

Strength modelling for real-world automatic continuous affect recognition from audiovisual signals[☆]

Jing Han^a, Zixing Zhang^{a,*}, Nicholas Cummins^a, Fabien Ringeval^{a,b}, Björn Schuller^{a,c}

^aChair of Complex and Intelligent Systems, University of Passau, Innstr. 41, Passau 94032, Germany

^bLaboratoire d'Informatique de Grenoble, Université Grenoble Alpes, 700 Avenue Centrale, Grenoble 38058, France

^cDepartment of Computing, Imperial College London, 180 Queens' Gate, London SW7 2AZ, UK

1. Introduction

Automatic affect recognition plays an essential role in smart conversational agent systems that aim to enable natural, intuitive, and friendly human-machine interaction. Early works in this field have focused on the recognition of prototypic expressions in terms of basic emotional states, and on the data collected in laboratory settings, where speakers either act or are induced with predefined emotional categories and content [9,29,30,47]. Recently, an increasing amount of research efforts have converged into dimensional approaches for rating naturalistic affective behaviours by continuous dimensions (e. g., arousal and valence) along the time continuum from audio, video, and music signals [8,10,16,24,32,33,39,46]. This trend is partially due to the benefits of being able to encode small difference in

affect over time and distinguish the subtle and complex spontaneous affective states. Furthermore, the affective computing community is moving toward combining multiple modalities (e. g., audio and video) for the analysis and recognition of human emotion [19,23,34,43,49], owing to (i) the easy access to various sensors like camera and microphone, and (ii) the complementary information that can be given from different modalities.

In this regard, this paper focuses on the realistic time- and value-continuous affect (emotion) recognition from audiovisual signals in the arousal and valence dimensional space. To handle this regression task, a variety of models have been investigated. For instance, *Support Vector Machine for Regression* (SVR) is arguably the most frequently employed approach owing to its mature theoretical foundation. Further, SVR is regarded as a baseline regression approach for many continuous affective computing tasks [27,31,36]. More recently, memory-enhanced *Recurrent Neural Networks* (RNNs), namely *Long Short-Term Memory RNNs* (LSTM-RNNs) [14], have started to receive greater attention in the sequential pattern recognition community [7,26,48,50]. A particular advantage offered by

[☆] This paper has been recommended for acceptance by Mohammad Soleymani.

* Corresponding author.

E-mail address: zixing.zhang@uni-passau.de (Z. Zhang).

LSTM-RNNs is a powerful capability to learn longer-term contextual information through the implementation of three memory gates in the hidden neurons. Wöllmer et al. [41] was among the first to apply LSTM-RNN on acoustic features for continuous affect recognition. This technique has also been successfully employed for other modalities (e. g., video, and physiological signals) [2,21,26].

Numerous studies have been performed to compare the advantages offered by a wide range of modelling techniques, including the aforementioned, for continuous affect recognition [21,27,35]. However, no clear observations can be drawn as to the superiority of any of them. For instance, the work in [21] compared the performance of SVR and *Bidirectional LSTM-RNNs* (BLSTM-RNNs) on the Sensitive Artificial Listener database [20], and the results indicate that the latter performed better on a reduced set of 15 acoustic *Low-Level-Descriptors* (LLD). However, the opposite conclusion was drawn in [35], where SVR was shown to be superior to LSTM-RNNs on the same database with functionals computed over a large ensemble of LLDs. Other results in the literature confirm this inconsistent performance observation between SVR and diverse neural networks like (B)LSTM-RNNs and *Feed-forward Neural Networks* (FNNs) [27]. A possible rationale behind this is the fact that each prediction model has its advantages and disadvantages. For example, SVRs cannot explicitly model contextual dependencies, whereas LSTM-RNNs are highly sensitive to overfitting.

The majority of previous studies have tended to explore the advantages (strength) of these models independently or in conventional early or late fusion strategies. However, recent results indicate that there may be significant benefits in fusing two, or more, models in hierarchical or ordered manner [15,18,22]. Motivated by these initial promising results, we propose a *Strength Modelling* approach, in which the strength of one model, as represented by its predictions, is concatenated with the original feature space which is then used as the basis for regression analysis in a subsequent model.

The major contributions of this study include: (1) proposing the novel machine learning framework of Strength Modelling specifically designed to take advantage of the benefits offered by various regression models namely SVR and LSTM-RNNs; (2) investigating the effectiveness of Strength Modelling for value- and time-continuous emotion regression on two *spontaneous multimodal* affective databases (RECOLA and SEMAINE); and (3) comprehensively analysing the robustness of Strength Modelling by integrating the proposed framework into frequently used multimodal fusion techniques namely early and late fusion.

The remainder of the present article is organised as follows: [Section 2](#) first discusses related works; [Section 3](#) then presents Strength Modelling in details and briefly reviews both the SVR and memory-enhanced RNNs; [Section 4](#) describes the selected spontaneous affective multimodal databases and corresponding audio and video feature sets; [Section 5](#) offers an extensive set of experiments conducted to exemplify the effectiveness and the robustness of our proposed approach; finally, [Section 6](#) concludes this work and discusses potential avenues for future work.

2. Related work

In the literature for multimodal affect recognition, a number of fusion approaches have been proposed and studied [45], with the majority of them relevant to *early* (aka *feature-level*) or *late* (aka *decision-level*) fusion. Early fusion is implemented by concatenating all the features from multiple modalities into one combined feature vector, which will then be used as the input for a machine learning technique. The benefit of early fusion is that, it allows a classifier to take advantage of the complementarity that exists between, for example, the audio and video feature spaces. The empirical experiments offered in [2,15,26] have shown that the early fusion strategy can deliver better results than the strategies without feature fusion.

Late fusion involves combining predictions obtained from individual learners (models) to come up with a final prediction. They normally consist of two steps: 1) generating different learners; and 2) combining the predictions of multiple learners. To generate different learners, there are two primary ways which are separately based on different *modalities* and *models*. Modality-based ways combines the output from learners trained on *different* modalities. Examples of this learner generation in the literature include [12,15,22,37], where multiple SVRs or LSTM-RNNs are trained separately for different modalities (e.g. audio and video). Model-based ways, on the other hand, aims to exploit information gained from multiple learners trained on a single modality. For example in [25], predictions obtained by 20 different topology structures of Deep Belief Networks (DBNs). However, due to the similarity of characteristics of different DBNs, the predictions cannot provide many variations that could be mutually complemented and improve the system performance. To combine the predictions of multiple learners, a straightforward way is to apply simple or weighted averaging (or voting) approach, such as Simple Linear Regression (SLR) [15,36]. Another common approach is to perform stacking [44]. In doing this, all the predictions from different learners are stacked and used as inputs of a subsequent non-linear model (e.g., SVR, LSTM-RNN) trained to make a final decision [12,25,37].

Different from these fusion strategies, our proposed Strength Modelling paradigm operates on a *single* feature space. Using an initial model, it gains a set of predictions which are then fused with the original feature set for use as a new feature space in a subsequent model. This offers the framework a vital important advantage as the single modality setting is often faced in affect recognition tasks, for example, if when either face or voice samples are missing in a particular recording.

