

Ensemble Learning to Assess Dynamics of Affective Experience Ratings and Physiological Change

Felix Dollack¹, Kiyoshi Kiyokawa¹, Huakun Liu¹, Monica Perusquia-Hernandez¹,
Chirag Raman², Hideaki Uchiyama¹, Xin Wei¹

¹Nara Institute of Science and Technology

²Delft University of Technology

{felix.d, kiyo, liu.huakun.li0, m.perusquia, hideaki.uchiyama, wei.xin.wy0}@is.naist.jp, c.a.raman@tudelft.nl

Abstract—The congruence between affective experiences and physiological changes has been a debated topic for centuries. Recent technological advances in measurement and data analysis provide hope to solve this epic challenge. Open science and open data practices, together with data analysis challenges open to the academic community, are also promising tools for solving this problem. In this entry to the Emotion Physiology and Experience Collaboration (EPiC) challenge, we propose a data analysis solution that combines theoretical assumptions with data-driven methodologies. We used feature engineering and ensemble selection. Each predictor was trained on subsets of the training data that would maximize the information available for training. Late fusion was used with an averaging step. We chose to average considering a “wisdom of crowds” strategy. This strategy yielded an overall RMSE of 1.19 in the test set. Future work should carefully explore if our assumptions are correct and the potential of weighted fusion.

Index Terms—affective computing, continuous ratings, biosignal processing, machine learning, data analysis challenge

I. INTRODUCTION

Understanding human emotion is instrumental for applications in mental healthcare, education, and communication [6]. These applications aim to automatically assess and generate affective cues by relying on an assumed relationship between affective experience and physiological changes. However, the debate on the precedence of body changes or subjective experience started in the previous century and remains current [7]. Recent research has discussed whether demand characteristics affect bias in our understanding of the relationship between facial expression and affective experience [9, 2]; and explored the relationship between dynamic Autonomic Nervous System responses and affective experiences [12, 22]. Physiological sensing technologies have been popular in studying the physiological changes correlating with affective experiences [5, 13]. Each physiological measurement type gives a different piece of information regarding the functioning of the sympathetic and parasympathetic nervous systems [1], leading to a popular multidimensional dataset collection. Traditional data analysis techniques require extensive knowledge about physiology characteristics, signal processing, and domain knowledge in affective sciences. This domain knowledge gave birth to hand-crafted feature engineering that improves data interpretability and reduces the number of comparisons

to be made when analyzing the data. Recent advances in Machine Learning (ML) and data-driven analyses have brought a new perspective. Purely data-driven analyses with end-to-end automated processing have become popular [16]. In end-to-end approaches, a machine learning network learns an intermediate representation of the input, thereby reducing manual work, and potentially enhancing the results [15]. However, the evidence does not always support this claim. A previous study showed that convolutional and recurrent neural networks yielded better results than other state-of-the-art methods [15]. Another study used a deep-learning approach to estimate momentary emotional states from multi-modal physiological data; and reported a higher correlation than traditional methods. Still, their mean absolute error (MAE) was higher (a lower MAE is better) [14]. Finally, another study found that end-to-end processes are suitable for predicting stress states with abrupt changes, but not as good when assessing subtle affective states like enjoyment [10]. Hence, end-to-end learning only provides a marginal improvement over feature engineering for physiological signal-based affect recognition. This is different from camera-based recognition, where performance is radically improved. One possible explanation is the limited amount of physiological data publicly available. Therefore, public data sets and multi-laboratory collaborations are necessary to assess the effectiveness of different training methods and cross-validation strategies.

The EPiC challenge aims to overcome the limitations in data availability and motivates researchers to work on the affect-embodiment coherence problem. Our team used theory-driven analysis and data engineering techniques to address the EPiC challenge. We opted for feature engineering and ensemble learning for our final submission. Also, we report an exploratory analysis validating our assumptions for the challenge submission.

II. RELATED WORK

ML has been used to model emotion recognition mechanisms from data following the public release of benchmark databases. The basic procedure of classical ML-based methods consists of four steps: physiological signal collection while eliciting participants’ emotions, feature extraction from the signals, training a classification model with the features, and emotion recognition based on the trained model [16]. Research

All authors contributed equally to this work.

issues include the design of discriminative features and the selection of the optimal classification technique. For instance, hypothesis testing is performed over some features followed by a predictive model that makes feature selection to see if the tested features were still relevant when all features were considered together [29]. It has been suggested that group synchrony improved arousal and valence classification [3]. Electrocardiograms (ECG) and Electrodermal Activity (EDA) have been used as tabular data with AutoGluon-Tabular to arousal and valence across individuals and datasets [8] with similar accuracy (around 56 – 62% respectively) to previous works. Nevertheless, the subject-independent classification remains only slightly above chance level.

A crucial challenge surrounding continuous-time annotation of emotions is the lag between observed features and the reported emotion measures [26, 24]. This lag arises from the time the rater requires to provide feedback about the experienced emotion. Such a temporal misalignment between features and labels has consequences for ML methods. Consequently, several compensation techniques have been investigated [20, 21, 18, 19]. These methods involve estimating the reaction lag from the data, by maximizing the correlation coefficient [20, 21, 18, 26] or the mutual information between the multimodal features and emotional ratings [19]. Others have used a recurrent neural network to handle the asynchronous dependencies [24].

