

# AMuCS: Affective Multimodal Counter-Strike video game dataset

Marios Fanourakis and Guillaume Chanel

**Abstract**—Video games are a versatile and multi-faceted stimulus which can elicit complex player experiences. As a consequence, several datasets have been curated or created for studying human cognition, behaviours, and physiological responses where video games are the primary stimulus. Many of these datasets have a low number of participants or do not have a rich set of modalities and are always recorded in a laboratory setting. To address these issues, we have recorded 256 participants at LAN events while they played the first person shooter, Counter-Strike: Global Offensive. Our dataset consists of several complementary modalities: physiological signals (ECG, EDA, Respiration), behavioural signals (facial expressions, eyetracking, depth images, seat pressure), computer interaction (keyboard and mouse events, game actions), and stimulus information (gameplay video, game logs). We show that the number of participants in our dataset and the variety of modalities recorded is advantageous for training machine learning models.

**Index Terms**—Dataset, affective computing, affective gaming, video games, emotion, multimodal, physiological, arousal, valence, ECG, EDA.

## I. INTRODUCTION

Interactive media experiences such as video games are a versatile and multi-faceted stimulus. Video games are becoming increasingly more realistic with detailed graphics, accurate physics, convincing emotional characters, and immersive virtual reality. As such, video games elicit complex player experiences which makes them an attractive subject for the research community.

Over the past several years the collection of player data has become a serious consideration and necessity for both researchers and developers[1]. Player telemetry allows designers to observe and obtain an accurate representation of several in-game behaviours, which in turn can be used to make important game design decisions and modifications. Telemetry is also often used for matchmaking and ranking players in competitive multiplayer games[2], as a way of providing a more enjoyable experience for players at different skill ranges. However, telemetry often only measures in-game behaviours and thus fails to capture the player behaviours required for getting a full representation of the player state.

Due to the ability of video games to immerse players and to elicit complex emotions, it has been widely studied in the domain of psychology and affective computing where, in addition to in-game measures, various physiological data (ex. ECG, EDA), facial expressions, and body posture are collected[3], [4], [5], [6], [7], [8].

In affective computing, the focus is in the emotional content of the user experience with the goal of understanding user behaviours, or making the experience more enjoyable or engaging by analyzing and adapting the experience to the user state. User emotions can be estimated continuously and in real-time by proxy through the physiological responses which are linked to the autonomous nervous system (ANS) and other behavioural information like facial expressions. A common finding is that it is necessary to perform multimodal recordings in order to capture as much as possible of the user experience.

As such, multimodal recordings of video game play have been researched in various contexts. To study the behavioural effects of video games, for example, to determine if video game play trains the response time or visual search patterns, or if violent video games increase aggressive behaviours. To study if different players (gender, personality, age, etc.) show differences in their affective responses to games. To implement dynamic difficulty adjustment based on the player's affective state and study if this increases the enjoyment of the game. To study the effect of game elements on player experience in order to improve game mechanics and the general player experience. Developers can also benefit from this research by creating emotionally intelligent games and peripherals to better engage players.

Although the usage of physiology to study player behaviour is quite common in the literature, it is rarely done within a multiplayer setting where each player's data streams are collected concurrently. Furthermore, although multimodal data collection has been achieved, we are not aware of a study which has collected all the modalities presented in this paper in a multiplayer scenario.

## II. AFFECTIVE GAMING DATASETS

To facilitate research in the fields of game analysis and psychology, several datasets have been made publicly available (see Table I). The datasets utilize a wide variety of games as a stimulus, from 2D platformers like Super Mario Bros to 3D action adventure games like Assassin's Creed. Very few datasets use multiplayer games and even fewer collect data from all players concurrently. Multiplayer data could be useful for analyzing synchronous behaviours, interaction dynamics, and more.

Typically, the data is collected in a laboratory environment where physiological signals are recorded such as the electrocardiogram (ECG), electrodermal activity (EDA), electroencephalogram (EEG), etc. and contextual data such as the game logs.

Annotations vary between the datasets as well. They can be collected through questionnaires that are administered after the

experience and ask participants to summarize the experienced intensity of a sentiment or other information of interest. A common type of questionnaire for indicating emotions is the self-assessment manikin (SAM) where participants are aided by a pictorial representation of the intensity of the emotions. Another type of annotation is based on a graphical element that represents a one or two-dimensional space of interest and participants can indicate a position anywhere in that space (continuous-space annotation). The arousal-valence space is commonly used to annotate emotional experiences in a continuous-space. This may be accompanied by visual aids of where some categorical emotions lie in this space. Annotations can also be collected in a continuous-time manner for a specific sentiment. Annotators are asked to watch a recording and at the same time use a software to indicate the absolute or relative intensity of a sentiment at each moment. These types of annotations can be used to compare the reactions to specific moments or events throughout the experience. The datasets listed in Table I are annotated in one or multiple of: fun, frustration, engagement, arousal, valence, and others.

Yannakakis et al.[9] collected their dataset about Maze-Ball, a prey/predator game, in which player experience is directly linked to camera settings, to investigate the impact of camera viewpoints on psychophysiology of players through preference surveys collected. Their data was annotated in the affective states of fun, challenge, boredom, frustration, excitement, anxiety and relaxation and they aim to predict the pairwise self-reported emotional preferences of the players. They achieved performances between 63% and 70% to predict the affective state preference using only the physiological signals. By including information about the camera viewpoints they achieved performances between 71% and 85% to predict the affective state references.

Karpouzis et al.[10] collected the Player Experience (PE) dataset about Infinite Mario Bros, a platformer game, along with face videos, self-reports on engagement, frustration, and challenge to model players and player experience. Analysis of facial expressions alongside game events proved to be a good predictor of player experience.

Alchalabi et al.[11] collected the FOCUS dataset about a custom serious game where players can control the movements of an avatar using an EEG device. They showed that they were able to detect players attention state using the EEG data (96% accuracy) and also able to classify players according to if they had ADHD or not (98% accuracy).

Song et al.[12] collected the Multimodal Game Frustration Database (MGFD) about a custom game, Crazy Trophy, a voice controlled 2D maze game, along with face videos and voice. Their goal was to manipulate the game mechanics (control and/or score feedback inconsistencies) to investigate their impact on the level of frustration and predict frustration from the facial expressions and voice prosodic features. Using an RNN-LSTM model and both audio and video features, they achieved an unweighted average recall of 60% for predicting frustration.

Blom et al.[13] collected data about the League of Legends (LoL) video game, a multiplayer open battle arena game, alongside physiological, input device, and game data to study

the relationship between player stress responses and in-game behaviour. They do not present any analysis in that publication.

Kutt et al.[14], [15] collected the BIRAFFE and BIRAFFE2 datasets about several games (affective spaceshooter 2, Freud me out 2 in BIRAFFE and Room of the Ghosts, Jump!, and Labyrith in BIRAFFE2) alongside physiological, input device, and game data aimed at the development of computer models for emotion classification and recognition. They do not present any analysis in these publications.

Beaudoin-Gagnon et al.[16] collected the FUNii database about Assassin's Creed : Unity and Assassin's Creed : Syndicate, third person action/adventure game, alongside physiological, eyetracking, face video, input device data to study physio-behavioral responses in a gaming context, player profiling, player experience modeling or adaptation strategies for affective adaptive games. They do not present any analysis in that publication.

Granato et al.[17] collected the RAGA dataset about Project Cars and RedOut, VR racing games, alongside physiological data to predict player emotions (valence and arousal) during gameplay using physiological data. They achieved a normalized RMSE of approximately 0.21 and 0.2 for arousal and valence respectively using a Gaussian Process Regression (GPR) model.

Alakus et al.[18] collected the multimodal eSports dataset about League of Legends alongside physiological, input device, game, and environmental data to study eSports athletes' psychophysiological data. The data of all players in the team are synchronized allowing for the analysis of data on the team level. They demonstrate that stress and concentration levels for professional players are less correlated compared to amateur players, meaning that professionals displayed a more independent playstyle. They also use the data for skill prediction and player re-identification using 3-minute sessions of sensor data. Best models achieved 0.856 and 0.521 (0.10 for a chance level) accuracy scores on a validation set for skill prediction (using SVM models) and player re-id (using random forest models) problems, respectively.

