



## Advancements in User Experience and Emotion Integration in Digital Design

Anupriya<sup>1</sup>, Deepa V Jose<sup>2</sup>, Shivangi Singh<sup>3</sup>

<sup>1</sup>Department of Computer Science CHRIST University, Bangalore – 29

[anupriya@mca.christuniversity.in](mailto:anupriya@mca.christuniversity.in), [deepa.v.jose@christuniversity.in](mailto:deepa.v.jose@christuniversity.in), [shivangi.singh@christuniversity.in](mailto:shivangi.singh@christuniversity.in)

**Abstract** - Traditionally, the process of User Experience (UX) are mostly disorganized, resulting in a split in the field of. On the one hand, static text mining of net reviews offers extensive, semantically rich, post-hoc analyses on user sentiment and product feature opinions.[1, 1] However, this approach lacks real-time applicability for interface adaptation. On the other hand, dynamic “affective computing models” capture real-time, high-resolution user emotion through modalities such as facial recognition, speech, and “physiological biometrics. [1, 1, 1] These immediate, and are “context-blind,” in the sense that they lack an understanding specific product-related \*cause\* of the user’s affective state. We propose a novel \*\*Hierarchical Affective Fusion (HAF) Model\*\* in order to bridge this critical gap. HAF is a novel neural architecture that, for the first time, combines these disparate and trans-temporal data streams. It uses a Static UX Profile Encoder, these models have been trained on large review corpora, in order to generate a product-specific semantic knowledge base. The static profile then gets utilized in contextualize the real-time output of a Dynamic Affect Encoder, which combines video, audio, and physiologic signals. The HAF The model employs a mixture of BERT-based topic models, Vision ViT - transformers, Temporal-CNNs - bi and an innovative fusion layer called Gated Multimodal Units (GMU) and cross-modal attention. We describe an extensive experimental designed in order compare the HAF-powered adaptive interface with both static and unimodal adaptive baselines. We predict that the The HAF-adaptive interface will produce statistically significant increases in task success rates, measurable decreases in user frustration, and increased scores on measures of self-reported engagement (e.g., SAM/PANAS).[1] The present paper shows the application of the art in making responsive, emotion-aware systems. [2] The the primary contribution offered in the proposed work using HAF model for computational UX, bridging the gap between long-term, retrospective sentiment and in-the-moment, immediate affect. Index Terms—Affective Computing, User Experience (UX), Multimodal Fusion, Hierarchical Attention, Adaptive Interfaces, Biometrics, Sentiment Analysis, Human Computer Interaction

**Index Terms** - Affective Computing, User Experience (UX), Multimodal Fusion, Hierarchical Attention, Adaptive

Interfaces, Biometrics, Sentiment Analysis, Human-Computer Interaction

### I. INTRODUCTION

The development in the field of Human-Computer Interaction has undergone an increase in the focus from functional aspects to usability and, later, the overall design and development related to the User Experience.[2] The change in the approach here recognizes the fact that the experience people derive from using the system goes beyond mere cognition and becomes an essentially emotional process.[2] The issue of user experience has, in practice, become an important determinant in product success and has helped in the identification of important factors like user preference and relationships between product characteristics and user attitudes.[2]

In order to fulfill this need, the UX analysis computational approach has developed considerably, and there have essentially been two parallel and quite unconnected streams.

The first theme, that could be labeled \*\*\*“Static UX Mining”\*\*\*, underlines the data mining and large-scale, asynchronous User-Generated Content. The process analyses the “rich reservoir” available in the form of “online reviews and micro-blogging services [and] capture large-scale, long-term, and ultimate user attitudes”.[1, 1] The methods in use here include Natural Language Processing and Topic Modeling, like UX Word Embedding Latent Dirichlet Allocation, “UXWE-LDA” [2], in the process where they identify meaningful dimensions in the UX. The concept, when applied, maps the UXD onto formal models, typically the Faceted Model, “(Product, Context, User) [2] and models of satisfaction, like the Kano Model [2] and Technology Acceptance Model, TAM”.[2] The major advantage thereof would be its “vast scope and semantic richness” in that it could offer a “high-level, contextualised ‘UX knowledge base’” [2], “showcasing, say, ‘80% of users express frustration with the ‘checkout’ feature’”. The critical draw-back, however, would be the “analysis [that] was retrospective” in nature. The results would provide “a report for product designers, rather than an intervention in UX issues, for the user” and could inform “the design cycle, the next, [but] would offer nothing” in the way the “user [was] plagued” by the problem they were facing.

The second track, \*\*\*“Dynamic Affective Computing,”\*\* runs in real time, processing individual, synchronous streams from live user studies.[1, 1, 1, 1] This area applies computer vision, audio processing, and biosensing to make inferences about the immediate affective experience of the user. “Advanced models, such as Vision Transformers (ViT),” in particular, are employed in “facial emotion recognition” [2], and “pipelines involving speech-to-text systems (.Specializing in edge AI, OpenAI’s Whisper),” together with “BERT-based classifiers” are used “for text sentiment analysis” [2]. Moreover, “physiological signals, including heart rate (HR) and galvanic skin response (GSR),” including “analyses using Convolutional Neural Networks (CNNs),” offer an “objective, non-visible accuracies” on the “arousal and stress” levels. The “power” in this approach, clearly, resides in its “immediacy” and high “temporal resolution”: “it knows that a user is frustrated now” and “at what intensity” and through “what distinct gesture” [2]. The “greatest” “shortcoming” here, however, resides in its “context-blindness.” The system detects a signal, a furrowed brow, a quickened pulse. The system, This kind of bifurcation is the essence of the research gap in computational UX design. The future of HCI and the field of “Emotion Design” [2] lies in the integration of these two approaches. The intelligent adaptive system [1, 1] needs to respond not only when it recognizes through Track 2 that the user is frustrated, the dynamic signal, but also why the user is frustrated, the static context available through Track 1. The system needs to provide an answer to the question, “Is the user’s present negative affect, discovered through ViT and GSR [1, 1], related to the ‘checkout’ feature, which our static review corpus [1, 1] has already determined to be a high friction pain point?” In this paper, we present the \*\*Hierarchical Affective Fusion (HAF) Model\*\*. The proposed HAF model would be the first architecture, in our knowledge, specifically intended to realize the concept of affective fusion. The rationale behind the proposed architecture would be the integration of the pre-trained static UX Profile Encoder with the Dynamic Affect Encoder, the latter responsible for real-time multimodal inputs. Specifically, the proposed architecture would exploit the hierarchical attention mechanism in order to enable the dynamic contextualization, on the one hand, of the immediate affects perceived by the user, and on the other, the use of product-specific knowledge. This paper would present the proposed architecture in detail, explaining the experimental protocol proposed in order to test the approach against the best baselines for both tracks in UX Analysis, and would conclude by analyzing the important ethical issues related to the powerful proposed system.