Deep learning (DL) has also been used in affective computing. Handcrafted feature design is not always necessary in DL [16], also, less modalities seem to be required to achieve equal performance. For example, Hssayeni and Ghoraani [14] presented a DL approach to estimate momentary emotional states from multi-modal physiological data. Used modalities included respiration, ECG, electromyography (EMG), EDA, and acceleration. The best emotion classification was achieved by a traditional method with 79% F1-score when all four physiological modalities were used. In contrast, using only two modalities, DL achieved 78% F1-score. Furthermore, there are several fusion strategies for multi-modal data: feature-level fusion, decision-level fusion, model-level fusion, and hybrid-level fusion. Late fusion by averaging class probabilities has performed well in the past [14].

In the case of continuous detection of valence and arousal, a previous work obtained 0.43 and 0.59 RMSE for valence and arousal, respectively [28], in the WESAD dataset. When dividing the continuous annotation into binary valence-arousal categories (high-low), another group of researchers reported subject-independent accuracy of 76.37% and 74.03% for valence and arousal, respectively [30], on the Continuously Annotated Signals of Emotion (CASE) dataset [27].

III. CHALLENGE CORPORA

The challenge corpora is an open dataset that collected six physiological signals while 30 participants (15 female, age range: 22 – 37 years) watched eight videos [27]. The videos aimed to elicit a range of emotions and were rated with continuous self-reported valence and arousal in the range

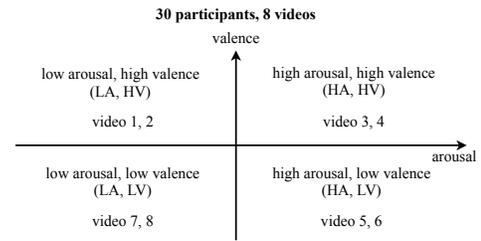


Fig. 1. Diagram depicting the dataset structure utilized in this competition.

of 0.5 to 9.5 using a joystick. Visual feedback was provided using the Self-Assessment Manikin [4]. Two videos were chosen per quadrant in the valence-arousal space, often called affect grid [25], and were presented in pseudo-random order to the participants. Figure 1 shows the video distribution. The physiological sensing was logged at 1000 Hz, and the continuous annotation was done at 20 Hz. The physiological sensors included are:

- 1) Cardiac activity as measured from Electrocardiography (ECG) and Photoplethysmography (BVP).
- 2) Muscle activity (EMG) recorded from three muscles: the corrugator supercilii (emg_coru), the zygomaticus major (emg_zygo), and the trapezius (emg_trap).
- 3) Electrodermal activity (EDA) measured from the non-dominant hand.
- 4) Respiration (RSP) recorded from the chest.
- 5) Skin Temperature (SKT) recorded from the little finger of the non-dominant hand.

For the EPiC Challenge, the Dataset was arranged in four scenarios to test four assumptions about the relationship between affective experiences and embodied cues. Each scenario is divided into training and test sets, as prepared by the organizers ¹.

A. Across-time scenario

This scenario evaluates subject-dependent and affective context-dependent model performance. The model is trained and tested over different durations of one data file (sub_vid). In this scenario, each of the 240 data files, that is, 30 participants watching eight videos, is divided into training and test parts based on the time series. For each data file, the earlier part is the training data, and the latter is the test data. The training and test data are not consecutive but are spaced by an unknown length of time. The length of the training data ranged from 48 s to 127 s depending on the size of each video, with an average of 88 s. All test data files were 50 s in length.

B. Across-subject scenario

This scenario evaluates subject-independent model performance. The model is trained on data from some participants and tested on data from another set of different, unseen, participants to verify the model’s generalization ability to new people. In this scenario, the data of 30 participants were

¹See: <https://github.com/Emognition/EPiC-2023-competition>

divided into five groups. Each group contains six participants watching eight videos for a total of 48 data files. This scenario consisted of five folds to use the cross-validation strategy. In each fold, the data of four groups of participants were set as the training data (192 data files in total). The 48 data files from the remaining six participants were set as the test data. The length of the training data files ranged from 50 s to 128 s, with an average of 90 s. All test data files lasted 50 s.

C. Across-elicitor scenario

This scenario evaluates affective context-independent model performance. The model is trained on data from several affective contexts and then tested on data from a different affective context. Each affective context represents one quadrant in the valence-arousal affect grid. There were two elicitors (i.e., videos) per affective context. This verifies whether the model can infer from the physiological signals triggered by one affective context to the physiological features triggered by another affectivity. In this scenario, the data from eight videos were divided into four groups, according to the video’s affective context. This resulted in four categories: low valence, high arousal; high valence, high arousal; high valence, low arousal; and low valence, low arousal. Each group contains 60 data files, which corresponds to 30 participants with two videos each. The two videos in one group are considered to trigger the same type of affectivity. By adopting a cross-validation strategy, this scenario contains four folds. In each fold, three groups, totaling 180 data files, were chosen as the training data for the challenge. The remaining two videos for a total of 60 data were chosen as the test data.

