmie C Holland, Brian J Kelly, and Mark I Weinberger. 2010. Why psychosocial care is difficult to integrate into routine cancer care: stigma is the elephant in the room. *Journal of the National Comprehensive Cancer Network* 8, 4 (2010), 362–366.
- [31] Linda A Jacobs and Lawrence N Shulman. 2017. Follow-up care of cancer survivors: challenges and solutions. *The Lancet Oncology* 18, 1 (2017), e19–e29.
- [32] Gyeyoung Jung, Soyeon Choi, Yuju Kang, and Jaejeung Kim. 2024. MyListener: An AI-Mediated Journaling Mobile Application for Alleviating Depression and Loneliness Using Contextual Data. In *Companion of the 2024 on ACM International Joint Conference on Pervasive and Ubiquitous Computing* (Melbourne VIC, Australia) (*UbiComp '24*). Association for Computing Machinery, New York, NY, USA, 137–141. doi:10.1145/3675094.3677601
- [33] Inyeop Kim, Hwarang Goh, Nematjon Narziev, Youngtae Noh, and Uichin Lee. 2020. Understanding User Contexts and Coping Strategies for Context-aware Phone Distraction Management System Design. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 4, 4, Article 134 (Dec. 2020), 134 pages.

- 33 pages. doi:10.1145/3432213
- [34] Taewan Kim, Seolyeong Bae, Hyun Ah Kim, Su-Woo Lee, Hwajung Hong, Chanmo Yang, and Young-Ho Kim. 2024. MindfulDiary: Harnessing Large Language Model to Support Psychiatric Patients' Journaling. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) (CHI '24). Association for Computing Machinery, New York, NY, USA, Article 701, 20 pages. doi:10.1145/3613904.3642937
 - [35] Predrag Klasnja and Wanda Pratt. 2012. Healthcare in the pocket: mapping the space of mobile-phone health interventions. *Journal of biomedical informatics* 45, 1 (2012), 184–198.
 - [36] Evan M Kleiman, Brianna J Turner, Szymon Fedor, Eleanor E Beale, Rosalind W Picard, Jeff C Huffman, and Matthew K Nock. 2018. Digital phenotyping of suicidal thoughts. *Depression and anxiety* 35, 7 (2018), 601–608.
 - [37] Lian Van Der Krieke, Bertus F Jeronimus, Frank J Blaauw, Rob BK Wanders, Ando C Emerencia, Hendrika M Schenk, Stijn De Vos, Evelien Snippe, Marieke Wichers, Johanna TW Wigman, et al. 2016. HowNutsAreTheDutch (HoeGekIsNL): A crowdsourcing study of mental symptoms and strengths. *International journal of methods in psychiatric research* 25, 2 (2016), 123–144.
 - [38] Jong Ho Lee, Sunghoon Ivan Lee, and Eun Kyoung Choe. 2024. GoalTrack: Supporting Personalized Goal-Setting in Stroke Rehabilitation with Multimodal Activity Journaling. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 8, 4, Article 167 (Nov. 2024), 29 pages. doi:10.1145/3699723
 - [39] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems* 33 (2020), 9459–9474.
 - [40] Jixin Li, Aditya Ponnada, Wei-Lin Wang, Genevieve Dunton, and Stephen Intille. 2024. Ask Less, Learn More: Adapting Ecological Momentary Assessment Survey Length by Modeling Question-Answer Information Gain. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 8, 4, Article 166 (Nov. 2024), 32 pages. doi:10.1145/3699735
 - [41] Myles-Jay Anthony Linton, Sarah Jelbert, Judi Kidger, Richard Morris, Lucy Biddle, and Bruce Hood. 2021. Investigating the use of electronic well-being diaries completed within a psychoeducation program for university students: Longitudinal text analysis study. *Journal of medical Internet research* 23, 4 (2021), e25279.
 - [42] Yinhan Liu, Mye Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019).
