

Calculating people's attraction based on their facial expressions, body movement, voice tone and words spoken

Mitreia Ioan David

I. ABSTRACT

Abstract—*Romantic attraction arises from dynamic interactions between facial expression, body movement, voice, and spoken language, yet most computational approaches examine these cues in isolation or as static summaries [1]. This article addresses the problem of predicting romantic interest from short segments of naturalistic dyadic interactions using temporally aligned multimodal behavioral signals.*

We present a pipeline that extracts facial Action Units, hand and upper-body movement dynamics, acoustic prosodic features, and sentiment cues from video recordings, synchronized at the second level. These features are modeled as short behavioral sequences using a stacked Long Short-Term Memory (LSTM) network to capture temporal dependencies across modalities. Unlike conventional affect recognition tasks, our target is an explicitly interpersonal outcome: binary romantic attraction toward a specific interaction partner.

Our contributions include (i) a practical multimodal pre-processing framework that aligns visual, vocal, and linguistic signals via participant visibility, (ii) a compact sequence-based representation of expressive behavior, and (iii) a temporal neural model for attraction prediction. Empirical results show that multimodal sequence modeling substantially outperforms static and unimodal baselines, supporting the view that attraction is communicated through dynamic, multimodal behavior rather than isolated cues.

II. INTRODUCTION

Human romantic attraction plays a central role in social behavior, relationship formation, and well-being. It emerges rapidly during interaction and is shaped by a complex interplay of facial expressions, body movement, vocal characteristics, and spoken language [2, 3, 4]. Psychological research has established that attraction is not determined by any single cue, but rather by the integration of multiple behavioral signals across these modalities [3, 4]. Recent work further demonstrates that attraction is sensitive to dynamic temporal patterns, such as synchronized body movement during brief interactions [5]. Despite this understanding, computational approaches to modeling attraction have remained limited, often focusing on static features, unimodal cues, or population-level attractiveness ratings rather than dyadic, temporally grounded interaction data [6, 7].

Recent advances in affective computing and multimodal machine learning have made it possible to extract rich behavioral signals from audiovisual recordings, including facial expressions, gestures, prosody, and language. These methods have been successfully applied to tasks such as emotion

recognition, sentiment analysis, and social signal processing. However, predicting romantic attraction poses additional challenges beyond standard affect recognition. Attraction is inherently interpersonal, context-dependent, and only partially observable through behavior. Moreover, it unfolds over short time scales, where subtle temporal patterns and coordination between modalities may carry more information than isolated snapshots [5].

The problem addressed in this article is the following: *can short sequences of multimodal behavioral signals be used to reliably predict whether a person experiences romantic attraction toward a specific interaction partner?* Formally, given synchronized facial, bodily, vocal, and linguistic cues extracted from brief segments of dyadic interaction, we aim to classify binary romantic interest at the individual–partner level. Addressing this problem requires not only robust multimodal feature extraction, but also temporal modeling capable of capturing how expressive behaviors evolve and co-occur over time.

The importance of this problem is twofold. From a scientific perspective, computational models can complement traditional psychological methods based on self-report and manual coding [8]. Such models enable precise analysis of how dynamic multimodal behavioral cues—including facial expressions, body movement, voice, and language—jointly shape interpersonal attraction judgments [2, 3]. From a methodological standpoint, attraction prediction represents a challenging testbed for multimodal sequence modeling. Because attraction is shaped by multiple interacting factors that vary across individuals and contexts [2], the outcome is not perfectly determined by observable behavior alone, involving noisy, partially missing signals and strong individual differences. Progress in this area may therefore inform the design of more general models for social and affective inference.

III. DATA COLLECTION

The quality and reliability of affective computing research depend critically on rigorous and multimodal data collection pipelines capable of capturing complex human emotional expressions. Traditional laboratory acquisition frameworks prioritize controlled environments and high-fidelity measurements; however, they often suffer from limited ecological validity and reduced generalizability [7]. Recent research has therefore shifted toward hybrid approaches that integrate audiovisual, physiological, and behavioral modalities collected in more naturalistic settings, thereby enabling richer affective modeling

⁰ACM CCS: Computing methodologies → Machine learning; Human-centered computing → Affective computing. AMS MSC: 68T07, 68T10, 62M45.

[8, 7]. Following this methodological direction, our study implements a multimodal extraction pipeline designed to capture fine-grained facial, acoustic, and bodily behavioral cues from video recordings. The complete pipeline is summarized below.

