# Exploring and Analyzing the Effect of Avatar's Realism on Anxiety of English as Second Language (ESL) Speakers

Tian-Qi Liu
liu-tq20@mails.tsinghua.edu.cn
Department of Computer Science and
Technology, Tsinghua University
Beijing, China

Joshua Rafael Sanchez
joshsan@uw.edu
University of Washington
Seattle, Washington, United States

Yun-Tao Wang
Department of Computer Science and
Technology, Tsinghua University
Beijing, China

Xin Yi
Institute for Network Sciences and
Cyberspace, Tsinghua University
Beijing, China

Yuan-Chun Shi
Department of Computer science and
Technology, Tsinghua University
Beijing, China

## ABSTRACT

The advent of virtual avatars has opened promising avenues for applications in remote conferencing, education, and beyond. This study investigates the impact of avatar realism on anxiety levels among English as a Second Language (ESL) speakers during interactions with native English speakers. Our findings reveal that cartoonish avatars or direct video interactions with real individuals tend to mitigate such anxieties. Conversely, avatars with higher degrees of human-like realism exacerbate anxiety levels among ESL users. This research carries significant implications for the design and deployment of virtual avatars, specifically in alleviating cross-cultural communication barriers and optimizing user experience.

## CCS CONCEPTS

• **Human-centered computing → Interaction paradigms**; **HCI design and evaluation methods**; **Visualization techniques**.

## KEYWORDS

Digital avatar, anxiety, English as second language (ESL), avatar realism

## 1 INTRODUCTION

As digital representations permeate various sectors and real-world applications, it has become evident that virtual systems offer distinct advantages over their physical counterparts. These advantages range from enhanced control, reduced susceptibility to errors, and the unique capability for synchronous multi-user interactions [54, 55]. One of the most promising facets of these virtual systems is the use of avatars as self-representations, a notable shift from conventional live video. Avatars, especially in remote communication, act as essential conduits, facilitating richer and more nuanced cross-cultural exchanges [9, 52]. Their relevance in contemporary communication cannot be overstated, with the prominence and efficacy of avatar-based communication witnessing a consistent surge [24, 36]. An extensive body of research further attests to the superior effectiveness of avatar-driven communication in varied scenarios, from online meetings [12] to emote education [21, 49].

Nevertheless, a prevalent challenge that often surfaces in cross-cultural communication is the anxiety that English as a Second Language (ESL) speakers frequently encounter due to linguistic barriers and apprehensions about others' evaluation [45]. Such anxiety doesn't merely reduce their communicative efficiency but may also inhibit their confidence and willingness to engage in interactions with native English speakers. Avatar-based communications change the presentation of the interlocutor in the eyes of the ESL speakers, which might present an avenue to alleviate the communicative anxiety of ESL speakers, optimize their experience, and enhance communication efficiency.

However, the question of how the visual representation of avatars impacts communication remains largely unexplored, especially among English as a Second Language (ESL) speakers. Given the diversity of communication styles and cultural backgrounds, it is impractical to propose a one-size-fits-all approach to avatar realism. The effectiveness of specific visual environments is influenced by a complex interplay of social history and learning objectives [35]. Previous research has shown that video calling technologies can significantly aid ESL speakers in improving their public speaking and communication skills [1, 11, 38, 56]. Building on these studies, our research aims to explore the impact of avatar realism on ESL speakers' anxiety levels.

To this end, we have formulated the following research questions:

(1) RQ1: How does the level of realism in an avatar influence the anxiety levels, as indicated by both self-reported psychological questionnaires and physiological signals, of ESL speakers during English communication?
(2) RQ2: Which specific visual characteristics of an avatar contribute to modulating the anxiety levels, as measured by psychological questionnaires and physiological indicators, of ESL speakers when communicating in English?

We initially hypothesized two core premises:

Firstly, with increasing levels of avatar realism (i.e., from cartoonish avatars to realistic-like avatars, culminating in live video), both physiological and psychological anxiety in users would proportionally escalate. Secondly, features that contribute to heightened realism, such as authentic facial features, body forms, and expressions, would intensify user anxiety.

To capture a comprehensive understanding of the effects, our methodology integrates biometric measurements, participant surveys, and a blend of quantitative and qualitative analyses.

## 2 RELATED WORK

### 2.1 Significance of Anxiety Treatment in Public Scenarios, and ESL Speakers' Needs

People of varying ages and social contexts must maintain emotional regulation concerning their emotions. On one end of the socioeconomic spectrum, professionals in high-pressure environments (E.g. social work and residential care) find emotional self-regulation to be essential, according to a study reviewing social workers in a residential unit in Norway [44]. In the educational field context, students often need to be taught emotional intelligence and self-regulation to learn effectively, especially given the rising usage of smartphones [19]. Successful emotion regulation may be achieved through various methods, such as self-realization or external observation. This paper will aim to help understand how novel, technology-forward methods can be advantageous in regulating our emotions, particularly with ESL speakers. Research suggests that addressing the needs of ESL speakers is quite complex; when trying to acquire a second language, learners should be approached with notions of investment and mutual understanding rather than motivation to capture the complexity of the context of a language [27]. Increased use of online learning environments for collaboration between scholars has been shown to improve knowledge sharing [5].

### 2.2 Applicable Signal Processing Methods on Detecting Anxiety

*2.2.1 Motivation For Analyzing Physiological Signals.* The enhancement of human-computer interaction through a computer's capability to recognize and express human emotions is well-documented [33]. Central to this is the collection and analysis of physiological signals, which play a pivotal role in automatic emotion recognition [50]. Various physiological signals, including electrocardiogram (ECG), electrodermal activity (EDA) and photoplethysmography (PPG), have proven effective in detecting fluctuations in emotional anxiety [2, 34]. Our study predominantly leverages ECG, EDA, and PPG, given their established efficacy in indicating anxiety.

*2.2.2 ECG.* The ECG signal captures the heart's electrical activity, allowing for the extraction of vital cardiac metrics such as heart rate (HR) and heart rate variability (HRV). Both HR and HRV are influenced by the sympathetic and parasympathetic branches of the autonomic nervous system, which are intricately linked to an individual's mental anxiety levels [20, 23, 32]. ECG has been employed in various studies, not only for cardiac diagnostics but also as a tool to discern emotional states and anxiety, especially with the advent of wearable ECG devices [4, 27, 48].

*2.2.3 EDA.* EDA measures the changes in conductance at the skin surface due to eccrine sweat production. EDA is a promising biomarker of autonomic system response, being a "Window on the arousal dimension of emotion" [39]. An EDA signal can be decomposed into two quantitative parts: The tonic component-Skin Conductance Level (SCL) and the phasic component-Skin Conductance Responses (SCR), and EDA can be extensively utilized as a stress and anxiety indicator [2, 15].

*2.2.4 PPG.* PPG captures blood volume changes beneath the skin, serves as an indicator of arousal levels linked to stress and anxiety. Its non-invasive nature, ease of use, and cost-effectiveness make it a preferred alternative to measure HR or HRV [2]. PPG's versatility in emotional cognition and anxiety monitoring has been highlighted in various studies [10, 29].

### 2.3 Definitions of Realism and Acknowledgement of Realism in Avatars

The avatar is generally defined as a user's graphical persona inside a virtual world, and photo realism is generally defined as the level of detail and texture used to create realistic, appealing images [18]. Decades of technological advances in virtual reality have made avatars increasingly realistic and life-like. In the meantime. as completely realistic avatars are still being developed, many avatar renditions range in realism. Research reveals that the realism of an avatar does not necessarily correlate with an increase in affinity for such an avatar and vice versa. For example, in user experience, the uncanny valley helps describe a profoundly negative response when a user observes a near-realistic character [22]. According to Mori's paper, the uncanny valley could also be depicted graphically by comparing the human likeness of an entity with the perceiver's affinity for that same entity. An exceptionally high level, around 80-90% of human likeness, could trend in the opposite direction.

Previous research suggests that such visual features of an avatar can affect the appeal of a character. Zibrek et. al [57] interviewed and studied over one thousand people from the general population to determine whether an avatar's rendering styles directly influenced appeal, or a character's personality was more determinant of a person's appeal in virtual reality. The "Realistic" style is designated the most realistic avatar due to its completed facial structure, including skin surface modification, eye refraction, and hair transparency. The remaining styles (I.e., "Toon CG", "Toon Shaded", "Creepy", "Zombie") add more non-human-like components to help reduce human likeness to induce the uncanny effect. Findings, in this case, showed that affinity towards a character is probably determined by a complex interaction between the character's appearance and personality.