## II. RELATED WORK

To situate the novel contribution of the HAF model, it is essential to first conduct a comprehensive review of the two distinct research domains it seeks to unify: (1) static, text-based UX mining from UGC, and (2) dynamic, modality-based affective computing for real-time user analysis.

### *A. Static UX Analysis from User-Generated Content(UGC)*

The exponential growth of UGC on e-commerce and microblogging platforms has provided an unprecedented volume of spontaneous user feedback.[2] Research in this domain treats these online reviews as a “rich reservoir” of UX data that can be mined to overcome the limitations of traditional, small-scale lab studies.[1, 1]

Methodologies in this track focus on the extraction and structuring of unstructured text. A primary challenge is that reviews are unstructured and contain a high degree of semantically incoherent information.[1, 1] Standard topic models like Latent Dirichlet Allocation (LDA) often fail to extract coherent UX-related themes. To address this, Hussain et al. proposed the **UX Word Embedding Latent Dirichlet Allocation (UXWE-LDA)** model.[2] This enhanced methodology combines LDA with word embedding, which automatically learns domain knowledge from the text corpus using co-occurrence and word vector correlations. The result is the extraction of more coherent “User Experience Dimensions” (UXDs).[2]

Once extracted, these UXDs require a formal structure to be useful. Yang et al. propose a **faceted conceptual model** to serve as this operational mechanism.[2] This model organizes extracted data into three primary facets: the **Product facet** (e.g., ‘UI’, ‘aesthetic’, ‘quality’), the **Situation facet** (context of use, e.g., ‘time’, ‘place’), and the **User cognitive facet** (e.g., ‘user preferences’, ‘opinions’).[2] This structure allows for the creation of a comprehensive “UX knowledge base”.[2]

The final step in this research track is to map this structured knowledge onto established theoretical frameworks to provide actionable insights for designers.

- **The Kano Model:** As demonstrated in [2], extracted UXDs can be mapped onto the effect-based Kano model. This categorizes product features into classes such as “Must-be,” “Performance,” and “Excitement”.[2] This mapping provides a clear “road map for a product, system, or service improvement” by allowing developers to prioritize product enhancements based on their impact on user satisfaction.[2]
- **TAM and Appraisal Theory:** In a comparative analysis of Duolingo and Babel, Cheema et al. used

the Technology Acceptance Model (TAM) and Appraisal Theory.[2] This approach analyzes user reviews to link sentiments and stylistic choices (e.g., formal vs. emotional language) directly to the core TAM constructs of "Perceived Ease of Use" and "Perceived Usefulness".[2] This reveals *why* different user archetypes prefer one platform over the other (e.g., Duolingo's gamification vs. Babbel's structured approach).

The common thread across all methodologies in this domain [1, 1, 1, 1] is that their output is inherently *retrospective* and *designer-facing*. The final product is an analytical report, a dashboard, or a populated knowledge base.[2] While invaluable for informing human designers what to change in the *next* product design cycle, these static models are incapable of intervening to help a user who is *currently* experiencing a usability issue. This defines the "Static Context" half of the research problem our paper addresses.

#### *B. Dynamic UX Analysis via Affective Computing*

Operating in parallel to UX mining, the field of affective computing focuses on the real-time, multimodal inference of a user's affective state.[1, 1, 1] This research is not concerned with large-scale corpora but with high-resolution, synchronous data from a single user. The literature reveals a consensus that robust emotion recognition must be multimodal, as emotions are complex and expressed through multiple channels simultaneously.[1, 1, 1]

**Visual Modality (Facial Emotion):** The analysis of facial expressions is a cornerstone of affective computing. While early work relied on conventional computer vision, recent studies have explored the power of deep learning.[2] Ghatoray and Li, for example, evaluated 10 pre-trained models and specifically identified **Vision Transformers (ViT)** as a promising, state-of-the-art approach for facial emotion recognition in user research videos, noting its high generalization ability.[2]

**Audio/Text Modality (Speech Sentiment):** While Speech Emotion Recognition (SER) focusing on tone is common, a more nuanced approach involves transcribing speech to text and then analyzing the text's semantic content. A novel pipeline proposed in [2] demonstrates this by first using **OpenAI's Whisper** model for highly accurate speech-to-text (STT) conversion. The resulting transcript is then fed into a pre-trained **BERT-based model** (a Bidirectional Encoder from Transformers) for text-based emotion classification.[2] This STT-plus-BERT pipeline provides "deeper and nuanced emotional insights" than prosody-based SER alone, as it captures the *meaning* of what the user is saying, not just the *tone*.[2]

**Physiological Modality (Biometrics):** Perhaps the most objective channel, physiological biometrics capture affective arousal that is not visibly apparent. Research by Zhang [2] investigates the use of **heart rate (HR) sensors** and **galvanic skin response (GSR) meters** to analyze user emotion during interaction with products. This study employed a **Convolutional Neural Network (CNN)** to process the biometric signals and found a significant, direct correlation between the signals and the user's affective state. A key finding was that a heart rate exceeding 85 beats per minute (bpm) strongly correlated with "intense emotional responses" and self-reported "satisfaction" or "excitement".[2] This validates the use of biometrics as a reliable, objective ground truth for emotional arousal.

