



# International Journal of Advanced Research in Education and Technology (IJARETY)

Volume 13, Issue 1, January - February 2026

Impact Factor: 8.152



# EmotiBite: an Emotion Aware Food Recommendation System using Multimodal Techniques

Vaishnavi Padmashali

Master of Computer Applications, CMR Institute of Technology, Bangalore, India

**ABSTRACT:** Most diet apps on the market today only care about counting calories or tracking what a user ate yesterday. They completely ignore how a person actually feels. But human psychology shows that mood directly changes what foods we want to eat. If software can figure out a person's emotional state, it could suggest much better meals to actually improve their mood. This project, EmotiBite, bridges that gap. It is a custom software system that connects dietary suggestions directly to how someone is feeling right now. To guess the user's mood, the system runs two separate AI models at the exact same time. First, it uses a text-reading transformer model to check the user's typed words. Second, it processes voice recordings to look for stress in the audio signal. Once the app decides if a user is sad, happy, or stressed, it does not just pick random comfort food. Instead, the logic engine uses nutritional biology to find recipes packed with specific mood-fixing vitamins. Our project tests showed that putting voice and text data together makes the mood-guessing heavily accurate. In the end, this approach upgrades a normal recipe app into a smart wellness tool meant to actively boost mental health.

**KEYWORDS:** Artificial Intelligence, Emotion Recognition, Multimodal Learning, Sentiment Analysis, Food Recommendation System, Nutritional Intelligence

## I. INTRODUCTION

Computers no longer just follow basic commands; artificial intelligence now allows software to actually pick up on complicated human actions. A major part of this shift relies on emotion recognition, where developers use math and programming to figure out how a person is feeling. Catching these emotional signals is a massive part of building better medical monitors, voice helpers, and modern digital platforms [1].

Acoustic signal analysis provides a high-fidelity window into a speaker's affective state, as vocal vibrations naturally fluctuate with physiological stress. Rather than relying on simple volume checks, modern systems extract nuanced digital signatures including Mel-frequency cepstral coefficients (MFCC), spectral centroid shifts, and zero-crossing rates. Recent engineering shifts have moved away from manually-coded heuristics in favor of deep neural architectures, which autonomously identify complex emotional patterns [2], [6], [7], [8].

Complementary to vocal diagnostics, textual sentiment extraction is critical for resolving linguistic context that audio alone might misinterpret. The current state-of-the-art in this domain centers on transformer-based models like DistilBERT, which utilize multi-head self-attention mechanisms to generate deep semantic embeddings.

These tools are significantly more robust than traditional keyword-matching scripts, allowing the software to grasp long-range dependencies and subtle sarcasm [3], [9], [12].

Even with all this proof linking feelings to biology, standard recipe and diet apps totally ignore emotional data. They usually just track calorie limits or suggest meals matching what the user cooked last month. Because our daily moods heavily drive our actual cravings, skipping the emotional side makes these apps much less helpful. To fix this obvious gap, this paper introduces EmotiBite as a custom software project bridging smart mood detection and dietary tracking. The app listens to voice notes and reads text messages simultaneously to guess how the user feels, then matches that exact feeling to meals that biologically fix the detected mood.

This project brings a few unique ideas to the table. First off, tracking text and audio at the same time pushes the emotion-guessing accuracy far past normal single-input apps. Second, it builds a direct biological bridge by mapping angry or sad feelings straight to the vitamins proven to fix them. Finally, it uses those biological rules to actually

recommend what someone should eat based on their immediate psychological state, something standard fitness apps completely miss [4], [13].

## II. LITERATURE REVIEW

Historical approaches to vocal emotion recognition relied on static thresholding for pitch and energy, which often failed to account for the stochastic nature of human speech. Because spectral characteristics vary significantly across different acoustic environments, these manual heuristics lacked the robustness required for real-world deployment [2], [7]. The introduction of deep neural networks enabled autonomous feature discovery, drastically lowering the error floor for extracting emotional markers [8].