## 3 EXPERIMENT SETTINGS AND METHODS

### 3.1 Overview

A user study was conducted on participants who were determined to be English language learners or people who were using English as a second language. For this study, participants engaged in three (3) verbal conversations, each with a different native English speaker. The native speaker either used the cartoon-like avatar, the realistic-like avatar, or the live video camera without using an avatar or any other visual embellishments. Observational data was collected while the study was being conducted and after the user study in the form of self-administered surveys. In the data analysis stage of the study, participants' anxiety was observed using a variety of methods in order to determine if there was any significant statistical difference between anxiety and stress levels between the three scenarios.

*3.1.1 User Study Workflow.* The following user study workflow illustrates the planning, execution, analysis and storing of research data throughout the study process (Figure 1). The study is divided into four phases: Participant Selection (Phase 1), Pre-Intervention Assessment (Phase 2), Intervention Phase (Phase 3), and Post-Intervention Assessment (Phase 4):

- Phase 1: We conducted the on-campus recruitment of participants. Potential candidates who showed interest in our study were instructed to fill out an evaluation form to certify their ability to hold a 5-minute conversation with a native speaker without major difficulty. The purpose of this evaluation was to exclude any candidate who exhibits physical/mental disabilities or similar extenuating circumstances that mainly prevent them from completing the study. Qualifying candidates were made participants.
- Phase 2: This phase involves a Pre-Intervention Assessment, where demographic information, command of and comfortability with the English language and technologies were shared by the participants. The resulting data from this assessment helped compare the characteristics of our sample of ESL speakers to that of the general population in the data analysis.
- Phase 3: In the Intervention Phase, participants were instructed to travel in-person to our laboratory setup, where they were tasked to engage in three conversations– one with a cartoon-like avatar, one with a realistic-like avatar, and one with no avatar. Throughout these exercises, we tracked the participants' physical and mental activity mainly through a MP160 physiological recorder and interviewer observation.
- Phase 4: Participants were prompted to answer the Post-Intervention Assessment shortly after the exercise. Data was collected through self-assessment surveys to analyze each participant's level of stress/anxiety for a particular scenario. We used qualitative and quantitative questions to determine overall anxiety level change between scenarios for our given population sample.

Each of the following sections below will explain in more detail the scenarios, virtual reality technologies, surveys, and data analysis and management methods, respectively, that our study will utilize to implement each phase.

*3.1.2 Details of Intervention: User Interaction Scenarios.* This section will explain the Intervention stage (Phase 3) in more detail. This Intervention stage involves the participant undergoing an in-person user study. This user study is a prospective study, and will resemble a crossover design. We will follow a controlled experimental design, in which the independent variable is the visual depiction of the native English speaker (the interviewer), and all other potential confounding variables (E.g., Interview script and tone of delivery) are kept constant across all 3 scenarios.

Each participant had three 3-6 minute 1-on-1 conversations, each with a different native English speaker and a different exposure (the English speaker uses either the cartoon-like avatar, the realistic-like avatar, or the live video stream). Each participant was exposed to each scenario once. When explaining instructions and logistics of the study not related to the exposure, the camera was turned off. Figure 2 illustrates the timeline of the 45-minute in-person user study.

The English-speaking conversation is meant to simulate a real 1-on-1 conversation between a native English speaker (interviewer) and an ESL speaker (interviewee) who meet for the first time. Interviewers were asked to keep a calm and friendly demeanor when delivering and answering questions. After each exposure, the interviewee was given a within-exposure survey to fill out (See Section 3.3 for more details). Native English speakers are given a transcript to follow, which includes 5 basic questions to ask to the ESL speaker to prompt a conversation and prompts the interviewee to ask a question in English.

The scenarios are as follows:

- Scenario 1: Exposure to 2D User Avatar w/ Cartoon-Like Features. Participants could converse with and view the cartoon-styled avatar rendition of the interviewer.
- Scenario 2: Exposure to 2D User Avatar w/ Realistic-Like Features. Participants were instructed to engage in the same activities and conversations with the interviewer, except that the interviewee (participant) would have been interacting with a realistic-styled avatar rendition of the interviewer.
- Scenario 3: Exposure to Video Camera / Live Video Feed. Participants were instructed to engage in the same activities and conversations with the interviewer, except that the interviewee (participant) would have been interacting with the live video of the interviewer.

It's crucial to note that the order of these conditions was counterbalanced to control for order effects, ensuring that each participant encountered the conditions in a varied sequence.

During each session of the study, the interviewer made personal notes and observations about the participant. The primary focus of these observations was on the participant's verbal response delivery, although other aspects were also considered.

### 3.2 Equipment and Their Implementation

*3.2.1 Avatar Generation.* A realistic avatar may be defined as having facial features akin to a human, while a cartoon avatar might be more characteristic of variation of skin color, glare reflection of the eyes, and transparency of the hair. [57]. Based on these descriptions of realistic and cartoon avatars, our group has developed
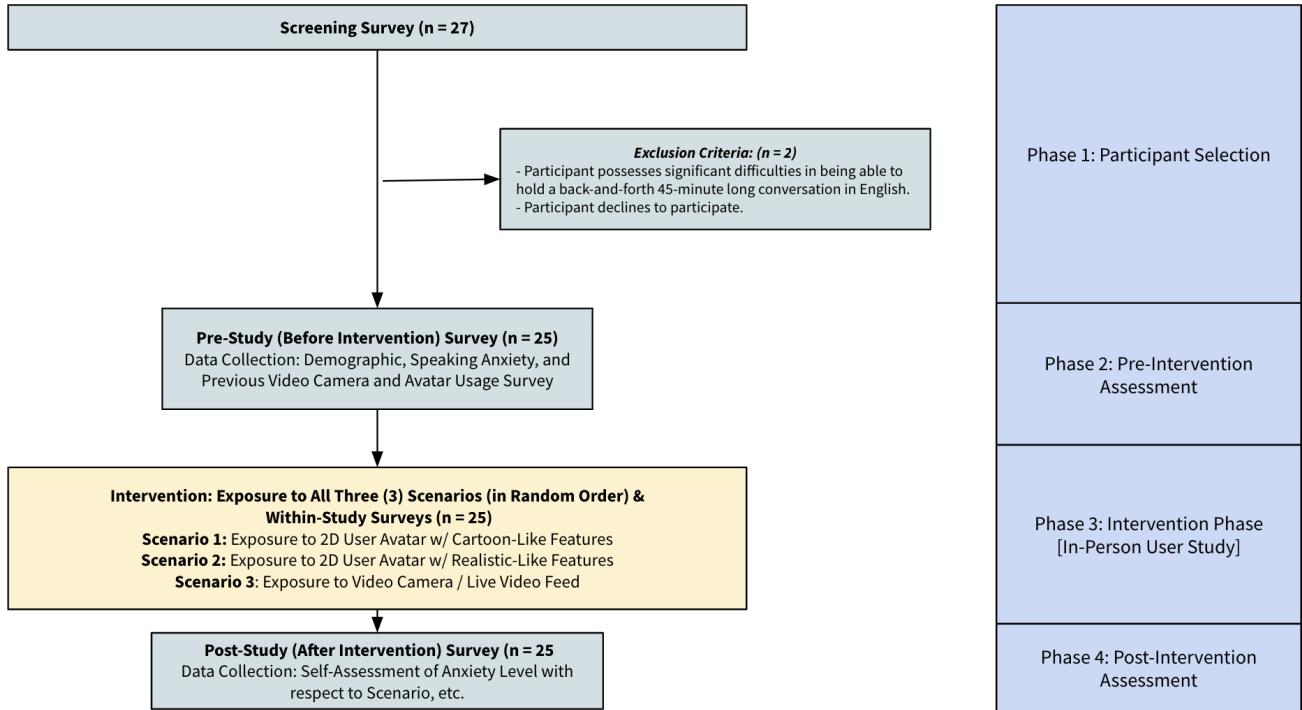
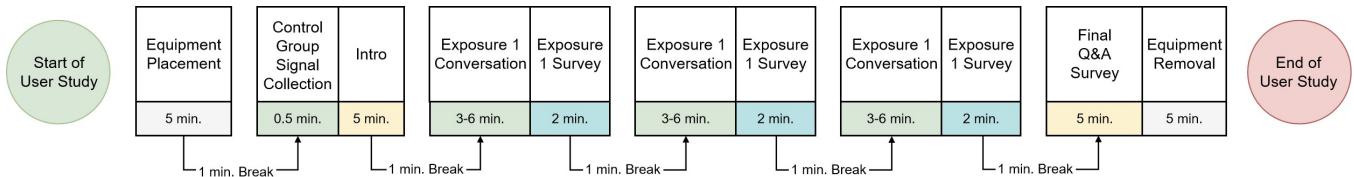**Figure 1: Participant Flowchart for the Entire User Study**



**Figure 2: Detailed Participant Flow During the Intervention Phase**

the following two types of avatars to develop the following facial features:

- Cartoon-Like Avatar (for Scenario 1): This avatar will feature different variations to lower human likeness, including facial feature propositions and varying sizes of facial appendages (E.g. arms, legs).
- Realistic-Like Avatar (for Scenario 2): This avatar will feature the most realistic rendition of the human avatar regarding facial features and body proportions.