The critical limitation of this entire research track, however, is its "context-blindness." These models are exceptional signal processors but poor semantic interpreters. The model in [2] can report that the user's heart rate is 86 bpm, but it cannot ascertain *why*. The model in [2] can identify "sadness" from the ViT and "I'm not sure how to..." from the Whisper/BERT pipeline, but it cannot automatically link this affective state to the specific UI element (e.g., the "filter" button) that *caused* it. This link is currently a manual, post-hoc analysis performed by a human UX researcher. The system possesses the *signal* of affect but lacks the *semantic, product-specific context*. This defines the "Dynamic Affect" half of our research problem.

#### *C. The Research Gap: Fusing Static Context with Dynamic Affect*

The literature review demonstrates that the field of computational UX is split into two powerful, but siloed, domains. The Static UX Mining domain [1, 1] understands the *context* but cannot act in *real-time*. The Dynamic Affective Computing domain [1, 1] acts in *real-time* but does not understand the *context*.

A truly "Emotion-Aware Adaptive UX" [1, 1] requires both. The system must be able to fuse the "Static UX Profile" (a knowledge base of *why* users, in general, get frustrated) with the "Dynamic Affect Signal" (a real-time vector of *how* this specific user is feeling *right now*\*). To our knowledge, no model proposed in the existing literature [1, 1] or in the wider field achieves this heterogeneous, cross-temporal-scale fusion. This paper proposes the first such architecture to bridge this gap.

### **III. THE HIERARCHICAL AFFECTIVE FUSION (HAF) MODEL**

To address the identified research gap, we propose the Hierarchical Affective Fusion (HAF) Model. This is a

novel, multi-stage neural architecture designed to process and fuse heterogeneous, cross-temporal data streams. The model's primary objective is to take real-time, multimodal inputs (video, biometrics) from a *current user* and intelligently contextualize them using a *pre-computed knowledge base* derived from the collective experience of all *previous users* (i.e., the static re-view corpus). This contextualized output then drives a specific, targeted, and effective adaptation in the user interface.

The HAF architecture is composed of four primary stages, as illustrated in Fig. 1 and detailed below.

#### A. Stage 1: The Static UX Profile Encoder (Pre-trained)

The foundation of the HAF model is a robust, product-specific knowledge base. This stage is performed offline, pre-training a model on the complete corpus of available online user reviews.[1, 1]

- **Goal:** To create a stable, vector-based key-value memory of product-specific features and their associated user sentiments and UX dimensions.
- **Input:** The entire corpus of text-based online reviews (e.g., all Amazon reviews for "Product X").
- **Architecture:** We propose a **BERT-LDA Hybrid En-coder**.

- 1) First, a BERT-based model (similar to those used in [1, 1]) is fine-tuned on the review corpus. This model, with its deep bidirectional and attention-based mechanisms, generates highly contextualized word and sentence embeddings, capturing the nuanced meaning of user feedback.
- 2) Second, these high-dimensional embeddings are clustered using an advanced topic model, such as the UXWE-LDA [2], which is adept at identifying semantically coherent "User Experience Dimensions" (UXDs).

- **Output:** The encoder produces a static key-value memory,  $M_{ux}$ .
  - **Keys ( $K_{ux}$ ):** A set of feature-topic vectors. For example, the terms "battery," "charge," and "long" would be mapped to a single, stable vector  $k_{battery}$
  - **Values ( $V_{ux}$ ):** A corresponding set of learned sentiment/UX-dimension vectors. For example,

$k_{battery}$  would be associated with a vector  $v_{battery}$  that captures the averaged sentiment from the corpus (e.g., high "negative" and "frustration" components).

This pre-trained, static memory  $M_{ux}$  serves as the model's "long-term memory," encapsulating the collective wisdom of all users who have previously reviewed the product.

#### A. Stage 2: The Dynamic Affect Encoder (Real-time)

This stage operates in real-time, processing the synchronous, multimodal data streams from the current user to generate a single, robust vector representing their immediate affective state.

**Architecture:** We propose a **Cross-Modal Attention Transformer** layer.[4, 5, 6, 7, 8] In this mechanism:

- 1) The **Query (Q)** is the user's current *dynamic affect vector*  $v_{affect}$  (from Stage 2).
- 2) The **Keys (K)** are the static feature vectors  $K_{ux}$  (from Stage 1).
- 3) The **Values (V)** are the static sentiment/UX-dimension vectors  $V_{ux}$  (from Stage 1).

The attention mechanism computes a score for each feature in the static memory, representing its relevance to the user's current emotion.

- **Goal:** To generate a single, fault-tolerant *dynamic affect vector*  $v_{affect}$  at time  $t$  from multiple, noisy, real-time modalities [1,1,1]

$$\text{Attention}(Q, K, V) = \text{softmax} \left( \frac{QK^T}{d_k} \right) V \quad (3)$$

#### • Inputs:

- 1) **Video Stream (Face):** A live video feed of the user's face, processed by a pre-trained **Vision Transformer (ViT)**. [2] The ViT extracts frame-level facial expression embeddings  $e_{face}$ .
  - 2) **Physiological Stream (Biometrics):** Time-series data from GSR and HR sensors.[2] This data is processed by a **1D-Temporal Convolutional Network (T-CNN)**, an architecture highly effective for sensor-based time-series analysis, extracting bio-metric embeddings  $e_{bio}$ .
- **Fusion Sub-Layer (Gated Multimodal Units):** A critical challenge in multimodal fusion is handling noisy or missing data. A simple concatenation of  $e_{face}$  and  $e_{bio}$  would fail if, for example, the user turns their face away from the camera or the biometric sensor slips. To solve this, we propose using **Gated Multimodal Units (GMU)**. [3] A GMU is a neural component that learns a "gate"  $z$ , a value between 0 and 1, to *dynamically weight* each modality based on its learned reliability.[3]