D. Across-version scenario

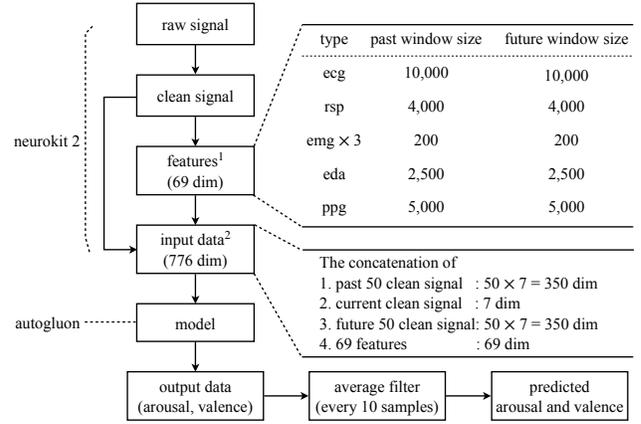
This scenario evaluates affective context-dependent model performance. The model is trained on data from a specific affective state instantiation and then tested on the other version of the same affective contexts to verify the model’s generalizability across similar affective contexts. The data from eight videos were divided into two groups. Each group contains 30 users watching four video types covering all quadrants in the affect grid, for a total of 120 data files. This scenario consists of two folds. In each fold, one group is the training data, while the other is the test data.

IV. TACKLING THE CHALLENGE

A. Preprocessing

We used NeuroKit2 to preprocess the physiological signals and feature extraction [17]. In the case of EMG, we opted to write a custom preprocessing pipeline based on [23] in combination with Neurokit.

The EMG pipeline to clean the signal consisted of a series of notch filters at frequencies of 60, 120, 180, and 240 Hz with a notch width of 3 Hz, followed by a bandpass filter with cutoff frequencies of 5 Hz and 250 Hz, and a detrending step. A z-transform is applied to the clean signal, before calculating the root mean square (RMS) over windows of 100 ms. To further smoothen the envelope, a Savitzky-Golay filter of third order



¹: Features are extracted per each valence and arousal data point.

In this case, the step is 50 physiology samples because the annotation frequency is 20 Hz.

²: The input data is constructed based on the timestamp at which each valence and arousal is recorded.

Fig. 2. The pipeline used to process each signal, and the post-processing after the machine learning model. Window sizes are represented in samples at 1 kHz.

with a length of 1 s was applied. This cleaned signal was then input to Neurokit’s amplitude function to keep the returned data structure consistent for Neurokit’s analyze function. The Neurokit’s analyze function was used to extract the features listed in Table I.

Feature extraction was performed over the whole signal using windows of different sizes as shown in Figure 2. These windows are described in samples at 1 KHz. ECG and PPG’s window size is longer to sample heart-rate variability. A similar reason applies to respiration. Regarding EDA, a medium-sized window is recommended to decompose the signal into tonic and phasic components. In contrast, the EMG window is shorter, because changes in facial expressions’ EMG can happen in the order of milliseconds. As a curious fact, we found a heart-rate artifact in the trapezius EMG. This artifact might be misinterpreted as an EMG feature, but we decided to leave it in because we did not formally distinguish between data types when modeling the data.

The feature vectors were sampled at 20 Hz to match the sampling frequency of the annotations. Additionally, the pre-processed (clean) data temporally surrounding each annotation datapoint were flattened and input as additional features to the modeling block.

B. Model training

We used a consistent architecture across all four scenarios, but adopted unique strategies for designing the input of model training and generating predictions to accommodate the distinct requirements of various scenarios.

For the architecture, we employed AutoGluon, an open-source AutoML framework developed by AWS, to train our model [11]. This framework expedites the development of machine learning models by automating model training, hyperparameter optimization, and model selection and ensembling. The AutoGluon-Tabular fits a total of 11 models that

includes gradient boosting methods (CatBoost, LightGBM, LightGBMLarge, LightGBMXT, XGBoost), extra trees (ExtraTreesMSE), K-nearest neighbors algorithm (KNeighborsDist, KNeighborsUnif), neural networks (NeuralNetFastAI, NeuralNetTorch), and random forests (RandomForestMSE). Furthermore, a weighted ensemble model (WeightedEnsemble_L2) is fitted and employed to combine the previously-trained models for generating predictions. To achieve optimal performance with AutoGluon, we designate the parameter presets as ‘best_quality’, allowing AutoGluon to automatically construct robust model ensembles while allocating sufficient training time.

1) *Across-time scenario*: We aimed to capture the unique characteristics and nuances of each subject’s emotional responses and the specific stimuli embedded within the emotional context. Therefore, we trained the model on discrete datasets originating per participant and per video. In this scenario, both training and test sets comprise 240 subsets. The training-test pairs were all collected from the same pool of participants and emotion elicitors that exhibit a one-to-one correspondence. We trained 240 models on each training dataset and subsequently selected the corresponding model to yield predictions on the test sets.