Indeed, Strength Modelling can be viewed as an *intermediate* fusion technology, which lies in the middle of the early and late fusion stages. Strength Modelling can therefore not only work independently of, but also be simply integrated into early and late fusion approaches. To the best of our knowledge, intermediate fusion techniques are not widely used in the machine learning community. Hermansky et al. [13] introduced a tandem structure that combines the output of a discriminative trained neural nets with dynamic classifiers such as *Hidden Markov Models* (HMMs), and applied it efficiently for speech recognition. This structure was further extended into a BLSTM-HMM [40,42]. In this approach the BLSTM networks provides a discrete phoneme prediction feature, together with continuous *Mel-Frequency Cepstral Coefficients* (MFCCs), for the HMMs that recognise speech.

For multimodal affect recognition, a relevant approach – Parallel Interacting Multiview Learning (PIML) – was proposed in [17] for the prediction of protein sub-nuclear locations. The approach exploits different modalities that are mutually learned in a parallel and hierarchical way to make a final decision. Reported results show that this approach is more suitable than the use of early fusion (merging all features). Compared to our approach, that aims at taking advantages of different models from a same modality, the focus of PIML is rather on exploiting the benefit from different modalities. Further, similar to early fusion approaches, PIML operates under a concurrence assumption of multiple modalities.

Strength Modelling is similar to the *Output Associative Relevance Vector Machine* (OA-RVM) regression framework originally proposed in [22]. The OA-RVM framework attempts to incorporate the contextual relationships that exist within and between different affective dimensions and various multimodal feature spaces, by training a secondary RVM with an initial set of multi-dimensional output predictions (learnt using any prediction scheme) concatenated with the original input features spaces. Additionally, the OA-RVM framework also attempts to capture the temporal dynamics by employing a sliding window framework that incorporates both past and future initial

outputs into the new feature space. Results presented in [15] indicate that the OA-RVM framework, is better suited to affect recognition problems than both conventional early and late fusion. Recently the OA-RVM model was extended in [18] to be multivariate, i.e., predicting multiple continuous output variables simultaneously.

Similar to Strength Modelling, OA-RVM systems take input features and output predictions into consideration to train a subsequent regression model to perform the final affective predictions. However, the strength of the OA-RVM framework is that it is underpinned by the RVM. Results in [15] indicate that, the framework is not as successful when using either a SVR or a SLR as the secondary model. Further, the OA-RVM is non-casual and requires careful tuning to find suitable window lengths in which to combine the initial outputs; this can take considerable time and effort. The proposed Strength Modelling framework, however, is designed to work with *any* combination of learning paradigms. Furthermore, Strength Modelling is casual; it combines input features and predictions on a frame-by-frame basis. This is a strong advantage over the OA-RVM in terms of employment in real-time scenarios (beyond the scope of this paper).

3. Strength Modelling

3.1. Strength Modelling

The proposed Strength Modelling framework for affect prediction is depicted in Fig. 1. As can be seen, the first regression model ($Model_1$) generates the original estimate \hat{y}_t based on the feature vector \mathbf{x}_t . Then, \hat{y}_t is concatenated with \mathbf{x}_t pair-wise as the input of the second model ($Model_2$) to learn the expected prediction y_t .

To implement the Strength Modelling for these suitable combination of individual models, $Model_1$ and $Model_2$ are trained subsequently, in other words, $Model_2$ takes the predictive ability of $Model_1$ into account for training. The procedure is given as follows:

- First, $Model_1$ is trained with \mathbf{x}_t to obtain the prediction \hat{y}_t .
- Then, $Model_2$ is trained with $[\mathbf{x}_t, \hat{y}_t]$ to learn the expected prediction y_t .

Whilst the framework should work with any arbitrary modelling technique we have selected two commonly used, in the context of affect recognition, for our initial investigations, namely the SVR and BLSTM-RNNs which are briefly reviewed in the subsequent subsection.

3.2. Regression models

SVR is extended from Support Vector Machine (SVM) to solve regression problems. It was first introduced in [4] and is one of the most dominant methods in the context of machine learning, particularly in emotion recognition [1,27]. Applying the SVR for a regression task, the target is to optimise the generalisation bounds for regression in the high-dimension feature space by using a ϵ -insensitive loss function which is used to measure the cost of the errors of the prediction. At the same time, a predefined hyperparameter C is set accordingly for different cases to balance the emphasis on the errors and the generalisation performance.

Normally, the high-dimension feature space is mapped from the initial feature space with a non-linear kernel function. However, in

our study, we use a linear kernel function, as the features in our cases (cf. Section 4.2) perform quite well for affect prediction in the original feature space, similar to [36].

One of the most important advantages of SVR is the convex optimisation function, the characteristics of which gives the benefit that the global optimal solution can be obtained. Moreover, SVR is learned by minimising an upper bound on the expected risk, as opposed to the neural networks trained by minimising the errors on all training data, which equips SVR a superior ability to generalise [11]. For a more in-depth explanation of the SVR paradigm the reader is referred to [4].

The other model utilised in our study is BLSTM-RNN which has been successfully applied to continuous emotion prediction [26] as well as for other regression tasks, such as speech dereverberation [48] and non-linguistic vocalisations classification [24]. In general, it is composed of one input layer, one or multiple hidden layers, and one output layer [14]. The bidirectional hidden layers separately process the input sequences in a forward and a backward order and connect to the same output layer which fuses them.

Compared with traditional RNNs, it introduces recurrently connected memory blocks to replace the network neurons in the hidden layers. Each block consists of a self-connected memory cell and three gate units, namely input, output, and forget gate. These three gates allow the network to learn when to write, read, or reset the value in the memory cell. Such a structure grants BLSTM-RNN to learn past and future context in both short and long range. For a more in-depth explanation of BLSTM-RNNs the reader is referred to [14].

It is worth noting that these paradigms bring distinct sets of advantages and disadvantages to the framework:

- The SVR model is more likely to achieve the global optimal solution, but it is not context-sensitive [21];
- The BLSTM-RNN model is easily trapped in a local minimum which can be hardly avoided and has a risk of over-fitting [7], while it is good at capturing the correlation between the past and the future information [21].

In this paper, $Model_1$ and $Model_2$ in Fig. 1 could be either an SVR model or a BLSTM-RNN model, resulting in four possible permutations, i.e., SVR-SVR (S-S), SVR-BLSTM (S-B), BLSTM-SVR (B-S), BLSTM-BLSTM (B-B). It is worth noting that the B-B structure can be regarded as a variation of the neural networks in a deep structure. Note, the S-S structure is not considered, because SVR training is achieved by solving a large margin separator. Therefore, it is unlikely to get any advantage in concatenating a set of SVR predictions with its feature space for subsequent SVR based regression analysis.

3.3. Strength Modelling with early and late fusion strategies

As previously discussed (Section 2), the Strength Modelling framework can be applied in both early and late fusion strategies. Traditional early fusion combines multiple feature spaces into one single set. When integrating Strength Modelling with early fusion, the initial predictions gained from models trained on the different feature sets are also concatenated to form a new feature vector. The new feature vector is then used as the basis for the final regression analysis via a subsequent model (Fig. 2).