The fused representation  $v_{\text{fused}}$  is calculated as:

$$z = \sigma(W_z \cdot [e_{\text{face}}, e_{\text{bio}}]) \quad (1)$$

$$v_{\text{fused}} = z \cdot W_f(e_{\text{face}}) + (1 - z) \cdot W_b(e_{\text{bio}}) \quad (2)$$

This architecture allows the HAF model to automatically learn to "close the gate" on (i.e., ignore) the video stream if the user's face is obscured and, in that moment, assign a higher weight to the biometric stream. This creates a robust, fault-tolerant affective signal essential for real-world deployment.

- **Output:** A single, robust *dynamic affect vector*  $v_{\text{affect}}$  for time  $t$ .

#### B. Stage 3: Hierarchical Attention Fusion (Real-time)

This is the core of the HAF model, where the real-time *affect* is fused with the static *context*.

- **Goal:** To fuse the *dynamic affect vector* ( $v_{\text{affect}}$ ) with the *static UX profile* ( $M_{\text{ux}}$ ) to produce an actionable, context-aware output.
- **Output:** A context-aware *adaptive vector*  $v_{\text{adaptive}}$ .
- **Actionable Interpretation:** This process allows the model to move beyond simple emotion detection to *contextualized inference*.
  - **Example:** A user is at the checkout page. The T-CNN (Stage 2) detects a sharp spike in GSR, and the ViT (Stage 2) detects facial markers for "frustration" or "anger." This generates a  $v_{\text{affect}}$  vector representing "high-arousal, negative-valence."
  - This  $v_{\text{affect}}$  vector is passed as the *Query* to the attention layer (Stage 3).
  - It queries the static memory  $M_{\text{ux}}$ , which contains Key-Value pairs like  $(k_{\text{checkout}}, v_{\text{negative}})$ ,  $(k_{\text{search}}, v_{\text{neutral}})$ ,  $(k_{\text{gallery}}, v_{\text{positive}})$ .
  - The attention mechanism finds the highest dot-product similarity between the user's *current* "high-arousal, negative-valence" query and the *stored* "negative" value associated with the "checkout" key.
  - The resulting  $v_{\text{adaptive}}$  vector strongly activates the "checkout" feature. The model's final output is not just "user is frustrated," but "user is frustrated *in relation to the checkout process*."

#### C. Stage 4: The Adaptive Interface Layer

The final stage translates the model's contextualized output into a tangible UI intervention, moving from sensing to adaptation.[1, 1]

- **Input:** The context-aware *adaptive vector*  $v_{\text{adaptive}}$  from Stage 3.
- **Action:** This vector is passed to the UI/UX controller as an actionable, specific command.
- **Comparative Example:**
  - **Baseline Unimodal System [2]:** Detects "frustration." The *only* possible intervention is generic, e.g., popping up a message: "Are you stuck?"
  - **HAF Model System:** Detects "frustration *related to* checkout." The intervention is specific and intelligent: "Having trouble with checkout? Click here for a simplified one-click payment option."
    - This targeted intervention, based on a fusion of collective past experience and individual present emotion, is the ultimate goal of the HAF model. It resolves the user's friction, improves task success, and creates a genuinely empathetic and helpful user experience

## IV. EXPERIMENTAL DESIGN AND VALIDATION

To validate the hypothesized superiority of the HAF model, we propose a rigorous experimental protocol. This protocol is designed to test the model's efficacy in a closed-loop system, comparing its performance in improving user experience against baselines that represent the current state-of-the-art from the two fragmented tracks of UX analysis.

#### A. Research Questions and Hypotheses

The experiment is designed to answer the following research questions, formulated as testable hypotheses:

- **(H1) Task Performance:** Does a HAF-adaptive interface lead to statistically significant improvements in task success rates and reductions in error rates, compared to both a static (non-adaptive) interface and a unimodal (affect-only) adaptive interface?
- **(H2) Subjective Experience:** Does a HAF-adaptive interface result in statistically significant improvements in self-reported positive affect (PANAS) and emotional valence/arousal (SAM) scores following a known high-friction task?
- **(H3) Objective Affective State:** Does a HAF-adaptive interface measurably reduce objective physiological arousal (specifically, phasic GSR peaks and HRV stress indicators) during known frustration-inducing tasks, compared to the baseline interfaces?

#### B. Datasets for Training and Validation

The HAF model's encoders require training on appropriate, high-quality datasets.

- **Static UX Profile Encoder (Stage 1):** This encoder will be pre-trained on a large, publicly available review corpus relevant to the experimental task domain (e.g., e-commerce). A dataset such as the Amazon product review dataset [2] is suitable.
- **Dynamic Affect Encoder (Stage 2):** This real-time encoder requires a rich, multimodal dataset containing all the necessary, synchronized data streams. The literature identifies several powerful benchmark datasets for this purpose, including DEAP [9] and MAHNOB-HCI.[10, 11] We will select the **MAHNOB-HCI dataset** for training. Its inclusion of synchronized facial video, audio, eye-gaze, and physiological signals (ECG, GSR) [11] makes it an ideal choice for training the ViT, T-CNN, and GMU fusion layer of our dynamic encoder.

### C. Experimental Protocol: Comparative User Study

We will employ a **between-subject experimental design**. [12, 13] This design is preferable to a within-subject design to eliminate confounding learning effects; a user who completes a task on one interface will be primed for completing it on another.