Smerdov et al.[19] collected the Toadstool dataset about Super Mario Bros alongside physiological, accelerometer and face video data to support research on emotionally aware machine learning algorithms, focusing on reinforcement learning and multimodal data fusion. They achieve an average RMSE of 0.126 using a CNN to predict blood volume pressure (BVP) from game video and face video data.

Alakus et al.[18] collected the GAMEEMO database about Train Sim World, Unravel, Slender-The arrival, and Goat simulator games (each selected to represent boring, calm, horror, and funny types) alongside EEG data to provide an emotion dataset based on computer games and to determine the success of the portable EEG device and compare the success of this device with classical EEG devices. They perform pattern recognition and signal-processing methods to observe the performance of the dataset and to classify EEG signals based on the arousal-valence emotion dimension and positive/negative emotions.

Melhart et al.[20] collected the AGAIN dataset about several custom games (3 racing, 3 shooters, 3 platformers) alongside

physiological, face video, input device, and game data to investigate the generality of affective computing across dissimilar tasks, and for affect modeling within each of its 9 specific interactive games. They achieve accuracies between 58% and 82% for predicting arousal using a random forest classifier with the physiological features.

Dresvyanskiy et al.[21] collected the DyCoDa dataset about a multiplayer survival game alongside audio, video, and depth data to investigate how humans interact with each other online via a video conferencing tool in a cooperation setting and which audio-visual cues are conveyed and perceived during this interaction.

Our dataset (AMuCS) is the largest *multiplayer* dataset by an order of magnitude (245 vs 30 participants) and includes 11 recorded modalities (average in Table I is 5 modalities). It also remains competitive with single player datasets like BIRAFFE, FUNii, and BIRAFFE2 which have a similar number of participants (206, 190, and 103 respectively).

### III. THE AMUCS DATASET

We recorded a total of 256 participants playing video games in realistic conditions (7.66% female, 0.81% non-binary, 0.40% no answer). The participants mean age was 22.68 years old (standard deviation of 3.71, ranging from 18 to 36 years old). The languages spoken by the participants varied between Swiss German, French, and English. The data was collected in 71 experimental sessions where groups of 2 or 4 participants played the *Counter-Strike: Global Offensive* (CS:GO) first person shooter (FPS) video game on a computer. Several modalities were recorded during the game using custom data acquisition software modules:

- mouse/keyboard button presses - recorded at an irregular rate (as button presses occurred).
- game data (health, armor, position, damage taken, damage received, etc.) - recorded at 64Hz using a custom game plugin.
- gameplay video - recorded at 30Hz using *Open Broadcasting Studio* (OBS).
- color and depth video of the face - recorded at up to 30Hz using an *Intel RealSense D435* camera.
- Seat pressure - recorded at 10Hz with a *Sensing Tex* seat pressure mat.
- Physiological data (electrocardiogram, electrodermal activity, respiration) - recorded at 100Hz using a *Bitalino (r)evolution* device
- Eyetracker data (gaze, pupil diameter) - recorded at 60Hz using a *Tobii pro nano*.

A detailed table of the recorded data types is listed in the Appendix A. The gameplay video also includes the gameplay audio and the microphone recordings on the same audio track although occasionally the microphone was not recorded due to technical issues.

All data was synchronized using the Lab Streaming Layer (LSL) software library<sup>1</sup>. More information about the data acquisition software and architecture can be found in [23]

where we measured the synchronization delays to be within 50ms on average after offset corrections. It is important to note that video frames were not sent over LSL. Video such as from the RealSense sensor (color and depth video) and from the gameplay (screen recording) was saved directly on the client PC (the participant's PC) and by means of custom modules we sent the frame number and timestamp of each recorded video frame through LSL. This significantly reduced the bandwidth usage on the local network.

The data was collected on-site at several video game LAN events in Switzerland over the course of two years: SwitzerLAN<sup>2</sup> 2020, SwitzerLAN 2021, PolyLAN<sup>3</sup> 36, and PolyLAN 37. The experimental area was setup in an approximately 5 square meter area within the event and included 4 gaming PCs each equipped with the sensors mentioned earlier and 1 server PC where the game server and LSL Lab Recorder were installed.

The SwitzerLAN 2020 event welcomed approximately 300 gamers in the BernExpo exposition area, a large open space in the city of Bern, Switzerland. This event took place in October of 2020 during the COVID 19 pandemic and various restrictions were in place such as wearing masks when not seated and temperature checks upon entering the LAN area. The physical setup of the experiment is illustrated in Figure 1. The experiment was in a corner of the same area as the LAN (in the big open space), consequently, the environmental lighting and noise was not controlled. Furthermore, the participants were informed that they could keep wearing their mask for the experiment and several did so. Sessions 1 to 14 (14 sessions totaling 41 participants) were recorded at this event.

The SwitzerLAN 2021 event welcomed approximately 1200 gamers in the BernExpo exposition area. This event took place in October 2021 and gamers had to show proof of vaccination, recovery, or a PCR test before entering the LAN area. The physical setup of the experiment was similar to the SwitzerLAN 2020 setup. Mask-wearing was not enforced in this event. Sessions 15 to 33 (18 sessions totaling 67 participants) were recorded at this event.

The PolyLAN 36 event welcomed approximately 200 gamers in a large auditorium of the EPFL campus in Ecublens, Switzerland. This event took place in November of 2021 and had similar restrictions as the SwitzerLAN 2021 event. The physical setup of the experiment was slightly different in this event and is illustrated in Figure 2. The experiment was located at the last row of the auditorium. As with the other events, the environmental lighting and noise was not controlled. Sessions 34 to 41 (7 sessions totaling 32 participants) were recorded at this event.

The PolyLAN 37 event welcomed approximately 1200 gamers in the SwissTech Convention Center in the EPFL campus. This event took place in April 2022 and had similar restrictions as the previous two events. The physical setup is illustrated in Figure 1 (same as in SwitzerLAN 2020 and 2021). The experiment was located in an office adjacent to the LAN area, this resulted in much less environmental noise from