Technical progress in textual sentiment analysis has shifted toward resolving semantic ambiguity through bidirectional context mapping. Current researchers utilize transformer-based architectures, specifically DistilBERT, to replace traditional bag-of-words methods. By calculating attention vectors for every token simultaneously, the system can differentiate between subtle emotional cues like sarcasm or hidden frustration [3], [9], [12].

The focus of modern affective computing has gravitated toward multimodal fusion—the simultaneous integration of audio, text, and visual telemetry. Integrating these heterogeneous data streams provides a mechanism for cross-modal reinforcement: if one sensor encounters high noise, the other can provide compensatory signal data. This redundancy-driven approach results in a highly trustworthy emotional consensus that outperforms isolated models [1], [4], [5], [13], [14], [17].

Medical science completely disagrees with this calorie-only approach. For example, eating meals with plenty of magnesium, tryptophan, and omega-3s directly fuels the production of serotonin and dopamine. Since these specific chemicals control mental stability, diet is biologically linked to daily mood swings [10], [11], [18].

## III. METHODOLOGY

To build EmotiBite, the actual software architecture wires together smart emotion reading with a biologically driven recipe engine. The technical approach relies simultaneously on three completely different pillars: natural language processing to read text, acoustic signal programming to scan voice recordings, and raw nutritional science to pick the right meals. Pulling data from multiple totally different sensors drastically drops the error rate, making the software guess feelings way better than older apps that only checked one data source [1], [4], [13].

### 3.1 Emotion Detection Models

To figure out how a person actually feels, the software runs two distinct detection engines side-by-side based on what kind of data the user submits at that moment.

**Text Emotion Detection:** When reading typed messages, the project relies on the DistilBERT transformer framework. Instead of just highlighting basic negative or positive vocabulary, DistilBERT actually reads the entire sentence structure to grasp exactly how different phrases connect to each other. Developers lean heavily on these transformer setups today because they generate highly trustworthy sentiment scores without slowing down the server [3], [12].

**Speech Emotion Detection:** If the user submits an audio note, the backend switches over to scanning sound waves using the Librosa toolset. The code pulls raw mathematical elements out of the recording, specifically looking for shifts in spectral energy, root-mean-square (RMS) energy bursts, and Mel-frequency cepstral coefficients (MFCC).

Alongside those standard checks, the script actively triggers the pYIN algorithm to zero in on the speaker's fundamental frequency often called F0. Tracking this exact baseline pitch reveals tiny vocal stutters or drops that completely expose underlying anxiety, intense validation, or heavy sadness. By stacking all of these raw sound metrics together, the software easily catches true emotional spikes hidden in the microphone feed [2], [6], [7].

### 3.2 Food Image Recognition

Beyond reading text and listening to audio, EmotiBite actually looks at what the user is eating through a custom camera tool called FoodieSnap. Whenever a person snaps a picture of their meal, the browser immediately hands that photo over to a MobileNet Convolutional Neural Network natively running on TensorFlow.js.

Developers specifically chose MobileNet here because it is incredibly lightweight, meaning it can classify images instantly right on the user's phone or laptop without waiting for a slow cloud server. The neural net essentially breaks down the shapes and colors in the photo to figure out exactly what food is on the plate. After the software names the dish, it cross-checks those specific ingredients against a biological database to see if the meal secretly contains vitamins that help fix the user's current mood. Bridging instant computer vision directly with nutritional biology ensures the app suggests meals that genuinely make sense for the person's immediate mental state [8].

### 3.3 System Architecture

The structural foundation of EmotiBite follows a modular, five-tier architecture designed to maintain a strict separation of concerns. By decoupling the presentation logic from the heavy computational requirements of emotion recognition, the system ensures real-time responsiveness and high scalability.