 - [43] Zechun Liu, Changsheng Zhao, Forrest Iandola, Chen Lai, Yuandong Tian, Igor Fedorov, Yunyang Xiong, Ernie Chang, Yangyang Shi, Raghuraman Krishnamoorthi, et al. 2024. Mobilellm: Optimizing sub-billion parameter language models for on-device use cases. In *Forty-first International Conference on Machine Learning*.
 - [44] Steven Loria et al. 2018. textblob Documentation. *Release 0.15 2*, 8 (2018), 269.
 - [45] Zilin Ma, Yiyang Mei, and Zhaoyuan Su. 2024. 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. 1105.
 - [46] Deborah K Mayer, Shelly Fuld Nasso, and Jo Anne Earp. 2017. Defining cancer survivors, their needs, and perspectives on survivorship health care in the USA. *The Lancet Oncology* 18, 1 (2017), e11–e18.
 - [47] Kiki Metsäranta, Marjo Kurki, Maritta Valimäki, and Minna Anttila. 2019. How do adolescents use electronic diaries? A mixed-methods study among adolescents with depressive symptoms. *Journal of Medical Internet Research* 21, 2 (2019), e11711.
 - [48] David C Mohr, Mi Zhang, and Stephen M Schueller. 2017. Personal sensing: understanding mental health using ubiquitous sensors and machine learning. *Annual review of clinical psychology* 13, 1 (2017), 23–47.
 - [49] Margaret E Morris, Qusai Kathawala, Todd K Leen, Ethan E Gorenstein, Farzin Guilak, Michael Labhard, and William Deleeuw. 2010. Mobile therapy: case study evaluations of a cell phone application for emotional self-awareness. *Journal of medical Internet research* 12, 2 (2010), e10.
 - [50] Abdulsalam Salihu Mustafa, Nor'ashikin Ali, Jaspaljeet Singh Dhillon, Gamal Alkaws, and Yahia Baashar. 2022. User engagement and abandonment of mHealth: a cross-sectional survey. In *Healthcare*, Vol. 10. MDPI, 221.
 - [51] 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.
 - [52] Subigya Nepal, Arvind Pillai, William Campbell, Talie Massachi, Eunsol Soul Choi, Xuhai Xu, Joanna Kuc, Jeremy F Huckins, Jason Holden, Colin Depp, Nicholas Jacobson, Mary P Czerwinski, Eric Granholm, and Andrew Campbell. 2024. Contextual AI Journaling: Integrating LLM and Time Series Behavioral Sensing Technology to Promote Self-Reflection and Well-being using the MindScape App. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) (CHI EA '24). Association for Computing Machinery, New York, NY, USA, Article 86, 8 pages. doi:10.1145/3613905.3650767
 - [53] Subigya Nepal, Arvind Pillai, William Campbell, Talie Massachi, Michael V. Heinz, Ashmita Kunwar, Eunsol Soul Choi, Xuhai Xu, Joanna Kuc, Jeremy F. Huckins, Jason Holden, Sarah M. Preum, Colin Depp, Nicholas Jacobson, Mary P. Czerwinski, Eric Granholm, and Andrew T. Campbell. 2024. MindScape Study: Integrating LLM and Behavioral Sensing for Personalized AI-Driven Journaling Experiences. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 8, 4, Article 186 (Nov. 2024), 44 pages. doi:10.1145/3699761
 - [54] Shehzad Niazi, Emily Vargas, Aaron Spaulding, Elaine Gustetic, Nancy Ford, David Paly, Kelsey Tatum, Matthew M Clark, and Teresa Rummans. 2020. Barriers to accepting mental health care in cancer patients with depression. *Social Work in Health Care* 59, 6 (2020), 351–364.
 - [55] Alicia L. Nobles, Jeffrey J. Glenn, Kamran Kowsari, Bethany A. Teachman, and Laura E. Barnes. 2018. Identification of Imminent Suicide Risk Among Young Adults using Text Messages. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems* (Montreal QC, Canada) (CHI