<sup>1</sup><https://labstreaminglayer.org>

<sup>2</sup><https://switzerlan.ch/>

<sup>3</sup><https://polylan.ch/>

Dataset	#Part.	Game	Modalities	Annotations	Notes
Maze-Ball[9] 2013	36	Maze-Ball (3D prey/predator)	BVP, SC	pairwise, comparing conditions more/less anxious, exciting, frustrating, fun and relaxing	environment: in lab; conditions: control style, camera angle (8x conditions); gameplay duration: 90s per condition (12m total)
PED[10] 2015	58	Infinite Mario Bros	face video, game logs	self-reported level of engagement, frustration and challenge (scale 0-4)	environment: in lab; conditions: level "A" and level "B"; gameplay duration average: 1m per level (6h for 380 game sessions)
FOCUS[11] 2018	9	FOCUS (custom)	EEG (Emotiv EPOC+)	-	environment: in lab; conditions: control character via keyboard vs EEG, ADHD vs non-ADHD subjects; gameplay duration average: 1m (keyboard), 2.5m (EEG, non-ADHD), 4m (EEG, ADHD)
MGFD[12] 2019	67	Crazy Trophy (custom voice controlled game)	face video, speech	self-reported frustration 4 point Likert scale; 4 external continuous annotations	environment: in lab; conditions: game control and feedback inconsistency; gameplay duration average: 45s per level (4.5m total for 6 levels)
[13] 2019	8	League of Legends (multiplayer)	keyboard, mouse, HR, GSR, PPG, game logs	-	environment: in lab
BIRAFFE[14] 2019	206	Affective SpaceShooter 2, Freud me out 2	ECG, GSR, face video, joystick input, game logs	self-reported valence and arousal	environment: in lab; conditions: the different games
FUNii[16] 2019	190	Assassin's Creed: Unity, Assassin's Creed: Syndicate	ECG, EDA, Respiration, EMG, eyetracking, face video, controller inputs	continuous self-reported fun	environment: in lab; gameplay duration max: 35m
RAGA[17] 2020	36	Project Cars, RedOut	ECG, fEMG, GSR, Respiration	continuous self-reported arousal and valence	environment: in lab; conditions: VR vs standard monitor
Multimodal eSports[19] 2020	10 (2x5)	League of Legends (multiplayer)	HR, IMU, EMG, GSR, eyetracker, EEG (Emotiv), pulse oximeter, keyboard, mouse, infrared (FLIR), game logs, environmental: temperature, humidity, and CO <sub>2</sub>	self-evaluation of own performance and teammates performance	environment: in lab; conditions: professional vs amateur team, vs bots or humans, with vs without team communication
Toadstool[22] 2020	10	Super Mario Bros	accelerometer, temperature, EDA, BVP, IBI, HR, face video, gameplay video	questionnaire	environment: in lab; gameplay duration: 35m
GAMEEMO[18] 2020	28	3x games: Train Sim World, Unravel, Slender - The Arrival, Goat Simulator	EEG (Emotiv)	SAM (arousal and valence)	environment: in lab; conditions: the different game types (boring, calm, horror, funny); gameplay duration: 5m per game (20m total)
AGAIN[20] 2021	124	custom games: 3x racing, 3x shooters, 3x platformers	gameplay video, game logs	continuous and unbounded self-reported arousal	environment: in lab; conditions: the different games; gameplay duration: 2m per game (18m total)
BIRAFFE2[15] 2022	103	3x 2D games: Room of the Ghosts, Jump!, Labyrinth	ECG, GSR, face video, gamepad input, game logs	self-reported valence and arousal	environment: in lab; conditions: the different games; gameplay duration max: 5m per game (15m total)
DyCoDa[21] 2022	30	Survival Game (multiplayer)	infrared video, depth video, gameplay video and audio, speech audio	self-evaluation questionnaires	environment: in lab; total 10h of recorded data; collaborative problem solving
AMuCS (ours)	245*	Counter-Strike: Global Offensive (multiplayer)	keyboard, mouse, ECG, EDA, Respiration, eyetracker, face features, depth video, seat pressure, gameplay video and audio, game logs	continuous and unbounded self-reported arousal or valence	environment: in-the-wild; gameplay duration: 10m

TABLE I

OPEN VIDEO GAME DATASETS WITH PHYSIOLOGICAL OR AFFECTIVE MODALITIES

the rest of the LAN as well as consistent lighting. Sessions 42 to 71 (29 sessions totaling 116 participants) were recorded at this event.

The study and data collection was approved by the ethical committee of the University of Geneva and conforms to all ethical guidelines set forth by the institution.

245 participants (out of 256 recorded) accepted to share their *anonymized* data with other research institutions. The AMuCS dataset consists of these 245 participants and access can be requested at <https://doi.org/10.26037/yareta:gvvoc4wfsfhupm4ygge26wupnm>.

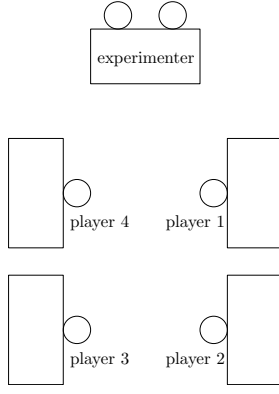


Fig. 1. Physical setup for SwitzerLAN 2020/2021 and PolyLAN 37.

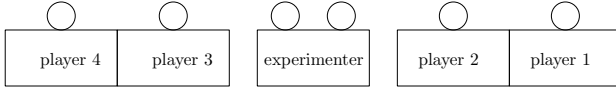


Fig. 2. Physical setup PolyLAN 36.

#### A. The Counter-Strike: Global Offensive game

In our experiments we utilized Valve’s Counter-Strike: Global Offensive (CS:GO). It is a free and modable multi-player first person shooter (FPS) developed in the Half-Life 2 game engine. It is also popular in the e-sport community. The game includes several games modes: demolition, hostage, deathmatch, and team deathmatch.

The experiment used the *team deathmatch* game mode where two teams try to eliminate each other. Players started with a 2 minute warmup round, where they could explore the game map and test their opponents without counting the score. After the warmup players were respawned to random locations and frozen in place for 1 minute. Once the freeze time was over the main round round started and had a duration of 10 minutes. Each player started with 100 health points and 100 armor points and were randomly placed on the game map. They were equipped with a random set of weapons from an assortment of assault rifles, long range rifles, pistols, light and heavy machine guns, and a knife. The goal of the game was to kill the players in the enemy team as many times as they can while avoiding to get killed. A player was killed once their health points reached 0, and they were subsequently revived (respawned) at a random location in the map after 2 seconds. If a player managed to get 2 kills in a row without dying they were rewarded with an item (healthshot) which restored 50 health points when used.

The game data was recorded using our custom *sourcmod*<sup>4</sup> plugin which enabled us to send the data on an LSL stream to be synchronized with the other experimental data. OBS was also used to record the screen and sound of both the game and the participant.

#### B. Experimental protocol

Groups of 2 or 4 participants were spontaneously recruited at the LAN events to play a single round of one versus one

(2 participants per experimental session) or two versus two (4 participants per experimental session) team deathmatch. The participants first read and signed a consent form which describes the experiment and the data that will be recorded. Immediately after, they answered a questionnaire with demographic questions (age, sex) as well as questions about their experience playing various types of video games (Bavelier lab Video Game questionnaire<sup>5</sup>), their fatigue level[24], and their closeness of relationship with the other participants in the experiment group[25].

Before attaching the ECG electrodes, the participants were given cotton soaked in alcohol in order to clean the electrode locations on their body. The 3-lead ECG sensor consisted of adhesive wet electrodes (with conductive gel) and was attached in the Einthoven triangle pattern. Specifically, the positive and negative leads were attached half way between the shoulder and the sternum on the left and right collar bone respectively, and the reference lead was attached on the right side just below the rib cage. The 2-lead EDA sensor consisted of dry electrodes secured by velcro straps around the proximal phalanx of the index and middle finger of the “keyboard” hand (the hand used to move the player character). The specific fingers and electrode locations were chosen to be the least obstructing for the participants while still having a signal of adequate quality. The “keyboard” hand does not move as much as the “mouse” hand and we found there were fewer movement artifacts at this location. The respiration belt (a stretch sensor) was attached over the shirt of the participant around their body and over the diaphragm and just under the breast area.

The participants were asked to sit approximately as they would be sitting during gameplay and then we used the Tobii eyetracker manager software<sup>6</sup> to adjust the screen distance and angle such that the eyetracker could reliably detect the eyes. The Tobii eyetracker was mounted on the bottom bezel of the screen. It was calibrated using a five point calibration procedure using the same software and the results were verified by asking the participant to look at various points on the screen. If there were significant errors in the gaze position, the calibration procedure was repeated.

Players adjusted their game settings (ex. mouse sensitivity, UI elements, keyboard bindings, screen resolution and aspect ratio) before joining the experiment’s game server where they started with a 2 minute warmup round. After the warmup round, there was a 1 minute period of time reserved for baselining where the game characters were frozen in place. After baselining, the main round started and had a duration of 10 minutes.

After the game was finished, the sensors were removed and the players used the PAGAN tool[26] to self-annotate their own recorded gameplay video according to the arousal or valence emotional dimensions using RankTrace[27], a relative and unbounded method for continuous annotations. The gameplay video served as a recall aid and we did not show the video of their own face so as not to bias the annotations towards facial expressions. Participants only annotated one of

<sup>4</sup><https://www.sourcmod.net/>

<sup>5</sup><https://www.unige.ch/fapse/brainlearning/vgq/>

<sup>6</sup><https://www.tobii.com/product-listing/eye-tracker-manager/>

arousal or valence, not both. This was done to reduce the cognitive load of the annotation process (in the case when a two dimensional annotation is used) with the goal of producing higher quality annotations. We also chose not to have the participant annotate the video twice (once for each dimension) due to time constraints for the experiment.