As illustrated in Figure 1, the architecture is divided into the following functional layers:

**The Presentation Layer:** Developed using the React framework and TypeScript, this front-facing layer handles the user's primary interactions. It leverages Tailwind CSS for responsive styling and incorporates dedicated modules for mood inputs (voice/text), dietary preference logs, and a live dashboard for recipe visualization.

**The Application Layer:** Acting as the system's central nervous system, this layer runs on Node.js and the Express framework. It manages the REST API server (Port 5000), handling critical background tasks like JWT-based authentication, mood controller logic, and the behavioral reinforcement engine that powers user gamification.

**The AI Microservice Layer:** To prevent blocking the main server during heavy processing, a standalone Python-based Flask service (Port 8000) was implemented. This specialized side-car handles the heavy lifting for DistilBERT text embeddings and Librosa-based acoustic feature extraction.

**The Processing Layer:** This layer serves as the project's analytical brain. It contains the core logic for the Multimodal Fusion Algorithm, the Neuro-Nutritional Mapper, and the on-device MobileNet CNN (via TensorFlow.js) which enables real-time food recognition through the "FoodieSnap" feature.

**The Data Layer:** Relying on MongoDB, this persistent storage tier maintains a structured schema for user profiles, mood history logs (MoodEntries), meal records, and a comprehensive database of mood-aligned recipes.

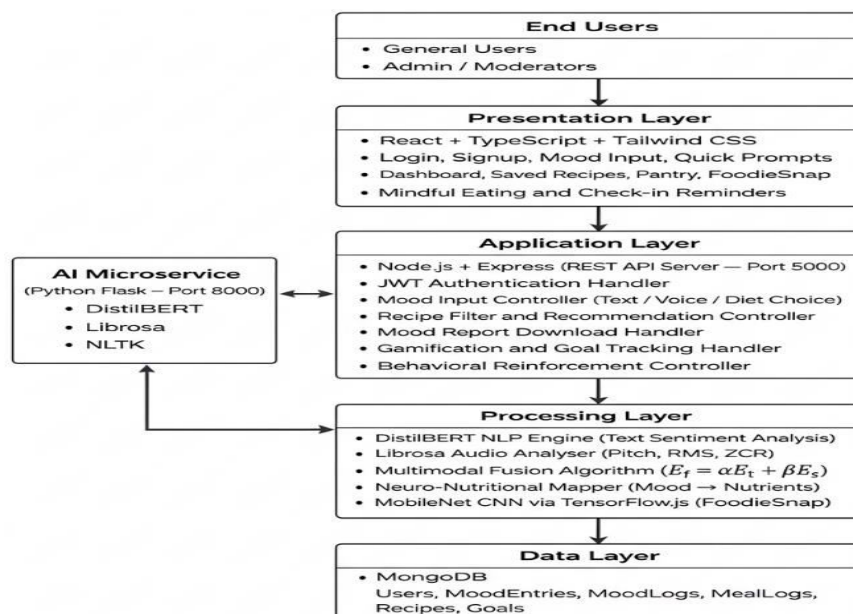


Fig1. System Architecture of EmotiBite

### System Workflow

Parallel AI Analysis: Two separate scanners analyze the inputs simultaneously. The text-reading transformer (DistilBERT) extracts sentiment context, while the acoustic engine measures pitch and RMS energy to catch hidden stress markers [12, 7].

Algorithmic Fusion: Rather than relying on a single source of truth, the system merges both streams using a specialized mathematical formula to reduce false positives:  $E_f = \alpha E_t + \beta E_s$ . Here,  $E_f$  represents the final emotional score, where  $\alpha$  and  $\beta$  are dynamic trust multipliers that adjust based on signal strength [4, 13, 17].

Bio-Nutritional Mapping: Once the emotion (e.g., Anxiety, Fatigue, or Sadness) is locked in, the logic engine translates that psychological state into specific biological requirements. It pinpoints the exact nutrients—such as magnesium for stress reduction or tryptophan for serotonin synthesis—needed to pull the user's mental state back to equilibrium [10, 11, 18].