A. Person Video Trimming and Facial Action Unit Extraction

To ensure that downstream feature extraction was performed only when the target participant was visible, we preprocessed the videos by temporally segmenting and retaining only frames in which the participant of interest appeared on screen. This yielded a video containing black frames during intervals when the participant was not visible. These videos functioned as the structural backbone for subsequent behavioral feature extraction.

From these preprocessed videos, we extracted frame-level facial Action Units (AUs) to quantify expressive behavior. Videos were sampled at 1 frame per second (every 30th frame at typical 30 fps playback), and each frame was classified into one of three categories: (1) black frames (mean pixel brightness < 15), indicating off-screen intervals; (2) frames without detectable faces, identified using Haar Cascade face detection with a minimum face size of 50×50 pixels; or (3) frames containing valid faces suitable for AU extraction.

For frames with valid faces, we performed emotion recognition using DeepFace and converted the resulting emotion probability distributions (happy, sad, angry, surprise, fear, disgust) into pseudo-Action Unit intensities based on established mappings from the Facial Action Coding System (FACS). For example, happiness was mapped primarily to AU06 (cheek raiser) and AU12 (lip corner puller), while anger was associated with AU04 (brow lowerer) and AU23 (lip tightener). The following 17 Action Units were estimated: AU01, AU02, AU04, AU05, AU06, AU07, AU09, AU10, AU12, AU14, AU15, AU17, AU20, AU23, AU25, AU26, and AU28. AU intensities were normalized to the range [0, 1].

Frames without faces or with black content were assigned neutral AU values (all zeros) to maintain temporal continuity. The final output consisted of time-indexed CSV files containing AU intensities at 1 Hz sampling rate, enabling alignment with other modalities for multimodal fusion.

B. Audio Feature Extraction

To analyze vocal affect and speaking patterns, we extracted frame-level acoustic features from each participant’s audio track using `librosa`. Audio files were loaded at a sampling rate of 22,050 Hz, and features were computed with a hop length of 512 samples (approximately 23 ms per frame).

The following features were extracted to capture prosodic, spectral, and temporal characteristics of speech:

- **Energy:** Root mean square (RMS) energy converted to decibels, reflecting vocal intensity and amplitude modulation.
- **Pitch (F0):** Fundamental frequency estimated using the probabilistic YIN (PYIN) algorithm [], with a valid range of C2 to C7 (approximately 65–2093 Hz). Pitch

confidence scores were retained to account for voicing probability.

- **Speaking rate proxy:** Onset strength (spectral flux) was extracted as a frame-level feature reflecting the rate of acoustic change. While true speaking rate requires phonetic transcription, onset strength provides a continuous signal-processing approximation that captures temporal articulation dynamics without requiring speech-to-text alignment.

All features were temporally aligned to ensure consistent frame counts, and frame timestamps were computed using `librosa.frames_to_time`. The resulting feature set was exported as a time-indexed CSV file for each pair.

C. Body, Hand, and Upper Body Movement Features

To capture nonverbal expressive behaviors beyond the face, we employed MediaPipe’s FaceMesh, Hands, and Pose modules to extract kinematic and gestural information. For each frame, we tracked 3D keypoints for the hands, face contour, and upper body. Time aligned features included:

- **Hand movement dynamics:** frame-to-frame velocity, acceleration, and movement smoothness.
- **Gesture frequency:** counts of meaningful movements exceeding a motion threshold.
- **Face-touch gestures:** detection of hand to face contact events using spatial proximity between hand landmarks and facial meshes.
- **Upper-body motion:** displacement and temporal derivatives of shoulder, torso, and arm keypoints.

By integrating these multimodal cues—facial expression, vocal behavior, and gesture/movement dynamics—we obtain a rich representation of affective behavior capable of supporting downstream analyses such as affect recognition, behavioral profiling, and interpersonal dynamics modeling. This multimodal pipeline aligns with contemporary trends in affective computing that emphasize ecological validity, multimodal integration, and fine-grained temporal annotation [8, 7].

D. Linguistic Sentiment Analysis

To capture the affective content of spoken language, we extracted second-level sentiment scores from timestamped transcriptions of each conversation. Transcripts were structured in blocks containing speaker identity, temporal boundaries (start and end timestamps in HH:MM:SS.mmm format), and utterance text.

Sentiment analysis was performed using VADER (Valence Aware Dictionary and sEntiment Reasoner), a lexicon-based sentiment analysis tool optimized for social media and conversational text [9]. For each utterance, VADER computed four sentiment scores: positive (pos), neutral (neu), negative (neg), and compound (a normalized aggregate score ranging from -1 to +1).