### C. Data Quality

All data modalities were visually inspected for determining their quality. We grouped data into four categories: not usable, partial, usable, and good. The first level of inspection was simply to determine how much of the data was missing or outside of an acceptable SNR during the main round of the game (10 minutes of data). The SNR was mainly a concern with the EDA, ECG, and respiration signals. For the EDA signal, we accounted for the presence and visibility of phasic responses versus other signal artifacts. For the ECG signal, we accounted for the visibility of the QRS complex or R-peaks versus other signal artifacts. The signals was tentatively labeled as:

- good - if less than 10% of the signal was missing or outside of SNR range
- usable - if less than 33% of the signal was missing or outside of SNR range
- partial - if less than 50% of the signal was missing or outside of SNR range
- not usable - if more than 50% of the signal was missing or outside of SNR range

Next we looked more closely at the quality of some of the signals. For EDA, we accounted for the resolution of the signal. Participants' base skin conductivity varied significantly and several had intrinsically low skin conductivity. This meant that their phasic responses had very low amplitude and could be more challenging to analyze but still usable, hence we labeled such EDA signals as usable. In several cases the resolution was too low to discern any phasic responses and these were labeled as not usable. At the other extreme, there were several participants whose skin perspiration was too high to be measured by the instrument, these were also labeled as not usable. If artifacts were visible such as those that may occur if the electrodes were influenced by the movement of the fingers or while pressing buttons on the keyboard, the signal was labeled as usable.

For ECG, we accounted for the visibility of the QRS complex. If all parts of the QRS complex were visible, then we labeled the data as good, if only the R-peaks were visible we labeled the data as usable, otherwise we labeled it as not usable or partial.

For respiration, we mainly focused on artifacts from the heart rate. If these artifacts had amplitude larger than 10% of the amplitude changes caused by breathing then the signal was labeled as usable, otherwise as good.

For the face video, if the participant's face was fully visible throughout the session it was labeled as good, if the face was partially covered (ex. wearing a mask) or otherwise not fully visible it was labeled as partial.

In Table II we summarize the number of usable data for some combinations of modalities.

signals	#sessions	#participants
ECG	67 (94%)	192 (78%)
EDA	65 (92%)	144 (58%)
Respiration	69 (97%)	210 (85%)
eyetracker	70 (98%)	233 (95%)
face video	64 (91%)	218 (89%)
gameplay video	69 (97%)	227 (92%)
keyboard buttons	67 (94%)	216 (88%)
mouse buttons	67 (94%)	221 (90%)
game logs	68 (95%)	234 (95%)
valence annotations	65 (91%)	111 (45%)
arousal annotations	68 (95%)	117 (47%)
face video + valence annotations	59 (84%)	102 (41%)
face video + arousal annotations	64 (90%)	108 (44%)
ECG + EDA + game logs	62 (87%)	132 (53%)
ECG + EDA + game logs + arousal annotations	50 (70%)	64 (26%)
ECG + EDA + game logs + face video + eyetracker	60 (84%)	126 (51%)
eyetracker + gameplay video + game logs	66 (92%)	218 (88%)
eyetracker + gameplay video + game logs + arousal annotations	63 (88%)	109 (44%)
eyetracker + gameplay video + game logs + ECG + EDA	59 (83%)	123 (50%)

TABLE II  
NUMBER OF SESSIONS (OUT OF 71) AND PARTICIPANTS (OUT OF 245)  
WITH USABLE DATA IN THE AMUCS DATASET

### D. Data pre-processing

The LSL LabRecorder software records data in the extensible data format (XDF). The files contain all local timestamps as well as timing information that facilitates the accurate synchronization of the data streams. We used the python xdf module (pyxdf<sup>7</sup>) to read the files and automatically apply the timestamp synchronization. We then used the pandas<sup>8</sup> module in python to format the data into easily queried data structures and to ultimately convert the xdf files to parquet<sup>9</sup> or CSV files. Since the annotations of the gameplay were performed after data collection they had to be aligned with the rest of the data. This was achieved by using the synchronized frame timing information of the gameplay that was recorded with LSL.

For convenience and ethical concerns, we derived some additional features that may be relevant and are either computational complex or rely on data which will not be made public such as the face videos. These features include the luminance of the screen, in-game combat and danger level indicator, and facial features like action unit (AU) activations.

1) *Screen Luminance*: To compute the perceived lightness (luminance) of each pixel on the screen the Lstar from CIELAB [28] was computed using the RGB values. First, the RGB values were converted from gamma encoding to linear encoding, then the standard coefficients for sRGB (0.2126, 0.7152, 0.0722 for R, G, and B respectively) were applied to compute the RGB luminance. Finally, the RGB luminance was converted to the perceived lightness, Lstar, which closely matches human light perception. It is important to note that Lstar does not take the Helmholtz–Kohlrausch effect [29] into account wherein the intense saturation of spectral hue is perceived as part of the color's luminance.

<sup>7</sup><https://github.com/xdf-modules/pyxdf>

<sup>8</sup><https://pandas.pydata.org/>

<sup>9</sup><https://parquet.apache.org/>

Having the Lstar value for each pixel, we then averaged the Lstar pixel values within an 8 degree horizontal foveal area of the screen centered at the gaze target of the participant. We used a rectangular area instead of circular since it simplified our computations. Under our experimental conditions this foveal area was approximately a 16cm by 9cm rectangular region on the screen (same aspect ratio as the screen). We did this for each frame of the video recording, always centering on the gaze target at each frame using the eye-tracking data.

In the dataset, we provide the mean luminance of the entire screen, the mean luminance of the gaze region and the central region of the screen as well as the mean luminance of the screen excluding the gaze region and excluding the central region.

Note that the Lstar luminance measures the pixel activations and does not correspond to the absolute luminance as measured by a luminance meter in units of candela per square meter or *blondel*. However, this information can still be useful when analyzing the pupil size and to attenuate the pupil light response from the pupil data like in Fanourakis et al. [30].

2) *Combat and other special game events*: From the game event data we computed some special indicators/events such as the number of enemies that are: in the field of view, in close range ( $< 500$  game distance units), in mid range ( $< 1000$  game distance units). We also computed a health danger indicator indicating if the health is below 70%, 50%, or 30%.

We labeled gameplay as combat if the player had received or dealt damage from/to another player within a 5 second window. Once there was a death event (either the player was killed or the enemy was killed), the combat state was reset even if it was within a 5 second window of the combat events previously mentioned.

The game event combinations mentioned above were selected among other recorded game events based on their relevance towards the game's two main goals: staying alive and killing the enemy. Combat is directly relevant to the two goals since the most common outcome of combat will either be a goal success (enemy killed) or a goal failure (player death). The "health danger" indicator gives an indication of the probability of achieving or failing the goals. All else being equal, a player with lower health will be killed more quickly during combat. The "number of enemies in field of view" event puts the player in the position to seek goal resolution by taking action to stay alive and/or kill the enemy. The "number of enemies in close/mid range" can also give an indicator of the amount of danger that a player might be in. In conjunction with these combined events some other simple events can be useful such as "reloading weapon" and "jumping" which prevent the player from making a fight or flight decision since in the former, the player is not able to fire their weapon for some time (cannot fight) and in the latter the player cannot move to take cover and save themselves from harm.

This type of information can be useful when analyzing the game context and summarizing the various events into different phases of the game like in Weber et al. [31].

3) *Face features*: We applied Baltrušaitis' OpenFace[32] feature extraction on the color video of the face to extract the following features: gaze, facial landmarks, head pose, and

continuous facial action unit activations. Although gaze is tracked by the Tobii eyetracker, this additional gaze estimate can be useful in the rare cases when the eyetracker failed to track the participant's eyes due to bad placement or movement outside of the eyetracker's operating range.

#### IV. ANALYSIS

In this section we will show the benefits of the large number of modalities and participants in our dataset. We will see how the number of modalities and number of participants in the training set influences machine learning prediction performance. This also gives a baseline of the predictive capacity of our dataset. Note that we did not aim to achieve competitive results and used classical machine learning methods such as gradient boost.

##### A. Multimodal prediction of game events

FPS games are typically fast paced with several different game events happening in rapid succession and in bursts. These game events can elicit physiological and behavioural responses from the players. This relationship allows us to use our dataset to predict game events from physiological signals.