2) *Across-subject scenario*: The substantial inter-subject variability in physiological responses to stimuli and the inherent limitations of self-report labels, such as subjectivity, bias, and emotional granularity, presents a considerable challenge in developing one-size-fits-all-subjects effective models for affective computing. To generate plausible predictions, we assume that a single video should elicit similar emotions in most participants. Guided by this assumption, we trained a dedicated model for each video. In this scenario, every fold consists of 24 subjects in the training dataset and six subjects in the test dataset, with all participants having watched the same eight videos. We combined data from different subjects of each video as input when training the model. In total, we trained eight models and employed the respective models

TABLE I
FEATURE LIST

Signal	Features
PPG (10 dim)	Rate : baseline, max, min, mean, SD, max time, min time, trend linear, trend quadratic, trend R2;
ECG (15 dim)	Rate : baseline, max, min, mean, SD, max time, min time, trend linear, trend quadratic, trend R2; Phase : atrial, completion atrial, ventricular, completion ventricular; quality mean;
RSP (20 dim)	Rate : baseline, max, min, mean, SD, max time, min time, trend linear, trend quadratic, trend R2 Amplitude : baseline, max, min, meanraw, mean, SD; Phase , phase completion, RVT baseline, RVT mean;
EDA (6 dim)	peak amplitude, SCR, SCR peak amplitude, SCR peak amplitude time, SCR RiseTime, SCR RecoveryTime;
EMG × 3 (6 × 3 dim)	Activation, Amplitude Mean, Amplitude Max, Amplitude SD, Amplitude Max Time, Bursts

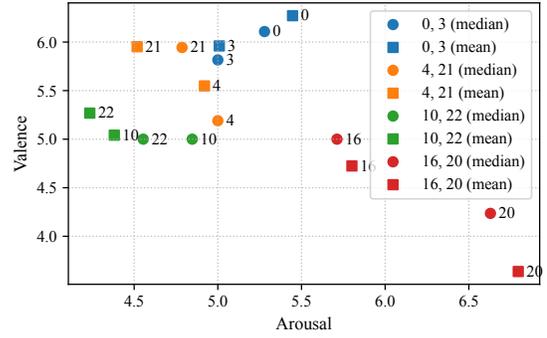


Fig. 3. The meta-analysis for determining the quadrant affiliation of videos by the mean and median of video ratings.

for affective state estimation when generating predictions for corresponding test set files.

3) *Across-elicitor scenario*: In this scenario, using only the data from the three quadrants available for estimating arousal and valence could compromise our ability to predict affective states associated with the missing quadrant accurately. To mitigate this concern, we performed a meta-analysis of the training data. We first calculated the mean value of user ratings per video file in the training set to categorize videos into the four affect grid quadrants systematically. By doing so, we could make well-founded assumptions regarding each video’s quadrant affiliation.

Within this scenario, a total of eight videos were provided. By analyzing the composition of the test dataset across four-folds, we categorized the eight videos into four groups: (0, 3), (4, 21), (10, 22), and (16, 20), see Figure 3. Among these ratings, it is evident that videos (0, 3) and (16, 20) belong to the high valence high arousal (HV, HA) and low valence high arousal (LV, HA) quadrants, respectively. The categorization of the other two groups is less apparent; therefore, our hypothesis relies on the video with the more prominent rating within each of the two video groups. Consequently, we assumed that video (0, 3) belongs to the (HV, HA) quadrant, (16, 20) belongs to the (LV, HA) quadrant, (10, 22) belongs to the (LV, LA) quadrant, and (4, 21) belongs to the (HV, LA) quadrant. Based on this assumption, we employed only two relevant quadrants to achieve sample balance and maximize the variance along the valence and arousal axes. For instance, when the videos in the test set belong to the (HV, HA) quadrant, we train the valence predictor on the dataset comprising videos from the (LV, LA) and (HV, LA) quadrants, and the arousal predictor on the dataset containing videos from the (LV, HA) and (LV, LA) quadrants. This approach effectively minimizes input bias and ensures more accurate emotional state estimations for the missing quadrant.

4) *Across-version scenario*: Given that instances of all the affective quadrants are available, albeit in only one version, we aimed to develop a general model that yields robust results by harnessing collective intelligence by applying late fusion. In other words, we assumed there is a “wisdom of crowds” effect when combining multiple weak classifiers into one. To

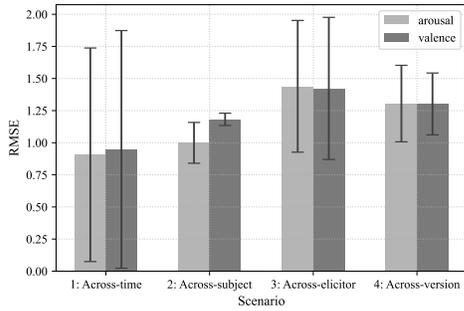


Fig. 4. Scenarios-level RMSE. Error bars represent standard deviation.