Strength Modelling can also be integrated with late fusion using three different approaches, i.e., (i) modality-based, (ii) model-based, and (iii) modality- and model-based (Fig. 3). *Modality*-based fusion combines the decisions from multiple independent modalities (i.e., audio and video in our case) with the same regression model; whilst *model*-based approach fuses the decisions from multiple different models (i.e., SVR and BLSTM-RNN in our case) within the same modality; and *modality- and model*-based approach is the combination of the above two approaches, regardless of which modality or

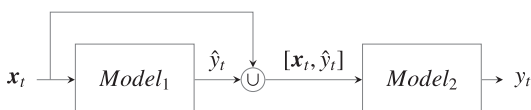


Fig. 1. Overview of the Strength Modelling framework.

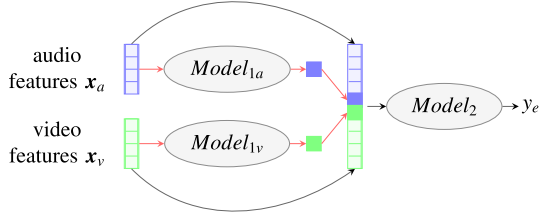


Fig. 2. Strength Modelling with early fusion strategy.

model is employed. For all three techniques the fusion weights are learnt using a linear regression model:

$$y_l = \epsilon + \sum_{i=1}^N \gamma_i \cdot y_i, \quad (1)$$

where y_i denotes the original prediction of the model i from N available ones; ϵ and γ_i are the bias and weights estimated on the development partition; and y_l is the final prediction.

4. Selected databases and features

For the transparency of experiments, we utilised the widely used multimodal continuously labelled affective databases – RECOLA [28] and SEMAINE [20], which have been adopted as standard databases for the AudioVisual Emotion Challenges (AVEC) in 2015/2016 [27,36] and in 2012 [31], respectively. Both databases were designed to study socio-affective behaviours from multimodal data.

4.1. Databases

4.1.1. RECOLA

The RECOLA database was recorded in the context of remote collaborative work. Spontaneous interactions were collected during resolving of a collaborative task that was performed in dyads and remotely through video conference. The corpus consists of multimodal signals, i.e., audio, video, Electro-CardioGram (ECG), and Electro-Dermal Activity (EDA), which were recorded continuously and synchronously from 27 French-speaking participants. It is worth to mention that, these subjects have different mother tongues (French, Italian, and German), which provides further diversity in the encoding of affect. In order to ensure speaker-independence, the corpus was equally divided into three partitions (training, development /validation, and test), with each partition containing nine unique recording approximately balanced for gender, age, and mother tongue of the participants.

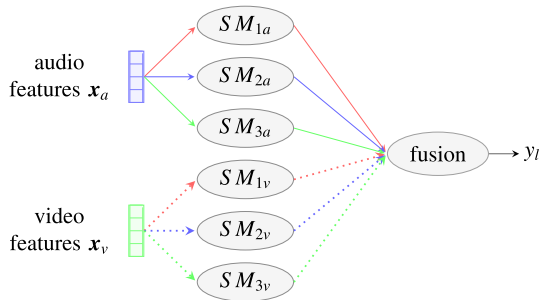


Fig. 3. Strength Modelling (SM) with late fusion strategy. Fused predictions are from multiple independent modalities with the same model (denoted by the red, green, or blue lines), multiple independent models within the same modality (denoted by the solid or dotted lines), or the combination.

To annotate the corpus, value- and time-continuous dimensional affect ratings in terms of arousal and valence were performed by six French-speaking raters (three males and three females) for the first five minutes of all recording sequences. The obtained labels were then resampled at a constant frame rate of 40 ms, and averaged over all raters by considering inter-evaluator agreement, to provide a ‘gold standard’ [28].

4.1.2. SEMAINE

The SEMAINE database was recorded in conversations between humans and artificially intelligent agents. In the recording scenario, a user was asked to talk with four emotionally stereotyped characters, which are even-tempered and sensible, happy and out-going, angry and confrontational, and sad and depressive, respectively.

For our experiments, the 24 recordings of the Solid-Sensitive Artificial Listener (Solid-SAL) part of the database were used, in which the characters were role-played. Each recording contains approximately four character conversation sessions. This Solid-SAL part was then equally split into three partitions: a training, development, and test partition, resulting in 8 recordings and 32 sessions per partition except for the training partition that contains 31 sessions. For more information on this database, the readers are referred to [31].

All sessions were annotated in continuous time and continuous value in terms of arousal and valence by two to eight raters, with the majority annotated by six raters. Different from RECOLA, the simple mean over the obtained labels was then taken to provide a single label as ‘gold standard’ for each dimension.

4.2. Audiovisual feature sets

For the acoustic features, we used the openSMILE toolkit [5] to generate 13 LLDs, i.e., 1 log energy and 12 MFCCs, with a frame window size of 25 ms at a step size of 10 ms. Rather than the official acoustic features, MFCCs were chosen as the LLDs since preliminary testing (results not given) indicated that they were more effective in association with both RECOLA [27,36] and SEMAINE [31]. The arithmetic mean and the coefficient of variance were then computed over the sequential LLDs with a window size of 8 s at a step size of 40 ms, resulting in 26 raw features for each functional window. Note that, for SEMAINE the window step size was set to 400 ms in order to reduce the computational workload in the machine learning process. Thus, the total numbers of the extracted segments of the training, development, and test partitions were 67.5 k, 67.5 k, 67.5 k for RECOLA, and were, respectively, 24.4 k, 21.8 k, and 19.4 k for SEMAINE.

For the visual features, we retained the official features for both RECOLA and SEMAINE. As to RECOLA, 49 facial landmarks were tracked firstly, as illustrated in Fig. 4. The detected face regions included left and right eyebrows (five points respectively), the nose (nine points), the left and right eyes (six points respectively), the outer mouth (12 points), and the inner mouth (six points). Then, the landmarks were aligned with a mean shape from stable points (located on the eye corners and on the nose region).

As features for each frame, 316 features were extracted, consisting of 196 features by computing the difference between the coordinates of the aligned landmarks and those from the mean shape and between the aligned landmark locations in the previous and the current frame, 71 ones by calculating the Euclidean distances (L2-norm) and the angles (in radians) between the points in three different groups, and another 49 ones by computing the Euclidean distance between the median of the stable landmarks and each aligned landmark in a video frame. For more details on the feature extraction process the reader is referred to [27].

Again, the functionals (arithmetic mean and coefficient of variance) were computed over the sequential 316 features within a fixed length window (8 s) that shifted forward at a rate of 40 ms. As a



Fig. 4. Illustration of the facial landmark features extraction from RECOLA database.

result, 632 raw features for each functional window were included in the geometric set. Feature reduction was also conducted by applying a *Principal Component Analysis* (PCA) to reduce the dimensionality of the geometric features, retaining 95% of the variance in the original data. The final dimensionality of the reduced video feature set is 49. It should be noted that a facial activity detector was used in conjunction with the video feature extraction; video features were not extracted for the frames where no face was detected, resulting in the number of video segments somewhat less than that of audio segments.

As to SEMAINE, 5908 frame-level features were provided as the video baseline features. In this feature set, eight features describes the position and pose of the face and eyes, and the rest are dense local appearance descriptors. For appearance descriptors, the uniform Local Binary Patterns (LBP) were used. Specifically, the registered face region was divided into 10×10 blocks, and the LBP operator was then applied to each block (59 features per block) followed by concatenating features of all blocks, resulting to another 5900 features.