The main challenge is that the physiological responses are not at the same cadence as game events. The elicited fluctuations in physiology through the ANS are generally at much lower frequencies (generally below  $0.5Hz$ )[33], [34] than game events from a fast paced game (bursts of events can be well above  $2Hz$  in our dataset). To facilitate the process we may analyze the game events to derive general game states. Weber et al. did so by defining several game phases based on game micro-events for the game "Tactical Ops: Assault on Terror" [31]. They defined a total of 6 phases: use of in-game menu, safe, danger (enemy in field of view), combat (player uses weapon), under attack, ghost mode (player has died and is viewing the arena in 3rd person). They then defined events in the context of these phases and found significant differences of heart rate responses to these game events.

We will proceed in a similar strategy but define only two phases: safe and danger. We added some complexity to the detection of the danger phase of Weber et al. by not only taking into account the enemies in the field of view but also the distance to the enemy, the current health of the player, and if the player is reloading their weapon or jumping. We also merged combat, and under attack phases of Weber et al. with the danger phase. We discarded the ghost phase since this was not enabled in our game, and we also discarded the use of in-game menu phase since the menu was used relatively rarely. Details on how we computed combat and danger can be found in the previous Section III-D2.

We wanted to verify that there is a perceived difference of emotional arousal between the game phases so we compared the participants' arousal annotations during portions of the game when the players were safe versus when they were in danger according to the previously computed game phases. The results are shown in Figure 3 where the mean arousal (z-score) is  $-0.36$  and  $0.12$  during the safe phase and danger phase respectively. A one-sided Mann-Whitney U test shows

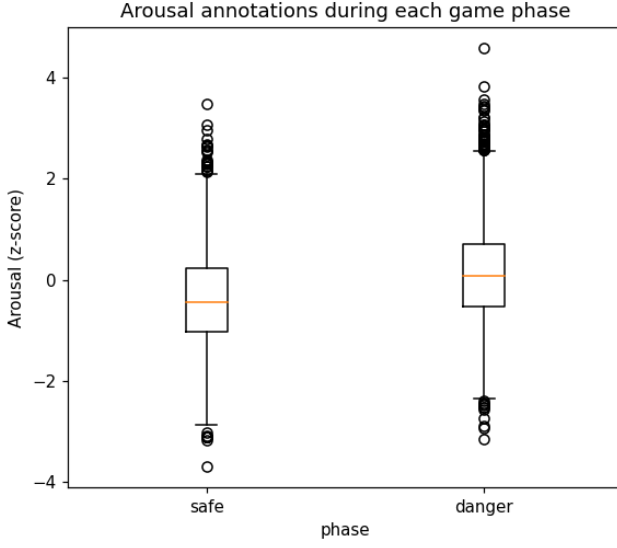


Fig. 3. Participant annotations during each of the game phases. Arousal annotations during the danger phase tend to be higher compared to the safe phase. Orange markings indicate the median value.

that the increase of the arousal annotation values during the danger phase is statistically significant (p-value smaller than 0.001).

The physiological modalities that we used for predicting the two phases were: the EDA, heart rate (HR), respiration, facial action units, on-screen gaze speed. The modalities were normalized per participant with a z-score. We used the signals of 121 participants who had usable data for these modalities.

A low pass filter with a cutoff of  $5Hz$  was applied to the EDA signal to remove high frequency noise. The same filter was also used for the respiration signal.

The instantaneous heart rate was computed from the ECG signal after a low pass filter of  $45Hz$  (to remove high frequency noise) by applying Hamilton method of R-peak detection of the biosppy Python package<sup>10</sup>, then computing the peak rates and smoothing them using a boxcar smoother of length 10.

The on-screen gaze speed was computed from the gaze data by measuring the distance between consecutive gaze points on the screen. This feature gives an indication of saccade and fixation behaviours without information about where on the screen the player is looking at.

We extracted several features from each of these signals by using a 15 second rolling window with a step size of 10 seconds. The features computed within each window are summarized in Table III.

We similarly windowed the game phase signals and extracted only the maximum value in the windows. Due to the fast paced nature of the game the players find themselves more frequently in the danger phase than the safe phase, resulting in imbalanced classes. Across all the windowed game phase signals and participants we had a total of 1050 instances of the safe phase class and 6762 instances of the danger phase class.

Feature	Description
mean	the mean value of the signal
variance	the variance of the signal
range	the difference between the maximum and minimum values of the signal
minimum	the minimum value of the signal
maximum	the maximum value of the signal
first value	the first value of the signal
last value	the last value of the signal
centroid	1D center of mass of the signal, indicative of the slope of the signal but less computationally expensive
number of peaks	the number of peaks in the signal detected from the inflection points of the derivative
peaks mean	the mean of the peak amplitudes
peaks variance	the variance of the peak amplitudes
number of valleys	the number of valleys in the signal detected from the inflection points of the derivative
valleys mean	the mean of the valley amplitudes
valleys variance	the variance of the valley amplitudes

TABLE III

FEATURES DERIVED FOR EACH WINDOWED SIGNAL.

We then used a gradient boost classifier<sup>11</sup> with leave-one-participant-out cross validation to predict the game phase from the physiological modalities. We used the default parameters of the model: log loss, learning rate of 0.1, 100 estimators, Friedman MSE criterion, maximum depth of 3. The results are summarized in Table IV where we report the mean f1 score and Cohen's kappa of the test sets.

We performed one-sided paired Wilcoxon statistical tests to determine if the increase in performance (f1 score) of the models was statistically significant with p-value smaller than 0.01 (\*) or 0.001 (\*\*). The results of the statistical tests are summarized in Table V.

On their own, EDA, gaze speed, and the facial action units have the best performance and when combined they result in significant performance gains for game phase prediction. With all the modalities together we reach an f1 score of 0.75 (Cohen's Kappa of 0.52). The HR signal and respiration perform poorly on their own, and when combined with other signals the performance is not affected significantly. A more diverse set of features or a different machine learning model may be able to utilize these signals more effectively.

group	modalities	f1 score	Cohen's kappa
A	HR	0.50	0.06
B	EDA	0.68	0.38
C	Respiration	0.47	0.02
D	on-screen gaze speed	0.63	0.30
E	Facial action units	0.67	0.36
F	HR, EDA	0.69	0.40
G	HR, EDA, respiration	0.69	0.41
H	HR, EDA, respiration, gaze speed	0.72	0.45
I	HR, EDA, respiration, facial action units	0.73	0.48
J	HR, EDA, respiration, facial action units, on-screen gaze speed	0.75	0.52

TABLE IV

GAME PHASE PREDICTION RESULTS, LEAVE-ONE-PARTICIPANT-OUT CROSS VALIDATION, 121 PARTICIPANTS

<sup>10</sup><https://biosppy.readthedocs.io/en/stable>

<sup>11</sup><https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.GradientBoostingClassifier.html>



↓better than→	A	B	C	D	E	F	G	H	I	J
A	-									
B	**	-	**	*						
C			-							
D	**		**	-						
E	**		**	*	-					
F	**		**	**		-				
G	**		**	**	*		-			
H	**	**	**	**	**	*	*	-		
I	**	**	**	**	**	*	*		-	
J	**	**	**	**	**	**	**	*	*	-

TABLE V

ONE-SIDED PAIRED WILCOXON STATISTICAL TEST OF GAME PHASE MODEL PERFORMANCE WITH DIFFERENT TRAINING MODALITIES. \*\* P-VALUE SMALLER THAN 0.001, \* P-VALUE SMALLER THAN .01.

Although there is extensive literature showing that heart rate is correlated with emotional arousal, our models were not able to distinguish between the safe phase and the danger phase despite that their arousal annotations had a statistical difference. One potential reason could be that the set of features we extracted from the heart rate were not appropriate. Indeed, heart rate variability (HRV) features are more often used in the literature and the lack of such features in our case could be the reason why our models performed poorly for this modality. Another potential reason could be that the heart rate fluctuations can be induced by both the sympathetic system and the parasympathetic system thus making the heart rate response origin uncertain between emotional stimuli or other functions. In combination with the complexity of the game stimuli, it could be the case that heart rate on its own is not enough to determine the game phase. In the literature, experiments showing the effects of arousal on the heart rate typically use a single unambiguous emotional stimulus and enough data is recorded (more than 40s) to capture the low frequency fluctuations ( $0.05Hz$  to  $0.15Hz$ ) of the heart rate which are related to the sympathetic nervous system. On the other hand, the EDA is directly linked to the sympathetic nervous system and emotional stimuli have a more direct influence[35].

