

# Adaptive XAI in High Stakes Environments: Modeling Swift Trust with Multimodal Feedback in Human AI Teams

Nishani Fernando<sup>1,\*</sup>, Bahareh Nakisa<sup>1,\*</sup>, Adnan Ahmad<sup>1,\*</sup> and Mohammad Naim Rastgoo<sup>2,\*</sup>

<sup>1</sup>Deakin University, Geelong, Victoria, Australia

<sup>2</sup>Monash University, Melbourne, Victoria, Australia

## Abstract

Effective human-AI teaming heavily depends on swift trust, particularly in high-stakes scenarios such as emergency response, where timely and accurate decision-making is critical. In these time-sensitive and cognitively demanding settings, adaptive explainability is essential for fostering trust between human operators and AI systems. However, existing explainable AI (XAI) approaches typically offer uniform explanations and rely heavily on explicit feedback mechanisms, which are often impractical in such high-pressure scenarios. To address this gap, we propose a conceptual framework for adaptive XAI that operates non-intrusively by responding to users' real-time cognitive and emotional states through implicit feedback, thereby enhancing swift trust in high-stakes environments. The proposed adaptive explainability trust framework (AXTF) leverages physiological and behavioral signals, such as EEG, ECG, and eye tracking, to infer user states and support explanation adaptation. At its core is a multi-objective, personalized trust estimation model that maps workload, stress, and emotion to dynamic trust estimates. These estimates guide the modulation of explanation features enabling responsive and personalized support that promotes swift trust in human-AI collaboration. This conceptual framework establishes a foundation for developing adaptive, non-intrusive XAI systems tailored to the rigorous demands of high-pressure, time-sensitive environments.

## Keywords

Adaptive Explainability, Human-Machine Teams, Swift Trust, Implicit Feedback, Affective Interaction, Dynamic Environments

## 1. Introduction

In high stakes domains such as emergency response [1] and military operations [2], human AI teams are often formed on the fly and operate under extreme time pressure, high cognitive workload, and rapidly evolving situational demands. These environments are characterized by rapid decision-making, elevated emotional intensity, and limited opportunities for explicit communication or coordination. Failures in such contexts can lead to significant safety, ethical, or operational consequences [3]. As a result, effective human-AI teaming in such scenarios hinges on the development of swift trust and the ability to support human operators through adaptive, context-sensitive system behavior. Swift trust, originally introduced in the context of temporary human teams [4], describes a form of trust that arises rapidly out of necessity, without the benefit of prolonged interaction or prior history. In high-stakes environments, humans are often compelled to place immediate trust in AI systems simply because there is no time to build it gradually. However, this initial trust is fragile and often vulnerable to performance errors, lack of transparency and high cognitive load [5]. Sustaining trust in such conditions demands that AI systems must be capable of communicating effectively and adapting responsively to the human's evolving cognitive and emotional state.

Explainability has emerged as a central mechanism for cultivating trust in AI, enabling humans to understand and anticipate AI behavior [6]. However, existing explainable AI (XAI) approaches [7] are often static and uniform, providing generic explanations that overlook situational awareness and fail to adapt to the dynamic nature of the environment, which directly influences the user's real-time cognitive and emotional state. Furthermore, these approaches typically rely on explicit human feedback,

---

MAI-XAI'25: Workshop on Multimodal, Affective and Interactive Explainable AI October 25–26, 2025, Bologna, Italy

✉ nlfernando11@gmail.com (N. Fernando); bahar.nakisa@deakin.edu.au (B. Nakisa); adnan.a@deakin.edu.au (A. Ahmad); naim.rastgoo@monash.edu (M. N. Rastgoo)



© 2025 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

such as verbal queries or stated preferences, which are often impractical in high-pressure, cognitively demanding environments where users are overloaded and time is constrained. Therefore, more advanced explainable systems are essential for high-stakes environments. Such systems must be capable of rapid adaptation not only to the human operator’s state but also to contextual variables like task urgency, system reliability, and environmental uncertainty.

To overcome these gaps, incorporating implicit human feedback is essential for advancing explainable AI (XAI) systems, especially in high-stakes environments. Non-invasive technologies, such as wearable sensors [8], provide a promising means of capturing physiological and behavioral signals that reflect a user’s internal state. These signals may include Electroencephalography (EEG) [8, 9], Electrocardiography (ECG) [10], and eye tracking [11, 12], serving as real-time proxies for trust, cognitive workload, and emotional state. However, effective explanation adaptation must also account for AI system performance and situational context, which collectively shape human-AI trust dynamics. A comprehensive, adaptive XAI framework must therefore integrate these diverse signals to provide personalized, context-aware support.

This work presents adaptive explainability trust framework (AXTF), a conceptual framework designed to advance human-AI teaming in high-stakes, time-sensitive domains by enabling adaptive, non-intrusive explainability driven by multi-objective trust estimation model. It outlines a foundational approach that combines implicit human feedback, AI performance metrics, and situational awareness to infer the evolving trust state of the user. At the core of this framework is a personalized trust inference model that integrates the user’s cognitive and emotional state along with situational awareness to infer dynamic trust levels. These estimates guide the adaptive modulation of explanation features such as timing, granularity, content, and presentation mode, enabling dynamic, context-aware explanation strategies that foster swift trust. This conceptual approach lays the foundation for future research and the practical development of trust-sensitive, non-intrusive XAI systems tailored to the demands of time-critical, high-pressure domains.

Unlike task specific or opaque AI models [13], our framework is generalizable across high stakes domains, interpretable by design, and adaptable in real time. It supports collaboration by recognizing the cognitive and affective constraints of human operators and responding accordingly closing the loop between trust inference, explanation adaptation, and mission performance. In the following sections, we review prior work, detail our conceptual model, and outline its application to real-world high-pressure settings such as emergency response.

## **2. Background and Related Work**

This section introduces key foundations for our proposed framework by synthesizing prior work across four core areas. First, we examine the role of implicit feedback in high-stakes human-AI collaboration, emphasizing the need for implicit, real-time cues such as physiological and behavioral signals to support decision-making under stress and cognitive load. Next, we explore the construct of swift trust, outlining its determinants, reliability, predictability, competence, transparency, and adaptability and their sensitivity to fluctuating user states. We then review the limitations of adaptive explainability, highlighting gaps in existing XAI systems that fail to adjust explanations to changing user conditions. Finally, we discuss user-centric and affect-aware XAI, emphasizing emerging evidence for modeling trust through real-time physiological inference, and outlining the need for integrated models that dynamically adapt explanation features to maintain trust and cognitive efficiency in time-sensitive contexts. This background frames the motivation and design of our proposed conceptual model.