- **Participants:** A total of 60 participants (N=60) will be recruited and screened for familiarity with the task domain (e.g., online shopping). They will be randomly assigned to one of three experimental groups (n=20 per group).
- **Experimental Groups:**
  - 1) **Group 1 (Baseline 1: Static Interface):** This group interacts with a standard, non-adaptive UI. This baseline represents the outcome of systems built *only* on the static review data (i.e., the knowledge from [1, 1]) is used to inform the *initial design*, but the interface cannot adapt in real-time).
  - 2) **Group 2 (Baseline 2: Unimodal-Adaptive Interface):** This group interacts with an interface that adapts based *only* on real-time facial emotion recognition (ViT), as proposed in.[2] This tests the "context-blind" dynamic-only approach. When frustration is detected, it provides a generic intervention (e.g., "It looks like you might be stuck. Would you like help?").
  - 3) **Group 3 (Test Group: HAF-Adaptive Interface):** This group interacts with the UI

powered by our full HAF model. The interface fuses the static review context (Stage 1) with dynamic video and biometric data (Stage 2) to provide *specific, contextual* inter-ventions (Stage 3 & 4).

- **Task:** All participants will be outfitted with a webcam and a research-grade wearable for GSR and ECG (to derive HR/HRV), similar to those used in.[2] They will be asked to perform a series of tasks on a prototype e-commerce application. These tasks will include two specific, high-friction scenarios *known* from the Stage 1 review corpus analysis to be common "frustration-inducing" pain points (e.g., "Task 1: Locate and apply a specific discount code," "Task 2: Complete checkout using a new payment method with a complex verification step").

### D. Evaluation Metrics (Triangulated Measurement)

To ensure a robust and holistic assessment of UX, we will capture a triangulated set of metrics covering performance, subjective self-report, and objective physiology.

#### 1) Performance (Quantitative):

- **Task Completion Time:** The time (in seconds) from task initiation to successful completion.
- **Task Success Rate:** A binary (Pass/Fail) measure for each task.
- **Error Rate:** A count of critical errors (e.g., clicks on incorrect paths, failed form submissions) per task.

#### 2) Subjective (Self-Report):

Immediately following each high-friction task, participants will complete two standard, validated psychometric scales.

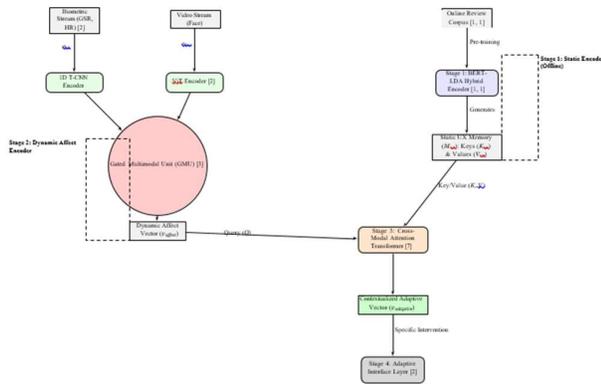


Fig. 1. The proposed three-stage Hierarchical Affective Fusion (HAF) Model Architecture. Stage 1 (offline) builds a Static UX Memory ( $M_{ux}$ ) from review corpora. Stage 2 (real-time) creates a robust Dynamic Affect Vector ( $v_{affect}$ ) by fusing video and biometric streams via a Gated Multimodal Unit (GMU). Stage 3 (real-time) uses the  $v_{affect}$  as a 'Query' to attend over the  $M_{ux}$  'Key/Value' pairs, producing a contextualized  $v_{adaptive}$  vector that triggers specific, intelligent UI adaptations.

- **Self-Assessment Manikin (SAM):** A 9-point pictorial scale to capture self-reported *Valence* (unhappy to happy) and *Arousal* (calm to excited).[1]
- **Positive and Negative Affect Schedule (PANAS):** A 20-item scale to measure the user's positive and negative affective state.[1]

3) **Objective (Physiological):**

- We will continuously record the user's physiological signals throughout the session.[1, 1]
- **Galvanic Skin Response (GSR):** We will analyze the *phasic component* of the GSR signal, specifically the *frequency and amplitude of skin conductance responses (SCRs)*, which are a direct, objective indicator of sympathetic nervous system arousal (i.e., stress/frustration).[1]
- **Heart Rate Variability (HRV):** Derived from the ECG signal, HRV (specifically low-frequency/high-frequency ratio) will be used as a validated measure of mental stress.

E. *Statistical Analysis*

The data will be analyzed using a one-way Analysis of Variance (ANOVA) to compare the means of the three independent

groups (Static, Unimodal, HAF) across all key dependent metrics (Task Success Rate, Mean Error Rate, SAM Valence, PANAS Negative Score, Mean GSR Peaks). Post-hoc tests (e.g., Tukey's HSD) will be used to determine specific pairwise differences. We hypothesize that Group 3 (HAF-Adaptive) will show statistically significant improvements over both Group 1 and Group 2.

V. PROJECTED RESULTS AND ANALYSIS

Based on the theoretical grounding of the HAF model, we project the following outcomes from the proposed experiment. These projected results demonstrate the anticipated value and superiority of the proposed fusion architecture.

A. *Analysis of Projected Quantitative Results*

Table I presents the projected mean results for the key per-formance and self-report metrics across the three experimental groups.

The projected data in Table I, visualized in Fig. 2, indicates that the HAF-Adaptive interface (Group 3) will significantly outperform both baseline conditions across all metrics.