Despite these shortcomings, it is evident from our results that there are statistically significant improvements in the performance of models when including multiple diverse modalities. The multimodal aspect of this dataset can therefore be a valuable attribute which can be utilized to train state of the art models.

### B. Prediction of emotional arousal annotations

Machine learning models often generalize better when trained with larger datasets. AMuCS is the largest multimodal dataset with continuous affect annotations and we will use it to show that model performance improves significantly as we increase the amount of data that is available for training. We will focus on emotional arousal as the target for our machine learning models to validate the continuous emotional annotations at the same time. Due to the high number of iterations necessary for statistical tests we must make sure the model is simple and will only use EDA and HR as input signals which further limits the target to arousal.

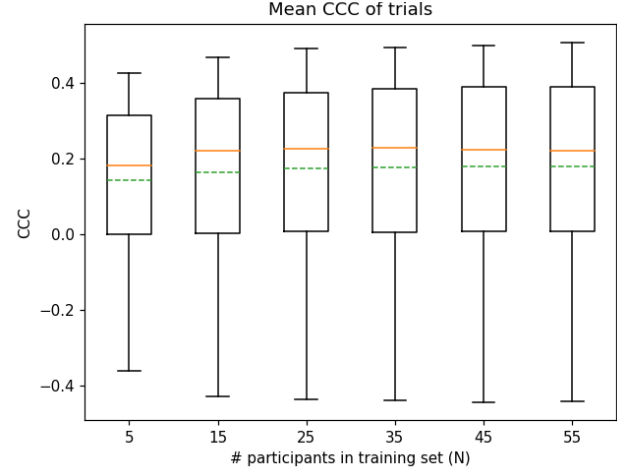


Fig. 4. Boxplots for each  $N$  of participants' random trials mean CCC. Orange markings indicate the median value, green markings indicate the mean value.

We fit a gradient boost regressor<sup>12</sup> in Python using an increasing amount of training data  $N$ , ranging from 5 training participants to 64 training participants in increments of 10. We have access to one more participant compared to what is published and summarized in Table II for the relevant modalities since some participants did not give their consent to share their data with other research institutions and are not included in AMuCS. We used the default parameters of the model: squared error loss, learning rate of 0.1, 100 estimators, Friedman MSE criterion, maximum depth of 3. We performed leave-one-out cross validation. For each left-out participant and for each total training data size  $N$ , we randomly selected  $N$  participants from the remaining 64 participants. We then fit the regressor using this random selection, and tested on the test participant. We repeated the random selection of  $N$  participants (with replacement) to fit and test new models until the mean of the CCC between all random trials was stable (i.e. the measured mean was within  $\pm 0.005$  of the real mean with 99% confidence) or a maximum of 1000 trials was reached.

We used the EDA and HR (computed from ECG as in Section IV-A) as input and the arousal annotation as target. The input and target modalities were normalized per participant (z-score). We used a window of size 7 seconds and step size 5 seconds for extracting the features.

In Table VI and Figure 4 we report the mean CCC of the random trials between all left-out participants for each training data size  $N$ . We observe that as we increase  $N$ , the mean CCC is increased. We performed one-sided paired Wilcoxon statistical tests to confirm that the improvement of the mean CCC as we increase  $N$  is statistically significant with a p-value smaller than 0.001 (\*\*) except between values of  $N = 35$  and  $N = 45$  where the p-value is smaller than 0.01 (\*).

One-sided paired Wilcoxon statistical tests show that the performance improvement as we increase the number of participants in the training set is significant. The results are summarized in Table VII.

<sup>12</sup><https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.GradientBoostingRegressor.html>

N	CCC mean
5 participants	0.141
15 participants	0.164
25 participants	0.172
35 participants	0.175
45 participants	0.177
55 participants	0.180
64 participants	0.181

TABLE VI

AROUSAL PREDICTION RESULTS, LEAVE-ONE-PARTICIPANT-OUT CROSS VALIDATION

↓better than→	5	15	25	35	45	55
5	-	**				
15	**	-				
25	**	**	-			
35	**	**	**	-		
45	**	**	**	*	-	
55	**	**	**	**	**	-

TABLE VII

ONE-SIDED PAIRED WILCOXON STATISTICAL TEST OF AROUSAL MODEL PERFORMANCE WITH DIFFERENT TRAINING SET SIZE  $N$ . \*\* P-VALUE SMALLER THAN 0.001, \* P-VALUE SMALLER THAN .01.

In Figure 4 we also observe an increase in the variance between the participants' results. The reason for this is that, although the CCC increases as we increase  $N$ , it does not increase equally across the participants. To verify this, we plot the change in CCC of each participant in Figure 5. We can clearly see that the participants have very different changes in performance. Most improve (green curves), some have no statistically significant change (blue curves), and for a few, the performance decreases (red curves), this explains the increase in variance that we saw in Figure 4 even though the results improve on average as we increase  $N$ .

In Figure 5, we observed that for some participants (8 participants out of 65), increasing  $N$  from 5 to 55 results in a statistically significant **decrease** in performance. Upon further investigation we have found that for 6 of those participants, the performance was generally bad (CCC below 0.1). This could be caused by poor quality annotations, for example if the participant did not understand the task. The other 2 participants had an above average CCC and we have not determined the precise cause of this performance decrease.

## V. CONCLUSION

We have presented our dataset, AMuCS, consisting of several modalities (EDA, ECG, Respiration, eyetracking, video, game data, continuous annotations of arousal or valence, and more) and collected from 256 participants over 4 LAN events in Switzerland. With data from 245 participants published, it is the largest dataset of its kind and has the following advantages over other existing datasets: more recorded modalities, a large number of participants, recorded on location at LAN events (in-the-wild).

We reported the number of usable data for different combinations of modalities so that researchers can have a better idea what to expect when using the AMuCS dataset. Depending on the desired modalities, researchers can expect to utilize data consisting of at least 64 participants to over 230 participants.

We showed that machine learning models which utilize more input modalities tend to perform better. We used a

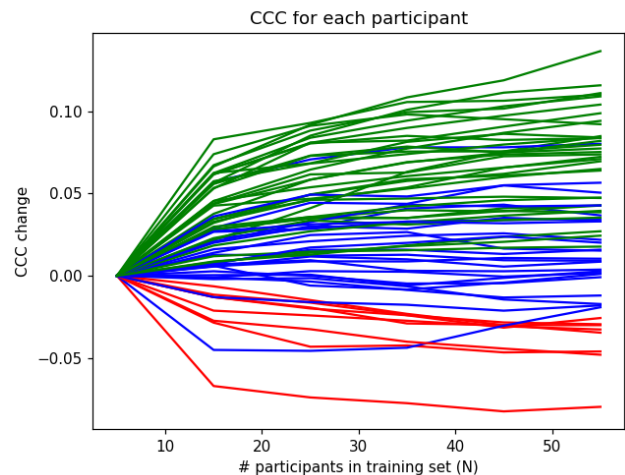


Fig. 5. Change in CCC for each participant vs  $N$ . Green curves show statistically significant increase (33 participants), red curves show statistically significant decrease (8 participants), blue curves do not have statistically significant changes (24 participants)

gradient boost model to classify the game phases (safe and danger). When using just one modality, the best f1-score we achieved was 0.68 with EDA. When using combinations of modalities, the best f1-score was 0.75 and was achieved when we combined all the modalities that we considered for this analysis (HR, EDA, respiration, facial action units, and on-screen gaze speed).

We also showed that the amount of training data (participants) can have a positive impact on machine learning performance. We used a gradient boost regressor to estimate continuous arousal from HR and EDA. We achieved a mean increase in CCC of 0.017 when we trained the model using data from 15 participants (CCC of 0.164) versus 64 participants (CCC of 0.181).

The use of simple features and models, allowed us to train the models with a sufficient number of iterations to achieve statistically significant results for the purposes of comparing the number of modalities and number of training data without the need for excessive resources (computational and time). We expect that more interesting results can be achieved with more complex features and models.

AMuCS includes some unique modalities. A depth video of the participant which can be used to analyze their movement, and a seat pressure mat which can be used to analyze their posture.