- '18). Association for Computing Machinery, New York, NY, USA, 1–11. doi:10.1145/3173574.3173987
- [56] Zizi Papacharissi. 2010. *A networked self: Identity, community, and culture on social network sites*. Routledge.
- [57] Soujanya Poria, Navonil Majumder, Rada Mihalcea, and Eduard Hovy. 2019. Emotion recognition in conversation: Research challenges, datasets, and recent advances. *IEEE access* 7 (2019), 100943–100953.
- [58] Narelle Powers, Judith Gullifer, and Rhonda Shaw. 2016. When the treatment stops: A qualitative study of life post breast cancer treatment. *Journal of Health Psychology* 21, 7 (2016), 1371–1382.
- [59] Varun Reddy, Zhiyuan Wang, Emma R Toner, Maria A Larrazabal, Mehdi Boukhechba, Bethany A Teachman, and Laura E Barnes. 2024. Audioinsight: Detecting social contexts relevant to social anxiety from speech. In *2024 12th International Conference on Affective Computing and Intelligent Interaction (ACII)*. IEEE, 55–62.
- [60] Protection Regulation. 2018. General data protection regulation. *Intouch* 25 (2018), 1–5.
- [61] N Reimers. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. *arXiv preprint arXiv:1908.10084* (2019).
- [62] Misha Sadeghi, Bernhard Egger, Reza Agahi, Robert Richer, Klara Capito, Lydia Helene Rupp, Lena Schindler-Gmelch, Matthias Berking, and Bjoern M Eskofier. 2023. Exploring the capabilities of a language model-only approach for depression detection in text data. In *2023 IEEE EMBS International Conference on Biomedical and Health Informatics (BHI)*. IEEE, 1–5.
- [63] John M Salsman, Laurie E McLouth, Janet A Tooze, Denisha Little-Greene, Michael Cohn, Mia Sorkin Kehoe, and Judith T Moskowitz. 2023. An eHealth, positive emotion skills intervention for enhancing psychological well-being in young adult Cancer survivors: results from a multi-site, pilot feasibility trial. *International Journal of Behavioral Medicine* 30, 5 (2023), 639–650.
- [64] Gerard Salton and Christopher Buckley. 1988. Term-weighting approaches in automatic text retrieval. *Information processing & management* 24, 5 (1988), 513–523.
- [65] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108* (2019).
- [66] Elizabeth M Seabrook, Margaret L Kern, Ben D Fulcher, and Nikki S Rickard. 2018. Predicting depression from language-based emotion dynamics: longitudinal analysis of Facebook and Twitter status updates. *Journal of medical Internet research* 20, 5 (2018), e168.
- [67] 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.
- [68] Ashish Sharma, Kevin Rushton, Inna Wanyin Lin, David Wadden, Khendra G Lucas, Adam S Miner, Theresa Nguyen, and Tim Althoff. 2023. Cognitive reframing of negative thoughts through human-language model interaction. *arXiv preprint arXiv:2305.02466* (2023).
- [69] Dai Shi, Dan Tao, Jiangtao Wang, Muyan Yao, Zhibo Wang, Houjin Chen, and Sumi Helal. 2021. Fine-Grained and Context-Aware Behavioral Biometrics for Pattern Lock on Smartphones. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 5, 1, Article 33 (March 2021), 30 pages. doi:10.1145/3448080
- [70] Saul Shiffman, Arthur A Stone, and Michael R Hufford. 2008. Ecological momentary assessment. *Annu. Rev. Clin. Psychol.* 4, 1 (2008), 1–32.
- [71] Timothy Sohn, Kevin A Li, William G Griswold, and James D Hollan. 2008. A diary study of mobile information needs. In *Proceedings of the sigchi conference on human factors in computing systems*. 433–442.
- [72] Elizabeth C Stade, Shannon Wiltsey Stirman, Lyle H Ungar, Cody L Boland, H Andrew Schwartz, David B Yaden, João Sedoc, Robert J DeRubeis, Robb Willer, and Johannes C Eichstaedt. 2024. Large language models could change the future of behavioral healthcare: a proposal for responsible development and evaluation. *NPJ Mental Health Research* 3, 1 (2024), 12.
- [73] Lesleigh Stinson, Yunchao Liu, and Jesse Dallery. 2022. Ecological momentary assessment: a systematic review of validity research. *Perspectives on Behavior Science* 45, 2 (2022), 469–493.
- [74] Arthur A Stone, Stefan Schneider, and Joshua M Smyth. 2023. Evaluation of pressing issues in ecological momentary assessment. *Annual Review of Clinical Psychology* 19, 1 (2023), 107–131.
- [75] Arthur A Stone, Cheng K Fred Wen, Stefan Schneider, and Doerte U Junghaenel. 2020. Evaluating the effect of daily diary instructional phrases on respondents' recall time frames: survey experiment. *Journal of Medical Internet Research* 22, 2 (2020), e16105.
- [76] Marcin Strackiewicz, Peter James, and Jukka-Pekka Onnela. 2021. A systematic review of smartphone-based human activity recognition methods for health research. *NPJ Digital Medicine* 4, 1 (2021), 148.
- [77] Maya Tamir. 2016. Why do people regulate their emotions? A taxonomy of motives in emotion regulation. *Personality and social psychology review* 20, 3 (2016), 199–222.
- [78] Daniela Teodorescu, Tiffany Cheng, Alona Fyshe, and Saif M Mohammad. 2023. Language and mental health: Measures of emotion dynamics from text as linguistic biosocial markers. *arXiv preprint arXiv:2310.17369* (2023).
- [79] Elsbeth Turcan and Kathleen McKeown. 2019. Dreddit: A reddit dataset for stress analysis in social media. *arXiv preprint arXiv:1911.00133* (2019).
- [80] Jin Wang, Liang-Chih Yu, K Robert Lai, and Xuejie Zhang. 2016. Dimensional sentiment analysis using a regional CNN-LSTM model. In *Proceedings of the 54th annual meeting of the association for computational linguistics (volume 2: Short papers)*. 225–230.
- [81] Weichen Wang, Weizhe Xu, Ayesha Chander, Subigya Nepal, Benjamin Buck, Serguei Pakhomov, Trevor Cohen, Dror Ben-Zeev, and Andrew Campbell. 2023. The Power of Speech in the Wild: Discriminative Power of Daily Voice Diaries in Understanding Auditory Verbal Hallucinations Using Deep Learning. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 7, 3 (2023), 1–29.
- [82] Zhiyuan Wang, Maria A Larrazabal, Mark Rucker, Emma R Toner, Katharine E Daniel, Shashwat Kumar, Mehdi Boukhechba, Bethany A Teachman, and Laura E Barnes. 2023. Detecting social contexts from mobile sensing indicators in virtual interactions with socially anxious individuals.