Further, to generate features on window-level, in this paper we used the method based on max-pooling. Specifically, the maximum of features were calculated with a window size of 8 s at a step size of 400 ms, to keep consistent with the audio features. We applied PCA for feature reduction on these window-level representations and generated 112 features, retaining 95% of the variance in the original data. To keep in line with RECOLA, we selected the first 49 principal components as the final video features.

5. Experiments and results

This section empirically evaluates the proposed Strength Modelling by large-scale experiments. We first perform Strength Modelling for the continuous affect recognition in the unimodal settings (cf. Section 5.2), i.e., audio or video. We then incorporate it with the early (cf. Section 5.3) and late (cf. Section 5.4) fusion strategies so as to investigate its robustness in the bimodal settings.

5.1. Experimental set-ups and evaluation metrics

Before the learning process, mean and variance standardisation was applied to features of all partitions. Specifically, the global means and variances were calculated from the training set, which were then applied over the development and test sets for online standardisation.

To demonstrate the effectiveness of the strength learning, we first carried out the baseline experiments, where the SVR or BLSTM-RNNs models were individually trained on the modalities of audio, video, or the combination, respectively. Specifically, the SVR was implemented in the LIBLINEAR toolkit [6] with linear kernel, and trained with L2-regularised L2-loss dual solver. The tolerance value of ϵ was set to be 0.1, and complexity (C) of the SVR was optimised by the best performance of the development set among [0.00001, 0.00002, 0.00005, 0.0001, ..., 0.2, 0.5, 1] for each modality and task.

For the BLSTM-RNNs, two bidirectional LSTM hidden layers were chosen, with each layer consisting of the same number of memory blocks (nodes). The number was optimised as well by the development set for each modality and task among [20, 40, 60, 80, 100, 120]. During network training, gradient descent was implemented with a learning rate of 10^{-5} and a momentum of 0.9. Zero mean Gaussian noise with standard deviation 0.2 was added to the input activations in the training phase so as to improve generalisation. All weights were randomly initialised in the range from -0.1 to 0.1 . Finally, the early stopping strategy was used as no improvement of the mean square error on the validation set has been observed during 20 epochs or the predefined maximum number of training epochs (150 in our case) has been executed. Furthermore, to accelerate the training process, we updated the network weights after running every mini batch of 8 sequences for computation in parallel. The training procedure was performed with our CURRENNT toolkit [38].

Herein we adapted the following naming conventions, the models trained with baseline approaches are referred to as *individual* models, whereas the ones associated with the proposed approaches are denoted as *strength* models. For the sake of a more even performance comparison the optimised parameters of individual models (i.e., SVR or BLSTM-RNN) were used in the corresponding strength models (i.e., S-B, B-S, or B-B models).

Annotation delay compensation was also performed to compensate for the temporal delay between the observable cues, as shown by the participants, and the corresponding emotion reported by the annotators [19]. Similar to [15,36], this delay was estimated in the preliminary experiments using SVR and by maximising the performance on the development partition, while shifting the gold standard annotations back in time. As in [15,36] we identified this delay to be 4 s which was duly compensated, by shifting the gold standard back in time with respect to the features, in all experiments presented.

Note that all fusion experiments require concurrent initial predictions from audio and visual modalities. However, as discussed in (Section 4.2), visual prediction cannot occur where a face has not been detected. For all fusion experiments where this occurred we replicated the initial corresponding audio prediction to fill the missing video slot.

Unless otherwise stated we report the accuracy of our systems in terms of the *Concordance Correlation Coefficient* (CCC) [27] metric:

$$\rho_c = \frac{2\rho\sigma_x\sigma_y}{\sigma_x^2 + \sigma_y^2 + (\mu_x - \mu_y)^2}, \quad (2)$$

where ρ is the *Pearson's Correlation Coefficient* (PCC) between two time series (e. g., prediction and gold-standard); μ_x and μ_y are the means of each time series; and σ_x^2 and σ_y^2 are the corresponding variances. In contrast to the PCC, CCC takes not only the linear correlation, but also the bias and variance between the two compared series into account. As a consequence, whereas PCC is insensitive to

bias and scaling issues, CCC reflects those two variations. The value of CCC is in the range of $[-1, 1]$, where $+1$ represents total concordance, -1 total discordance, and 0 no concordance at all. One may further note that, it has also been successfully used as objective function to train discriminative neural networks [39], and has been used as the official scoring metric in the last two editions of the AVEC. We further intuitively compared the difference between PCC and CCC by Fig. 5. From the figure, the obtained PCC of the two series (black and blue) is 1.000, while the obtained CCC is only 0.467 as it takes the bias of the mean and variance of the two series into account. For continuous emotion recognition, ones are often interested in not only the variation trend but also the absolute value/degree of personal emotional state. Therefore, the metric of CCC fits better for continuous emotion recognition than PCC.

In addition to CCC, results are also given in all tables in terms of *Root Mean Square Error* (RMSE), a popular metric for regression tasks. To further access the significance level of performance improvement, a statistical evaluation was carried out over the whole predictions between the proposed and the baseline approaches by means of Fisher’s r -to- z transformation [3]. Unless stated otherwise, a p value less than 0.05 indicates significance.

5.2. Affect recognition with Strength Modelling

Table 1 displays the results (RMSE and CCC) obtained from the strength models and the individual models of SVR and BLSTM-RNN on the development and test partitions of RECOLA and SEMAINE databases from the *audio*. As can be seen, the three Strength Modelling set-ups either matched or outperformed their corresponding individual models in most cases. This observation implies that the advantages of each model (i.e., SVR and BLSTM-RNN) are enhanced via Strength Modelling. In particular the performance of the BLSTM model, for both arousal and valence, was significantly boosted by the inclusion of SVR predictions (S-B) on the development and test sets. We speculate this improvement could be due to the initial SVR predictions helping the subsequent RNN avoid local minima.

Similarly, the B-S combination brought additional performance improvement for the SVR model (except the valence case of SEMAINE), although not as obvious as for the S-B model. Again, we speculate that the temporal information leveraged by the BLSTM-RNN is being exploited by the successive SVR model. The best results for both arousal and valence dimensions were achieved with the framework of B-B for RECOLA, which achieved relative gains of 6.5 % and 29.1 % for arousal and valence respectively on the test set when compared to the single BLSTM-RNN model (B). This indicates there are potential benefits for audio based affect recognition by the deep structure formed by combining two BLSTM-RNNs using the Strength Modelling framework. Additionally, one can observe that there is no much performance improvement by applying Strength Modelling in the case of the valence recognition of SEMAINE. This might be attribute to the poor performance of the baseline systems, which

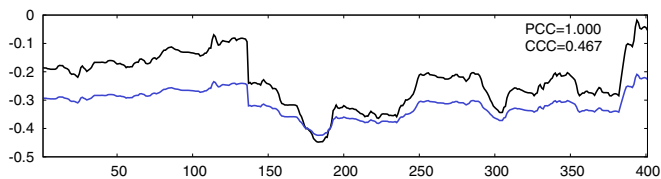


Fig. 5. Comparison of PCC and CCC between two series. The black line is gold standard from RECOLA database test partition, and the blue line is generated by shifting and scaling the gold standard.

can be regarded as noise and possibly not able to provide useful information for the other models.

