

A Framework for ECA-Based Psychotherapy

Santosh V. Patapati
Dept. of HCI
Cyrion Labs
santosh@cyrionlabs.org

Trisanth Srinivasan
Dept. of HCI
Cyrion Labs
trisanth@cyrionlabs.org

Himani Musku
School of Computer Science
Carnegie Mellon University
hmusku@andrew.cmu.edu

Amith Adiraju
Dept. of HCI
Cyrion Labs
aadiraju@cyrionlabs.org

Abstract—We present Tessa, an autonomous framework for structured bi-weekly psychotherapy intervention. Tessa alternates between screening sessions that use motivational interviewing and follow-up CBT sessions to help users identify, challenge, and reframe negative thoughts. Using text analysis with Large Language Models (LLMs), an Embodied Conversational Agent (ECA) adapts its dialogue and synchronizes non-verbal cues in real time. The framework incorporates a therapist-matching pipeline and clinical report generation to enable seamless hand-off to professional care. Taken together, Tessa provides accessible, stigma-free, and data-driven mental health support that improves personalized care and reduces barriers to therapy.

Index Terms—Human-Computer Interaction, Psychotherapy

I. INTRODUCTION

Many with mental health challenges are underserved due to limited access and social stigma. In response, digital therapy solutions are emerging as scalable alternatives. ECAs are of particular interest due to their ability to simulate human empathy. They mitigate stigma’s impact by providing a non-judgmental, accessible interface. However, existing systems struggle to deliver longitudinal support, and past autonomous ECAs have not combined context-aware speech with coordinated non-verbal cues, limiting empathetic effectiveness.

We introduce Tessa, a therapeutic system that delivers structured bi-weekly psychotherapy. Tessa uses a dual-session cycle: an initial screening with Motivational Interviewing (MI) and a follow-up Cognitive Behavioral Therapy (CBT) session to identify, challenge, and reframe negative thoughts. Leveraging LLMs for real-time text analysis, Tessa adapts its dialogue and monitors 25 mental health dimensions to improve well-being. Its autonomous ECA delivers empathetic responses with synchronized gestures. Tessa also offers a streamlined pathway to professional care through therapist matching and clinical report generation. Our demonstration showcases Tessa’s scalable approach to long-term mental health support.

II. METHODOLOGY

A. Bi-Weekly Interaction Structure

Tessa follows a bi-weekly cycle of Motivational Interviewing (MI) followed by a Cognitive Behavioral Therapy (CBT) session. MI and CBT complement each other well [1].

Before the initial MI screening session, the user completes a 25-item questionnaire based on the Diagnostic and Statistical Manual of Mental Disorders-5 (DSM-5), the clinical standard that outlines diagnostic criteria for mental disorders. The questionnaire combines the PHQ-8, GAD-7, and PDS [2]

to cover a wide range of mood disorders and symptoms. Questionnaire responses are passed as context to Tessa during future sessions via its running state.

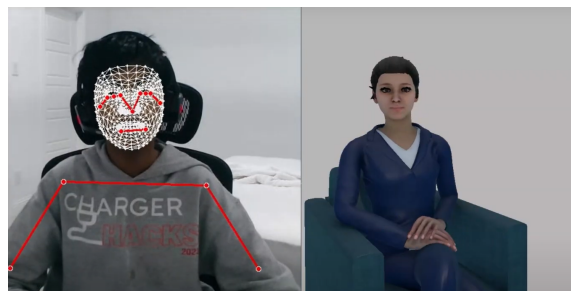


Fig. 1. Sample frame from a Tessa session. The ECA maintains a neutral listening position as the user speaks. Visual features are tracked in real-time.

In the first session (Screening), the ECA conducts a semi-structured interview using MI (Figure 2) to assess the user’s mental state and boost motivation for change. To motivate users, the ECA must build rapport, and it does this by engaging in small talk and providing a transparent introduction.

Audio and video recordings, the session transcript, and prior questionnaire responses are input into a quad-modal transformer model [3] that predicts the user’s well-being across 25 mental health dimensions, such as ”Restless or Unsatisfying Sleep” and ”Significant Weight Changes.”

In the second bi-weekly session (Intervention), held a few days after screening, the ECA and user conduct a CBT session as a follow-up to MI. This session targets dimensions flagged during screening or by the ML model. CBT follows three steps: identify, challenge, and replace negative thoughts with constructive ones (Figure 3). The goal is for users to adopt more constructive thoughts in day-to-day life and thus face fewer mental health challenges. For this reason, we expect users’ well-being across the 25 dimensions to improve over time. By session end, the user formulates a balanced perspective and actionable coping strategy. Each session thus builds on the last, creating a continuous therapeutic loop.

B. Multimodal Interaction Pipeline

Modeled after GenECA [4], each time the user presses the speaking button, an ”Utterance Segment” is recorded. Audio and video signals are analyzed by lightweight analyzers to detect distress, while Automatic Speech Recognition (ASR) transcribes speech for text validation. Combined