- Proceedings of the ACM on interactive, mobile, wearable and ubiquitous technologies* 7, 3 (2023), 1–26.
- [83] Zhiyuan Wang, Varun Reddy, Karen Ingersoll, Tabor Flickinger, and Laura E Barnes. 2024. Rapport Matters: Enhancing HIV mHealth Communication through Linguistic Analysis and Large Language Models. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*. 1–8.
 - [84] Zhiyuan Wang, Fangxu Yuan, Virginia LeBaron, Tabor Flickinger, and Laura E Barnes. 2024. PALLM: Evaluating and Enhancing PALLiative Care Conversations with Large Language Models. *ACM Transactions on Computing for Healthcare* (2024).
 - [85] Ruolan Wu, Chun Yu, Xiaole Pan, Yujia Liu, Ningning Zhang, Yue Fu, Yuhan Wang, Zhi Zheng, Li Chen, Qiaolei Jiang, et al. 2024. MindShift: Leveraging Large Language Models for Mental-States-Based Problematic Smartphone Use Intervention. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. 1–24.
 - [86] Haoyi Xiong, Zhiyuan Wang, Xuhong Li, Jiang Bian, Zeke Xie, Shahid Mumtaz, Anwer Al-Dulaimi, and Laura E Barnes. 2024. Converging paradigms: The synergy of symbolic and connectionist ai in llm-empowered autonomous agents. *arXiv preprint arXiv:2407.08516* (2024).
 - [87] Xuhai Xu, Bingsheng Yao, Yuanzhe Dong, Saadia Gabriel, Hong Yu, James Hendler, Marzyeh Ghassemi, Anind K Dey, and Dakuo Wang. 2024. Mental-llm: Leveraging large language models for mental health prediction via online text data. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 8, 1 (2024), 1–32.
 - [88] Kailai Yang, Shaoxiong Ji, Tianlin Zhang, Qianqian Xie, Ziyan Kuang, and Sophia Ananiadou. 2023. Towards interpretable mental health analysis with large language models. *arXiv preprint arXiv:2304.03347* (2023).
 - [89] Xi Yang, Aokun Chen, Nima PourNejatian, Hoo Chang Shin, Kaleb E Smith, Christopher Parisien, Colin Compas, Cheryl Martin, Anthony B Costa, Mona G Flores, et al. 2022. A large language model for electronic health records. *NPJ digital medicine* 5, 1 (2022), 194.
 - [90] Zhilin Yang. 2019. XLNet: Generalized Autoregressive Pretraining for Language Understanding. *arXiv preprint arXiv:1906.08237* (2019).
 - [91] Tsungcheng Yao, Ernest Foo, and Sebastian Binnewies. 2024. Personalised Abusive Language Detection Using LLMs and Retrieval-Augmented Generation. In *Proceedings of the 7th International Conference on Natural Language and Speech Processing (ICNLSP 2024)*. 92–98.
 - [92] James Zabora, Karlynn BrintzenhofeSzoc, Barbara Curbow, Craig Hooker, and Steven Piantadosi. 2001. The prevalence of psychological distress by cancer site. *Psycho-oncology: journal of the psychological, social and behavioral dimensions of Cancer* 10, 1 (2001), 19–28.
 - [93] Brad J Zebrack. 2000. Cancer survivor identity and quality of life. *Cancer practice* 8, 5 (2000), 238–242.
 - [94] Zhihong Zeng, Yuxiao Hu, Yun Fu, Thomas S Huang, Glenn I Roisman, and Zhen Wen. 2006. Audio-visual emotion recognition in adult attachment interview. In *Proceedings of the 8th international conference on Multimodal interfaces*. 139–145.
 - [95] Penghao Zhao, Hailin Zhang, Qinhan Yu, Zhengren Wang, Yunteng Geng, Fangcheng Fu, Ling Yang, Wentao Zhang, and Bin Cui. 2024. Retrieval-augmented generation for ai-generated content: A survey. *arXiv preprint arXiv:2402.19473* (2024).