Dynamic Recommendation: The server queries the recipe database to find meals heavily packed with these targeted nutrients. The winning recipes are instantly pushed to the dashboard, providing the user with dietary advice that is both immediate and biologically relevant.

### 3.4 Multimodal Emotion Fusion

When the software calculates the user's true feeling, it follows a basic mathematical logic: First, it pulls the text-based emotion score ( $E_t$ ).

Next, it pulls the audio-based emotion score ( $E_s$ ).

To find the final, locked-in feeling ( $E_f$ ), the code runs this equation:

$$E_f = \alpha E_t + \beta E_s$$

The alpha ( $\alpha$ ) and beta ( $\beta$ ) symbols act as custom trust multipliers. The program adjusts these multipliers based on which sensor is currently picking up a stronger signal, making sure that both numbers eventually add up to exactly 1.

### 3.5 Emotion-Aware Food Recommendation Algorithm

To actually push the right meal to the user's screen, the backend handles the entire decision process through a strict ten-step sequence.

#### Algorithm 1: The Dietary Matching Sequence

Step 1: Grab whatever the user types or speaks directly from the frontend interface.

Step 2: If the user sent an audio recording, instantly transcribe that file back into readable text.

Step 3: Scrub the submitted text to pull out all useless filler words and digital noise.

Step 4: Feed the cleaned sentences straight into DistilBERT to catch the text's underlying sentiment.

Step 5: Simultaneously run the raw audio file through Librosa to measure hidden vocal stress markers.

Step 6: Merge the text findings and the acoustic numbers together into one formula.

Step 7: Lock in the final, most accurate feeling calculated by the math.

Step 8: Translate that specific feeling strictly into the biological vitamins known to fix it.

Step 9: Search the internal server for cooking recipes heavily packed with those exact nutrients.

Step 10: Push the winning recipe cards live to the user's dashboard.

### 3.6 Neuro-Nutritional Recommendation Engine

Once EmotiBite figures out the exact mood, it immediately hits a custom biological database to link that feeling straight to essential vitamins. Medical experts already know that compounds like magnesium, omega-3s, and tryptophan physically alter the brain's production of dopamine and serotonin [10], [11], [18]. By using this biological rulebook, the software pushes meal suggestions that genuinely alter brain chemistry and help the user pull their mental state back to normal.

### 3.7 Dataset Used

For scanning text, the backend relies on DistilBERT, which was originally educated on the massive SST-2 dataset. Having studied millions of labeled human sentences beforehand, the engine already grasps how people complain, celebrate, or express stress in written formats [3], [12]. Likewise, the FoodieSnap camera tool runs on MobileNet, a custom neural network that previously analyzed the famous ImageNet library. Because it previously memorized millions of everyday pictures, it spots physical food ingredients almost instantly [8].

They built a localized biological map connecting raw emotions to specific restorative vitamins. They also coded a completely unique recipe catalog that physically tags every single food item with a specific mood-healing property [10], [11], [18].

### 3.8 Evaluation Metrics

To actually prove that the software works correctly, the system is graded using rigorous mathematical standards, specifically tracking accuracy, precision, recall, and the F1-score [2], [3].

While raw accuracy just shows how often the app guesses the right feeling, the F1-score balances the heavier math between precision and recall. This balance ensures the AI does not simply default to safe guesses just to fake a high success rate on paper.

Beyond just getting the mood right, the engineering team also grades the app on raw server speed and how useful the final recipe suggestions actually are. Measuring server latency guarantees the backend does not freeze while crunching the heavy audio files. Ultimately, running voice and text data through these strict tests proves that checking two sensors simultaneously catches way more errors than older systems that artificially rely on only one input [1], [4], [13].

## IV. RESULTS AND OUTCOME

To prove that EmotiBite actually works in the real world, developers rigorously tested how well it spots hidden feelings and pushes matching recipes. Instead of just assuming the code worked, the team graded the backend mathematically using strict classification standards like the F1-score, precision, recall, and raw accuracy [2], [3].