The same set of experiments were also conducted on the video feature set (Table 2). As for valence, the highest CCC obtained on test set achieves at 0.477 using the S-B model for RECOLA and at 0.158 using the B-B model for SEMAINE. As expected, we observe that the models (individual or strength) trained using only acoustic features is more efficient for interpreting the dimension of arousal rather than valence. Whereas, the opposite observation is seen for models trained only on the visual features. This finding is in agreement with similar results in the literature [8,10,26].

Additionally, Strength Modelling achieved comparable or superior performance to other state-of-the-art methods applied on the RECOLA database. The OA-RVM model was used in [15,18], and the reported performance in terms of CCC, with audio features on the development set, was 0.689 for arousal [15], and 0.510 for valence using video features [18]. We achieved 0.755 with audio features for arousal, and 0.592 with video features for valence with the proposed Strength Modelling framework, showing the interest of our method.

To further highlight advantages of Strength Modelling, Fig. 6 illustrates the automatic predictions of arousal via audio signals (a) and valence via video signals (b) obtained with the best settings of the strength models and the individual models frame by frame for a single test subject from RECOLA. Note that, similar plots were observed for the other subjects in the test set. In general, the predictions generated by the proposed Strength Modelling approach are closer to the gold standard, which consequently contributes to better results in terms of CCC.

5.3. Strength Modelling integrated with early fusion

Table 3 shows the performance of both the individual and strength models integrated with the early fusion strategy. In most cases, the performance of the *individual* models of either SVR or BLSTM-RNN was significantly improved with the fused feature vector for both arousal and valence dimensions in comparison to the performance with the corresponding individual models trained only on the unimodal feature sets (Section 5.2) in most cases for both RECOLA and SEMAINE datasets.

For the strength model systems, the early fusion B-S model generally outperformed the equivalent SVR model, and the structure of S-B outperformed the equivalent BLSTM model. However, the gain obtained by Strength Modelling with the early fused features is not as obvious as that with individual models. This might be due to the higher dimensions of the fused feature sets which possibly reduce the weight of the predicted features.

5.4. Strength Modelling integrated with late fusion

This section aims to explore the feasibility of integrating Strength Modelling into three different late fusion strategies: modality-based, model-based, and the combination (see Section 3.3). A comparison of the performance of different fusion approaches, with or without Strength Modelling, is presented in Table 4. For the systems *without* Strength Modelling for RECOLA, one can observe that best individual model test set performances, 0.625 and 0.394, for arousal and valence respectively (Section 5.2) were boosted to 0.671 and 0.405 with the modality-based late fusion approach, and to 0.651 and 0.497 with the model-based late fusion approach. These results were further promoted to 0.664 and 0.549 when combining the modality- and model-based late fusion approaches. This result is in line with other results in the literature [15,27], and again confirms the importance of multimodal fusion for affect recognition. However, similar observation can only be seen on the validation set for SEMAINE, which might be due to the huge mismatch between the validation and test partitions.

Table 1
Results based on audio features only: performance comparison in terms of RMSE and CCC between the *strength*-involved models and the *individual* models of SVR (S) and BLSTM-RNN (B) on the development and test partitions of RECOLA and SEMAINE databases from the *audio* signals. The best achieved CCC is in bold. The symbol of * indicates the significance of the performance improvement over the related individual systems.

Audio based method	RECOLA				SEMAINE			
	Arousal		Valence		Arousal		Valence	
	RMSE	CCC	RMSE	CCC	RMSE	CCC	RMSE	CCC
a. On the <i>development</i> set								
S	0.126	0.714	0.149	0.331	0.218	0.399	0.262	0.172
B	0.142	0.692	0.117	0.286	0.209	0.387	0.261	0.117
B-S	0.127	0.713	0.144	0.348*	0.206	0.417*	0.255	0.179
S-B	0.122	0.753*	0.113	0.413*	0.210	0.434 *	0.262	0.172
B-B	0.122	0.755 *	0.112	0.476 *	0.206	0.417*	0.255	0.178*
b. On the <i>test</i> set								
S	0.133	0.605	0.165	0.248	0.216	0.397	0.263	0.017
B	0.155	0.625	0.119	0.282	0.202	0.317	0.256	0.008
B-S	0.133	0.606	0.160	0.264	0.205	0.332	0.258	0.006
S-B	0.133	0.665*	0.117	0.319*	0.203	0.423 *	0.262	0.017
B-B	0.133	0.666 *	0.123	0.364 *	0.205	0.332*	0.258	0.006

Interestingly when incorporating Strength Modelling into late fusion we can observe significant improvements over the corresponding non-strength set-ups. This finding confirms the effectiveness and the robustness of the proposed method for multimodal continuous affect recognition. In particular, the best test results of RECOLA, 0.685 and 0.554, were obtained by the strength models integrated with the modality- and model-based late fusion approach. This arousal result matches the performance with the AVEC 2016 affect recognition subchallenge baseline system, 0.682, which was obtained using a late fusion strategy involving eight feature sets [36].

As for SEMAINE, although obvious performance improvement can be seen on the development set, a similar observation can not be observed on the test set. This finding is possibly attributed to the mismatch between the development set and the test set, since all parameters of the training models were optimised on the development set. However, these parameters are not fit for the test set anymore.

Further, for a comparison with the OA-RVM system, we applied the same fusion system as used in [15], with only audio and video

features. The results are shown in Tables 4 and 5 for the RECOLA and SEMAINE database, respectively. It can be seen that, for both databases, the proposed methods outperform the OA-RVM technique, which further confirms the efficiency of the proposed Strength Modelling method.

In general, to provide an overview of the contributions of Strength Modelling to the continuous emotion recognition, we averaged the relative performance improvement of Strength Modelling over RECOLA and SEMAINE for arousal and valence recognition. The corresponding results from four cases (i.e., audio only, video only, early fusion, and late fusion) are displayed in Fig. 7. From the figure, one can observe an obvious performance improvement gained by Strength Modelling, except for the late fusion framework. This particular case is highly attributed to the mismatch between validation and test sets of SEMAINE as aforementioned, as all parameters of the training models were optimised on the development set. Employing some state-of-the-art generation techniques like dropout for training neural networks might help to tackle this problem in the future.

Table 2
Results based on visual features only: performance comparison in terms of RMSE and CCC between the *strength*-involved models and the *individual* models of SVR (S) and BLSTM-RNN (B) on the development and test partitions of RECOLA and SEMAINE databases from the *video* signals. The best achieved CCC is in bold. The symbol of * indicates the significance of the performance improvement over the related individual systems.