### **2.1. Human AI Collaboration and the Role of Implicit Feedback**

In high stakes domains, human AI collaboration is often task critical, requiring both agents to operate in close coordination under intense cognitive and temporal pressure. The effectiveness of this teaming rests heavily on the human operator’s ability to trust the AI system to understand its role, predict its behavior, and rely on its outputs in moments of uncertainty. Studies across domains such as emergency

response and autonomous operations show that trust in AI systems enhances decision making efficiency, reduces cognitive burden, and improves overall team performance [5, 14, 15]. However, traditional models of trust formation often assume explicit communication between humans and AI agents, such as requests for clarification, preference adjustment, or corrective feedback. In practice, such explicit feedback is limited or infeasible in high stakes situations [16]. Human operators are typically focused on the task at hand, operating under cognitive overload, and have minimal capacity to verbally assess or tune their interaction with the AI system. This constraint necessitates an alternative trust support mechanism, one that can function implicitly, adaptively, and in real time.

Implicit feedback, measured through physiological and behavioral signals, offers a promising foundation for adaptive support in human-machine teams (HMTs). Unlike explicit feedback, implicit indicators are passively observable and continuous, providing a non-intrusive method for assessing the user's cognitive and emotional state. A growing body of research supports the viability of using such signals for trust estimation. For instance, EEG has been shown to correlate with cognitive load and attention [3]. Similarly, ECG and galvanic skin response (GSR) are reliable indicators of physiological arousal and stress, which have been linked to trust erosion under pressure [17, 8, 18, 19, 20, 21]. Moreover, gaze patterns, facial expressions, and vocal features reflect emotional valence and engagement, serving as real-time proxies for user affect and trust [22]. These findings suggest that trust-relevant mental states can be inferred in real time from sensor data, enabling AI systems to detect when a user is confused, overloaded, disengaged, or stressed without requiring explicit articulation.

By grounding our model in these implicit signals, we enable AI systems to dynamically assess the operator's cognitive and emotional state and adapt their behavior accordingly. In the context of XAI, this means tailoring explanations to be more concise, timely, or expressive depending on the inferred user state. For example, a spike in physiological arousal following an AI action may indicate confusion or concern, prompting the system to proactively issue a clarifying explanation. Similarly, indicators of high trust and low load might invite more detailed, exploratory explanations to support learning or calibration. In this way, implicit feedback enables real time, user sensitive adaptation, forming a critical bridge between human trust dynamics and machine explainability. It allows the AI system to act as an responsive teammate not just explaining what it does, but choosing when and how to explain based on the operator's needs. This perspective forms the foundation for the next section, where we examine the elements of swift trust, and how they intersect with user state and explainability in high stake team settings.

## **2.2. Swift Trust and Its Determinants**

In high stakes human AI collaboration, trust must be formed quickly often in the absence of prolonged interaction or past performance history. This phenomenon, referred to as swift trust [4], is essential for enabling rapid coordination in dynamic environments such as disaster response or critical medical care. Unlike traditional trust, which emerges gradually through relationship building, swift trust is assumed provisionally based on contextual cues like system role, professionalism, and perceived competence. However, swift trust is inherently fragile. It can erode rapidly when system behavior is unclear, inconsistent, or perceived as unreliable. Maintaining and calibrating this trust is a non-trivial challenge especially under conditions of high workload and emotional strain, where human perception of system behavior becomes volatile. A breakdown in trust can lead to over reliance (complacency) or under reliance (disuse), both of which are detrimental to team performance.

A large body of empirical work (e.g.[14, 23, 15]) confirms that trust in AI systems is closely linked to perceived performance, system reliability and predictability, as well as the user's workload, stress, and emotional state. Hancock et al. [5] found that performance was the strongest predictor of trust, with workload and environmental risk also contributing significantly. More recent work shows that stress and cognitive overload reduce perceived trust, while positive emotional states such as engagement promote confidence and trust [3, 24, 15]. These variables are not only correlational but causal, influencing how operators interpret and respond to AI recommendations in real time. For example, when workload is high and system behavior is unclear, trust may drop even if the AI is performing correctly. Conversely,

under calm conditions with transparent AI behavior, trust can remain stable even after minor failures. Trust, therefore, operates as a feedback variable, modulated by both system performance and the human’s cognitive and emotional state. Hoff and Bashir’s [25] three-level trust model highlights that initial trust or disposition to trust is shaped by individual traits, prior experiences, and cultural factors, serving as a baseline for interaction with automation systems. While these dispositional influences are important, our work focuses on the dynamic adaptation of trust during interaction, particularly how AI systems can respond to evolving cognitive and emotional states to support trust formation and calibration in high-stakes environments.

### Key Elements of Swift Trust

To model and support swift trust effectively, it is useful to decompose it into key elements, as outlined in Table 1 and commonly cited in the literature [25, 26, 16]. Among these, adaptability plays a critical role in high-stakes settings where user state and task demand shift rapidly. It reflects the AI system’s responsiveness to physiological, cognitive, and contextual signals, allowing it to adjust its explanations (e.g., simplifying content during stress) and behaviors (e.g., increasing feedback frequency during uncertainty) to maintain trust. As Cho et al. [26] and Seong and Bisantz [27] note, adaptive systems promote more accurate and timely trust calibration, allowing the user to rapidly align trust with the actual performance of AI.

**Table 1**  
Key elements influencing trust in AI systems

Element	Definition
Reliability	Perceived consistency and dependability of the AI actions
Competence	Perceived skill or ability of the AI to complete its task
Predictability	Operator’s ability to anticipate AI behavior based on context
Transparency	Clarity in how and why the AI makes decisions
Adaptability	The ability of AI to adjust its behavior and output based on evolving user states and task demands

These elements are not static; rather, they are dynamically influenced by both user state and system behavior. An effective trust-supporting system must monitor changes in stress, workload, and emotional valence and adjust its communication strategy accordingly. Further, these elements can be modulated by explainability features, but only if the system is responsive to the underlying user state. For instance, low predictability can be improved by proactive explanations. Low transparency can be mitigated with “why” explanations about intent. Low reliability perception under stress may be best addressed with brief confidence statements (e.g., “High certainty: obstacle detected”) [16]. In the next section, we explore how explainability mechanisms can be leveraged to modulate these elements and support the dynamic formation and maintenance of swift trust.

### 2.3. Adaptive Explainability for Trust Formation

Explainability has long been recognized as a mechanism for fostering trust in AI systems. However, most existing approaches are static and context agnostic, offering fixed explanations that do not adjust to the user’s state or task environment. Recent advances have begun to explore adaptive explanation strategies for instance, using reinforcement learning or partially observable Markov decision processes (POMDPs) to tailor explanations based on user type or task progression [28]. However, these approaches often rely on predefined user profiles or require explicit feedback, making them difficult to apply in high stakes, real time environments. Model reconciliation approaches [28] align AI explanations with human mental models, but typically assume static trust misalignment and do not account for fluctuating physiological or affective states. Floyd et al. [29] introduced trust guided transparency, where the AI

modulates its behavior based on estimated user trust, but their system relied on explicit interaction logs and performance scores, rather than physiological signals.