Avatars were created using Adobe Character Animator (2020 Version). Figure 3 shows screenshots of the 6 avatars used during the study. Each set of 2 avatars, one realistic-like 2D avatar, and one cartoon-like 2D avatar, was designed in the likeness of 1 of the 3 volunteers native English speakers. During each conversation, each avatar's movement corresponded with the natural person's facial movement as tracked by the video feed (Figure 4). If the speaker was assigned an avatar for their specific discussion, only the live movements of the 2D avatar were shown to the ESL speaker.

*3.2.2   Biometric Equipment.* In our study, we focused on three specific biomarkers to gauge the physiological responses related to anxiety in ESL speakers: electrodermal response (EDA), electrocardiogram (ECG), and photoplethysmography (PPG). To efficiently monitor and capture these parameters throughout the experiment, we utilized a combination of three devices: MacBook Pro 2020, served as the main interface for the conversations between the participants and the avatars, providing a clear and engaging visual platform. Alienware x15 R2, responsible for logging the physiological signal data. MP160 Multi-Channel Physiological Recorder, responsible for detecting the participants' physiological signals, offering a comprehensive and accurate insight into their real-time responses.

### 3.3   Surveys

Surveys are a reliable way of determining an emotional health interest concerning factoring for education, age, and other related demographics. [37, 53] We created surveys for participants to fill out

**Figure 3: Avatars Deployed in the Study: Each native English speaker was represented through one of three visual modalities: a cartoon-like avatar, a realistic-like avatar, or a direct live video camera stream. Both the cartoonish and realistic-like avatars were custom-designed based on the distinct facial features of the corresponding native English speaker.**
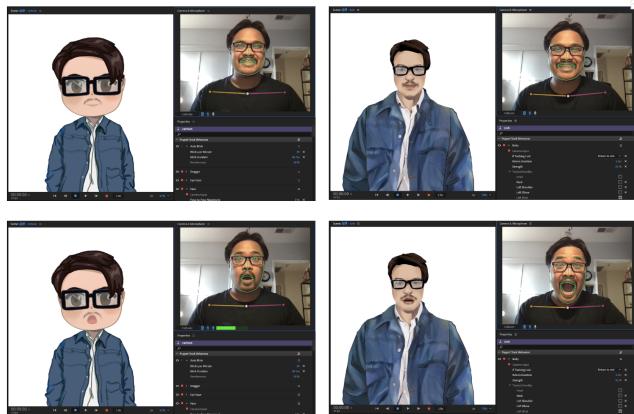


**Figure 4: Screenshots of avatar usage during study. (Left images w/ cartoon-like, Right images w/ realistic-like). In avatar exposures, only the left screen was shown to the participant. The right screen (live video) remained hidden.**

to provide a first-hand perspective on their own perceived anxiety and to see whether survey responses correlate with the results of the biometric data collected. The surveys, and their content and purpose, are as follows:

- 1. Screening Survey: Asks for participant's contact information and English speaking level. The purpose is to select participants willing and able to participate in the study to completion.
- 2. Pre-Study Survey: Asks for participant's demographic information (age, gender, nationality), native language, spoken command of native and English language, anxiety experienced when speaking native versus the English language,

and preference for using a video camera and 2D-avatars when speaking online. The main goals of this survey are to understand our sample of the population pool better and to determine the English-speaking ability of each participant for further analysis.

- 3. Within-Exposure Surveys: These consist of the following: Three (3) in-study surveys. Each survey requires a score from 1 to 7 (1 being "Highly Disagree," to 7 being "Highly Agree") for (10) statements inquiring about their anxiety with a particular exposure. One survey was given after each exposure. One (1) free question-and-answer survey. 4-7 questions were verbally asked in Chinese about the participant's preference for avatars and usage of avatars in different situations.
- 4. Post-Study Survey: Asks for participants' attitude change towards video and avatar usage after in-person study and overall opinions about the study's usefulness in addressing their speaking anxiety.

## 3.4 Data Collection and Management Procedures

*3.4.1 Participant Recruitment.* We have designated selection criteria for participants as follows:

- Participants must be identified as an ESL speaker or an English language learner. Such speakers will be identified through self-assessment via a screening questionnaire.
- Participants must not harbor any physical or mental difficulties that significantly affect the proper usage of online video calling.
- Participants must complete the entire user study process: Completing the pre-use surveys, participating in each scenario in the user study (intervention), participating in heart rate data monitoring during intervention, and completing the post-use surveys.

We recruited 25 candidates based on this selection criteria. Out of this population pool, our sample reflects the following statistics:

- All the participants speak Mandarin as their first language, and speak English as their second language.
- The age range of participants was from 18-26.
- The study covered 11 males and 14 females.

From these 25 participants, a majority (21) reported learning English in the classroom– all are presumed to have taken state-mandated English language courses from primary school to high school. Several students reported having studied abroad in a predominantly English-speaking country. Our sample also ranged in English speaking ability outside of the classroom. We inquired about how often they use the English language by the situation. Most of the ESL participants never use English casually, sometimes use English in education scenario, sometimes or never use it for professional usage, the result can be seen in Table 1. We also asked them to self-report their command of English through Listening, Spoken Interaction, and Spoken Production– definitions of "Beginner", "Intermediate", and "Advanced" comprehension were drawn from the Common European Framework of References for Languages [25]. The majority of ESL users self-assessed their proficiency in listening, spoken interaction, and spoken production as 'intermediate'

**Table 1: Number of Participants by Self-Reported Usage of English by Situation.**

| Field | Never | Sometimes | 1/2 of the time | Most of the time | Always |
|---|---|---|---|---|---|
| Casual | 16 | 7 | 2 | 0 | 0 |
| Educational | 0 | 15 | 7 | 3 | 0 |
| Professional | 11 | 13 | 0 | 0 | 1 |

**Table 2: Number of Participants by Self-Reported Command of English Via Listening, Spoken Interaction, and Spoken Production.**

| Field | Beginning | Intermediate | Advanced |
|---|---|---|---|
| Listening | 6 | 18 | 1 |
| Spoken Interaction | 9 | 15 | 1 |
| Spoken Production | 10 | 14 | 1 |

when categorized among 'beginning', 'intermediate', and 'advanced' levels, as detailed in Table 2.

*3.4.2 Biometric Collection Procedure.* Before the experiment, ECG electrodes were attached correctly to the right clavicle and under the ribs on both sides. The EDA electrodes were attached to participants' pulp of the index and middle fingers in the left hand of subjects, and the PPG sensor was attached to the pulp of the participants' ring finger in the left hand. After placement of all the electrodes, participants were asked to sit in a relaxed posture for a minute before recording a 30 seconds' baseline signal. The MP160 recorder acquired all three signals at a sampling frequency of 2000Hz.

# 4 DATA ANALYSIS AND RESULTS

We will be analyzing data through multimodal analysis. To determine the participant's anxiety measurement, we will consider both quantitative and qualitative dimensions of analysis. We will qualitatively ask open-response and multiple-choice questions to further contextualize personal information. We will ask Likert-scale questions to determine self-assessed anxiety levels under each scenario and analyze heart rate data recorded during the study.

The three physiological signals (EDA, ECG and PPG) were recorded in real-time during the user's entire experiment by BIOPAC Acqknowledge software, and after the experiment, the time-dependent data of the three signals were exported in a CSV format.

The physiological signals were segmented into four distinct categories: control, cartoon, realistic, and person. For the control group, signals were captured for a span of 30 seconds prior to the experiment, reflecting the baseline or relaxed state of the ESL participants. The subsequent groups were delineated based on the avatar used by English native speakers: the cartoon group engaged in conversations with a cartoon avatar, the realistic group interacted with a realistic avatar, and the person group conversed via a live video feed of a native English speaker.