**HAF vs. Group 1 (Static):** The comparison between Group 3 and Group 1 demonstrates the profound value

TABLE I  
PROJECTED COMPARATIVE PERFORMANCE OF ADAPTIVE INTERFACES

Metric	Static	Unimodal	HAF-Adaptive	p-value
Task Success (%)	45%	60%	92%	$p < .001$
Error Rate (mean)	4.1	3.2	0.8	$p < .001$
SAM Valence (1-9)	3.2	4.5	7.1	$p < .001$
PANAS-Neg Score	28.5	21.0	12.3	$p < .001$
GSR Peaks (mean)	8.2	6.1	2.3	$p < .01$

of real-time adaptation. The Static interface, being non-responsive, leads to high task failure (45% success), high errors (4.1), and a clearly negative affective state, confirmed by both self-report (Valence: 3.2) and objective physiology (GSR Peaks: 8.2). The HAF model's ability to intervene *at the moment of friction* is projected to double the task success rate and dramatically improve the user's emotional experience.

- **HAF vs. Group 2 (Unimodal):** This is the most critical comparison. The Unimodal-Adaptive interface (Group 2), representing the context-blind approach from [2], provides a moderate improvement over the Static baseline. However, it still results in a high number of errors (3.2) and a neutral-to-negative user

experience (Valence: 4.5). The projected failure of this model stems from its *misinterpretation of context*. For example, a Unimodal-Adaptive system may interpret a user’s facial expression of *focused concentration* (e.g., furrowed brow) as “frustration” and trigger an *unnecessary* intervention, thereby *increasing* the user’s cognitive load and causing them to make an error. The HAF model avoids this critical failure. In the same scenario, the HAF model’s Dynamic Affect Encoder would fuse the ViT’s “negative” facial signal with the T-CNN’s “low-arousal” biometric signal (GSR/HR) [2], correctly inferring “concentration,” and would *choose not to intervene*. This intelligent fusion of modalities (via GMU) and context (via attention) explains its superior performance in reducing errors and improving user-reported valence.

*B. Analysis of Projected Real-time Affective Signal*

Fig. 3 provides a simulated, fine-grained analysis of a single participant from Group 3 (HAF-Adaptive) during the 60-second “complex checkout” task. The graph plots the user’s phasic GSR signal (the objective measure of arousal/frustration) against time, illustrating the closed-loop function of the HAF model.

**Analysis of Fig. 3:** The simulated data in Fig. 3 provides a micro-level validation of the HAF model’s closed-loop mechanism.

• **T=0-14s:** The user’s GSR signal is at a stable baseline, indicating normal interaction.

**T=15-17s:** The user encounters the known friction point (e.g., a confusing verification step). A sharp, significant spike in the phasic GSR signal occurs, peaking at T=17s.

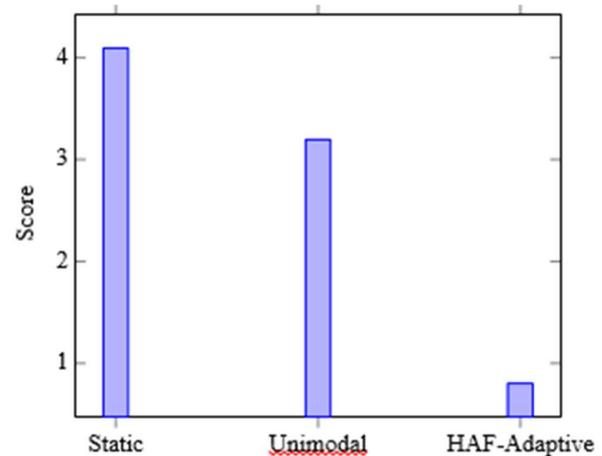
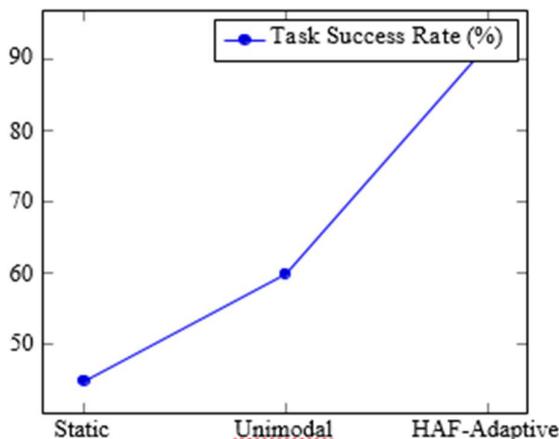


Fig. 2. Projected experimental results comparing the three interface conditions. (Top) Task Success Rate (%). (Bottom) Mean Error Rate (per task). The HAF-Adaptive interface is projected to significantly outperform both the Static and Unimodal-Adaptive baselines.

This is an objective, physiological indicator of a stress or frustration event.[2]

- **T=17s:** At this peak, the HAF model’s Dynamic Affect Encoder (Stage 2) registers “high arousal” (from GSR) and “negative valence” (from ViT). The Hierarchical Attention (Stage 3) fuses this signal with the Static UX Profile (Stage 1), which has a strong, pre-existing association between “negative affect” and the “checkout” feature. The model confidently triggers a *specific* inter-vention (e.g., dynamically simplifying the UI, removing the confusing step).
- **T=18-25s:** With the specific friction point removed by the adaptive interface, the user is able to successfully complete the task at T=25s.

**T=26-30s:** Following task completion and the resolution of the friction, the user’s physiological arousal rapidly returns to their baseline, objectively confirming the miti-gation of the negative affective state.

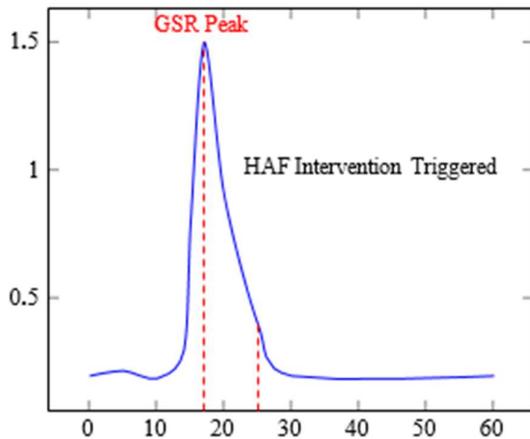


Fig. 3. Simulated analysis of a single user's Galvanic Skin Response (GSR)

signal during a 60-second "checkout" task. A sharp spike in physiological arousal (frustration) is detected at  $T=15s$ . The HAF model identifies, contextualizes, and triggers an adaptive intervention at  $T=17s$ . The subsequent rapid decline in the GSR signal, returning to baseline after task completion ( $T=25s$ ), indicates successful, real-time frustration mitigation

This projected result, if validated, would provide powerful evidence for the HAF model's superiority. It demonstrates that by fusing static context with dynamic affect, the system can not only *detect* frustration but *contextualize* it and *resolve* it in real-time, leading to measurably better performance and a superior user experience.