Recent research has demonstrated that trust related cognitive and emotional states can be inferred using physiological and behavioral signals such as heart rate variability (HRV), electrodermal activity (EDA), facial expressions, and gaze [22, 3]. Fuzzy and neuro fuzzy models have been employed to classify trust states in real time, providing interpretable trust metrics for adaptive systems. However, these models have rarely been connected to explanation generation, leaving a gap in integrating trust estimation with communication behavior.

Moreover, trust calibration aligning user trust with system capability is especially critical in high risk or time sensitive domains. Studies in aviation, medicine, and robotics have shown that miscalibrated trust leads to automation bias or disuse [5, 30]. Meanwhile, cognitive load plays a central role in explainability: too much detail can overwhelm the user, while too little can induce confusion or mistrust. Paleja et al. [15] found that tailoring explanation granularity benefits novice users under load but may frustrate experts, reinforcing the need for adaptive strategies that consider real time workload and user expertise.

## 2.4. User Centric and Affective XAI

Early work on user aware XAI [13] focused on clustering users by behavior patterns to personalize explanations. However, most approaches lack real time responsiveness, and few incorporate affective signals to dynamically adjust content. Ali et al. [7] emphasize that explanations tailored to user context and emotional state enhance perceived competence and empathy, but current systems lack the infrastructure to connect affective inference with explanation logic. As summarized in Table 2, existing work highlights the individual importance of trust, explainability, and physiological modeling. However, few frameworks bring these together into a cohesive, real time trust adaptive explainability model that operates effectively in high stakes human AI teams.

**Table 2**

Trust-relevant user states and corresponding explanation adaptation strategies

Trust-Relevant Factor	Empirical Insight	Explanation Adaptation Rule
High cognitive load	High workload reduces trust and task performance [5, 15]	Short, high-level summaries preferred.
High stress	Increased stress is negatively correlated with trust in automated systems[31]	Proactive, calming explanations delivered early.
Low emotional valence	Negative valence correlates with reduced trust [22]	Use reactive, empathetic tone or voice modality.
Low AI performance	Lowers reliability and competence perceptions [14, 25]	Provide corrective, fallback or reassuring explanations to recover trust. Reinforce competence with confidence indicators.
Low predictability	Hinders user ability to anticipate system behavior [14]	Provide rationale (why/how) behind actions.
Unfamiliar task (Expertise)	Increases mental effort; risks trust erosion [15]	Simplify content, emphasize goal relevance with proactive and guiding explanations.
Disengagement	Reduced engagement can impact performance and trust levels [5, 16, 15]	Switch to visual/interactive formats to regain attention. Mode of delivery can influence information absorption and processing cost.

### 2.4.1. Trust Modeling through Explainability Cues

In Human Machine Teams (HMTs), trust and explainability are dynamically interlinked, unfolding as a sequence of cause-and-effect interactions. AI behaviors whether task related actions, feedback, or navigation decisions directly influence the human operator's physiological state, emotional response, and cognitive processing. These changes, in turn, shape the operator's perception of trust, impacting collaboration quality and task performance [6]. These suggest how real-time adaptation of AI explanations, grounded in these causal pathways, can mitigate cognitive and emotional strain while reinforcing trust in high stakes and time sensitive environments.

**AI Behavior and Trust Perception.** The observable behavior of AI such as decision making, task execution, or error handling shapes the operator's perception of its competence, reliability, and intent. When the AI behaves transparently and contextually, users are more likely to perceive it as trustworthy. In contrast, opaque or inconsistent behaviors introduce uncertainty and distrust. Shin [6] reports that causability and explainability account for over 58% of variance in trust perceptions, underscoring the importance of clarity and communicative alignment in fostering trust.

**Emotional and Physiological Responses.** Trust perceptions are further mediated by affective responses, which manifest physiologically through changes in heart rate (HR), heart rate variability (HRV), electrodermal activity, or neural activity (EEG) [32]. For example, unexpected or ambiguous AI actions may trigger stress or frustration, while cooperative and predictable behavior fosters engagement and calm. These physiological markers serve as real time, implicit indicators of trust state [16], enabling continuous user monitoring without explicit intervention.

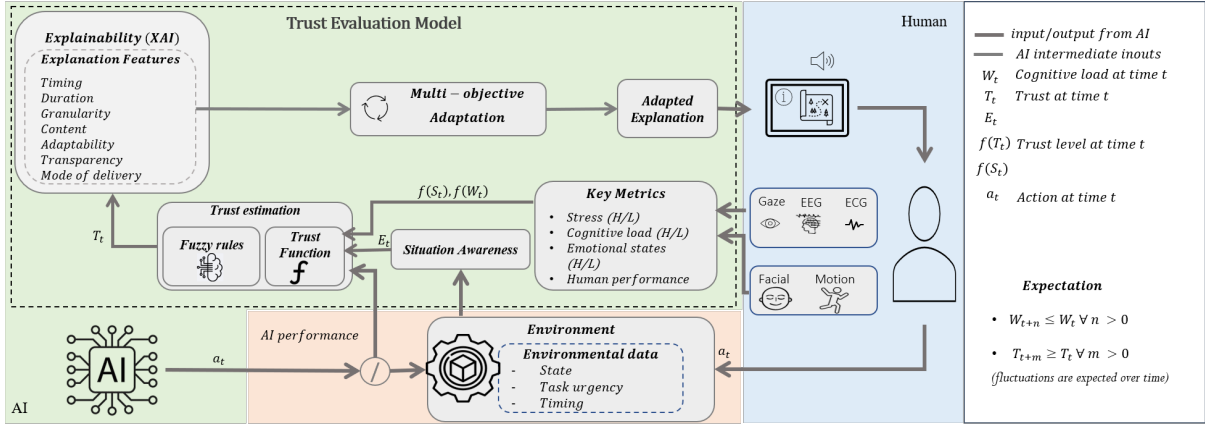
**Cognitive Load and Information Processing.** AI outputs that are complex, ambiguous, or mistimed can impose a high cognitive burden, impairing the user's ability to process information and make timely decisions. This effect is especially pronounced in emergency response scenarios, where attention is divided, time is limited, and errors are costly. Prior studies [15, 16] show that cognitive overload negatively affects both trust and performance in collaborative human AI systems. Adaptive explainability can address this by modulating explanation complexity, timing, and content helping reduce overload, maintain attention, and recalibrate trust in real time.

However, most XAI systems do not modulate their explanations based on implicit human signals such as stress, workload, or emotion, nor do they account for rapidly changing contextual cues. This limitation is particularly consequential in high-stakes environments, where cognitive overload and uncertainty diminish the operator's ability to process static or overly generic explanations. There remains a significant gap in designing closed-loop adaptive XAI systems that leverage real-time physiological and contextual data to both interpret user states and dynamically tailor explanation strategies. Such systems would not only reflect the user's cognitive and emotional state but also influence it over time supporting trust formation, cognitive efficiency, and collaborative resilience under pressure.