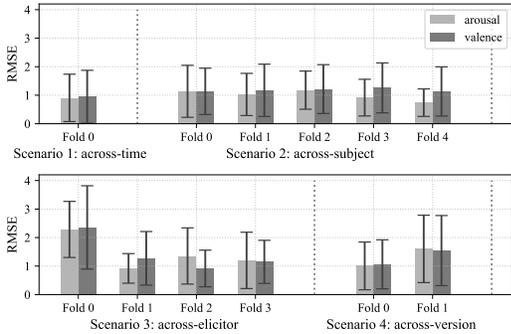


Fig. 5. Folds-level RMSE. Error bars represent standard deviation.

this aim, we developed four separate models, each trained on a distinct set of four videos from the training dataset. During the testing process, we input the preprocessed and feature-extracted data from the test set into each of these four models to generate predictions. Then we applied a late fusion strategy to obtain the final estimation. The four predictions from each model were fused by calculating their mean predicted rating values.

V. VALIDATION RESULTS

The models were assessed using the root mean square error (RMSE) metric. A lower RMSE is better. It has the same units as the valence and arousal annotations. The final score for the EPiC challenge for our team was 1.19. Additionally, we report the detailed performance for each scenario on the test set, as reported by the workshop organizers. Our training was done on the full train set to maximize data availability. For each test data, i.e., the data of one subject and one video, the RMSE is calculated for arousal and valence, respectively. Then in each fold, the performance is assessed by averaging the RMSE values among all test data. Similarly, the model performance in each scenario is evaluated by averaging all RMSE values within the scenario. The final RMSE result was obtained by calculating the mean score on all scenarios and two prediction targets, i.e., arousal and valence. The scenarios-level RMSE and folds-level RMSE are shown in Figure 4, Figure 5, and Table II.

A. Across-time scenario

The RMSE of predicted arousal and valence are 0.91 and 0.95, with a standard deviation of 0.83 and 0.93, respectively. Among the 240 test data, the lowest RMSE of arousal and valence is 0 and 0, and the highest is 3.87 and 4.33. For the predicted results of arousal and valence, the RMSE of 69% and 66% of the test results were below 1. Overall, the prediction error for arousal is slightly lower than that for valence.

B. Across-subject scenario

The RMSE for the predicted arousal and valence are 1 and 1.1, respectively, with standard deviations of 0.16 and 0.04 across the five folds. In each fold, the predicted arousal RMSE is consistently lower than the predicted valence RMSE. This difference is especially noticeable in fold 4, where the RMSE and standard deviation for predicted arousal are 0.74 and 0.48, respectively, in contrast to the valence RMSE and standard deviation, which are 1.13 and 0.86, respectively.

Despite the seemingly positive results in this scenario, there is a limitation on how we predicted the final ratings. We assumed that training and testing across elicitors would help to predict better the outcomes in the ratings done by other people not seen in the dataset. However, in the real world, we do not have information about the type of stimuli used. Therefore, the performance will probably be reduced, as exemplified in the following section.

C. Across-elicitor scenario

Our model yielded the highest RMSE values, with the RMSE for arousal and valence being 1.44 and 1.42, respectively, and standard deviations of 0.51 and 0.55. We note that the high RMSE was due to poor prediction results in fold 0, where the RMSE values were twice those observed in other folds, reaching 2.29 and 2.36 for arousal and valence, respectively. By analyzing the data for this scenario, we noticed that a potential cause for the high RMSE lies in the

TABLE II
FOLDS-LEVEL AND SCENARIOS-LEVEL RMSE

Scenario	Fold	Arousal		Valence	
		RMSE	STD	RMSE	STD
Across-time	0	0.91	0.83	0.95	0.93
	Scenario level	1.00	0.16	1.18	0.05
Across-subject	0	1.14	0.91	1.14	0.81
	1	1.03	0.74	1.17	0.92
	2	1.18	0.67	1.21	0.86
	3	0.92	0.64	1.26	0.87
	4	0.74	0.48	1.13	0.86
Scenario level		1.00	0.16	1.18	0.05
Across-elicitor	0	2.29	0.98	2.36	1.46
	1	0.92	0.52	1.27	0.94
	2	1.35	0.99	0.92	0.64
	3	1.20	0.99	1.15	0.76
Scenario level		1.44	0.51	1.42	0.55
Across-elicitor	0	1.00	0.84	1.06	0.86
	1	1.60	1.18	1.54	1.23
Scenario level		1.30	0.30	1.30	0.24

significant deviation between the test data and the training data in fold 0. The training data does not include similar patterns in the test data. As shown in Figure 6, the test data in fold 0 contains videos 16 and 20 in the upper left quadrant of the affective grid. The rating patterns of arousal and valence differ from the data in the other three quadrants, which were used as training data. This suggests that the larger variations in ratings characteristic of negative, high-arousing emotions are not present in the other types of emotion. Furthermore, these results also demonstrate the reliance of our model on data similarity, indicating a weaker generalization capability for novel data patterns.

D. Across-version scenario

In this scenario, the prediction RMSE using our model is 1.30 for both arousal and valence, with standard deviations of 0.30 and 0.24, respectively. Although each of the two groups’ data in this scenario covered all four emotional states, the cross-validation results revealed that our model’s RMSE in fold 0 was 0.5 lower than in fold 1. This suggests that an appropriately balanced training set, encompassing all four emotional states, can significantly enhance the model’s generalization capabilities.