This dataset is rich and can be exploited in many ways: player performance analysis, player modeling, affect recognition, attention analysis, game event recognition, e-sport analytics, behaviour synchronicity in collaborative/competitive scenarios, and more.

## ACKNOWLEDGMENTS

This work was funded by Innosuisse with grant number 34316.1 IP.ICT.

The authors would like to thank the organizers of Switzerland and PolyLAN for accommodating our study in their events.

## APPENDIX A

signal	sensor	Hz	notes	raw file	processed file	published
ECG	Bitalino (r)evolution BT	100	10-bit resolution	xdf	csv	Y
EDA	Bitalino (r)evolution BT	100	10-bit resolution	xdf	csv	Y
respiration	Bitalino (r)evolution BT	100	10-bit resolution	xdf	csv	Y
left eye x gaze	Tobii pro nano	60		xdf	csv	Y
left eye y gaze	Tobii pro nano	60		xdf	csv	Y
left eye pupil diameter	Tobii pro nano	60		xdf	csv	Y
right eye x gaze	Tobii pro nano	60		xdf	csv	Y
right eye y gaze	Tobii pro nano	60		xdf	csv	Y
right eye pupil diameter	Tobii pro nano	60		xdf	csv	Y
seat pressure	Sensing Tex seat pressure mat	10	10-bit resolution; 16x16 sensor grid	xdf	csv	Y
face color video	Intel RealSense D435	15/30 <sup>13</sup>	RGB24; 640x480 pixels	rosbag + xdf	mkv + csv	N
face depth video	Intel RealSense D435	30	gray16le; 640x480 pixels	rosbag + xdf	mkv + csv	Y
openFace features	openFace[32]	15/30		N/A	csv	Y
gameplay video	OBS	30	same as game reso- lution	mp4 + xdf	mp4 + csv	Y
keyboard buttons	keyboard	N/A		xdf	csv	Y
mouse buttons	mouse	N/A		xdf	csv	Y
game - isDucking	CS:GO plugin	64		xdf	csv	Y
game - isJumping	CS:GO plugin	64		xdf	csv	Y
game - isReloading	CS:GO plugin	64		xdf	csv	Y
game - health	CS:GO plugin	64		xdf	csv	Y
game - armor	CS:GO plugin	64		xdf	csv	Y
game - bulletShots	CS:GO plugin	64		xdf	csv	Y
game - damageToEnemy	CS:GO plugin	64		xdf	csv	Y
game - damageFromEnemy	CS:GO plugin	64		xdf	csv	Y
game - inFOV1/2/3/4	CS:GO plugin	64		xdf	csv	Y
game - aimTarget	CS:GO plugin	64		xdf	csv	Y
game - position X/Y/Z	CS:GO plugin	64		xdf	csv	Y
game - velocity X/Y/Z	CS:GO plugin	64		xdf	csv	Y
game - eyeVector X/Y/Z	CS:GO plugin	64		xdf	csv	Y
valence/arousal annotations	PAGAN (gameplay vid.)	N/A		csv	csv	Y

TABLE VIII

LIST OF RECORDED DATA. ONLY A SUBSET OF THE GAME DATA IS LISTED

## REFERENCES

- [1] A. Drachen, A. Canossa, and G. N. Yannakakis, "Player modeling using self-organization in Tomb Raider: Underworld," in *2009 IEEE Symposium on Computational Intelligence and Games*. IEEE, sep 2009, pp. 1–8. [Online]. Available: <http://ieeexplore.ieee.org/document/5286500/>
- [2] O. Delalleau, E. Contal, E. Thibodeau-Laufer, R. C. Ferrari, Y. Bengio, and F. Zhang, "Beyond Skill Rating: Advanced Matchmaking in Ghost Recon Online," *IEEE Transactions on Computational Intelligence and AI in Games*, vol. 4, no. 3, pp. 167–177, sep 2012. [Online]. Available: <http://ieeexplore.ieee.org/document/6156756/>
- [3] J. M. Kivikangas, G. Chanele, B. Cowley, I. Ekman, M. Salminen, S. Järvelä, and N. Ravaja, "A review of the use of psychophysiological methods in game research," *Journal of Gaming & Virtual Worlds*, vol. 3, no. 3, pp. 181–199, sep 2011. [Online]. Available: [http://www.ingentaconnect.com/content/10.1386/jgvw.3.3.181\\_{\\_}1](http://www.ingentaconnect.com/content/10.1386/jgvw.3.3.181_{_}1)
- [4] T. Christy and L. I. Kuncheva, "Technological Advancements in Affective Gaming: A Historical Survey," *GSTF Journal on Computing (JoC)*, vol. 3, no. 4, p. 38, apr 2014. [Online]. Available: <http://www.globalsciencejournals.com/article/10.7603/s40601-013-0038-5>
- [5] G. N. Yannakakis, H. P. Martinez, and M. Garbarino, "Psychophysiology in Games," 2016, pp. 119–137. [Online]. Available: [http://link.springer.com/10.1007/978-3-319-41316-7\\_{\\_}7](http://link.springer.com/10.1007/978-3-319-41316-7_{_}7)
- [6] A. Clerico, C. Chamberland, M. Parent, P.-E. Michon, S. Tremblay, T. H. Falk, J.-C. Gagnon, and P. Jackson, "Biometrics and classifier fusion to predict the fun-factor in video gaming," in *2016 IEEE Conference on Computational Intelligence and Games (CIG)*. IEEE, sep 2016, pp. 1–8. [Online]. Available: <http://ieeexplore.ieee.org/document/7860418/>
- [7] R. V. Aranha, C. G. Correa, and F. L. S. Nunes, "Adapting software with Affective Computing: a systematic review," *IEEE Transactions on Affective Computing*, vol. 3045, no. c, pp. 1–1, 2019. [Online]. Available: <https://ieeexplore.ieee.org/document/8656550/>
- [8] G. Chanele and P. Lopes, "User Evaluation of Affective Dynamic Difficulty Adjustment Based on Physiological Deep Learning," in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 2020.
- [9] G. N. Yannakakis, H. P. Martínez, and A. Jhala, "Towards affective camera control in games," *User Modeling and User-Adapted Interaction*, vol. 20, no. 4, pp. 313–340, oct 2010. [Online]. Available: <http://link.springer.com/10.1007/s11257-010-9078-0>
- [10] K. Karpouzis, G. N. Yannakakis, N. Shaker, and S. Asteriadis, "The platformer experience dataset," in *2015 International Conference on Affective Computing and Intelligent Interaction (ACII)*. IEEE, sep 2015, pp. 712–718. [Online]. Available: <http://ieeexplore.ieee.org/document/7344647/>
- [11] A. E. Alchalabi, S. Shirmohammadi, A. N. Eddin, and M. Elsharnouby, "FOCUS: Detecting ADHD Patients by an EEG-Based Serious Game," *IEEE Transactions on Instrumentation and Measurement*, vol. 67, no. 7, pp. 1512–1520, jul 2018. [Online]. Available: <https://ieeexplore.ieee.org/document/8370717/>
- [12] M. Song, Z. Yang, A. Baird, E. Parada-Cabaleiro, Z. Zhang, Z. Zhao, and B. Schuller, "Audiovisual Analysis for Recognising Frustration during Game-Play: Introducing the Multimodal Game Frustration Database," in *2019 8th International Conference on Affective Computing and Intelligent Interaction (ACII)*. IEEE, sep 2019, pp. 517–523. [Online]. Available: <https://ieeexplore.ieee.org/document/8925464/>

<sup>13</sup>The first 108 participants (33 sessions) have a lower sampling rate of approximately 15Hz