## 3. Adaptive Explainability Trust Framework

The findings presented in Section 2 lead to a critical insight: explainability is not merely a communication tool, but a dynamic mechanism for shaping and stabilizing swift trust. To support this, we propose, Adaptive Explainability Trust Framework (AXTF), a conceptual framework designed to enhance human decision-making within human-AI teams by improving swift trust and team performance through dynamic adaptation of AI explanations based on real-time assessments of user states and environmental factors. Specifically, our framework enables AI explanations to be continuously tailored according to the user's cognitive load, stress, and emotional state, as well as the task context, including urgency, goals, and environmental complexity, particularly in high-stakes domains such as emergency response.

To support swift trust and effective human AI teaming in high stakes environments, our proposed conceptual model links real time physiological and behavioral indicators of human state to adaptive explainability mechanisms. The model integrates three interconnected components: (1) multimodal feedback sensing and inference, (2) multi-objective trust modeling, and (3) explanation feature adaptation, with environmental inputs to ensure contextual relevance.



**Figure 1:** Adaptive Explainability Trust Framework (AXTF). The framework supports swift trust formation by dynamically adjusting explanation features—*timing, duration, and granularity*—based on the user’s cognitive load, emotions, and performance. This adaptive explainability reduces cognitive overload, enhances trust, and improves decision-making.

The proposed framework (Fig. 1) forms a closed-loop pipeline that integrates real-time physiological and behavioral signals (e.g., EEG, ECG, heart rate variability—HRV) with environmental data such as task goals, urgency, and state to assess the user’s cognitive load ( $W$ ), stress ( $S$ ), and emotional valence ( $E$ ) levels [15, 33]. Based on these user states and environmental performance metrics (e.g., task errors, success rates), dynamic trust estimation is performed using a multi-objective neurofuzzy rule-based inference engine. The framework then adapts key explanation features, including timing, duration, granularity, and mode of delivery, by mapping these trust estimates and contextual knowledge to reduce cognitive overload, enhance trust, and guide the subsequent human action ( $a_{t+1}$ ). For example, in scenarios characterized by low trust and high cognitive load, short, moderate steps, and reactive explanations are more effective for fostering trust compared to lengthy, detailed ones. While some fluctuations in cognitive load ( $W_t$ ) and trust ( $T_t$ ) are expected, the ultimate goal of these adaptive explanations is to increase trust ( $T$ ) and decrease cognitive load ( $W$ ) over time, thereby promoting swift trust and improving situational awareness, decision-making, and overall team performance. This feedback-driven framework is designed to support temporal situational awareness, workload balancing, and trust resilience in high-pressure environments where explicit communication may be limited, but implicit signals provide actionable insight. The following subsections detail each component of the framework, including multimodal feedback sensing and inference, multi-objective trust modeling, and explanation feature adaptation.

### 3.1. Real-Time Multimodal Inference of Human State

The system continuously monitors a range of physiological and behavioral signals to infer latent cognitive and emotional states that are critical for trust assessment, including:

- EEG and pupillometry for evaluating cognitive workload,
- ECG, galvanic skin response (GSR), and heart rate variability (HRV) to detect stress and arousal levels,
- Facial expressions, gaze tracking, and voice features to infer emotional valence and user engagement.

These raw signals undergo preprocessing and feature extraction pipelines before being classified into meaningful, interpretable states such as “High Workload” or “Low Valence” through specialized physiological inference models. By transforming complex biosignals into actionable, high-level indicators, the system enables downstream reasoning components to adapt interactions dynamically based on the user’s real-time cognitive and emotional state.

**Table 3**

Fuzzy rules for trust inference based on workload, stress, emotion, and performance

Interpretation	Fuzzy Rule
Captures: low competence and predictability	IF $W = \text{High}$ AND $S = \text{High}$ AND $E = \text{Negative}$ THEN $T = \text{Low}$
Captures: reliability and transparency	IF $S = \text{Low}$ AND $E = \text{Positive}$ AND $P > 0.8$ THEN $T = \text{High}$
Captures: competence and predictability	IF $W = \text{Low}$ AND $S = \text{Low}$ AND $E = \text{Neutral}$ THEN $T = \text{High}$
Captures: adaptability	IF $W = \text{Medium}$ AND $E = \text{Positive}$ AND $P > 0.6$ THEN $T = \text{Medium}$
Captures: reliability under pressure	IF $S = \text{High}$ AND $P < 0.4$ THEN $T = \text{Low}$
Captures: transparency degradation due to overload/affect	IF $W = \text{High}$ OR $E = \text{Negative}$ THEN $T = \text{Low}$
Captures: balanced state supporting trust	IF $W = \text{Medium}$ AND $S = \text{Medium}$ AND $E = \text{Positive}$ THEN $T = \text{High}$

### 3.2. Fuzzy Logic Based Trust Inference

To translate these human state metrics into actionable trust estimates, we suggest a multiobjective neurofuzzy inference method. This approach allows the encoding of literature grounded, rule-based mappings (see Table 3) into an interpretable decision layer.

Trust is estimated as a categorical variable with three levels Low, Medium, or High based on a combination of user physiological and affective states and system performance metrics. The model enables interpretable reasoning that links observed human state and AI behavior to well-established trust dimensions such as reliability, competence, predictability, transparency, and adaptability. The model takes the following inputs:

- **Input Variables**

- **Workload (W)**: categorized as Low, Medium, or High; derived from EEG features, gaze data, or behavioral task-switching patterns.
- **Stress level (S)**: categorized as Low, Medium, or High; inferred from ECG, GSR, and HRV metrics.
- **Emotion valence (E)**: categorized as Negative, Neutral, or Positive; estimated from facial expressions, tone of voice, or affective models.
- **System performance score (P)**: a normalized score between 0 and 1; computed from task metrics such as success rate and error frequency.

- **Trust Output**

- **Trust (T)**: classified as Low, Medium, or High.

The fuzzy trust inference system maps normalized input variables into fuzzy linguistic categories using membership functions. Each input (e.g., workload, stress, emotion valence, performance) is associated with three fuzzy sets: *Low*, *Medium*, and *High*, except for emotion valence, which uses *Negative*, *Neutral*, and *Positive*.