A Appendix

A.1 Prompt Design

Our prompt design consists of five key components, as shown in Figure 12, which extends the simplified demonstration in Figure 7 with comprehensive prompt language details:

- **Introduction:** Provides essential context about the cancer survivor being analyzed:
 - Trait information extracted from participant profiles (e.g., age, gender, cancer type and stage)
 - Data collection context (mobile EMA-based diary entries)
 - Current diary entry timestamp and content
- **Metrics Description:** Defines the target metrics to predict:
 - PANAS scores for positive and negative affect (0-30 scale)
 - Emotion regulation desire score (0-10 scale) and its binary indicator derived from individual mean comparison
 - Intervention availability (binary)
 - Binary states for specific emotions and behaviors derived from individual mean comparisons (binary)
- **[Optional] Personal Trajectory:** Includes historical diary entries from the previous day up to the current moment:
 - Retrieves timestamps and contents of recent diary entries within the time range identified (e.g., since last or current day before the current diary)
 - Ordered chronologically to show temporal progression
- **[Optional] Similar Peer Diaries:** Presents comparable cases retrieved through RAG:
 - Top-K semantically similar diary entries selected using FAISS vector similarity search, serving as few-shot learning examples
 - Each case includes the complete diary text and all associated ground truth
- **Structure Restriction:** Enforces strict JSON formatting for predictions:
 - Requires probability scores in [0.0, 1.0] range for all binary states, with 0.5 specified as uncertain.
 - Mandates `< PREDICTIONS >` tags to ensure reliable parsing and post-processing

This structured prompt design ensures consistent input format while providing comprehensive context for emotional state prediction. The