VI. REVISITING ASSUMPTIONS

A. Lag between physiological signals and ratings

When preparing the data to train our models, we assumed that the physiological changes happen at different speeds depending on the measurement metrics used. Here we investigate the effect of the time delay in reporting emotions. In particular, we followed the general procedure described by Schmitt, Ringeval, and Schuller [26]. We shifted the features forward in time in steps of 0.005 s up to a maximum of 0.05 s and trained a Gated Recurrent Unit (GRU) model to predict arousal and valence. Note that the annotations were performed at 20 Hz while the physiological signals were sampled at 1000 Hz. Then, we experimented with using each individual signal as input to the model in isolation before combining all signals as input. The analysis results are in Table III. The results suggest that predictive performance generally improves when accounting for annotation delays. However, the delay yielding the most empirical gains varies for each biosignal, and we often found several minimum values. A further investigation of the timing relationships between physiological change, experience, and annotation is needed to understand when the differences significantly disrupt the predictions.

B. The gradual nature of changes in emotion

Qualitatively, we noticed that the predictions from our models were characterized by more high-frequency changes than the annotations provided in the training data, which change more gradually over time. Consequently, we utilized a moving average window comprising 10 samples, equivalent to a 2 Hz low-pass filter. The choice of the window length follows from the observation that the annotations were provided at 20 Hz, so the Nyquist frequency is 10 Hz – only changes at 10 Hz

can be measured with the joystick described in the dataset. Furthermore, we assumed people would not make more than two abrupt changes per second. Further investigation is required to establish whether the nature of emotion changes is a gradual process, or whether the relative smoothness of the ground-truth ratings is an annotation artifact.

C. Single- and multi-label predictors

Considering that participants simultaneously rated their emotions for valence (X-axis) and arousal (Y-axis) using a two-dimensional joystick, we speculated on the potential connection between these two values, despite the orthogonal nature of the valence-arousal model’s axes. To evaluate this hypothesis, we extracted 24 datasets from scenario 1, each containing data from the combination of six participants and four videos. We compared the following approaches: (1) independent prediction of valence and arousal, (2) predicting valence first and then using both physiological signals and predicted valence values for arousal prediction, and (3) predicting arousal first and subsequently using the predicted arousal values for valence prediction. This comparison investigated any potential connections between valence and arousal predictions. As demonstrated in Figure 7, the performance differences among these three strategies are minimal. As a result, owing to the slower training process associated with multi-label predictors compared to single-label cases, we ultimately chose to employ the independent prediction strategy for this competition.

VII. DISCUSSION AND FUTURE DIRECTIONS

We introduced an attempt to address the EPiC Challenge. We predicted continuous valence and arousal ratings from several biosignals, across four scenarios. We used readily available algorithms, with our novelty being (a) the window choices for feature calculation; (b) the use of data around each annotation point; and (c) our use of theoretical assumptions to maximize data variance in each scenario’s training. Our overall RMSE for the four validation scenarios was 1.19, as provided by the competition organizers. This result still has considerable room for improvement compared to previous work on

TABLE III
RMSE AT DIFFERENT RATING-BIOSIGNAL DELAYS. BOLD VALUES ARE THE MINIMUM.

Delay	ALL	BVP	ECG	RMSE (10^{-3})			GSR	RSP	SKT
				EMG <i>coru</i>	EMG <i>trap</i>	EMG <i>zygo</i>			
0	1.21	1.21	1.22	1.24	1.24	1.23	1.24	1.19	1.22
0.005	1.24	1.22	1.22	1.21	1.24	1.22	1.21	1.20	1.20
0.01	1.22	1.21	1.20	1.20	1.23	1.25	1.21	1.27	1.21
0.015	1.21	1.23	1.22	1.21	1.20	1.23	1.22	1.24	1.20
0.02	1.24	1.23	1.22	1.23	1.23	1.24	1.2	1.21	1.20
0.025	1.23	1.21	1.23	1.22	1.22	1.22	1.22	1.21	1.23
0.03	1.22	1.20	1.21	1.22	1.24	1.21	1.22	1.20	1.19
0.035	1.21	1.23	1.21	1.26	1.19	1.23	1.18	1.20	1.24
0.04	1.23	1.23	1.23	1.20	1.20	1.24	1.23	1.23	1.24
0.045	1.23	1.21	1.20	1.20	1.22	1.23	1.22	1.21	1.23
0.05	1.26	1.21	1.23	1.23	1.23	1.23	1.21	1.22	1.19