Each conversation's duration typically ranged between 180 to 360 seconds, with the length primarily influenced by the time participants took to answer five questions. Nevertheless, to ensure data uniformity and neutralize potential discrepancies stemming from

these varying durations, only the signals from the first 120 seconds of each interaction were subjected to preprocessing and analysis.
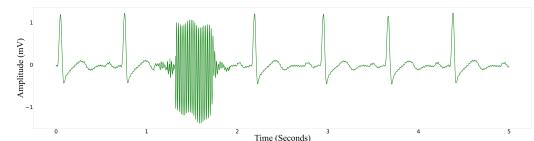
## 4.1 Data Pipeline of Bio-Signal Features

*4.1.1 EDA Signal Preprocess.* EDA signals can be noised by some factors, for example, thermal agitation of electrons in amplifier and power line noises [42]. The raw EDA signal is firstly filtered by a low-pass filter with upper cutoff frequency of 3Hz. After that, a powerline filter of powerline frequency 50Hz is used to further denoise the signal.

The following step is decomposition. Skin conductance data is characterized by the overlap of two components, namely SCL and SCR. SCL is a slowly and continuously varying tonic component, while SCR is a rapid component overlaps associated with a specific stimulus.

*4.1.2 ECG Signal Preprocess.* Raw ECG data undergoes preprocessing that involves artifact removal and signal filtering. Artifacts can notably distort features in both the time-domain and frequency-domain, as highlighted by [28]. For instance, a missed heartbeat might result in a pronounced elongation of the R-R interval and a subsequent drop in heart rate. Once these artifacts are visually identified, the corresponding RR intervals are manually edited, as detailed by [40]. Segments of the ECG that were removed can be seen in Figure 5.

Following the removal process, the refined ECG signal is subjected to a high-pass filter with a cutoff frequency of 0.5Hz and a powerline filtering where the powerline is calibrated to 50Hz. This filtering approach is instrumental in reducing noise present in the raw ECG data.



**Figure 5: Example of ECG Signal Segment Requiring Removal**

*4.1.3 PPG Signal Preprocess.* The PPG signal can be highly susceptible to disturbances from artifacts, especially those stemming from user finger movements. Such disturbances often manifest as erratic fluctuations within brief intervals. Initially, the signal underwent a visual inspection to identify these irregular variations. Portions of the signal affected by these artifacts were manually edited out; examples of which can be found in Figure 6. Subsequently, the PPG signal was processed through a high-pass filter with a lower cutoff frequency set at 0.5Hz, followed by a low-pass filter with its upper cutoff frequency designated at 8Hz.

*4.1.4 Feature Extraction.* Feature extraction is an essential step for anxiety comparison. However, although various features can be extracted from the signals, not all of them can properly and equally show the user's anxiety. Therefore, we need to determine relevant features for anxiety evaluation. Our study focuses on the features extracted from ECG, EDA, and PPG, respectively.
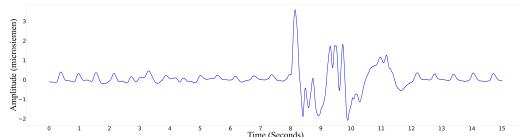
**Figure 6: Example of PPG Signal Segment Requiring Removal**

**mHR** is the mean value of the heart rate (HR) in a period of time. Heart rate (HR) is an important body parameter of healthy regulatory systems that can reflect some sudden environmental and psychological challenges [20, 31]. In particular, anxiety has been found to prompt increased HR in a lot of previous works, and mHR in an interval is often used to help classify a person's anxiety during the period [13].

mHR can be extracted from ECG signals based on the following steps:

After the signal is preprocessed, the R peaks of the ECG signal are then detected, based on the Pan-Tompkins algorithm [26].

mRR (mean R-R Interval) can be calculated after R peaks are detected by the following formula:

$$mRR = \frac{\sum_{i=1}^{N}(RR_i)}{N}$$

where $RR_i$ represents the time interval between the i-th R peak and its' subsequent (i+1)-th peak, N is number of adjacent R waves in ECG signal.

mHR can also be figured out based on the formula:

$$mHR = \frac{\sum_{i=1}^{N}(\frac{60000}{RR_i})}{N}$$

**RMSSD** is the root mean square of successive differences between normal heartbeats, reflecting the beat-to-beat variance in HR and is an essential feature to estimate the vagally mediated changes [41]. As one of the commonly used short-term heart rate variability (HRV) indices, RMSSD has been found that its reduction is related to users' serious anxiety increases [6, 7, 32, 47]

$$RMSSD = \sqrt{\frac{\sum_{i=1}^{N-1}(RR_i - RR_{i+1})^2}{N-1}}$$

**LF/HF** is the ratio of LF to HF, where LF means the absolute power in the low-frequency band (0.04 - 0.15Hz) and HF means the absolute power in the high-frequency band (0.15 - 0.4Hz) [40].

$$LF/HF = \frac{HR_{LF}}{HR_{HF}}$$

**NR** is the number of response amplitudes in the analyzed time segment, which is a typical EDA feature correlated to subjects' anxiety degree changes. The minimum amplitude threshold to detect SCR peaks is $0.1\mu S$. SCR Peaks are identified through the algorithm proposed by D. Makowski et al. [16].

$$NR = count(Amplitudes)$$

**maxR** is the maximum response amplitude in the analyzed time segment, which is also helpful to reflect the subjects' response to

events in the interested time period.

$$maxR = \max_{K=0...NR-1}(Amplitude_K)$$

where $Amplitude_K$ is the k-th amplitude during the analyzed period of time.

**MH** is the mean value of the SCR amplitude height include the Tonic component.

**MA** is the mean value of the SCR amplitude height exclude the Tonic component [31].

$$MA = \frac{1}{NR}\sum_{k=0}^{NR-1}|Amplitude_K|$$

**MPR** is the mean heart rate measured based on PPG peaks. The mPR is calculated from the similar steps like mHR, after the PPG peaks are detected based on Elgendi's method [8].

Research focusing on the relationship between physiological signal features and anxiety level suggests that, anxiety as a typical kind of anxiety, can result in increased mHR [2, 20, 31, 51], increased LF/HF [2, 30, 32], increased NR [2, 31], increased maxR [14, 31], increased MH [14, 31], increased MA [2, 14, 31], increased MPR [2, 14] and also decreased RMSSD [2, 17, 30, 43]. The selected biosignal features and their descriptions can be found in Table 3.

## 4.2 Results From Bio-signal Features

The analysis focused on features associated with user anxiety, specifically:mHR, RMSSD, LF/HF, NR, maxR, MH, MA and MPR. To account for differences due to the total communication time, we restricted our analysis to features derived from the initial 120 seconds of each interaction. Repeated measures ANOVA was utilized to discern statistical differences in these metrics.

Outliers were retained because they may represent important cases as they can embody more severe cases of anxiety, this decision did cause certain groups to breach the normality assumption.

**NR Analysis:** Initial checks revealed that NR was not normally distributed for all the groups, as confirmed by the Shapiro-Wilk's test (p < .05). Subsequently, a Friedman test showed significant differences in NR across the control and experimental groups, $\chi^2(3)$ = 9.452, p = .024 < .05. Post hoc analysis with Bonferroni correction identified a significant difference between the cartoonish avatar and the realistic avatar groups (p = .037 < .05).

**maxR Anlysis:** Shapiro-Wilk's test indicated that the maxR distribution approached approximately normality (p = .047). Mauchly's test of sphericity confirmed the sphericity assumption ($\chi^2(5)$ = 4.891, p = .430 > .05). Significant differences in maxR across groups were found (F(3, 72) = 10.788, p < .001). Bonferroni-adjusted post hoc analysis pinpointed significant differences between:

- Control and realistic avatar groups (-.842 (95% CI, -1.330 to -.353)$\mu s$, p < .001).
- Control and live video groups (-.464 (95% CI, -.908 to -.021)$\mu s$, p = .036 < .05).
- Cartoonish avatar and realistic avatar groups (-.540 (95% CI, -1.010 to -.070)$\mu s$, p = .018 < .05).

**MH Analysis:** Given the presence of outliers, a Friedman test was conducted and found differences in MH among the avatar groups, $\chi^2(3)$ = 8.856, p = .031 < .05. Bonferroni-adjusted pairwise

**Table 3: Notations and Descriptions of the Selected Biosignal Features**

| Signal | Domain | Parameter | Description | Unit |
|---|---|---|---|---|
| ECG | Time | mHR | mean Heart Rate | 1/min |
| | | RMSSD | Root mean square of successive N-N interval differences | ms |
| | Frequency | LF | Absolute power in the low-frequency band(0.04 - 0.15Hz) | $ms^2$ |
| | | HF | Absolute power in the high-frequency band(0.15 - 0.4Hz) | $ms^2$ |
| | | LF/HF | Ratio of LF to HF | - |
| EDA | Event | NR | number of the response amplitudes in the analyzed period | - |
| | | maxR | maximum of the response amplitudes in the analyzed period | $\mu s$ |
| | Time | MH | the mean height of SCR amplitude including the Tonic component | $\mu s$ |
| | | MA | the mean SCR amplitude excluding the Tonic component | $\mu s$ |
| PPG | Time | MPR | mean PPG rate | 1/min |

comparisons revealed a significant difference between the cartoon and realistic avatar groups (p = .019 < .05).