To enable temporal alignment with other modalities, sentiment scores were expanded from utterance-level to second-level resolution. Each whole second covered by an utterance was assigned the sentiment scores of that utterance. When

multiple utterances overlapped within a single second, the most recent utterance's sentiment was retained to avoid redundancy. The resulting time-indexed CSV files contained per-second sentiment scores, speaker labels, and the corresponding utterance text, facilitating multimodal fusion with acoustic and visual features at 1 Hz temporal resolution.

IV. METHODOLOGY

Our approach models short temporal sequences of multimodal behavioral signals to predict binary labels - False if the person is not attracted to the other person and True otherwise. The methodology consists of feature aggregation every second, sequence construction, preprocessing, and sequential LSTM-based modeling. This section provides a precise description of each step, including hyperparameters and architectural details.

A. Sequence Construction and Feature Representation

Each date's video recordings are preprocessed to produce **per-second feature vectors** $\mathbf{x}_t \in \mathbb{R}^{28}$ comprising:

- 17 **Facial Action Units (AUs) for each person**
- 4 **Hand and upper-body movement features for each person**
- 3 **Acoustic features for the date**
- 4 **Sentiment features for each person**

These per second vectors are concatenated to form sequences of length $T = 15$ seconds:

$$\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T] \in \mathbb{R}^{T \times 28}.$$

Sequences are generated with a sliding window (stride 1) to capture temporal dependencies, producing overlapping samples for model training. Missing modalities are handled deterministically via zero filling, allowing the model to learn from absent or occluded signals without introducing NaNs or imputed values.

B. Feature Normalization

All features are standardized using a **StandardScaler (z-score normalization)** fit on the training set. For efficiency, the training sequences are reshaped from (N_{train}, T, D) to $(N_{\text{train}} \cdot T, D)$ prior to fitting. Validation and test sequences are transformed using the same scaler.

C. Sequential LSTM Architecture

The model is a stacked LSTM followed by fully connected layers:

- **Input:** sequences of shape $(T = 15, D = 28)$.
- **LSTM layers:**
 - 1) LSTM(64), return_sequences=True, dropout=0.2
 - 2) BatchNormalization
 - 3) LSTM(32), return_sequences=True, dropout=0.2
 - 4) BatchNormalization
 - 5) LSTM(16), return_sequences=False, dropout=0.2
 - 6) BatchNormalization
- **Dense layers:** Dense(32) with ReLU + Dropout(0.3), Dense(16) with ReLU + Dropout(0.2)

- **Output:** Dense(1) with sigmoid activation, producing probability $\hat{y} \in [0, 1]$

Formally, the model learns a mapping:

$$f_\theta : \mathbb{R}^{T \times D} \rightarrow [0, 1], \quad \hat{y} = f_\theta(\mathbf{X}),$$

where θ denotes all learnable parameters.

D. Training Objective and Optimization

The model is trained to minimize **binary cross-entropy**:

$$\mathcal{L}(\theta) = -\frac{1}{N} \sum_{i=1}^N w_{y_i} [y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i)],$$

where w_{y_i} are optional class weights computed from the training data to address class imbalance. The optimizer is **Adam**, with default parameters (learning rate = 0.001).

E. Training Protocol

Training uses:

- **Batch size:** 16
- **Epochs:** up to 100 with EarlyStopping (monitor=val_loss, patience=20, restore_best_weights=True)
- **Learning rate reduction:** ReduceLROnPlateau (factor=0.5, patience=10, min_lr=1e-6)
- **Metrics:** accuracy, precision, recall

Per person sequences are split such that validation data contains the last 20% of each participant's sequences, preventing data leakage and ensuring evaluation reflects generalization to unseen individuals.

F. Summary

This methodology integrates multimodal behavioral signals into a sequential modeling framework. The approach adheres to best practices in affective computing and sequence-based behavioral modeling [7, 6].

V. RELATED WORK

Our work lies at the intersection of multimodal affective computing [8, 7, 6] and the psychology of romantic attraction [3, 4, 2]. In contrast to classical studies that examine isolated channels (e.g., only facial attractiveness [1] or only body movement [5]), we adopt an integrated perspective combining facial expressions, body movement, voice, and spoken language over time [8, 3]. This section reviews three main strands of prior research: (i) multimodal affective computing and behavioral signal processing, (ii) psychological and behavioral studies of attraction, and (iii) multimodal and multichannel approaches to human attractiveness.