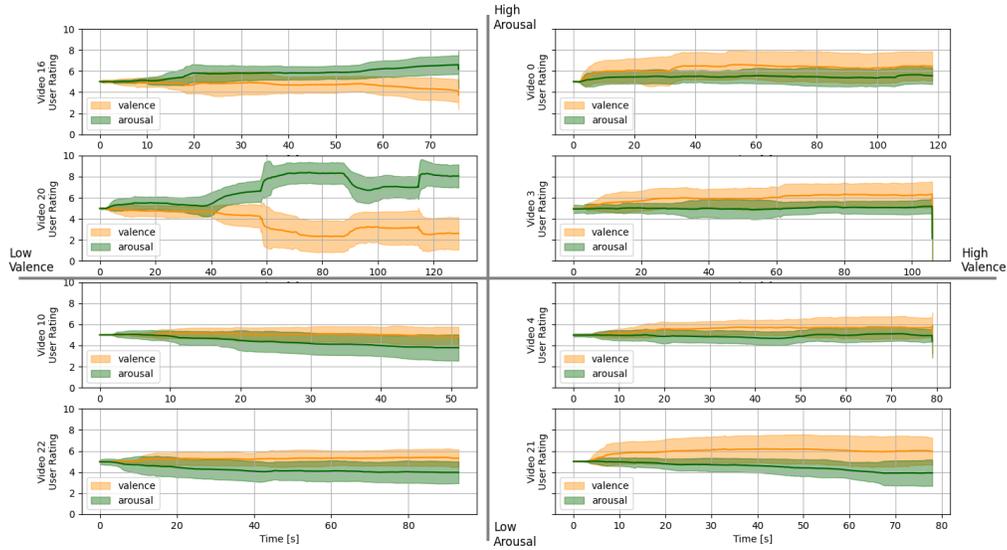


Fig. 6. Plot of the rating averages per each stimuli video. Shaded areas represent the standard deviation among participants.

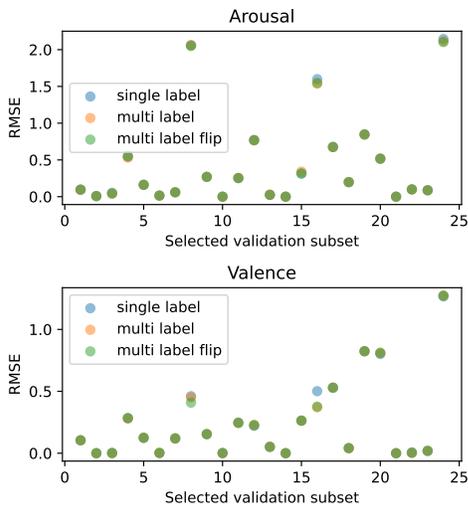


Fig. 7. The comparison between single and multi-label predictors. The “single label” represents independent prediction of valence and arousal. “Multi label” and “multi label flip” correspond to predicting arousal after predicting valence and vice versa.

predicting continuous valence and arousal annotations from physiological data, and can be used as a baseline for future regression studies on the CASE dataset. As expected from the literature, fitting a personalized model to predict a single person’s reaction at a future point (across-time) is an easier problem than in the other scenarios. To extend our model to other, unseen people, we used information about the stimuli, and capitalized in our knowledge of the affective context in which the data was collected. By training several models per stimuli, we reduced the RMSE in the across-subject scenario. However, this strategy is unlikely to work in the real world, as

we typically would not have information about the stimulus. A similar approach was used in the across-elicitor scenario. By examining the results (Figure 6), we hypothesize that high arousal and low valence emotions display abrupt physiological changes not present in other affective quadrants. Future work should explore whether this is consistently true, and devise a method to tackle the lack of information during training. Furthermore, the results of the across-version validation signal that expected affective messages depicted in a stimulus might not produce the same effect in different individuals. Future work should validate if this is the case, and assess if a weighted late fusion provides improvements with respect to an averaging function. This would also validate or refute whether our bet for a “wisdom of a crowd classifier” is suitable. Finally, future work should be done to formally assess if end-to-end methods outperform ensemble methods and feature engineering similar to those used in this work.

ETHICAL IMPACT STATEMENT

This research was a data analysis of the CASE dataset [27]. All data provided is anonymous, and obtained following the Declaration of Helsinki. Our results have several limitations. The sample size is only 30 people, and no mentions of the cultural background of the participants are made in the dataset. The interpretation of the stimuli, and therefore the ratings and physiological changes, might differ from person to person. Therefore, our results need to be replicated in other corpora. Moreover, we built our model based on certain assumptions that are described, but not formally validated. These assumptions might not yield the same performance in other datasets. Finally, our results are biased to better predict the situations in the dataset. Therefore, our model should be used with caution.