Video based method	RECOLA				SEMAINE			
	Arousal		Valence		Arousal		Valence	
	RMSE	CCC	RMSE	CCC	RMSE	CCC	RMSE	CCC
a. On the <i>development</i> set								
S	0.197	0.120	0.139	0.456	0.249	0.241	0.253	0.393
B	0.184	0.287	0.110	0.478	0.224	0.232	0.247	0.332
B-S	0.183	0.292	0.110	0.592 *	0.222	0.250	0.252	0.354
S-B	0.186	0.350 *	0.118	0.510*	0.231	0.291 *	0.242	0.405
B-B	0.185	0.344*	0.113	0.501*	0.222	0.249*	0.256	0.301
b. On the <i>test</i> set								
S	0.186	0.193	0.156	0.381	0.279	0.112	0.278	0.115
B	0.183	0.193	0.122	0.394	0.240	0.112	0.275	0.063
B-S	0.176	0.265 *	0.130	0.464*	0.235	0.072	0.285	0.043
S-B	0.186	0.196	0.121	0.477 *	0.249	0.125	0.284	0.068
B-B	0.197	0.184	0.120	0.459*	0.235	0.072	0.255	0.158 *

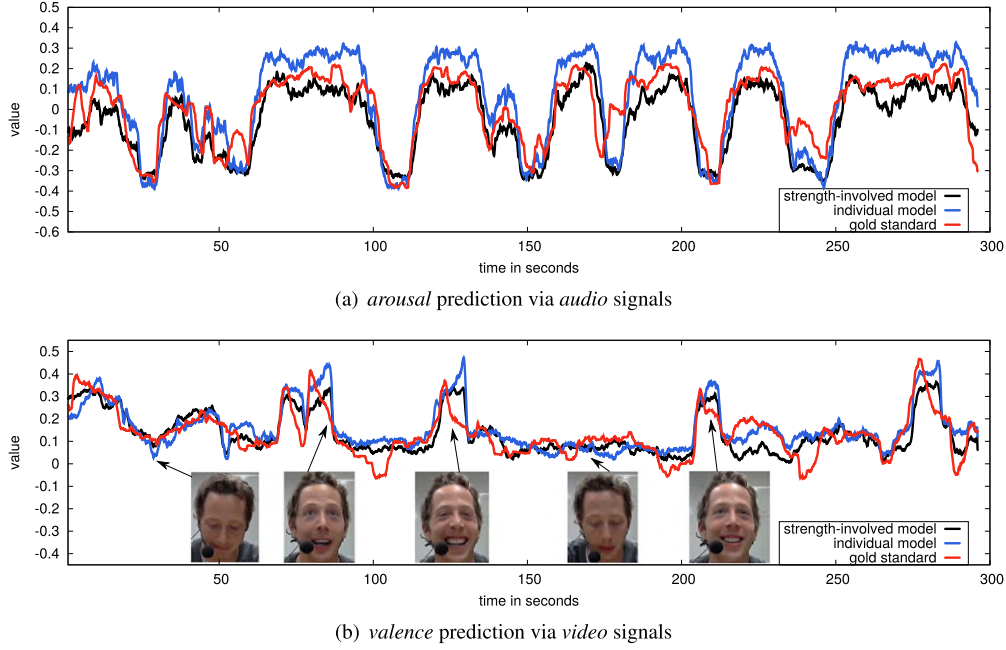


Fig. 6. Automatic prediction of arousal via audio signals (a) and valence via video signals (b) obtained with the best settings of the *strength-involved* models and *individual* models for a subject from the test partition on RECOLA database.

6. Conclusion and future work

This paper proposed and investigated a novel framework, *Strength Modelling*, for continuous audiovisual affect recognition. Strength Modelling concatenates the strength of an initial model, as represented by its predictions, with the original features to form a new feature set which is then used as the basis for regression analysis in a subsequent model.

To demonstrate the suitability of the framework, we jointly explored the benefits from two state-of-the-art regression models, i.e., Support Vector Regression (SVR) and Bidirectional Long Short-Term Memory Recurrent Neural Network (BLSTM-RNN), in three different Strength Modelling structures (SVR-BLSTM, BLSTM-SVR, BLSTM-BLSTM). Further, these three structures were evaluated in

both unimodal settings, using either audio or video signals, and the bimodal settings where early fusion and late fusion strategies were integrated.

Results gained on the widely used RECOLA and SEMAINE databases indicate that Strength Modelling can match or outperform the corresponding conventional individual models when performing affect recognition. An interesting observation was that, among our three different Strength Modelling set-ups no one case significantly outperformed the others. This demonstrates the flexibility of the proposed framework, in terms of being able to work in conjunction with different combination of regression strategies.

A further advantage of Strength Modelling is that, it can be implemented as a plug-in for use in both early and late fusion stages. Results gained from an exhaustive set of fusion experiments confirmed this

Table 3

Early fusion results on RECOLA and SEMAINE databases: performance comparison in terms of RMSE and CCC between the *strength-involved* models and the *individual* models of SVR (S) and BLSTM-RNN (B) with *early fusion strategy* on the development and test partitions of RECOLA and SEMAINE databases. The best achieved CCC is in bold. The symbol of * indicates the significance of the performance improvement over the related individual systems.

Early fusion method	RECOLA				SEMAINE			
	Arousal		Valence		Arousal		Valence	
	RMSE	CCC	RMSE	CCC	RMSE	CCC	RMSE	CCC
a. On the development set								
S	0.121	0.728	0.113	0.544	0.213	0.392	0.252	0.436
B	0.132	0.700	0.109	0.513	0.217	0.354	0.257	0.205
B-S	0.122	0.727	0.118	0.549	0.210	0.374	0.239	0.363
S-B	0.127	0.712	0.096	0.526	0.208	0.423 *	0.253	0.397
B-B	0.126	0.718*	0.095	0.542*	0.210	0.421*	0.241	0.361*
b. On the test set								
S	0.132	0.610	0.139	0.463	0.224	0.304	0.292	0.057
B	0.148	0.562	0.114	0.476	0.204	0.288	0.244	0.127
B-S	0.132	0.610	0.121	0.520 *	0.204	0.328*	0.264	0.063
S-B	0.144	0.616*	0.112	0.473	0.198	0.408 *	0.275	0.144 *
B-B	0.143	0.618 *	0.114	0.499*	0.220	0.307*	0.265	0.060

Table 4
Late fusion results on the RECOLA database: performance comparison in terms of RMSE and CCC between the *strength*-involved models and the *individual* models of SVR (S) and BLSTM-RNN (B) with *late fusion strategies* (i.e., modality-based, model-based, or the combination) on the development and test partitions of RECOLA database. The best achieved CCC is in bold. The symbol of * indicates the significance of the performance improvement over the related individual systems.