A. Multimodal Affective Computing and Sequence Modeling

Multimodal affective computing has made substantial progress in integrating heterogeneous signals such as facial expression, body posture, speech, and language for emotion recognition and related tasks. Survey work highlights the importance of combining modalities at different levels (early, late, and hybrid fusion) and representing their temporal dynamics to capture complex affective states [6, 7]. Architectures such as the Tensor Fusion Network explicitly model unimodal, bimodal, and trimodal interactions between vision, audio, and language to improve sentiment analysis [10], while more recent approaches introduce attention-based mechanisms and feature reuse strategies to enhance multimodal emotion recognition from audiovisual data [11].

Within this broader context, [8] proposes a multimodal emotion recognition system that integrates facial expressions, body movement, speech, and spoken language, demonstrating that jointly modeling these modalities yields more robust affect estimates than using any single channel alone. Methodologically, our data collection and feature-extraction pipeline is closely aligned with this line of work: we extract facial Action Units, hand and upper-body motion descriptors, acoustic prosodic features, and sentiment scores from transcripts, all temporally aligned at the per-second level. At the sequence-modeling level, we employ stacked recurrent neural networks to capture temporal dependencies, following the general recommendations and conceptual overview of recurrent architectures in [12].

Unlike most of these studies, however, our target variable is not an internal affective state such as valence or arousal, but an explicitly interpersonal outcome: a binary indicator of romantic interest toward another person. Thus, while our methodological foundation is typical of multimodal affective computing [13, 7, 8], our prediction task is shifted from intra-individual emotion classification to dyadic attraction estimation.

B. Psychological Perspectives on Romantic Attraction

Psychological research on romantic attraction has identified multiple contributing factors, including physical appearance, voice, movement, similarity, and context [2]. Recent work emphasizes that attraction is inherently multimodal and dynamic, emerging from the integration of cues across different sensory and behavioral channels during interaction. For example, recent work [3] synthesizes evidence showing that looks, voice, movement, and even scent jointly shape attraction to future lovers and friends, arguing that no single channel fully explains interpersonal appeal.

Experimental paradigms such as speed dating have been used to study the formation of initial romantic interest under controlled but ecologically valid conditions. [5] show that synchronized body sway and movement patterns during speed dates predict romantic interest, highlighting the importance of nonverbal motor coordination. [14] extend this paradigm to virtual environments, using online meeting platforms to investigate how initial attraction and relationship formation unfold

when interaction is mediated by video conferencing. These studies demonstrate that attraction judgments are sensitive to subtle temporal patterns in behavior, not just static snapshots of appearance or isolated self-report measures.

Our work is conceptually aligned with this literature in that we also focus on initial romantic interest following brief interactions. However, instead of analyzing a single behavioral channel (e.g., body sway alone) or relying primarily on self-report and static ratings, we attempt to predict attraction by jointly modeling dynamic facial, bodily, vocal, and linguistic cues within short temporal windows. This aligns with the theoretical view that attraction is a situated, context-dependent process unfolding over time [13, 2].

C. Multimodal Attractiveness and Communication of Appeal

Beyond emotion recognition and laboratory paradigms, several studies explicitly examine how different modalities communicate attractiveness and desirability. [4] analyze 881 judgments of men's and women's physical, vocal, and olfactory attractiveness, showing that these channels each contribute to perceived attractiveness, with partially overlapping but distinct information. Similarly, [3] argue for an "attraction in every sense" perspective, in which looks, voice, movement, and scent collectively inform evaluations of potential romantic partners and friends. These findings reinforce the idea that multimodal integration is necessary to approximate attraction judgments.

Our study extends this multimodal perspective in two important ways. First, we move from static or cross-sectional attractiveness ratings to temporally resolved behavioral sequences, incorporating not just what a person looks or sounds like, but how their expressions, gestures, and vocal prosody evolve over 15-second windows. Second, whereas prior work often focuses on global attractiveness or desirability ratings aggregated across many observers [4, 1], we aim to predict attraction at the dyadic level: whether a specific participant expresses romantic interest in a specific interaction partner [5, 14]. This shift from population-level attractiveness to pair attraction introduces additional variability but brings the modeling task closer to the phenomenon of romantic choice.