REFERENCES

- [1] Stephanie Balters and Martin Steinert. “Capturing emotion reactivity through physiology measurement as a foundation for affective engineering in engineering design science and engineering practices”. In: *Journal of Intelligent Manufacturing* 28.7 (2017), pp. 1585–1607.
- [2] Lisa Feldman Barrett et al. “Emotional Expressions Reconsidered: Challenges to Inferring Emotion From Human Facial Movements”. In: *Psychological Science in the Public Interest* 20.1 (2019), pp. 1–68.
- [3] Patrícia Bota et al. “Group Synchrony for Emotion Recognition using Physiological Signals”. In: *IEEE Trans. Affect* (2023), pp. 1–12.
- [4] Margaret M. Bradley and Peter J. Lang. “Measuring emotion: The self-assessment manikin and the semantic differential”. In: *Journal of Behavior Therapy and Experimental Psychiatry* 25.1 (1994), pp. 49–59.
- [5] John T. Cacioppo and Louis G. Tassinary. “Inferring psychological significance from physiological signals.” In: *American Psychologist* 45.1 (1990), pp. 16–28.
- [6] Rafael A. Calvo et al. *The Oxford Handbook of Affective Computing*. Oxford University Press, 2015.
- [7] Walter B. Cannon. “The James-Lange Theory of Emotions: A Critical Examination and an Alternative Theory”. In: *The American Journal of Psychology* (1987).
- [8] David Chhan and Vernon J Lawhern. *An Evaluation of Tabular Neural Network Approaches for Human Affective State Classification from Physiological Signals*. Tech. rep. 2022.
- [9] Nicholas A. Coles et al. “Fact or artifact? Demand characteristics and participants’ beliefs can moderate, but do not fully account for, the effects of facial feedback on emotional experience”. In: *Journal of Personality and Social Psychology* (2022).
- [10] Maciej Dzieżyc et al. “Can We Ditch Feature Engineering? End-to-End Deep Learning for Affect Recognition from Physiological Sensor Data”. In: *Sensors* (2020).
- [11] Nick Erickson et al. *AutoGluon-Tabular: Robust and Accurate AutoML for Structured Data*. Mar. 13, 2020.
- [12] Yulia Golland et al. “Studying the dynamics of autonomic activity during emotional experience”. In: *Psychophysiology* 51.11 (2014), pp. 1101–1111.
- [13] Hatice Gunes and Hayley Hung. “Is automatic facial expression recognition of emotions coming to a dead end? The rise of the new kids on the block”. In: *Image and Vision Computing* 55 (2016), pp. 6–8.
- [14] Murtadha D. Hssayeni and Behnanumbersz Ghoraani. “Multi-Modal Physiological Data Fusion for Affect Estimation Using Deep Learning”. In: *IEEE Access* 9 (2021), pp. 21642–21652.
- [15] Gil Keren et al. “End-to-end learning for dimensional emotion recognition from physiological signals”. In: *2017 ICME*. 2017, pp. 985–990.
- [16] Shan Li and Weihong Deng. “Deep Facial Expression Recognition: A Survey”. In: *IEEE Trans. Affect* 13.3 (2022), pp. 1195–1215.
- [17] Dominique Makowski et al. “NeuroKit2: A Python toolbox for neurophysiological signal processing”. In: *Behavior Research Methods* (2021).
- [18] Soroosh Mariooryad and Carlos Busso. “Analysis and compensation of the reaction lag of evaluators in continuous emotional annotations”. In: *2013 Humaine Association Conference on ACII*. IEEE. 2013, pp. 85–90.
- [19] Soroosh Mariooryad and Carlos Busso. “Correcting time-continuum emotional labels by modeling the reaction lag of evaluators”. In: *IEEE Trans. Affect* 6.2 (2014), pp. 97–108.
- [20] Mihalis A Nicolaou et al. “Automatic segmentation of spontaneous data using dimensional labels from multiple coders”. In: *Proc. of LREC Int. Workshop on Multimodal Corpora: Advances in Capturing, Coding and Analyzing Multimodality*. 2010, pp. 43–48.
- [21] Jérémie Nicolle et al. “Robust continuous prediction of human emotions using multiscale dynamic cues”. In: *Proceedings of the 14th ACM ICMI*. 2012, pp. 501–508.
- [22] Lorenzo Pasquini et al. “Dynamic autonomic nervous system states arise during emotions and manifest in basal physiology”. In: *Psychophysiology* (2023).
- [23] Monica Perusquía-Hernández et al. “Smile Action Unit detection from distal wearable Electromyography and Computer Vision”. In: *16th IEEE International Conference on Automatic Face and Gesture Recognition*. 2021.
- [24] Fabien Ringeval et al. “Prediction of Asynchronous Dimensional Emotion Ratings from Audiovisual and Physiological Data.” In: *Pattern Recognition Letters* 66 (2014).
- [25] James A Russell et al. “Affect Grid: A Single-Item Scale of Pleasure and Arousal”. In: *Journal of Personality and Social Psychology* 57.3 (1989), pp. 493–502.
- [26] Maximilian Schmitt et al. “At the Border of Acoustics and Linguistics: Bag-of-Audio-Words for the Recognition of Emotions in Speech”. In: *Interspeech 2016*. 2016, pp. 495–499.
- [27] Karan Sharma et al. “A dataset of continuous affect annotations and physiological signals for emotion analysis”. In: *Scientific Data* 6.1 (2019), p. 196.
- [28] Pekka Siirtola et al. “Predicting Emotion with Biosignals: A Comparison of Classification and Regression Models for Estimating Valence and Arousal Level Using Wearable Sensors”. In: *Sensors* 23.3 (2023), p. 1598.
- [29] Alexandria K. Vail et al. “Visual attention in schizophrenia: Eye contact and gaze aversion during clinical interactions”. In: *2017 Seventh International Conference on ACII*. 2017, pp. 490–497.
- [30] Tianyi Zhang et al. “CorrNet: Fine-Grained Emotion Recognition for Video Watching Using Wearable Physiological Sensors”. In: *Sensors* 21.1 (2021), p. 52.