#### 1. Workload (W), Stress (S)

Both workload and stress are defined over the domain  $[0, 1]$  and share the same triangular membership structure:

$$\mu_{\text{Low}}(x) = \begin{cases} 1, & x \leq 0.2 \\ \frac{0.5-x}{0.3}, & 0.2 < x \leq 0.5 \\ 0, & x > 0.5 \end{cases} \quad \mu_{\text{Medium}}(x) = \begin{cases} 0, & x \leq 0.2 \text{ or } x \geq 0.8 \\ \frac{x-0.2}{0.3}, & 0.2 < x \leq 0.5 \\ \frac{0.8-x}{0.3}, & 0.5 < x < 0.8 \end{cases} \quad \mu_{\text{High}}(x) = \begin{cases} 0, & x \leq 0.5 \\ \frac{x-0.5}{0.3}, & 0.5 < x \leq 0.8 \\ 1, & x > 0.8 \end{cases}$$

## 2. Emotion Valence (E)

Emotion valence is modeled over the range  $[-1, 1]$ :

$$\mu_{\text{Negative}}(e) = \begin{cases} 1, & e \leq -0.5 \\ \frac{-e-0.1}{0.4}, & -0.5 < e \leq -0.1 \\ 0, & e > -0.1 \end{cases} \quad \mu_{\text{Neutral}}(e) = \begin{cases} 0, & |e| > 0.5 \\ \frac{0.5-|e|}{0.5}, & |e| \leq 0.5 \end{cases} \quad \mu_{\text{Positive}}(e) = \begin{cases} 0, & e < 0.1 \\ \frac{e-0.1}{0.4}, & 0.1 \leq e < 0.5 \\ 1, & e \geq 0.5 \end{cases}$$

## 3. System Performance (P)

System performance is a normalized score  $p \in [0, 1]$ :

$$\mu_{\text{Low}}(p) = \begin{cases} 1, & p \leq 0.3 \\ \frac{0.5-p}{0.2}, & 0.3 < p \leq 0.5 \\ 0, & p > 0.5 \end{cases} \quad \mu_{\text{Medium}}(p) = \begin{cases} 0, & p \leq 0.3 \text{ or } p \geq 0.7 \\ \frac{p-0.3}{0.2}, & 0.3 < p \leq 0.5 \\ \frac{0.7-p}{0.2}, & 0.5 < p < 0.7 \end{cases} \quad \mu_{\text{High}}(p) = \begin{cases} 0, & p < 0.5 \\ \frac{p-0.5}{0.3}, & 0.5 \leq p < 0.8 \\ 1, & p \geq 0.8 \end{cases}$$

These functions are derived from literature indicating that increased workload and stress reduce trust in automation [5, 16], while positive valence and higher performance promote trust. Mapping continuous inputs into interpretable fuzzy categories supports transparent, adaptable trust modeling in real time. The fuzzy trust inference model estimates trust levels based on these real-time assessments of workload, stress, emotional valence, and system performance. Using a set of fuzzy rules (Table 3) derived from empirical research and domain knowledge (see Table 1), the model captures complex interactions among these factors to produce a unified trust estimate  $T$  classified as Low, Medium, or High. This approach enables interpretable reasoning and supports adaptive AI behavior aligned with the operator's current cognitive-affective state, improving human-AI collaboration.

### 3.3. Trust Sensitive Explanation Adaptation

The model focuses on seven core explainability features that shape user experience and collaboration. We define these key features as follows:

1. **Timing:** Timely delivery of explanations is vital in high-pressure situations. Explanations can be delivered proactively, ahead of an AI action, to minimize potential confusion (e.g., "Avoiding unstable debris ahead"), or reactively, triggered by user hesitation or unexpected AI behavior. In an emergency context, timing must align with both the task phase and the user's attentional bandwidth to avoid distraction or delay [15, 16, 6].
2. **Duration:** The length of explanations should be carefully adapted to the time sensitivity of the situation and the user's cognitive capacity. Under high cognitive load or time pressure, brief explanations, typically lasting 2–3 seconds, are more effective in preserving situational focus and preventing distraction. On the other hand, when cognitive load is lower, longer or layered explanations can offer deeper insight without overwhelming the user [15, 16]. For instance, during search and rescue triage, the system may initially provide short verbal alerts, then follow up with optional elaboration once the situation stabilizes.
3. **Granularity:** Explanation granularity refers to the level of detail provided in the explanation. High-level summaries, such as "Scanning lower level first," help reduce the user's information processing demands. In contrast, detailed step-by-step explanations, for example, "Entering sector B → mapping → thermal anomaly detected," are better suited for experienced users or situations where trust is high. The granularity should be adapted based on factors such as user familiarity, workload, and trust levels [15, 30, 34] to ensure explanations remain cognitively accessible while still informative.

**Table 4**

XAI features aligned with trust and cognitive principles, and corresponding explanation strategies

XAI Feature	Supports Swift Trust & Cognitive Element(s)	Explanation Strategy Example
Timing	Predictability, Reliability	Provide proactive clarification during decision spikes or delays
Granularity	Competence, Cognitive Efficiency	Provide brief summary under high workload, detailed explanation when user trust is higher
Duration	Cognitive Efficiency, Predictability	Provide short explanations during time-critical phases, extended versions in low-pressure phases
Delivery Mode	Transparency, Engagement, Clarity	Switch between text, voice, or visual modes based on user attentiveness, emotional state, or context. Combine multiple modalities such as voice and visual to deliver context for enhanced information processing.
Content	Reliability, Context Sensitivity	Use local (contextual) explanations in dynamic settings, hierarchical in stable scenarios
Transparency	Transparency, Competence	Explain “why” decisions are made to improve model clarity and intention transparency
Adaptability	Adaptability, Personalization, Trust Calibration	Tailor content and modality dynamically based on inferred trust and workload signals

4. **Content:** Explanation content is chosen based on task relevance and the user’s current focus. Contextual or local explanations, such as “Rerouting due to obstacle,” emphasize immediate actions or environmental conditions. In contrast, hierarchical explanations, for example, “Prioritizing lower floors due to heat signature density,” communicate broader planning strategies. In emergency scenarios, providing contextual content enhances temporal situational awareness and responsiveness [35, 27].
5. **Transparency:** Transparency shapes the user’s understanding of the AI’s decision-making process. “How” transparency, such as “Based on heatmap and terrain risk, path updated,” helps users evaluate the system’s methods. “Why” transparency, for example, “Avoiding risk to maximize coverage,” clarifies the AI’s intent and goal reasoning. Both types support mental model alignment and trust; however, the appropriate level and form of transparency should be adapted based on the user’s state and the context [36, 25, 29].
6. **Adaptability:** Adaptability acts as the central mechanism that dynamically modulates all other explainability features in real time. This capability allows AI systems to selectively tailor explanations to align with the user’s current state and task objectives. Under conditions of high workload or stress, the system simplifies explanations and employs lower-detail modes to reduce cognitive load. Conversely, when users are calm and trust levels are high, the system can provide more complex and interactive explanations. This adaptability ensures that explanations remain both informative and sustainable, especially in high-pressure environments [37, 25, 29, 13].
7. **Mode of Delivery:** The medium used to deliver explanations significantly influences user comprehension and cognitive load. Visual formats such as maps, trajectory overlays, or alert icons are ideal when the user’s auditory attention is available. Textual explanations, including on-screen summaries or status updates, are more suitable during quieter moments or post-task phases. Auditory delivery through spoken instructions works best when the user’s hands and eyes are occupied. Combining multiple channels in a multimodal approach—such as spoken plus visual alerts—enhances system resilience and inclusivity. Research by Adadi and Berrada [38] demonstrates that multimodal explanations improve user understanding, particularly when users are multitasking or experiencing stress.