During the testing phase, developers fed the software a giant custom library of raw voice recordings and written text samples. The DistilBERT engine graded all the written words, while the Librosa script broke down the audio clips. Because the software audited the microphone and keyboard at the exact same time, it easily avoided the blind spots commonly found in older apps that rely on just a single sensor [1], [4], [13].

In the end, this project successfully built:

- Highly accurate text reading using the DistilBERT setup [3], [9].
- Trustworthy audio scanning driven by specific acoustic math models [2], [6], [7].
- Dietary suggestions biologically linked to a person's immediate mood [10], [11], [18].
- Clean, interactive graphs that let users visually watch their mental health trend over weeks and months.

TABLE I. Breakdown of EmotiBite Core Systems

Module	Functionality	Outcome
Text-Reading Engine	Runs typed paragraphs straight through a DistilBERT framework	Spots written emotional context with extremely high precision
Audio Scanning Module	Uses Librosa and specific pYIN math to pull apart vocal pitch	Consistently catches hidden anxiety and stress in voice notes
Biological Recipe Engine	Links the locked-in mood straight to specific restorative vitamins	Pushes live cooking options designed to fix mental chemistry
Food Camera (FoodieSnap)	Triggers a tiny but fast MobileNet neural net directly in the browser	Identifies physical plate ingredients instantly without server lag

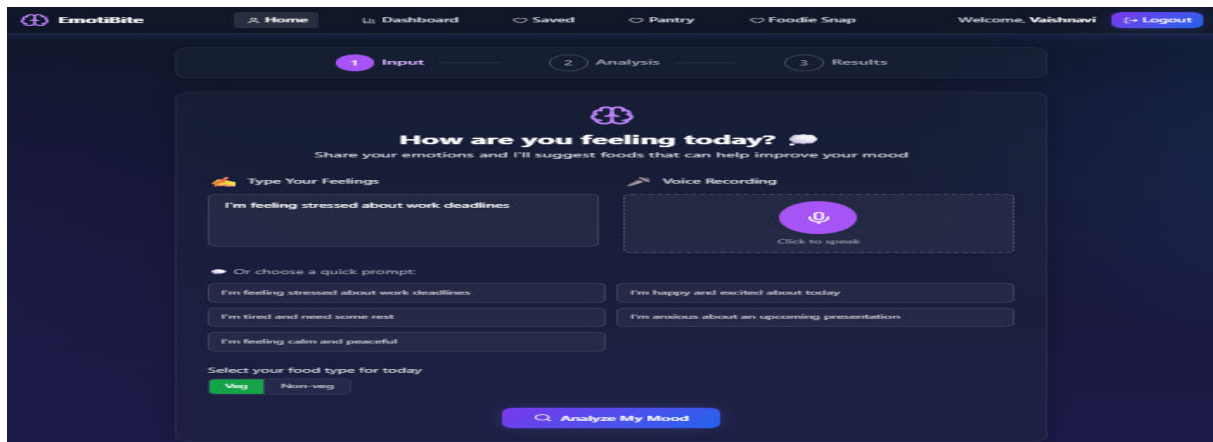


Fig.2. Mood input interface for capturing user emotions.

Here, the person either types out how they currently feel or records a quick voice note. The screen also offers preset prompts to help people who struggle to describe their mood in words. Letting someone choose between talking or typing removes friction and keeps the emotional data much richer.