Late fusion	RECOLA							
	Arousal				Valence			
	Dev.		Test		Dev.		Test	
	RMSE	CCC	RMSE	CCC	RMSE	CCC	RMSE	CCC
a. Modality-based								
A + V; S	0.117	0.777	0.134	0.654	0.128	0.493	0.149	0.386
A + V; B	0.126	0.736	0.134	0.671	0.104	0.475	0.113	0.405
A + V; B-S	0.114	0.791 *	0.130	0.668	0.090	0.664 *	0.105	0.542 *
A + V; S-B	0.117	0.778	0.128	0.681 *	0.096	0.586*	0.105	0.495*
A + V; B-B	0.117	0.779*	0.130	0.680*	0.095	0.601*	0.106	0.506*
b. Model-based								
A; S + B	0.119	0.771	0.132	0.651	0.112	0.335	0.117	0.284
V; S + B	0.179	0.230	0.172	0.184	0.096	0.588	0.110	0.497
A; (B-S) + (S-B) + (B-B)	0.117	0.778 *	0.132	0.664 *	0.108	0.409 *	0.120	0.303 *
V; (B-S) + (S-B) + (B-B)	0.171	0.344 *	0.171	0.222 *	0.095	0.599 *	0.111	0.477
c. Modality- and model-based								
A + V; S + B	0.113	0.795	0.130	0.664	0.089	0.670	0.107	0.549
A + V; (B-S) + (S-B) + (B-B)	0.110	0.808 *	0.127	0.685 *	0.088	0.671	0.103	0.554
State-of-the-art method								
OA-RVM	0.135	0.725	0.150	0.612	0.171	0.384	0.169	0.392

advantage. The best Strength Modelling test set results on the RECOLA dataset, 0.685 and 0.554, for arousal and valence respectively were obtained using Strength Modelling integrated into a modality- and model-based late fusion approach. These results are much higher than the ones obtained from other state-of-the-art systems. Moreover, on the SEMAINE dataset, competitive results can also be obtained.

There is a wide range of possible future research direction associated with Strength Modelling to build on this initial set of promising results. First, only two widely used regression model were investigated in the present article for affect recognition. Much of our future efforts will concentrate around assessing the suitability of more other regression approaches (e. g., Partial Least Squares Regression) for use in the framework. Investigating a more general rule of what kind of models can be implemented together in the framework help to expand

its application. In addition, it is interesting to extend the framework widely and deeply. Second, motivated by the work in [17], we will also combine the original features with the predictions from different modalities (integrating the predictions based on audio features with the original video features for a final arousal or valence prediction), rather than from different models only. Furthermore, we also plan to generalise the promising advantages offered by Strength Modelling, by evaluating its performance on other behavioural regression tasks.

Acknowledgments

This work was supported by the EU's Horizon 2020 Programme through the Innovative Action No. 645094 (SEWA) and the

Table 5
Late fusion results on the SEMAINE database: performance comparison in terms of RMSE and CCC between the *strength*-involved models and the *individual* models of SVR (S) and BLSTM-RNN (B) with *late fusion strategies* (i.e., modality-based, model-based, or the combination) on the development and test partitions of SEMAINE database. The best achieved CCC is in bold. The symbol of * indicates the significance of the performance improvement over the related individual systems.

Late fusion	SEMAINE							
	Arousal				Valence			
	Dev.		Test		Dev.		Test	
	RMSE	CCC	RMSE	CCC	RMSE	CCC	RMSE	CCC
a. Modality-based								
A + V; S	0.205	0.416	0.205	0.370	0.231	0.422	0.271	0.097
A + V; B	0.202	0.439	0.210	0.313	0.240	0.351	0.276	0.055
A + V; B-S	0.200	0.460 *	0.211	0.305	0.238	0.369	0.271	0.033
A + V; S-B	0.201	0.445	0.207	0.368	0.231	0.424	0.278	0.062
A + V; B-B	0.200	0.460 *	0.211	0.304	0.242	0.336	0.257	0.099 *
b. Model-based								
A; S + B	0.207	0.394	0.201	0.348	0.254	0.212	0.261	0.021
V; S + B	0.222	0.238	0.229	0.125	0.237	0.376	0.273	0.096
A; (B-S) + (S-B) + (B-B)	0.204	0.420 *	0.202	0.364 *	0.253	0.226	0.262	0.014
V; (B-S) + (S-B) + (B-B)	0.221	0.246	0.231	0.084	0.235	0.390 *	0.300	0.036
c. Modality- and model-based								
A + V; S + B	0.201	0.447	0.206	0.353	0.235	0.395	0.277	0.054
A + V; (B-S) + (S-B) + (B-B)	0.198	0.470 *	0.207	0.346	0.224	0.477 *	0.301	0.026
State-of-the-art method								
OA-RVM	0.253	0.433	0.247	0.346	0.312	0.315	0.351	0.021

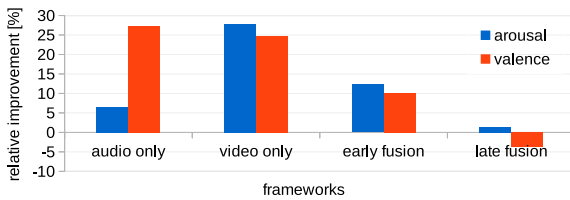


Fig. 7. Averaged relative performance improvement (in terms of CCC) cross RECOLA and SEMAINE for arousal and valence recognition. The performance of the Strength Modelling was compared with the best individual systems in the case of audio only, video only, early fusion, and late fusion frameworks.

EC's Seventh Framework Programme through the ERC Starting Grant No. 338164 (iHEARu). We further thank the Nvidia Corporation for their support of this research by Tesla K40-type GPU donation.