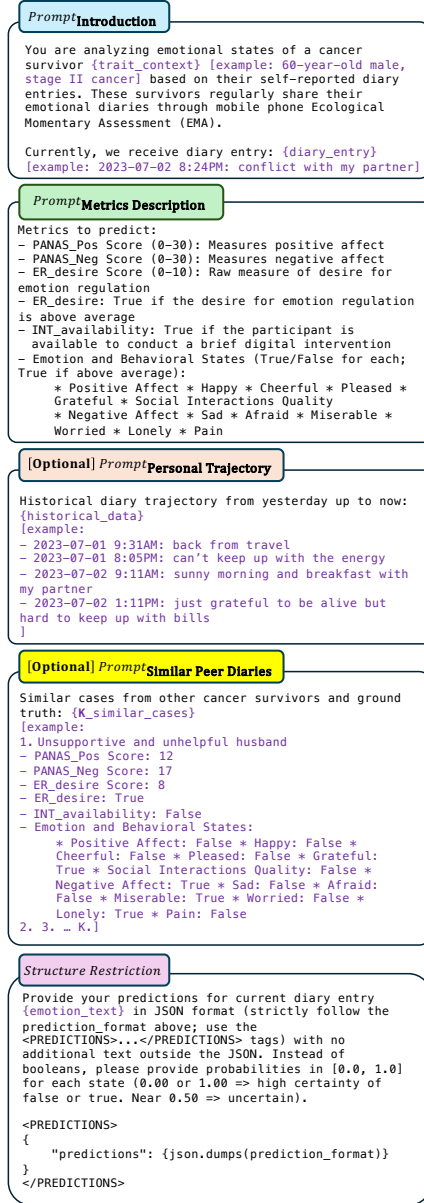


Fig. 12. Detailed CALLM prompt structure.

modular design allows for flexible configuration of optional components based on data availability and specific prediction needs.

B Baseline Hyperparameter Tuning

We conducted extensive hyperparameter tuning for all baseline models to ensure a fair comparison with the CALLM framework. The tuning process involved nested cross-validation to avoid data leakage, with parameters optimized for the multitask learning objective across all target variables. In our nested scheme, the outer loop consisted of the 5-fold grouped cross-validation (stratified by participant), providing unbiased performance estimates, while the inner loop performed hyperparameter selection using 5-fold cross-validation on the training portion only. This nested approach ensured that hyperparameter selection had no access to the test data, maintaining the integrity of our reported performance metrics.

For traditional machine learning approaches, we used Logistic Regression classifiers with various text representation methods. For Bag-of-Words (BoW) and TF-IDF feature extraction, we explored a comprehensive parameter space. Vectorization parameters included various document frequency thresholds (min_df values of 1, 2, and 5; max_df values of 0.9, 0.95, and 1.0) and n-gram ranges (unigrams only, and unigrams with bigrams). For TF-IDF specifically, we additionally tuned normalization methods (L1 and L2) and tested both with and without inverse document frequency weighting. The Logistic Regression classifier hyperparameters encompassed regularization strengths (C values of 0.1, 1.0, and 10.0), penalty types (L1 and L2), solver method (liblinear), and class weighting approaches (balanced versus none). This multi-output configuration of Logistic Regression classifiers enabled simultaneous prediction of all emotional and behavioral states.

For transformer-based models (BERT, RoBERTa, EmoBERT, etc.), we tuned several critical parameters affecting both computational efficiency and model performance. These included batch sizes (8, 16, and 32), learning rates (1e-5, 5e-5, 1e-4, and 5e-4, 1e-3), maximum sequence lengths (128 and 256), and training epochs (2, 3, 4, and 6). To mitigate overfitting, we implemented multiple regularization strategies: classifier dropout rates (0.1, 0.2, and 0.3), weight decay (0.01 and 0.1), and architecture-specific dropouts for attention mechanisms and hidden states (0.1 and 0.2). We also employed early stopping with a patience value of 2, halting training when validation performance ceased to improve, which both prevented overfitting and reduced computational costs. Due to computational constraints, we employed a two-stage tuning approach for transformer models. First, we conducted a coarse grid search on a subset of hyperparameter combinations using groups 1-3 for training and groups 4-5 for validation. Then, we fine-tuned the most promising configurations on the full dataset. For memory-intensive models like SentenceBERT and Emotion-Transformer, we used smaller batch sizes (8, 16) and implemented chunked processing to manage memory efficiently.

All tuning was performed by optimizing the mean balanced accuracy across all four target variables (positive affect, negative affect, emotion regulation desire, and intervention availability). The best parameters for each model were selected based on validation performance and then applied consistently across all five stratified groups of the final evaluation.