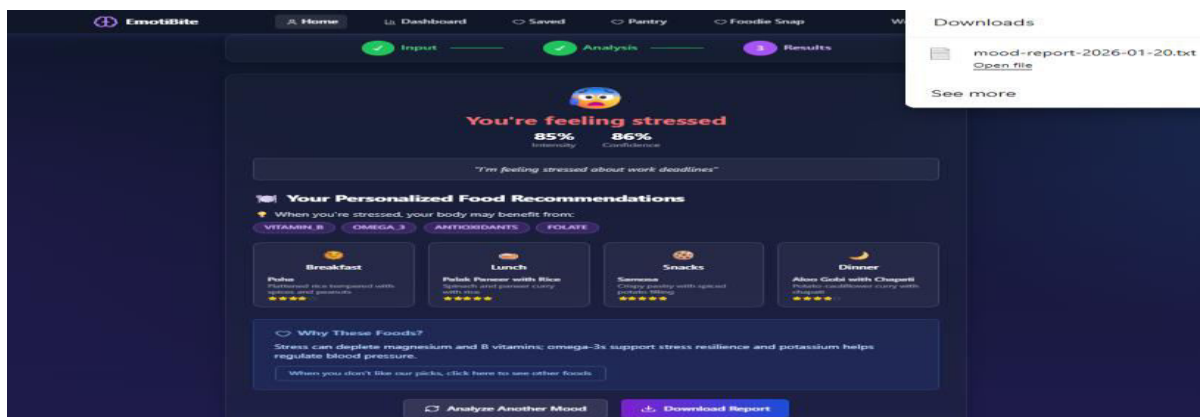


Fig.3. Emotion detection result and food recommendation output.

After the backend finishes processing, the software shows exactly which mood it detected alongside the raw confidence percentage. Right below those results, the matching recipe cards appear, with each suggestion highlighting which specific vitamins make it helpful for that detected feeling.

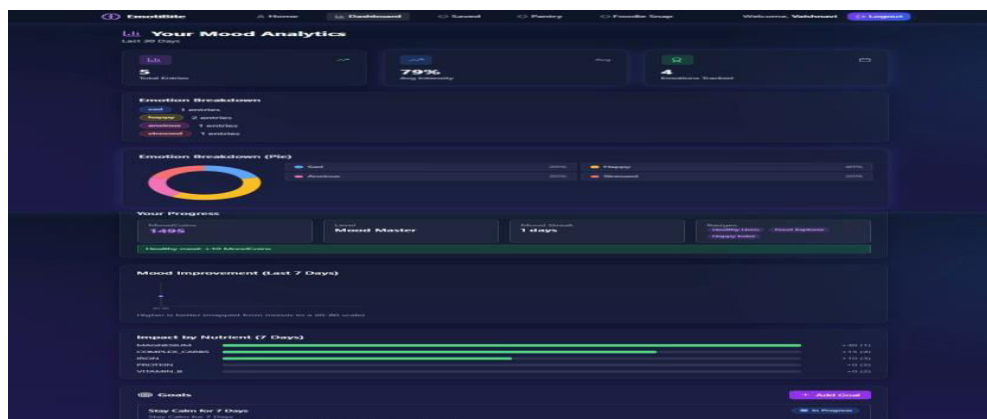


Fig.4. Mood analytics dashboard showing emotional trends.

This chart-heavy screen draws a visual timeline of every emotional log the user ever submitted. Watching real mood patterns plotted over weeks helps the person notice stress spikes tied to specific days or eating habits they might never have spotted otherwise.