**MA Analysis:** The MA from the control group did not adhere to a normal distribution. A Friedman test identified differences in MA across avatar interactions, $\chi^2(3)$ = 14.472, p = .002 < .05. Subsequent pairwise comparisons with Bonferroni correction highlighted significant differences between the realistic avatar group and both the control group (p = .006 < .05) and the cartoon avatar group (p = .006 < .05).

For the remaining features (mHR, RMSSD, LF/HF, MPR), no significant differences were found among the groups. The result is demonstrated in Figure 7. Compared to the control group, there was no significant increase in users' anxiety when interacting with a cartoon avatar or engaging in direct video calls. However, conversing with a realistic-like avatar did significantly elevate users' anxiety levels. Interestingly, according to the features extracted from the physiological signals, ESL users exhibited a slightly higher sense of anxiety during video calls compared to interactions with the cartoon avatar.

Our findings indicate that EDA-related features, particularly NR, maxR, MH and MA, exhibit heightened sensitivity to increased anxiety within the context of ESL-native speaker conversations. In these scenarios, physiological manifestations of anxiety did not result in significant alterations in other discussed markers, such as heart rate. Thus, EDA emerges as a more responsive indicator to this specific communication setting.

## 4.3 Data Analysis of Survey Data

*4.3.1 Quantitative Data Analysis.* To measure a survey construct, Cronbach's Alpha is used to determine consistency between items, and is the whole widely used objective measure of reliability [46]. When measuring statements within groups, an alpha value of 0.7 to 0.8 is regarded as satisfactory [3]. For this study, we created two groups of items (statements) and measured each alpha of each item within a group to determine consistency. The groups we created and the constructs are listed as follows.

- Group 1: For Measuring Participants' Language Speaking Anxiety in Pre-Study Survey (See Table 4 for selected statements):
  - C1: Measuring Native-Language Speaking Anxiety of Participants
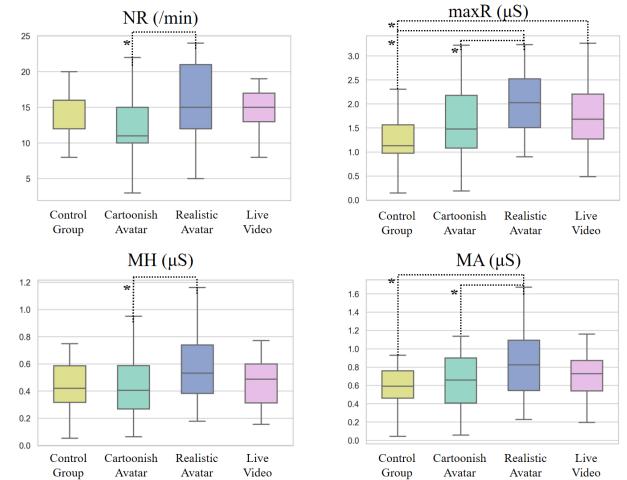


**Figure 7: Distributions of NR, maxR, MH, and MA features when participants were in the control group or interacting with cartoonish avatars, realistic-like avatars, or live videos. Dotted lines marked with an asterisk (*) highlight significant differences between paired groups (p < 0.05).**

  - C2: Measuring English-Language Speaking Anxiety of Participants
- Group 2: For Measuring Participants' Anxiety When Speaking to Avatar or Video Camera during User Study (See Table 5 for selected statements):
  - C3: Measuring Anxiety When Speaking to Cartoon-Like Avatar
  - C4: Measuring Anxiety When Speaking to Realistic-Like Avatar
  - C5: Measuring Anxiety When Speaking to Live Video of Real Person

Participants were asked in a written survey to choose on a Likert scale their opinions about each of these statements, from 1 (strongly disagree) to 7 (strongly agree). In the survey, we asked for responses for more than 6 statements per group (12 statements for Group 1, 10 statements for Group 2), in order to account for potential low Cronbach's Alpha values (< 0.60) from multiple statements. After

the user study finished, we used this selected group of items to measure each construct. Note we slightly changed the wording to fit the construct we inquired about. This means, for example, for Measuring English-Language Speaking Anxiety (C2) using the statement "I am not afraid to express myself at meetings" (L-S2), the statement was edited to the following to reflect the change: "I am not afraid to express myself at meetings in English." We edited each statement similarly to fit the meaning of each construct, but as little as possible to avoid distorting the meaning of the statement. Exact questions used are included in the supplementary material.

After validating consistency, we can firstly use the Language Speaking Items (Group 1) to determine the overall anxiety change when speaking the English language versus when speaking their native language. The responses were rated on a 7-point Likert scale (from 1 "Strongly Disagree" to 7 "Strongly Agree").

Based on the responses shown in Table 6, we could see slight deviations in scores depending on the statement. We interpret a 1-point decrease as only a slight change from either situation. Statements L-S2 and L-S3 address anxiety experienced in meetings; because of the negligible change in sentiment between situations, language barrier seems to have less of an impact. Statements L-S4 and L-S5 inquire about comfort and relaxation when speaking; from native to English mean scores, there is over a 1-point difference, meaning language may have a stronger impact on internal feelings experienced when speaking.

For almost every statement E-S1 through E-S6 (see Table 7), the scores increase in the following order: Realistic-Like Avatar, Cartoon-Like Avatar, then Speaking to a Live Person. The greatest score increases between Realistic-Like Avatar and Live Person exposures are shown in the statements E-S1 and E-S5, each a 1.2-point and 1.6-point increase, respectively. E-S1 asks how easy it is to look at the interviewer when speaking, and E-S5 asks whether the appearance distracts from the conversation. We could observe that looks and appearance may more strongly affect the preference between different exposures.

*4.3.2 Qualitative Data Analysis.* For qualitative data analysis, we want to understand which avatars users preferred and for what reasons, whether avatar-related or not. Near the end of the in-person user study, after three conversations, we asked users to rank the avatar or person they would rather speak to between the three options and their reasons for choosing so. The following table depicts the number of participants listing their preferences under each combination. Users' preferences are shown in Table 8. A majority of participants quoted to rank the "Video" exposure (Scenario 3) first, then the "Cartoon" exposure (Scenario 1), then the "Realistic" exposure (Scenario 2), in that order.

Through thematic analysis of participants' reasonings, we can decipher several themes which guided their decision to rank.

- Exposure 1: Participants were able to converse with and to view the cartoon-styled avatar rendition of the interviewer. Preference: Majority of participants preferred as 2nd choice, above the realistic-style avatar and below the live video, in rank.
  - Reasons for Preference: The cartoonish appearance conveyed a non-threatening demeanor, making the figure

seem more adorable, approachable and easy to get along with.
  - Reasons for Dislike: Movement of face does not give enough information (E.g. limited movement of mouth, automatic eye blinking)
  - Preferred features: A baby-faced appearance, large eyes.
  - Least Liked Features: Subtle or less pronounced mouth movements.
  - Scenarios of Appropriate Usage: Most participants would consider using in daily situations for entertainment, but would refrain from using in professional situations.
- Exposure 2: Participants were instructed to engage in the same activities and conversations with the interviewer, with the exception that the interviewee (participant) would have been interacting with a realistic-styled avatar rendition of the interviewer.
  - Preference: Majority of participants preferred as 3rd choice, below the two order exposures.
  - Reasons for Preference: A few participants noted higher expressiveness; closer in resemblance to a human.
  - Reasons for Dislike: Most participants perceived the avatar as robotic and cold. A few mentioned more distractions (E.g. clothing, coloration of person) as compared to the cartoon-styled avatar.
  - Preferred features: The realistic-like eyes.
  - Least Liked Features: Rigid facial expressions, realistic eyes accompanied by incongruous and unnatural blinking motions.
  - Scenarios of Appropriate Usage: Many participants preferred to never use the avatar, or to use it as seldom as possible.
- Exposure 3: Participants were instructed to engage in the same activities and conversations with the interviewer, with the exception that the interviewee would have been interacting with the live video of the interviewer.
  - Preference: Majority of participants preferred this exposure as 1st choice, regardless of order of exposure and native speaker assignment.
  - Reasons for Preference: Many participants expressed natural, comfortable and reassuring; they reported being accustomed to seeing a live video of a person. Many appreciated being able to read the facial expressions of the interviewer.
  - Reasons for Dislike: A handful of participants reported experiencing stage fright. Seeing another person's face directly made them feel as though their own facial expressions and appearance were entirely exposed, leading to feelings of vulnerability and a diminished sense of privacy.
  - Preferred features: The rich and natural facial dynamics, including engaging eye contact and a contagious smile.
  - Least Liked Features: Overly direct gaze, feeling of being intently stared at.
  - Scenarios of Appropriate Usage: Many expressed video usage as a necessity for professional usage, and even for casual usage.