The inferred trust level is used to modulate these key explainability features. This trust adaptive mechanism supports continuous calibration, allowing the system to respond not just to performance, but to how the human feels and functions during collaboration.

While explainability is often discussed as a means of improving user understanding or meeting regulatory requirements, in high stakes human AI teaming, it serves a deeper role: calibrating trust in real time. As trust is highly sensitive to changes in workload, stress, and affect, static or misaligned explanations may unintentionally erode confidence or overload the user. In contrast, adaptive explainability tuned to the user’s current cognitive and emotional state can serve as a powerful tool for trust repair, reinforcement, and regulation. Consider a search and rescue drone system operating in post-disaster environments. A human operator under high stress and workload may be overwhelmed by frequent decision updates. The system detects high HRV, elevated EEG load, and negative valence, infers low trust, and adapts by issuing short, confidence-framed explanations through audio (“Clear path detected. High certainty.”). If later states show reduced load and increased engagement, it shifts to detailed, interactive visualizations for planning and collaboration.

## 4. Future Directions

While this work establishes a conceptual foundation for adaptive explainability in high-stakes human-AI teaming, several important avenues remain for future investigation.

First, two-way communication between human and AI teammates should be more explicitly modeled not only to adapt AI explanations to user states, but also to enable reciprocal influence where the AI dynamically adjusts its behavior based on user reactions and evolving task demands. Such bi-directional interaction will allow the AI to both respond to human needs and evaluate and refine its own performance, closing the loop between perception, explanation, and behavior. Second, cultural and dispositional factors play a foundational role in trust formation but remain underexplored in adaptive XAI. Future research should investigate how traits such as uncertainty avoidance, communication preferences, and prior experience influence trust dynamics and explanation preferences, enabling more inclusive and culturally-aware explanation strategies. Third, implementation and evaluation in interactive, dynamic environments, such as simulation-based emergency response scenarios, will be essential to validate the framework. These settings offer controllable, high-fidelity contexts for assessing how real-time physiological and behavioral feedback impacts explanation effectiveness, trust calibration, and team performance under pressure. Fourth, advances in generative agent simulations [39] open opportunities for large-scale validation using synthetic populations embedded with memory, social reasoning, and behavioral diversity. These agent-based testbeds can be used to examine long-term trust trajectories and cross-profile adaptation strategies in simulated high-stakes team settings. Finally, to ensure AI safety in high-pressure decision environments, future work must incorporate safeguards that prevent explanation misuse or cognitive overload. Adaptive explanation systems must remain transparent, interpretable, and bounded by safety constraints that prevent miscalibration of trust—especially under uncertainty or stress. Embedding safety-aware logic into adaptation rules (e.g., thresholds on explanation complexity or delivery timing) will help maintain alignment with human cognitive capacity, trust boundaries, and ethical standards in mission-critical operations.

In conclusion, advancing adaptive explainability through bi-directional interaction, cultural awareness, scalable simulation, and embedded safety principles will help realize the next generation of trust-sensitive, cognitively aligned, and ethically grounded human-AI systems.

## 5. Conclusion and Contribution

This work contributes to affective, situated, and trustworthy AI by introducing a conceptual framework for real-time adaptation of explainability to support swift trust and effective teamwork in high-stakes environments. The framework integrates multimodal implicit feedback physiological, behavioral, environmental, and contextual signals to infer user states such as workload, stress, and emotional valence. These inferred states inform the dynamic adjustment of explanation features (e.g., timing, granularity, modality), enabling alignment with the user’s cognitive and affective demands in pursuit of time-critical task goals. Explainability is reframed as a multi-objective adaptive function balancing

transparency, cognitive efficiency, and trust calibration. The framework is designed to be model-agnostic and extensible, supporting the integration of diverse trust modeling and learning mechanisms. It enables AI systems to act as responsive teammates fostering trust, maintaining collaboration under pressure, and supporting decision-making when explicit communication is limited. This lays a foundation for generalizable real-world deployment of adaptive XAI in high-stakes domains such as emergency response, medical operations, and mission-critical decision support, where human-machine teaming must remain transparent, affect-aware, and cognitively efficient under uncertainty.

## **Declaration on Generative AI**

During the preparation of this work, the author(s) used GPT-4 in order to: Grammar and spelling check. After using these tool(s)/service(s), the author(s) reviewed and edited the content as needed and take(s) full responsibility for the publication's content.