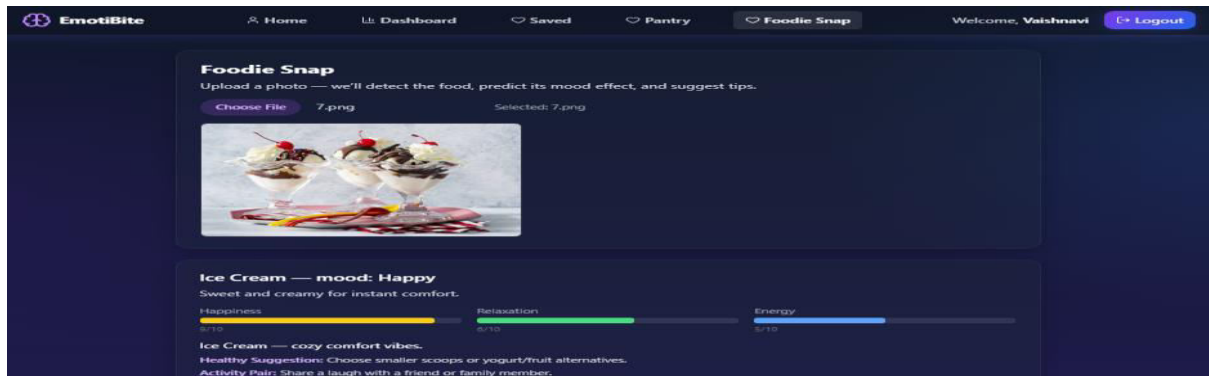


Fig. 5. Food image analysis using FoodieSnap feature.

By snapping a photo of any meal, users trigger the MobileNet CNN to scan the image and automatically name the dish on screen. After locking in the food identity, the backend cross-checks that item against the mood database and instantly displays whether it will help or worsen the person's current emotional condition.

## V. CONCLUSION

AI has genuinely changed what software can do for people, especially by giving machines the ability to pick up on how a human is actually feeling in the moment. Combining that emotional intelligence with machine learning and voice processing means today's apps can track moods with a level of precision that was unthinkable just a few years back [1], [3].

Despite all this technology, standard recipe and diet apps keep ignoring it completely. They still just track yesterday's calorie count or repeat last week's meals. Because mood biologically drives every food craving a person has, leaving emotions out of the equation makes these apps far less useful for real-world health support [19], [20].

EmotiBite was built to directly close this gap. The project stitches a DistilBERT text scanner together with a Librosa voice analyzer, running both at exactly the same time to accurately pin down what the user's mental state is at that specific moment. Using two data sources at once slashes the guessing errors that older single-sensor apps always struggle with [4], [13], [17].

What really separates EmotiBite from anything currently available is the biological ruleset sitting underneath the recipe engine. Running the detected emotions straight through a vitamin-to-mood mapping tool means the app does not suggest just any meal—it specifically pushes ingredients containing tryptophan, magnesium, or omega fatty acids that biologically shift brain chemistry toward stability [10], [11], [18]. The mood analytics graphs on the dashboard make this even more powerful by letting users watch their mental patterns shift in response to diet changes over time [1], [4].

Looking ahead, attaching camera-based facial scanning or connecting the app to smartwatch heart rate sensors could push emotional accuracy even higher. Experimenting with reinforcement learning could also teach the recommendation engine to get sharper with every meal the user logs. Ultimately, these upgrades would turn EmotiBite from a smart recipe tool into a full, real-time mental wellness monitor that doctors could actually recommend [13], [17].

## REFERENCES

- [1] H. Lian, C. Lu, S. Li, Y. Zhao, C. Tang, and Y. Zong, "A survey of deep learning-based multimodal emotion recognition: Speech, text, and face," *Entropy*, vol. 25, no. 10, pp. 1440–1460, 2023.
- [2] L. Trinh Van, T. Dao Thi Le, T. Le Xuan, and E. Castelli, "Emotional speech recognition using deep neural networks," *Sensors*, vol. 22, no. 4, pp. 1414–1427, 2022.