**Table 4: Group 1: Language Speaking Items/Statements to Measure Anxiety - Selected Statements From Survey**

| Code | Statement | Cronbach's Alpha for Item |
|---|---|---|
| L-S1 | Engaging in a group discussion with new people makes me loose and relaxed. | 0.828 |
| L-S2 | Generally, I am relaxed when I have to participate in a meeting. | 0.826 |
| L-S3 | I am not afraid to express myself at meetings. | 0.712 |
| L-S4 | Communicating at meetings usually makes me comfortable. | 0.774 |
| L-S5 | Ordinarily, I am very calm and relaxed in conversations. | 0.793 |
| L-S6 | My thoughts do not become confused and jumbled when I am giving a speech. | 0.639 |

**Table 5: Group 2: Avatar/Video Exposure Items/Statements - Selected Statements From Survey**

| Code | Statement | Cronbach's Alpha for Item |
|---|---|---|
| E-S1 | It is easy for me to look at his/her face when speaking. | 0.778 |
| E-S2 | Most of the time, I felt quite comfortable when looking at him/her. | 0.747 |
| E-S3 | My eyes did not wander when speaking to him/her. | 0.715 |
| E-S4 | I had no trouble speaking clearly when looking at him/her. | 0.701 |
| E-S5 | His/her appearance did not distract me at all during the conversation. | 0.768 |
| E-S6 | I was not nit-picky at all when looking at his/her face. | 0.692 |

**Table 6: Mean scores given to each statement in Group 1 (L-S1 through L-S6), rated on a 7-point Likert scale (from 1 "Strongly Disagree" to 7 "Strongly Agree")**

| Situation of Statement | L-S1 | L-S2 | L-S3 | L-S4 | L-S5 | L-S6 |
|---|---|---|---|---|---|---|
| When Speaking Native Language | 4.2 | 3.6 | 3.44 | 4.24 | 4.72 | 3.56 |
| When Speaking English Language | 3.44 | 3.4 | 3.88 | 3.32 | 3.44 | 4.04 |

**Table 7: Mean scores given to each statement in Group 2 (E-S1 through E-S6), rated on a 7-point Likert scale (from 1 "Strongly Disagree" to 7 "Strongly Agree")**

| Situation of Statement | E-S1 | E-S2 | E-S3 | E-S4 | E-S5 | E-S6 |
|---|---|---|---|---|---|---|
| When Speaking to Cartoon-Like Avatar | 5.12 | 4.76 | 3.88 | 4.76 | 4.44 | 4.52 |
| When Speaking to Realistic-Like Avatar | 4.72 | 4.4 | 3.44 | 4.76 | 3.92 | 4.44 |
| When Speaking to Live Person | 5.92 | 5.2 | 4.36 | 4.96 | 5.52 | 4.76 |

## 5 DISCUSSION

### 5.1 The Uncanny Valley and Our Realistic-Like Avatar

We had hypothesized that a cartoon-like avatar interpretation will result in lower anxiety levels than more realistic interpretations. While we do acknowledge the definition of "realistic" may vary, earlier studies suggest the importance of avoiding the uncanny valley when creating avatars. Participant feedback to our realistic-like avatar suggested that this avatar's stylings approached the uncanny valley. Biometric data shows an increase in anxiety when speaking with the realistic avatar as opposed to the cartoon-like avatar. In fact, we received both from self-reported data from participants and from verbally raised concerns during the user study about the perceived robotic movements of the avatar.

### 5.2 Live Videos and Cartoon Avatars Physiologically Relax Users, Yet the Former Garners More Subjective Preference

Interactions with realistic-like avatars were found to elevate anxiety levels in ESL users, a sentiment echoed both in their subjective responses and physiological data. From the feedback, it's evident that ESL users felt most at ease during live video dialogues with native English speakers, a sentiment that was closely paralleled by their interactions with cartoonish avatars. Physiologically speaking, the differences between interactions using cartoonish avatars and live video were minimal, suggesting low anxiety levels in both cases. However, when delving deeper into the physiological responses, cartoonish avatars seemed to have a marginal advantage over live video.

Several factors can be posited for this observed difference:

- Feedback from user surveys suggests that cartoon avatars were perceived to be more fitting for informal, laid-back conversations. In contrast, live video interactions were seen

**Table 8: Number of Participants by Preference of Exposures**

| 1st Preference: | 2nd Preference: | 3rd Preference: | Participants Who Prefer: |
|---|---|---|---|
| Live Video | Cartoon-Like Avatar | Realistic-Like Avatar | 17 |
| Live Video | Realistic-Like Avatar | Cartoon-Like Avatar | 2 |
| Cartoon-Like Avatar | Live Video | Realistic-Like Avatar | 2 |
| Cartoon-Like Avatar | (2nd/3rd Tie) Live Video | (2nd/3rd Tie) Realistic-Like Avatar | 1 |
| Cartoon-Like Avatar | Realistic-Like Avatar | Live Video | 3 |
| Realistic-Like Avatar | Cartoon-Like Avatar | Live Video | 0 |
| Realistic-Like Avatar | Live Video | Cartoon-Like Avatar | 0 |

as versatile, suitable for both informal chats and more serious discussions. This flexibility of live videos led to them being slightly more preferred. Given that our experiment revolved around light-hearted conversations, the cartoon avatars seemed apt, which is reflected in the reduced physiological anxiety indicators among participants.

- Elaborating on the physiological aspect, it was primarily the EDA-related signals that manifested significant variances among different avatars. The heightened sensitivity of EDA in this particular scenario – an ESL user conversing with a native English speaker – might hint at nuanced anxiety differentials that are subtle and perhaps harder for users to consciously recognize. This underscores the interesting observation that even when interaction feels most authentic, as in the case of live video dialogues, ESL users still experience underlying anxiety.

## 5.3 Optimal Avatar Usage Scenarios and Design Considerations for ESL Users

Based on physiological signal outcomes and subjective feedback, direct video dialogue appears to be the most ideal method in a majority of contexts, as it tends to lower user anxiety, psychologically and physiologically. Nonetheless, avatar-based communication also finds relevance across a variety of settings. Users tend to prefer cartoonish avatars in lighter, more entertainment-focused scenarios, while gravitating towards more realistic avatars in more serious contexts. Compared to direct video dialogue, engaging with an avatar can heighten a user's sense of security. As long as the avatar's features are designed appropriately, user anxiety can be maintained at a minimal level, making avatars a recommended medium for remote meetings or education with ESL users.

Users desire avatars that can circumvent the challenges of direct gaze in video, yet they also hope that avatars retain essential social interaction features. Aspects such as eye contact, expressive facial movements, and emotion conveyance are crucial. While the realism of these features can be enhanced as the avatar becomes more anthropomorphic, it's vital to avoid the pitfalls of the uncanny valley phenomenon to ensure a comfortable user experience.

## 5.4 The Importance of Associating

The most deciding factor of whether or not a participant felt more or less anxious with an exposure is whether or not they were able to associate with the avatar or video they were viewing. We designed the study so that every conversation that was delivered by

the native English speakers would be identical in content, potential responses to questions, demeanor, emotion, and any other attribute that would characterize a conversation in order to ensure results were dependent only on their perception of the avatar or live video. Our study suggests that visual representation is still an important component in determining how well a participant can feel comfortable. In the case of ESL speakers, being able to associate, perhaps if that stresses knowing how to anticipate emotion from the other person, is an essential component of learning a new language and its culture.

## 5.5 Optimal Opportunities for ESL Speakers to Enhance English Proficiency

Our research underscores the pivotal role of consistent interaction and immersive experiences with native English speakers in mitigating anxiety for ESL learners. Feedback gathered post-study illuminates a strong inclination among participants to re-engage in similar interactions, viewing them as valuable practice sessions. Notably, many expressed an eagerness to confront more intricate and contextually rich questions, emphasizing their drive for enhanced learning and engagement.

A noteworthy proposition from some participants revolved around the adaptive use of avatars. They contemplated starting their learning journey using avatars as a buffer, especially during the initial phases when their English proficiency might still be nascent. As their confidence and linguistic skills grow, they would gradually pivot towards direct camera interactions, harnessing the authenticity and immediacy it offers. This adaptive approach underscores the value of personalizing English learning trajectories, emphasizing the unique needs and comfort levels of each learner.