## References

- [1] Y. X. Lim, Cognitive Human-Machine Interfaces and Interactions for Avionics Systems, PhD thesis, RMIT University, 2021. URL: <https://doi.org/10.25439/rmt.27601791>. doi:10.25439/rmt.27601791.
- [2] F. Dehais, A. Duprès, S. Blum, N. Drougard, S. Scannella, R. N. Roy, F. Lotte, Monitoring pilot's mental workload using erps and spectral power with a six-dry-electrode eeg system in real flight conditions, *Sensors (Basel, Switzerland)* 19 (2019). URL: <https://api.semanticscholar.org/CorpusID:83462067>.
- [3] L. Rodriguez Rodriguez, C. E. Bustamante Orellana, E. K. Chiou, L. Huang, N. Cooke, Y. Kang, A review of mathematical models of human trust in automation, *Frontiers in Neuroergonomics Volume 4 - 2023* (2023). URL: <https://www.frontiersin.org/journals/neuroergonomics/articles/10.3389/fnrgo.2023.1171403>. doi:10.3389/fnrgo.2023.1171403.
- [4] D. Meyerson, K. E. Weick, R. M. Kramer, *Swift trust and temporary groups*, Sage Publications, Inc, Thousand Oaks, CA, US, 1996, pp. 166–195. URL: <https://doi.org/10.4135/9781452243610.n9>. doi:10.4135/9781452243610.n9.
- [5] P. A. Hancock, D. R. Billings, K. E. Schaefer, J. Y. C. Chen, E. J. de Visser, R. Parasuraman, A meta-analysis of factors affecting trust in human-robot interaction, *Human Factors* 53 (2011) 517–527. URL: <https://journals.sagepub.com/doi/abs/10.1177/0018720811417254>. doi:10.1177/0018720811417254.
- [6] D. Shin, The effects of explainability and causability on perception, trust, and acceptance: Implications for explainable ai, *International Journal of Human-Computer Studies* 146 (2021) 102551. URL: <https://www.sciencedirect.com/science/article/pii/S1071581920301531>. doi:<https://doi.org/10.1016/j.ijhcs.2020.102551>.
- [7] S. Ali, T. Abuhmed, S. El-Sappagh, K. Muhammad, J. M. Alonso-Moral, R. Confalonieri, R. Guidotti, J. Del Ser, N. Díaz-Rodríguez, F. Herrera, Explainable artificial intelligence (xai): What we know and what is left to attain trustworthy artificial intelligence, *Information Fusion* 99 (2023) 101805. URL: <https://www.sciencedirect.com/science/article/pii/S1566253523001148>. doi:<https://doi.org/10.1016/j.inffus.2023.101805>.
- [8] K. Akash, W.-L. Hu, N. Jain, T. Reid, A classification model for sensing human trust in machines using eeg and gsr, *ACM Trans. Interact. Intell. Syst.* 8 (2018) Article 27. URL: <https://doi.org/10.1145/3132743>. doi:10.1145/3132743.
- [9] S. Choo, C. Nam, Detecting human trust calibration in automation: A convolutional neural network approach, 2022. URL: <https://research.ebsco.com/linkprocessor/plink?id=adc76bfc-bd10-3828-ad61-8d1d02d56aab>. doi:10.1109/THMS.2021.3137015.
- [10] A. Ahmad, B. Nakisa, M. N. Rastgoo, Robust emotion recognition via bi-level self-supervised continual learning, 2025. URL: <https://arxiv.org/abs/2505.10575>.
- [11] C. Hu, S. Huang, Y. Zhou, S. Ge, B. Yi, X. Zhang, X. Wu, Dynamic and quantitative trust modeling and real-time estimation in human-machine co-driving process, *Transportation Research Part F: Traffic Psychology and Behaviour* 106 (2024) 306–327. URL: <https://www.sciencedirect.com/science/article/pii/S1369847824002006>. doi:<https://doi.org/10.1016/j.trf.2024.08.001>.
- [12] N. Hulle, S. Aroca-Ouellette, A. J. Ries, J. Brawer, A. Roncone, Eyes on the game: Deciphering implicit human signals to infer human proficiency, trust, and intent, *arXiv.org* (2024). URL: <https://arxiv.org/abs/2407.03298>.
- [13] U. Soni, S. Sreedharan, S. Kambhampati, Not all users are the same: Providing personalized explanations for sequential decision making problems, in: 2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), 2021, pp. 6240–6247. URL: <https://doi.org/10.1109/IROS51168.2021.9636331>. doi:10.1109/IROS51168.2021.9636331.
- [14] S. Milivojevic, M. Sobhani, N. Webb, Z. Madin, J. Ward, S. Yusuf, C. Baber, E. R. Hunt, Swift trust in mobile ad hoc human-robot teams, 2024. URL: <https://doi.org/10.1145/3686038.3686057>. doi:10.1145/3686038.3686057.
- [15] R. Paleja, M. Ghuy, N. R. Arachchige, R. Jensen, M. Gombolay, The utility of explainable ai

- in ad hoc human-machine teaming, 2021. URL: <https://dl.acm.org/doi/10.5555/3540261.3540308>, <https://github.com/CORE-Robotics-Lab/Utility-of-Explainable-AI-NeurIPS2021>.
- [16] M. R. Endsley, Supporting human-ai teams: transparency, explainability, and situation awareness, *Computers in Human Behavior* 140 (2023) 107574. URL: <https://www.sciencedirect.com/science/article/pii/S0747563222003946>. doi:<https://doi.org/10.1016/j.chb.2022.107574>.
  - [17] F. Shaffer, J. P. Ginsberg, An overview of heart rate variability metrics and norms, *Front Public Health* 5 (2017) 258. doi:[10.3389/fpubh.2017.00258](https://doi.org/10.3389/fpubh.2017.00258), 2296-2565 Shaffer, Fred Ginsberg, J P Journal Article Review Switzerland 2017/10/17 Front Public Health. 2017 Sep 28;5:258. doi: 10.3389/fpubh.2017.00258. eCollection 2017.
  - [18] H. N. Green, T. Iqbal, Using physiological measures, gaze, and facial expressions to model human trust in a robot partner, 2025. URL: <https://arxiv.org/abs/2504.05291>.
  - [19] A. Aygun, H. Ghasemzadeh, R. Jafari, Robust interbeat interval and heart rate variability estimation method from various morphological features using wearable sensors, *IEEE J Biomed Health Inform* 24 (2020) 2238–2250. doi:[10.1109/jbhi.2019.2962627](https://doi.org/10.1109/jbhi.2019.2962627), 2168-2208 Aygun, Ayca Ghasemzadeh, Hassan Jafari, Roozbeh R01 EB028106/EB/NIBIB NIH HHS/United States Journal Article Research Support, N.I.H., Extramural Research Support, U.S. Gov't, Non-P.H.S. United States 2020/01/04 IEEE J Biomed Health Inform. 2020 Aug;24(8):2238-2250. doi: 10.1109/JBHI.2019.2962627. Epub 2019 Dec 27.
  - [20] M. N. Rastgoo, N. Bahareh, R. Andry, M. Frederic, , V. Chandran, Driver stress levels detection system using hyperparameter optimization, *Journal of Intelligent Transportation Systems* 28 (2024) 443–458. URL: <https://doi.org/10.1080/15472450.2022.2140046>. doi:[10.1080/15472450.2022.2140046](https://doi.org/10.1080/15472450.2022.2140046), doi: 10.1080/15472450.2022.2140046.
  - [21] B. Nakisa, M. N. Rastgoo, D. Tjondronegoro, V. Chandran, Evolutionary computation algorithms for feature selection of eeg-based emotion recognition using mobile sensors, *Expert Systems with Applications* 93 (2018) 143–155. URL: <https://www.sciencedirect.com/science/article/pii/S0957417417306747>. doi:<https://doi.org/10.1016/j.eswa.2017.09.062>.
  - [22] H. M. Khalid, L. W. Shiung, P. Nooralishahi, Z. Rasool, M. G. Helander, L. C. Kiong, C. Ai-vyrn, Exploring psycho-physiological correlates to trust: implications for human-robot-human interaction, *Proceedings of the Human Factors and Ergonomics Society Annual Meeting* 60 (2016) 697–701. URL: <https://journals.sagepub.com/doi/abs/10.1177/1541931213601160>. doi:[10.1177/1541931213601160](https://doi.org/10.1177/1541931213601160).
  - [23] S. C. Kohn, E. J. de Visser, E. Wiese, Y. C. Lee, T. H. Shaw, Measurement of trust in automation: A narrative review and reference guide, *Front Psychol* 12 (2021) 604977. URL: [https://libkey.io/libraries/136/articles/502135231/content-location?utm\\_source=nomad](https://libkey.io/libraries/136/articles/502135231/content-location?utm_source=nomad). doi:[10.3389/fpsyg.2021.604977](https://doi.org/10.3389/fpsyg.2021.604977), 1664-1078 Kohn, Spencer C de Visser, Ewart J Wiese, Eva Lee, Yi-Ching Shaw, Tyler H Journal Article Review Switzerland 2021/11/06 Front Psychol. 2021 Oct 19;12:604977. doi: 10.3389/fpsyg.2021.604977. eCollection 2021.
  - [24] B. Sadrifaridpour, H. Saeidi, Y. Wang, J. Burke, Modeling and control of trust in human and robot collaborative manufacturing, *AAAI Spring Symposium - Technical Report* (2014) 64–70. doi:[10.1007/978-1-4899-7668-0\\_7](https://doi.org/10.1007/978-1-4899-7668-0_7).
  - [25] K. A. Hoff, M. Bashir, Trust in automation: integrating empirical evidence on factors that influence trust, *Human Factors* 57 (2015) 407–434. URL: <https://journals.sagepub.com/doi/abs/10.1177/0018720814547570>. doi:[10.1177/0018720814547570](https://doi.org/10.1177/0018720814547570).
  - [26] J.-H. Cho, K. Chan, S. Adali, A survey on trust modeling, *ACM Comput. Surv.* 48 (2015) Article 28. URL: <https://doi.org/10.1145/2815595>. doi:[10.1145/2815595](https://doi.org/10.1145/2815595).
  - [27] Y. Seong, A. M. Bisantz, The impact of cognitive feedback on judgment performance and trust with decision aids, *International Journal of Industrial Ergonomics* 38 (2008) 608–625. URL: <https://www.sciencedirect.com/science/article/pii/S0169814108000279>. doi:<https://doi.org/10.1016/j.ergon.2008.01.007>.
  - [28] S. Sreedharan, T. Chakraborti, S. Kambhampati, Handling model uncertainty and multiplicity in explanations via model reconciliation, *Proceedings of the International Conference on Automated Planning and Scheduling* 28 (2018) 518–526. URL: <https://ojs.aaai.org/index.php/ICAPS/article/>