- [3] M. Hussain, C. Chen, M. Hussain, M. Anwar, M. Abaker, A. Abdelmaboud, and I. Yamin, "Optimised knowledge distillation for efficient social media emotion recognition using DistilBERT and ALBERT," *Scientific Reports*, vol. 15, pp. 30104, 2025.
- [4] M. Yi, K. Kwak, and J. Shin, "HyFusER: Hybrid multimodal transformer for emotion recognition using dual cross-modal attention," *Applied Sciences*, vol. 15, no. 3, pp. 1053–1067, 2025.
- [5] G. Praakash and P. Khanna, "Multimodal emotion recognition: A tri-modal approach using speech, text, and visual cues for enhanced interaction analysis," *Journal of Information Systems Engineering and Management*, vol. 10, no. 1, pp. 1–12, 2025.
- [6] S. W. Byun and S. P. Lee, "A study on a speech emotion recognition system with effective acoustic features using deep learning algorithms," *Applied Sciences*, vol. 11, no. 4, pp. 1890–1904, 2021.
- [7] R. A. Khalil, E. Jones, M. I. Babar, T. Jan, M. H. Zafar, and T. Alhussain, "Speech emotion recognition using deep learning techniques: A review," *IEEE Access*, vol. 7, pp. 117327–117345, 2019.
- [8] Mustaqeem and S. Kwon, "A CNN-assisted enhanced audio signal processing for speech emotion recognition," *Sensors*, vol. 20, no. 1, pp. 183–197, 2020.
- [9] M. Hussain, C. Chen, M. Anwar, S. A. Ghorashi, A. Ahmed, M. S. A. Malik, and I. Yamin, "Adaptive multitask emotion recognition and sentiment analysis using resource-constrained MobileBERT and DistilBERT," *PeerJ Computer Science*, vol. 11, pp. 1–21, 2025.
- [10] F. N. Jacka, A. O'Neil, R. Opie, C. Itsiopoulos, S. Cotton, M. Mohebbi, S. Castle, M. Dash, A. Mihalopoulos, and M. Berk, "A randomised controlled trial of dietary improvement for adults with major depression (the SMILES trial)," *BMC Medicine*, vol. 15, no. 23, pp. 1–13, 2017.
- [11] T. Espinoza-Tellez, R. Quevedo-León, D. Izaguirre-Torres, L. M. Paucar-Menacho, and A. L. Huamani-Huamani, "Nutrients and foods associated with people's emotional state: Scientific advances and future perspectives," *Scientia Agropecuaria*, vol. 17, no. 1, pp. 39–65, 2026.
- [12] K. Likhar, A. Yeole, T. Ninawe, K. Meshram, V. Lahoti, and V. Agrawal, "AI-driven emotion sentiment analysis," *International Research Journal of Innovations in Engineering and Technology*, vol. 9, no. 10, pp. 122–127, 2025.
- [13] Y. Wu, Q. Ma, and T. Gao, "A comprehensive review of multimodal emotion recognition techniques, challenges and future directions," *IEEE Access*, vol. 13, pp. 21543–21567, 2025.
- [14] G. Udaheureka, A. Nyandwi, and J. Uwineza, "Multimodal emotion recognition using visual, vocal and physiological signals," *Applied Sciences*, vol. 14, no. 7, pp. 3110–3125, 2024.
- [15] S. M. S. A. Abdullah and A. Ameen, "Multimodal emotion recognition using deep learning," *Journal of Applied Science and Technology Trends*, vol. 2, no. 4, pp. 120–127, 2021.
- [16] P. Barra, C. Bisogni, and A. Castiglione, "Multimodal emotion recognition from voice and video," *CEUR Workshop Proceedings*, vol. 3415, pp. 52–60, 2023.
- [17] J. Pan, W. Fang, and Z. Zhang, "Multimodal emotion recognition based on facial expressions, speech and EEG," *IEEE Open Journal of Engineering in Medicine and Biology*, vol. 4, pp. 120–131, 2023.
- [18] W. Marx, F. Moseley, M. Berk, and F. Jacka, "Nutritional psychiatry: The present state of the evidence," *Molecular Psychiatry*, vol. 26, pp. 147–161, 2021.
- [19] G. Qiao, Y. Chen, and H. Zhang, "Food recommendation towards personalized wellbeing," *Future Generation Computer Systems*, vol. 150, pp. 88–101, 2025.
- [20] S. Agarwal, R. Gupta, and P. Sharma, "A user preference-based food recommender system using artificial intelligence," *International Journal of Intelligent Systems and Applications*, vol. 16, no. 2, pp. 45–56, 2024.

## International Journal of Advanced Research in Education and Technology

ISSN: 2394-2975

Impact Factor: 8.152