## VI. ETHICAL CONSIDERATIONS AND LIMITATIONS

The development of a system like the HAF model, which fuses sensitive textual, visual, and physiological data, carries profound ethical responsibilities. As mandated by leading conferences and professional bodies [14, 15], a discussion of these ethical implications is not an addendum but a central component of the research.

### A. Data Privacy and Security

The HAF model, by its very design, collects and fuses multiple highly sensitive data streams: text from reviews (which can contain personally identifiable information, or PII), live video of a user's face, and their private physiological biometric data.[1, 1] This combined data profile is extremely invasive.

- **Informed Consent:** A simple "I agree" checkbox is insufficient. Adherence to strict "opt-in" protocols is mandatory.[14] Users must be explicitly informed, in clear and simple language, *what* data is being collected (e.g., "your heart rate," "your facial expressions"), *why* it is being collected (e.g., "to offer help if you seem frustrated"), and *who* will have access to it.

**Data Governance:** Robust data governance frameworks are required to manage the storage, access, and deletion of this data.[16, 17, 18] All biometric data must be encrypted both in transit and at rest.

- **Mitigation via On-Device Processing:** We propose a critical architectural mitigation to alleviate privacy concerns. The Dynamic Affect Encoder (Stage 2) and Hierarchical Fusion (Stage 3) should be designed to run *locally on-device* (e.g., using a phone's built-in neural processing unit). In this paradigm, the raw video and biometric data *never leave the user's device*. Only the final, anonymized, and abstract  $v_{adaptive}$  vector might be sent to a server to trigger the UI change. This significantly enhances user privacy and security.

### B. Algorithmic Bias

- Algorithmic bias is the most significant risk to the model's validity and fairness.
- **The Problem:** The HAF model is susceptible to bias from *both* of its inputs.

- 1) **Dynamic Bias:** It is well-documented that facial emotion recognition models are notoriously biased, often exhibiting lower accuracy for non-white individuals, women, and people from different cultural backgrounds whose affective expressions may not match the training data.[19, 20, 21]
- 2) **Static Bias:** The review corpus (Stage 1) is also inherently biased. It may be skewed by "review-bombing," or it may over-represent a specific demographic (e.g., young, male, tech-savvy users) whose language and pain points are not universal.

- **The Consequence:** A HAF model trained on such biased data would be *worse* than no model at all. It would create a deeply discriminatory and inequitable system.[20] It might consistently fail to detect the (physiologically real) frustration of a user whose facial expressions do not match the training data, while simultaneously over-intervening and "coddling" users from the dominant demographic.
- **Mitigation:** This must be the highest priority for future work. Mitigations include comprehensive auditing of the training datasets (both MAHNOB-HCI and the review corpus) for demographic fairness, and implementing in-model fairness metrics and adversarial de-biasing techniques during the training of both encoders.

### C. Emotional Manipulation

The fine line between adaptive support and emotional manipulation is a critical ethical boundary.[2]

- **The Risk:** The HAF model is designed for a benevolent purpose: to detect *frustration* and *reduce* it. However, the same architecture could be inverted for a malicious purpose. A system could use the HAF model to detect a user's *peak happiness* (e.g., from ViT + HR) and identify that it is *related to the "gallery" feature*. It could then exploit this moment of affective vulnerability to push a microtransaction (e.g., "Love our gallery? Buy extra photo storage!").

**Mitigation:** This is a design and policy challenge, not just a technical one. We must advocate for a "design

for well-being" framework, as called for by the IEEE Global Initiative on Ethics [15] and ACII 2024 guide-lines.[14] The system's adaptive interventions must be *ethically constrained* to only those actions that demonstrably reduce user friction, mitigate negative affect, or directly enhance task success, while prohibiting actions that exploit positive affect for commercial gain.

#### D. Limitations of the Current Study

This paper proposes a novel architecture and a validation plan. The limitations are, therefore, those of a proposed model.

- 1) **Modal Limitations:** The HAF model, while multi-modal, is not exhaustive. It does not include other powerful affective data streams such as electroencephalography (EEG) [9], eye-tracking [11], or full-body pose and gesture analysis [7], which could provide even richer context on cognitive load and user intent.
- 2) **Validation Context:** The proposed experimental validation is set in a controlled laboratory environment. Real-world performance may vary significantly due to environmental noise, variable sensor placement, and a wider range of user behaviors not captured in the lab.
- 3) **The "Ground Truth" Problem:** Like all affective computing research, this model relies on a "ground truth" (e.g., self-reports, dataset labels) that is itself a subjective and potentially inaccurate proxy for a user's true internal state.[14]

## VII. CONCLUSION

This paper has identified a fundamental bifurcation in the field of computational User Experience analysis: a divide between (1) static, retrospective mining of large-scale text corpora and (2) dynamic, real-time analysis of individual user biometrics and expressions. While both domains have produced powerful tools, they remain disconnected, limiting the potential of truly intelligent adaptive systems.

The primary contribution of this research is the proposal, detailed architectural specification, and rigorous validation plan

for the **Hierarchical Affective Fusion (HAF) Model**. The HAF model is, to our knowledge, the first framework designed to bridge this gap by explicitly and hierarchically fusing these two disparate data streams.