VI. RESULTS AND CONCLUSIONS

The model achieved a validation accuracy of 92% with high precision (96.6%) and recall (92.1%) on the positive class (attraction). These results suggest that short multimodal behavioral sequences contain predictive information about binary romantic attraction, supporting the motivation for temporal and multimodal modeling [6, 8]. However, several limitations should be noted. First, the dataset exhibits class imbalance (71% attracted vs. 29% not attracted), which may influence model predictions. Second, the relatively small validation set (300 samples) limits the strength of generalization claims. Third, while the integrated feature set combining facial, bodily, vocal, and linguistic signals shows promise [2, 3], systematic ablation studies are needed to determine the relative contribution of each modality.

Several aspects of the experimental design support the validity of these conclusions. Temporal segmentation and per-second aggregation ensure that predictions are based on dynamic interaction patterns rather than isolated frames [15]. Participant-level data splitting prevents information leakage between training and validation sets, promoting generalization to unseen individuals. Furthermore, deterministic handling of missing modalities through zero filling avoids introducing artificial correlations via imputation. Together, these design choices strengthen confidence that the reported performance reflects learned multimodal behavioral patterns rather than artifacts of preprocessing or evaluation.

In conclusion, this work provides preliminary evidence that multimodal, sequence-based modeling offers a promising approach to predicting romantic attraction from short segments of dyadic interaction [8, 13]. The presented results demonstrate feasibility but also highlight the need for larger, more balanced datasets and more rigorous evaluation protocols. Future work should include systematic ablation studies, attention-based architectures for interpretability, cross-validation across multiple data splits, and testing on independent datasets [10, 11].

VII. REFERENCES

REFERENCES

- [1] T. Sano and H. Kawabata. A computational approach to investigating facial attractiveness factors using geometric morphometric analysis and deep learning. *Scientific Reports*, 13(1):19797, 2023.
- [2] Theresa DiDonato and Brett Jakubiak. Romantic attraction. In *The Science of Romantic Relationships*, pages 123–157. Cambridge University Press, 2023.
- [3] A. Schirmer, M. Franz, L. Krismann, V. Nöring, M. Große, M. Mahmut, and I. Croy. Attraction in every sense: How looks, voice, movement and scent draw us to future lovers and friends. *British Journal of Psychology*, 116(3):684–701, 2025.
- [4] Megan Nicole Williams and Coren Lee Apicella. A test of multimodal communication in humans using 881 judgements of men and women’s physical, vocal, and olfactory attractiveness. *Heliyon*, 9(6):e16895, 2023.
- [5] A. Chang, H. E. Kragness, W. Tsou, D. J. Bosnyak, A. Thiede, and L. J. Trainor. Body sway predicts romantic interest in speed dating. *Social Cognitive and Affective Neuroscience*, 16(1–2):185–192, 2021.
- [6] Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe Morency. Multimodal machine learning: A survey and taxonomy. 2017.
- [7] A. Sherstinsky. Facial expression and body gesture emotion recognition: A systematic review on the use of visual data in affective computing. *IEEE Transactions on Affective Computing*, 13(4):1648–1670, 2022.
- [8] K. Kraack. A multimodal emotion recognition system: Integrating facial expressions, body movement, speech, and spoken language. *arXiv preprint*, 2023.
- [9] C. J. Hutto and Eric Gilbert. Vader: A parsimonious rule-based model for sentiment analysis of social media text. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 8, pages 216–225, 2014.
- [10] Amir Zadeh, Minghai Chen, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency. Tensor fusion network for multimodal sentiment analysis. 2017.
- [11] Z. Zhao, T. Gao, H. Wang, and B. W. Schuller. Swrr: Feature map classifier based on sliding window attention and high-response feature reuse for multimodal emotion recognition. pages 2433–2437, 2023.
- [12] R. M. Schmidt. Recurrent neural networks (rnns): A gentle introduction and overview. *arXiv preprint arXiv:1912.05911*, 2019.
- [13] L. A. M. Lebois, C. D. Wilson-Mendenhall, W. K. Simmons, L. F. Barrett, and L. W. Barsalou. Learning situated emotions. *Neuropsychologia*, 145:106637, 2018.
- [14] J. E. French, L. J. Bolton, and A. L. Meltzer. Virtual speed dating: Utilizing online-meeting platforms to study initial attraction and relationship formation. *Personal Relationships*, 31(2):420–444, 2024.
- [15] K. K. Haidet, J. Tate, D. Divirgilio-Thomas, A. Kolanowski, and M. B. Happ. Methods to improve reliability of video-recorded behavioral data. *Research in Nursing & Health*, 32(4):465–474, 2009.