view/13930. doi:10.1609/icaps.v28i1.13930.

- [29] M. W. Floyd, D. W. Aha, Incorporating transparency during trust-guided behavior adaptation, in: A. Goel, M. B. Díaz-Agudo, T. Roth-Berghofer (Eds.), *Case-Based Reasoning Research and Development*, Springer International Publishing, 2022, pp. 124–138.
- [30] P. Bobko, L. Hirshfield, L. Eloy, C. Spencer, E. Doherty, J. Driscoll, H. Obolsky, Human-agent teaming and trust calibration: a theoretical framework, configurable testbed, empirical illustration, and implications for the development of adaptive systems, *Theoretical Issues in Ergonomics Science* 24 (2023) 310–334. URL: <https://doi.org/10.1080/1463922X.2022.2086644>. doi:10.1080/1463922X.2022.2086644.
- [31] S. C. Kohn, E. J. de Visser, E. Wiese, Y. C. Lee, T. H. Shaw, Measurement of trust in automation: A narrative review and reference guide, *Front Psychol* 12 (2021) 604977. URL: [https://libkey.io/libraries/136/articles/502135231/content-location?utm\\_source=nomad](https://libkey.io/libraries/136/articles/502135231/content-location?utm_source=nomad). doi:10.3389/fpsyg.2021.604977, 1664-1078 Kohn, Spencer C de Visser, Ewart J Wiese, Eva Lee, Yi-Ching Shaw, Tyler H Journal Article Review Switzerland 2021/11/06 Front Psychol. 2021 Oct 19;12:604977. doi: 10.3389/fpsyg.2021.604977. eCollection 2021.
- [32] L. Mou, Y. Zhao, C. Zhou, B. Nakisa, M. N. Rastgoo, L. Ma, T. Huang, B. Yin, R. Jain, W. Gao, Driver emotion recognition with a hybrid attentional multimodal fusion framework, *IEEE Transactions on Affective Computing* 14 (2023) 2970–2981. doi:10.1109/TAFFC.2023.3250460.
- [33] N. Fernando, B. Nakisa, M. N. Rastgoo, A. Ahmad, Adaptive explainability in human-machine teams: Enhancing swift trust and collaboration in dynamic environments, *CHI - Affective interaction and affective computing - past, present and future* (2025). URL: <https://sites.google.com/view/affectiveinteraction-chi25/call-for-participation?authuser=0>.
- [34] J. Zerilli, U. Bhatt, A. Weller, How transparency modulates trust in artificial intelligence, *Patterns* 3 (2022). URL: <https://doi.org/10.1016/j.patter.2022.100455>. doi:10.1016/j.patter.2022.100455, doi: 10.1016/j.patter.2022.100455.
- [35] M. Setzu, R. Guidotti, A. Monreale, F. Turini, D. Pedreschi, F. Giannotti, Glocalx - from local to global explanations of black box ai models, *Artificial Intelligence* 294 (2021) 103457. URL: <https://www.sciencedirect.com/science/article/pii/S0004370221000084>. doi:<https://doi.org/10.1016/j.artint.2021.103457>.
- [36] V. Chamola, V. Hassija, A. R. Sulthana, D. Ghosh, D. Dhingra, B. Sikdar, A review of trustworthy and explainable artificial intelligence (xai), *IEEE Access* 11 (2023) 78994–79015. URL: <https://doi-org.ezproxy-f.deakin.edu.au/10.1109/ACCESS.2023.3294569>. doi:10.1109/ACCESS.2023.3294569.
- [37] J. D. Lee, K. A. See, Trust in automation: Designing for appropriate reliance, *Human Factors* 46 (2004) 50–80. URL: [https://journals.sagepub.com/doi/abs/10.1518/hfes.46.1.50\\_30392](https://journals.sagepub.com/doi/abs/10.1518/hfes.46.1.50_30392). doi:10.1518/hfes.46.1.50\_30392.
- [38] A. Adadi, M. Berrada, Peeking inside the black-box: A survey on explainable artificial intelligence (xai), *IEEE Access* 6 (2018) 52138–52160. URL: <https://doi.org/10.1109/ACCESS.2018.2870052https://ieeexplore.ieee.org/ielx7/6287639/8274985/08466590.pdf?tp=&arnumber=8466590&isnumber=8274985&ref=>. doi:10.1109/ACCESS.2018.2870052.
- [39] J. S. Park, C. Q. Zou, A. Shaw, B. M. Hill, C. Cai, M. R. Morris, R. Willer, P. Liang, M. S. Bernstein, Generative agent simulations of 1,000 people, 2024. URL: <https://arxiv.org/abs/2411.10109>.