## 5.6 Limitations

In terms of study limitations, we acknowledge that we have used only 2 types of renditions (6 avatars total) of a 2D avatar. The definitions of what can be defined as "realistic-like" and as "cartoon-like" can constitute a wide spectrum of perspectives, so the exploration of these definitions– for example, how "realistic" an avatar should be for ESL speakers– could be an area of further research. Additionally, the biometric and survey results were obtained from a controlled environment. All participants were holding guided 1-on-1 conversations with native English speakers who were trained to guide the dialogue in a particular manner, and under a time constraint. We welcome research looking to explore ESL speakers' anxiety levels from a variety of different public and private scenarios, and perhaps

studies observing ESL speakers' improvement over a longer period of time.

## 6 CONCLUSION

Through meticulous data collection and comprehensive analysis, we have addressed our two primary research questions.

For RQ1: Our findings illuminate the impact of an avatar's realism on the anxiety levels of ESL speakers during English communication. Both psychological questionnaires and physiological indicators substantiate the influence. ESL speakers are more fluent and at ease when the visual representation feels more relatable. This is evident from users' subjective feedback, which displays a pronounced inclination towards live video. Notably, interactions with a cartoonish avatar also lead to a palpable decrease in anxiety, offering a distinct physiological benefit. This finding both complements and complicates our initial hypothesis: while an escalation in avatar realism generally corresponded with reduced anxiety, cartoonish avatars emerged as an exception, showcasing their particular charm. The realistic-like avatar, however, veered towards the unsettling domain of the uncanny valley.

For RQ2: Our research pinpoints particular visual characteristics of an avatar that modulate the anxiety levels of ESL speakers. Key among these is the avatar's perceived likability, augmented by its association with users, especially in vital aspects such as emotion conveyance and maintaining eye contact. While increased likability can derive from a multitude of visual elements worth future exploration, our findings emphasize the importance of avoiding the uncanny valley to maintain an avatar's appeal and effectiveness. Our results both corroborate and elaborate on our second hypothesis. While heightened realism in certain elements did amplify anxiety, the combined dynamics of familiarity, likability, and strategic avoidance of the uncanny valley emerged as vital determinants shaping the user experience.

## REFERENCES

[1] Dukhayel Aldukhayel. 2022. Remote Presentations: Making L2 Presentations Less Stressful. *Education Research International* 2022 (2022).
[2] Lou Ancillon, Mohamed Elgendi, and Carlo Menon. 2022. Machine Learning for Anxiety Detection Using Biosignals: A Review. *Diagnostics* 12, 8 (2022). https://www.mdpi.com/2075-4418/12/8/1794
[3] J Martin Bland and Douglas G Altman. 1997. Statistics notes: Cronbach's alpha. *Bmj* 314, 7080 (1997), 572.
[4] Jos F Brosschot and Julian F Thayer. 2003. Heart rate response is longer after negative emotions than after positive emotions. *International journal of psychophysiology* 50, 3 (2003), 181–187.
[5] Juanjuan Chen, Minhong Wang, Paul A Kirschner, and Chin-Chung Tsai. 2018. The role of collaboration, computer use, learning environments, and supporting strategies in CSCL: A meta-analysis. *Review of Educational Research* 88, 6 (2018), 799–843.
[6] Anthony B. Ciccone, Jacob A Siedlik, Jill M. Wecht, Jake Andrew Deckert, Nhuquynh D. Nguyen, and Joseph P. Weir. 2017. Reminder: RMSSD and SD1 are identical heart rate variability metrics. *Muscle & Nerve* 56 (2017).
[7] Dimitri Dimitriev, Elena Saperova, Aleksey D. Dimitriev, and Martin G. Frasch. 2016. State Anxiety and Nonlinear Dynamics of Heart Rate Variability in Students. *PLoS ONE* 11 (2016).
[8] Mohamed Elgendi, Ian Norton, Matt B Brearley, Derek Abbott, and Dale Schuurmans. 2013. Systolic Peak Detection in Acceleration Photoplethysmograms Measured from Emergency Responders in Tropical Conditions. *PLoS ONE* 8 (2013).
[9] Gwen D Erlam, Nick Garrett, Norina Gasteiger, Kelvin Lau, Kath Hoare, Shivani Agarwal, and Ailsa Haxell. 2021. What really matters: Experiences of emergency remote teaching in university teaching and learning during the COVID-19 pandemic. In *Frontiers in Education*, Vol. 6. Frontiers Media SA, 639842.
[10] Hao-Yu Jan, Mei-Fen Chen, Tieh-Cheng Fu, Wen-Chen Lin, Cheng-Lun Tsai, and Kang-Ping Lin. 2019. Evaluation of Coherence Between ECG and PPG Derived Parameters on Heart Rate Variability and Respiration in Healthy Volunteers With/Without Controlled Breathing. *Journal of Medical and Biological Engineering* 39 (2019), 783–795.
[11] Matthew K. Miller, Martin Johannes Dechant, and Regan L. Mandryk. 2021. Meeting You, Seeing Me: The Role of Social Anxiety, Visual Feedback, and Interface Layout in a Get-to-Know-You Task via Video Chat.. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–14.
[12] Yuka Koyano, Toshiaki Oguchi, Shingo Akagaki, Sayaka Doi, Hiroyuki Mitani, Takayuki Yasunaga, and Hiroaki Tobita. 2022. Development and evaluation of online meeting system to promote effective communication. In *Proceedings of the Future Technologies Conference (FTC) 2021, Volume 3*. Springer, 712–725.
[13] Sylvia D. Kreibig. 2010. Autonomic nervous system activity in emotion: A review. *Biological Psychology* 84 (2010), 394–421.
[14] Seungji Lee, Taejun Lee, Taeyang Yang, Changrak Yoon, and Sung-Phil Kim. 2020. Detection of Drivers' Anxiety Invoked by Driving Situations Using Multimodal Biosignals.
[15] Yun Liu and Siqing Du. 2018. Psychological stress level detection based on electrodermal activity. *Behavioural brain research* 341 (2018), 50–53.
[16] Dominique Makowski, Tam Pham, Zen Juen Lau, Jan C. Brammer, Françoise Lespinasse, Pham Tien Hung, Christopher Schölzel, and SH Annabel Chen. 2021. NeuroKit2: A Python toolbox for neurophysiological signal processing. *Behavior research methods* (2021).
[17] Marek Malik, John Thomas Bigger, A. John Camm, Robert E. Kleiger, Alberto Malliani, Arthur J. Moss, and Peter J. Schwartz. 1996. Heart rate variability. Standards of measurement, physiological interpretation, and clinical use. *European Heart Journal* 17 (1996), 354–381.
[18] Carol Manetta and Richard A Blade. 1995. Glossary of virtual reality terminology. *International Journal of Virtual Reality* 1, 2 (1995), 35–39.
[19] Maria Lidia Mascia, Mirian Agus, and Maria Pietronilla Penna. 2020. Emotional intelligence, self-regulation, smartphone addiction: which relationship with student well-being and quality of life? *Frontiers in psychology* 11 (2020), 375.
[20] Darien Miranda, Marco A. Calderon, and Jesús Favela. 2014. Anxiety detection using wearable monitoring. In *MexIHC '14*.
[21] Vardan Mkrttchian, Dina Kharicheva, Ekaterina Aleshina, Svetlana Panasenko, Yulia Vertakova, Leyla Ayvarovna Gamidullaeva, Mikhail Ivanov, and Vsevolod Chernyshenko. 2020. Avatar-based learning and teaching as a concept of new perspectives in online education in Post-Soviet Union countries: Theory, environment, approaches, tools. *International Journal of Virtual and Personal Learning Environments (IJVPLE)* 10, 2 (2020), 66–82.
[22] Masahiro Mori, Karl F MacDorman, and Norri Kageki. 2012. The uncanny valley [from the field]. *IEEE Robotics & automation magazine* 19, 2 (2012), 98–100.
[23] Isaac Moshe, Yannik Terhorst, Kennedy Opoku Asare, Lasse Bosse Sander, Denzil Ferreira, Harald Baumeister, David C Mohr, and Laura Pulkki-Råback. 2021. Predicting symptoms of depression and anxiety using smartphone and wearable data. *Frontiers in psychiatry* 12 (2021), 625247.
[24] Kristine L Nowak and Jesse Fox. 2018. Avatars and computer-mediated communication: a review of the definitions, uses, and effects of digital representations. *Review of Communication Research* 6 (2018), 30–53.
[25] Council of Europe. Council for Cultural Co-operation. Education Committee. Modern Languages Division. 2001. *Common European framework of reference for languages: Learning, teaching, assessment*. Cambridge University Press.
[26] Jiapu Pan and Willis J. Tompkins. 1985. A Real-Time QRS Detection Algorithm. *IEEE Transactions on Biomedical Engineering* BME-32 (1985), 230–236.
[27] Bonny Norton Peirce. 1995. Social identity, investment, and language learning. *TESOL quarterly* 29, 1 (1995), 9–31.
[28] Mirja A. Peltola. 2012. Role of Editing of R–R Intervals in the Analysis of Heart Rate Variability. *Frontiers in Physiology* 3 (2012).
[29] David Perpetuini, Antonio Maria Chiarelli, Daniela Cardone, Chiara Filippini, Sergio Rinella, Simona Massimino, Francesco Bianco, Valentina Bucciarelli, Vincenzo Vinciguerra, Piero Fallica, et al. 2021. Prediction of state anxiety by machine learning applied to photoplethysmography data. *PeerJ* 9 (2021), e10448.
[30] David Perpetuini, Antonio Maria Chiarelli, Daniela Cardone, Chiara Filippini, Sergio Rinella, Simona Massimino, Francesco Bianco, Valentina Bucciarelli, Vincenzo Vinciguerra, Pier Giorgio Fallica, Vincenzo Perciavalle, Sabina Gallina, Sabrina Conoci, and Arcangelo Merla. 2021. Prediction of state anxiety by machine learning applied to photoplethysmography data. *PeerJ* 9 (2021).
[31] Livia Petrescu, Catalina Petrescu, Oana Mitruţ, Gabriela Moise, Alin Dragos Bogdan Moldoveanu, Florica Moldoveanu, and Marius Leordeanu. 2020. Integrating Biosignals Measurement in Virtual Reality Environments for Anxiety Detection. *Sensors (Basel, Switzerland)* 20 (2020).
[32] Tam Pham, Zen Juen Lau, Shen-Hsing Annabel Chen, and Dominique Makowski. 2021. Heart Rate Variability in Psychology: A Review of HRV Indices and an Analysis Tutorial. *Sensors (Basel, Switzerland)* 21 (2021).
[33] Rosalind W Picard. 2000. *Affective computing*. MIT press.
[34] Andre Pittig, Joanna J. Arch, C. W. Lam, and Michelle G. Craske. 2013. Heart rate and heart rate variability in panic, social anxiety, obsessive-compulsive,