- [13] P. M. Blom, S. Bakkes, and P. Spronck, "Towards Multi-modal Stress Response Modelling in Competitive League of Legends," in *2019 IEEE Conference on Games (CoG)*, vol. 2019-Augus. IEEE, aug 2019, pp. 1–4. [Online]. Available: <https://ieeexplore.ieee.org/document/8848004/>
- [14] K. Kutt, D. Drążyk, P. Jemioło, S. Bobek, B. Giżycka, V. Rodriguez-Fernandez, and G. J. Nalepa, "BIRAFFE : bio-reactions and faces for emotion-based personalization," in *3rd Workshop on Affective Computing and Context Awareness in Ambient Intelligence (AfCAI 2019)*, Cartagena, Spain, 2019. [Online]. Available: [http://ceur-ws.org/Vol-2609/AfCAI2019\\_{\\_}paper\\_{\\_}6.pdf](http://ceur-ws.org/Vol-2609/AfCAI2019_{_}paper_{_}6.pdf)
- [15] K. Kutt, D. Drążyk, L. Zuchowska, M. Szelążek, S. Bobek, and G. J. Nalepa, "BIRAFFE2, a multimodal dataset for emotion-based personalization in rich affective game environments," *Scientific Data*, vol. 9, no. 1, p. 274, dec 2022. [Online]. Available: <https://www.nature.com/articles/s41597-022-01402-6>
- [16] N. Beaudoin-Gagnon, A. Fortin-Cote, C. Chamberland, L. Lefebvre, J. Bergeron-Boucher, A. Campeau-Lecours, S. Tremblay, and P. L. Jackson, "The FUNii Database: A Physiological, Behavioral, Demographic and Subjective Video Game Database for Affective Gaming and Player Experience Research," in *2019 8th International Conference on Affective Computing and Intelligent Interaction (ACII)*. IEEE, sep 2019, pp. 1–7. [Online]. Available: <https://ieeexplore.ieee.org/document/8925502/>
- [17] M. Granato, D. Gadia, D. Maggiorini, and L. A. Ripamonti, "An empirical study of players' emotions in VR racing games based on a dataset of physiological data," *Multimedia Tools and Applications*, vol. 79, no. 45–46, pp. 33 657–33 686, dec 2020. [Online]. Available: <http://link.springer.com/10.1007/s11042-019-08585-y>
- [18] T. B. Alakus, M. Gonen, and I. Turkoglu, "Database for an emotion recognition system based on EEG signals and various computer games – GAMEEMO," *Biomedical Signal Processing and Control*, vol. 60, p. 101951, jul 2020. [Online]. Available: <https://doi.org/10.1016/j.bspc.2020.101951https://linkinghub.elsevier.com/retrieve/pii/S1746809420301075>
- [19] A. Smerdov, B. Zhou, P. Lukowicz, and A. Somov, "Collection and Validation of Psychophysiological Data from Professional and Amateur Players: a Multimodal eSports Dataset," vol. XX, pp. 1–12, nov 2020. [Online]. Available: <http://arxiv.org/abs/2011.00958>
- [20] D. Melhart, A. Liapis, and G. N. Yannakakis, "The Affect Game AnnotatIoN (AGAIN) Dataset," apr 2021. [Online]. Available: <http://arxiv.org/abs/2104.02643>
- [21] D. Dresvyanskiy, Y. Sinha, M. Busch, I. Siegert, A. Karpov, and W. Minker, "Dycoda: A multi-modal data collection of multi-user remote survival game recordings," in *Speech and Computer*, S. R. M. Prasanna, A. Karpov, K. Samudravijaya, and S. S. Agrawal, Eds. Cham: Springer International Publishing, 2022, pp. 163–177.
- [22] H. Svoren, V. Thambawita, P. Halvorsen, P. Jakobsen, E. Garcia-Ceja, F. M. Noori, H. L. Hammer, M. Lux, M. A. Riegler, and S. A. Hicks, "Toadstool: A Dataset for Training Emotional Intelligent Machines Playing Super Mario Bros," in *Proceedings of the 11th ACM Multimedia Systems Conference*. New York, NY, USA: ACM, may 2020, pp. 309–314. [Online]. Available: <https://dl.acm.org/doi/10.1145/3339825.3394939>
- [23] M. Fanourakis, P. Lopes, and G. Chanel, "Remote Multi-Player Synchronization using the Labstreaming Layer System," in *Foundations of Digital Games Demos*, Malta, 2020. [Online]. Available: <https://archive-ouverte.unige.ch/unige:148594>
- [24] S. Greenberg, P. Aislinn, and D. Kirsten, "Development and Validation of the Fatigue State Questionnaire: Preliminary Findings," *The Open Psychology Journal*, vol. 9, no. 1, pp. 50–65, jun 2016. [Online]. Available: <https://openpsychologyjournal.com/VOLUME/9/PAGE/50/>
- [25] S. Gächter, C. Starmer, and F. Tufano, "Measuring the Closeness of Relationships: A Comprehensive Evaluation of the 'Inclusion of the Other in the Self' Scale," *PLOS ONE*, vol. 10, no. 6, p. e0129478, jun 2015. [Online]. Available: <https://dx.plos.org/10.1371/journal.pone.0129478>
- [26] D. Melhart, A. Liapis, and G. N. Yannakakis, "PAGAN: Video Affect Annotation Made Easy," in *2019 8th International Conference on Affective Computing and Intelligent Interaction (ACII)*. IEEE, sep 2019, pp. 130–136. [Online]. Available: <https://ieeexplore.ieee.org/document/8925434/>
- [27] P. Lopes, G. N. Yannakakis, and A. Liapis, "RankTrace: Relative and unbounded affect annotation," in *2017 Seventh International Conference on Affective Computing and Intelligent Interaction (ACII)*, vol. 2018-Janua. IEEE, oct 2017, pp. 158–163. [Online]. Available: <http://ieeexplore.ieee.org/document/8273594/>
- [28] A. R. Robertson, "The CIE 1976 Color-Difference Formulae," *Color Research & Application*, vol. 2, no. 1, pp. 7–11, mar 1977. [Online]. Available: <https://onlinelibrary.wiley.com/doi/10.1002/j.1520-6378.1977.tb00104.x>
- [29] R. L. Donofrio, "Review Paper: The Helmholtz-Kohlrausch effect," *Journal of the Society for Information Display*, vol. 19, no. 10, p. 658, 2011. [Online]. Available: <http://doi.wiley.com/10.1889/1.1828693>
- [30] M. Fanourakis and G. Chanel, "Attenuation of the dynamic pupil light response during screen viewing for arousal assessment," *Frontiers in Virtual Reality*, vol. 3, 2022. [Online]. Available: <https://www.frontiersin.org/articles/10.3389/frvir.2022.971613>
- [31] R. Weber, K.-M. Behr, R. Tamborini, U. Ritterfeld, and K. Mathiak, "What Do We Really Know About First-Person-Shooter Games? An Event-Related, High-Resolution Content Analysis," *Journal of Computer-Mediated Communication*, vol. 14, no. 4, pp. 1016–1037, jul 2009. [Online]. Available: <https://academic.oup.com/jcmc/article/14/4/1016-1037/4583562>
- [32] T. Baltrusaitis, A. Zadeh, Y. C. Lim, and L.-P. Morency, "OpenFace 2.0: Facial Behavior Analysis Toolkit," in *2018 13th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2018)*. IEEE, may 2018, pp. 59–66. [Online]. Available: <https://ieeexplore.ieee.org/document/8373812/>
- [33] G. G. BERNTSON, J. THOMAS BIGGER JR., D. L. ECKBERG, P. GROSSMAN, P. G. KAUFMANN, M. MALIK, H. N. NAGARAJA, S. W. PORGES, J. P. SAUL, P. H. STONE, and M. W. VAN DER MOLEN, "Heart rate variability: Origins, methods, and interpretive caveats," *Psychophysiology*, vol. 34, no. 6, pp. 623–648, 1997. [Online]. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1469-8986.1997.tb02140.x>
- [34] R. Freeman and M. W. Chapleau, "Testing the autonomic nervous system," in *Peripheral Nerve Disorders*, ser. Handbook of Clinical Neurology, G. Said and C. Krarup, Eds. Elsevier, 2013, vol. 115, pp. 115–136. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/B9780444529022000072>
- [35] M. E. Dawson, A. M. Schell, and D. L. Filon, "The Electrodermal System," in *Handbook of Psychophysiology*, J. T. Cacioppo, L. G. Tassinary, and G. Berntson, Eds. Cambridge: Cambridge University Press, 2017, pp. 217–243. [Online]. Available: <http://ebooks.cambridge.org/ref/id/CBO9780511546396A015>