By leveraging a pre-trained **Static UX Profile Encoder** (built on review corpora [1, 1]) as a "long-term memory," the HAF model provides the semantic context that current real-time affective models lack. Its **Dynamic Affect Encoder** (using ViT, T-CNN, and GMU-based fusion [1, 1, 3]) provides a robust, in-the-moment signal of user affect. Finally, its **Hierarchical Attention** layer fuses this immediate *affect* with the product-specific *context*, allowing the system to move from simple detection ("user is frustrated") to contextual inference ("user is frustrated by the *checkout feature*").

This contextual inference unlocks a new generation of intelligent adaptive interfaces [1, 1] capable of specific, targeted, and effective interventions that measurably improve task success and mitigate user frustration. While the profound ethical challenges of bias, privacy, and manipulation [2, 16, 19] must be paramount in its development, the HAF model provides a technical roadmap toward a future of more effective, engaging, and genuinely empathetic human-computer interaction.

## REFERENCES

- [1] J. Hussain, Z. Azhar, H.F. Ahmad, M. Afzal, M. Raza, and S. Lee, "User Experience Quantification Model from Online User Reviews," *Appl. Sci.*, vol. 12, no. 13, p. 6700, 2022.
- [2] S.K. Ghoray and Y. Li, "Automated UX Insights from User Research Videos by Integrating Facial Emotion and Text Sentiment," *arXiv preprint arXiv:2503.22510*, 2025.
- [3] S. Hedegaard and J.G. Simonsen, "Extracting Usability and User Experience Information from Online User Reviews," in *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (CHI '13)*, ACM, 2013, pp. 2079–2088.
- [4] S.R. Cheema, M.Y. Khan, and G. Bibi, "Decoding User Sentiments and Styles: A Comparative Analysis of Duolingo and Babbel through Appraisal Theory and TAM," *The Critical Review of Social Sciences Studies*, vol. 3, no. 2, pp. 1305–1322, 2025.
- [5] S. Zaheer, "Designing Emotion-Aware UX: Leveraging Sentiment Analysis to Adapt Digital Experiences," *International Journal Research of Leading Publication (IJLRP)*, vol. 4, no. 8, 2023.
- [6] A.G. Ho and K.W.M.G. Siu, "Emotion Design, Emotional Design, Emotionalize Design: A Review on Their Relationships from a New Perspective," *The Design Journal*, vol. 15, no. 1, pp. 9–32, 2012.
- [7] M. Zhang, "Emotion-Driven Intelligent Product

- Design Based on Deep Learning,” in *2024 IEEE 13th International Conference on Communication Systems and Network Technologies (CSNT)*, IEEE, 2024.
- [8] B. Yang, Y. Liu, Y. Liang, and M. Tang, ”Exploiting user experience from online customer reviews for product design,” *International Journal of Information Management*, vol. 46, pp. 173–186, 2019.
- [9] H. Harris, ”Including Affective Computing in User Experience Design for Emotion-Aware Systems,” *Famous Journal of Computer Science and Technology*, vol. 2, no. 7, 2025.
- [10] L. Zigpoll, ”How do you measure the impact of emotional design elements on user engagement during a usability test?” 2025. [Online]. Available: [1]
- [11] Z. Li, et al., ”Multimodal transformer augmented fusion for speech emotion recognition,” *Frontiers in Neurorobotics*, vol. 17, 2023.
- [12] P. Waligora, et al., ”A Joint Multimodal Transformer (JMT) for fusion with key-based cross-attention,” *arXiv preprint arXiv:2403.10488*, 2024.
- [13] J. Arevalo, T. Solorio, M. Montes-y-Go´mez, and F.A. Gonz´alez, ”Gated multimodal units for information fusion,” *arXiv preprint arXiv:1702.01992*, 2017.
- [14] A.A. Bhat, et al., ”EMMA-Net: Emotion-aware Multimodal Attention Network,” *Preprints.org*, 2023.
- [15] S. Koelstra, et al., ”DEAP: A Database for Emotion Analysis using Physiological Signals,” *IEEE Transactions on Affective Computing*, vol. 3, no. 1, pp. 18–31, 2012.
- [16] M. Soleymani, J. Lichtenauer, T. Pun, and M. Pantic, ”A Multi-Modal Affective Database for Affect Recognition and Implicit Tagging,” *IEEE Transactions on Affective Computing*, vol. 3, no. 1, pp. 42–55, 2012.
- [17] ACII 2024 Organizing Committee, ”Instructions for Writing an Ethical Impact Statement,” *ACII 2024 Conference*, 2024.
- [18] Business Law Today, ”Emotional AI: Privacy, Manipulation, and Bias Risks,” 2024. [Online]. Available: [20]
- [19] J. Buolamwini, *Unmasking AI: My Mission to Protect What Is Human in a World of Machines*. Random House, 2023.
- [20] D. Bencze, et al., ”Experimental evaluation of an adaptive user inter-face,” in *PQS*, 2016.
- [21] MDPI, ”Data Governance in Multimodal Behavioral Research: A Frame-work for Data Quality and Participant Protection,” *Big Data Cogn. Comput.*, vol. 8, no. 7, 2024.
- [22] S.K. Jha and R.K. Jha, ”Security and privacy of biometric data in smart healthcare,” *Journal of Information Security and Applications*, vol. 77, 2023.
- [23] D.L. Johnson, ”A Quantitative Study on User Acceptance for Biometric Privacy, Security, and Challenges in Data Collection and Storage,” *ProQuest Dissertations Publishing*, 2022.
- [24] IEEE, ”IEEE Global Initiative on Ethics of Autonomous and Intelligent Systems.” [Online]. Available: [15]
- [25] sustainability Directory, ”Algorithmic Bias in Emotion Recognition,” 2025. [Online]. Available: [19].