and generalized anxiety disorders at baseline and in response to relaxation and hyperventilation. *International journal of psychophysiology : official journal of the International Organization of Psychophysiology* 87 1 (2013), 19–27.

[35] Sally-Ann Robertson and Mellony Graven. 2015. Exploring South African mathematics teachers' experiences of learner migration. *Intercultural Education* 26, 4 (2015), 278–295.

[36] MD Roblyer, Marclyn Porter, Talbot Bielefeldt, and Martha B Donaldson. 2009. "Teaching online made me a better teacher" studying the impact of virtual course experiences on teachers' face-to-face practice. *Journal of Computing in Teacher Education* 25, 4 (2009), 121–126.

[37] Catherine E Ross and Chia-Ling Wu. 1996. Education, age, and the cumulative advantage in health. *Journal of health and social behavior* (1996), 104–120.

[38] Md Nazmus Sakib, Megha Yadav, Theodora Chaspari, and Amir H Behzadan. 2019. Coupling virtual reality and physiological markers to improve public speaking performance. In *International Conference on Construction Applications of Virtual Reality*.

[39] Henrique Sequeira, Pascal Hot, Laetitia Silvert, and Sylvain Delplanque. 2009. Electrical autonomic correlates of emotion. *International Journal of Psychophysiology* 71, 1 (2009), 50–56. https://doi.org/10.1016/j.ijpsycho.2008.07.009 Electrophysiology of Affect and Cognition.

[40] Fred Shaffer and Jay P. Ginsberg. 2017. An Overview of Heart Rate Variability Metrics and Norms. *Frontiers in Public Health* 5 (2017).

[41] Fred Shaffer, Rollin Mccraty, and Christopher L. Zerr. 2014. A healthy heart is not a metronome: an integrative review of the heart's anatomy and heart rate variability. *Frontiers in Psychology* 5 (2014).

[42] Vivek Sharma, Neelam Rup Prakash, and Parveen Kalra. 2016. EDA wavelet features as Social Anxiety Disorder (SAD) estimator in adolescent females. *2016 International Conference on Advances in Computing, Communications and Informatics (ICACCI)* (2016), 1843–1846.

[43] Phyllis K. Stein, Matthew S. Bosner, Robert E. Kleiger, and Brooke M. Conger. 1994. Heart rate variability: a measure of cardiac autonomic tone. *American heart journal* 127 5 (1994), 1376–81.

[44] Heine Steinkopf, Dag Nordanger, Anne Halvorsen, Brynjulf Stige, and Anne Marita Milde. 2021. Prerequisites for maintaining emotion self-regulation in social work with traumatized adolescents: A qualitative study among social workers in a Norwegian residential care unit. *Residential Treatment for Children & Youth* 38, 4 (2021), 346–361.

[45] Muhammad Tanveer. 2007. Investigation of the factors that cause language anxiety for ESL/EFL learners in learning speaking skills and the influence it casts on communication in the target language. *University of Glasgow, Scotland* (2007).

[46] Mohsen Tavakol and Reg Dennick. 2011. Making sense of Cronbach's alpha. *International journal of medical education* 2 (2011), 53.

[47] Julian F. Thayer, Bruce H. Friedman, and Thomas D. Borkovec. 1996. Autonomic characteristics of generalized anxiety disorder and worry. *Biological Psychiatry* 39, 4 (1996), 255–266. https://doi.org/10.1016/0006-3223(95)00136-0

[48] Somchanok Tivatansakul and Michiko Ohkura. 2015. Improvement of emotional healthcare system with stress detection from ECG signal. In *2015 37th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*. IEEE, 6792–6795.

[49] Jun-Jie Tseng, Ya-Hsun Tsai, and Rih-Chang Chao. 2013. Enhancing L2 interaction in avatar-based virtual worlds: Student teachers' perceptions. *Australasian Journal of Educational Technology* 29, 3 (2013).

[50] Johannes Wagner, Jonghwa Kim, and Elisabeth André. 2005. From physiological signals to emotions: Implementing and comparing selected methods for feature extraction and classification. In *2005 IEEE international conference on multimedia and expo*. IEEE, 940–943.

[51] Wanhui Wen, Guangyuan Liu, Zhi-Hong Mao, Wenjin Huang, Xu Zhang, Hui Hu, Jiemin Yang, and Wenyan Jia. 2020. Toward Constructing a Real-time Social Anxiety Evaluation System: Exploring Effective Heart Rate Features. *IEEE Transactions on Affective Computing* 11 (2020), 100–110.

[52] Andrea Stevenson Won, Jakki O. Bailey, and Siqi Yi. 2020. Work-in-Progress—Learning about Virtual Worlds in Virtual Worlds: How Remote Learning in a Pandemic Can Inform Future Teaching. In *2020 6th International Conference of the Immersive Learning Research Network (iLRN)*. 377–380. https://doi.org/10.23919/iLRN47897.2020.9155201

[53] Julie L Yee and Debbie Niemeier. 1996. Advantages and disadvantages: Longitudinal vs. repeated cross-section surveys. (1996).

[54] Mingshao Zhang, Zhou Zhang, Yizhe Chang, El-Sayed Aziz, Sven Esche, and Constantin Chassapis. 2018. Recent developments in game-based virtual reality educational laboratories using the microsoft kinect. *International Journal of Emerging Technologies in Learning (iJET)* 13, 1 (2018), 138–159.

[55] Yuanxing Zhang, Zhuqi Li, Chengliang Gao, Kaigui Bian, Lingyang Song, Shaoling Dong, and Xiaoming Li. 2018. Mobile social big data: Wechat moments dataset, network applications, and opportunities. *IEEE network* 32, 3 (2018), 146–153.

[56] Chunping Zheng, Lili Wang, and Ching Sing Chai. 2021. Self-assessment first or peer-assessment first: Effects of video-based formative practice on learners' English public speaking anxiety and performance. *Computer Assisted Language Learning* (2021), 1–34.

[57] Katja Zibrek, Elena Kokkinara, and Rachel McDonnell. 2018. The effect of realistic appearance of virtual characters in immersive environments-does the character's personality play a role? *IEEE transactions on visualization and computer graphics* 24, 4 (2018), 1681–1690.