References

- [1] C.-Y. Chang, C.-W. Chang, J.-Y. Zheng, P.-C. Chung, Physiological emotion analysis using support vector regression, *Neurocomputing* 122 (Dec. 2013) 79–87.
- [2] L. Chao, J. Tao, M. Yang, Y. Li, Z. Wen, Long short term memory recurrent neural network based multimodal dimensional emotion recognition, *Proc. the 5th International Workshop on Audio/Visual Emotion Challenge (AVEC)*, Brisbane, Australia, 2015, pp. 65–72.
- [3] J. Cohen, P. Cohen, S.G. West, L.S. Aiken, *Applied Multiple Regression/correlation Analysis for the Behavioral Sciences*, Routledge, Abingdon, UK, 2013.
- [4] H. Drucker, C.J. Burges, L. Kaufman, A.J. Smola, V. Vapnik, *Support Vector Regression Machines*, Denver, CO, 1997, 155–161.
- [5] F. Eyben, W. Wöllmer, B. Schuller, openSMILE - the Munich versatile and fast open-source audio feature extractor, *Proc. ACM International Conference on Multimedia (ACM MM)*, Florence, Italy, 2010, pp. 1459–1462.
- [6] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, C.-J. Lin, LIBLINEAR: a library for large linear classification, *J. Mach. Learn. Res.* 9 (Jun. 2008) 1871–1874.
- [7] A. Graves, J. Schmidhuber, Framewise phoneme classification with bidirectional LSTM and other neural network architectures, *Neural Netw.* 18 (5) (Jul. 2005) 602–610.
- [8] H. Gunes, M. Pantic, Automatic, dimensional and continuous emotion recognition, *Int. J. Synth. Emot.* 1 (1) (Jan. 2010) 68–99.
- [9] H. Gunes, M. Piccardi, Automatic temporal segment detection and affect recognition from face and body display, *IEEE Trans. Syst. Man Cybern. Part B Cybern.* 39 (1) (Feb. 2009) 64–84.
- [10] H. Gunes, B. Schuller, Categorical and dimensional affect analysis in continuous input: current trends and future directions, *Image Vis. Comput.* 31 (2) (Feb. 2013) 120–136.
- [11] S.R. Gunn, *Support vector machines for classification and regression*, Tech. Rep. 14 School of Electronics and Computer Science, University of Southampton, Southampton, England, May 1998.
- [12] L. He, D. Jiang, L. Yang, E. Pei, P. Wu, H. Sahli, Multimodal affective dimension prediction using deep bidirectional long short-term memory recurrent neural networks, *Proc. the 5th International Workshop on Audio/Visual Emotion Challenge (AVEC)*, Brisbane, Australia, 2015, pp. 73–80.
- [13] H. Hermansky, D.P.W. Ellis, S. Sharma, Tandem connectionist feature extraction for conventional HMM systems, *Proc. IEEE International Conference on Audio, Speech, and Signal Processing (ICASSP)*, Istanbul, Turkey, 2000, pp. 1635–1638.
- [14] S. Hochreiter, J. Schmidhuber, Long short-term memory, *Neural Comput.* 9 (8) (Nov. 1997) 1735–1780.
- [15] Z. Huang, T. Dang, N. Cummins, B. Stasak, P. Le, V. Sethu, J. Epps, An investigation of annotation delay compensation and output-associative fusion for multimodal continuous emotion prediction, *Proc. the 5th International Workshop on Audio/Visual Emotion Challenge (AVEC)*, Brisbane, Australia, 2015, pp. 41–48.
- [16] N. Kumar, R. Gupta, T. Guha, C. Vaz, M. Van Segbroeck, J. Kim, S.S. Narayanan, Affective feature design and predicting continuous affective dimensions from music, *Proc. MediaEval*, Barcelona, Spain, 2014.
- [17] O. Kursun, H. Seker, F. Gürgen, N. Aydin, O.V. Favorov, C.O. Sakar, Parallel interacting multiview learning: an application to prediction of protein sub-nuclear location, *Proc. 9th International Conference on Information Technology and Applications in Biomedicine (ITAB)*, Larnaca, Cyprus, 2009, pp. 1–4.
- [18] A. Manandhar, K.D. Morton, P.A. Torrione, L.M. Collins, Multivariate output-associative RVM for multi-dimensional affect predictions, *World Acad. Sci. Eng. Technol. Int. J. Comput. Electr. Autom. Control. Inf. Eng.* 10 (3) (Jan. 2016) 408–415.
- [19] S. Mariooryad, C. Busso, Correcting time-continuous emotional labels by modeling the reaction lag of evaluators, *IEEE Trans. Affect. Comput.* 6 (2) (April 2015) 97–108.
- [20] G. McKeown, M. Valstar, R. Cowie, M. Pantic, M. Schroder, The semaine database: annotated multimodal records of emotionally colored conversations between a person and a limited agent, *IEEE Trans. Affect. Comput.* 3 (1) (Jan. 2012) 5–17.
- [21] M.A. Nicolaou, H. Gunes, M. Pantic, Continuous prediction of spontaneous affect from multiple cues and modalities in valence-arousal space, *IEEE Trans. Affect. Comput.* 2 (2) (Apr. 2011) 92–105.
- [22] M.A. Nicolaou, H. Gunes, M. Pantic, Output-associative RVM regression for dimensional and continuous emotion prediction, *Image Vis. Comput.* 30 (3) (Mar. 2012) 186–196.
- [23] M. Pantic, L.J.M. Rothkrantz, Toward an affect-sensitive multimodal human-computer interaction, *Proc. IEEE* 91 (9) (Sep. 2003) 1370–1390.
- [24] S. Petridis, M. Pantic, Prediction-based audiovisual fusion for classification of non-linguistic vocalisations, *IEEE Trans. Affect. Comput.* 7 (1) (Jan. 2016) 45–58.
- [25] X. Qiu, L. Zhang, Y. Ren, P.N. Suganthan, G. Amaratunga, Ensemble deep learning for regression and time series forecasting, *Proc. Computational Intelligence in Ensemble Learning (CIEL)*, Orlando, FL, 2014, pp. 1–6.
- [26] F. Ringeval, F. Eyben, E. Kroupi, A. Yuce, J.-P. Thiran, T. Ebrahimi, D. Lalanne, B. Schuller, Prediction of asynchronous dimensional emotion ratings from audiovisual and physiological data, *Pattern Recogn. Lett.* 66 (Nov. 2015) 22–30.
- [27] F. Ringeval, B. Schuller, M. Valstar, S. Jaiswal, E. Marchi, D. Lalanne, R. Cowie, M. Pantic, AV+EC 2015: the first affect recognition challenge bridging across audio, video, and physiological data, *Proc. the 5th International Workshop on Audio/Visual Emotion Challenge (AVEC)*, Brisbane, Australia, 2015, pp. 3–8.
- [28] F. Ringeval, A. Sonderegger, J. Sauer, D. Lalanne, Introducing the RECOLA multimodal corpus of remote collaborative and affective interactions, *Proc. EmoSPACE (FG)*, Shanghai, China, 2013, pp. 1–8.
- [29] B. Schuller, S. Reiter, R. Muller, M. Al-Hames, M. Lang, G. Rigoll, Speaker independent speech emotion recognition by ensemble classification, *Proc. IEEE International Conference on Multimedia and Expo (ICME)*, Amsterdam, Netherlands, 2005, pp. 864–867.
- [30] B. Schuller, G. Rigoll, M. Lang, Hidden Markov model-based speech emotion recognition, *Proc. IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, vol. I. Hong Kong, China, 2003, pp. I401–I404.
- [31] B. Schuller, M. Valster, F. Eyben, R. Cowie, M. Pantic, AVEC 2012: the continuous audio/visual emotion challenge, *Proc. the 14th ACM International Conference on Multimodal Interaction (ICMI)*, Nara, Japan, 2012, pp. 449–456.
- [32] M. Soleymani, A. Aljanaki, Y.-H. Yang, M.N. Caro, F. Eyben, K. Markov, B.W. Schuller, R. Veltkamp, F. Wieringer, F. Wieringer, Emotional analysis of music: a comparison of methods, *Proc. ACM International Conference on Multimedia*, New York, NY, 2014, pp. 1161–1164.
- [33] M. Soleymani, S. Asghari-Esfeden, Y. Fu, M. Pantic, Analysis of EEG signals and facial expressions for continuous emotion detection, *IEEE Trans. Affect. Comput.* 7 (1) (Jan. 2016) 17–28.
- [34] M. Soleymani, S. Asghari-Esfeden, M. Pantic, Y. Fu, Continuous emotion detection using EEG signals and facial expressions, *Proc. IEEE International Conference on Multimedia and Expo (ICME)*, 2014, pp. 1–6.
- [35] L. Tian, J.D. Moore, C. Lai, Emotion recognition in spontaneous and acted dialogues, *Proc. Affective Computing and Intelligent Interaction (ACII)*, Xi'an, China, 2015, pp. 698–704.
- [36] M.F. Valstar, J. Gratch, B.W. Schuller, F. Ringeval, D. Lalanne, M. Torres, S. Scherer, G. Stratou, R. Cowie, M. Pantic, AVEC 2016 - depression, mood, and emotion recognition workshop and challenge, *Proc. the 6th International Workshop on Audio/Visual Emotion Challenge*, Amsterdam, The Netherlands, 2016, pp. 3–10.
- [37] J. Wei, E. Pei, D. Jiang, H. Sahli, L. Xie, Z. Fu, Multimodal continuous affect recognition based on LSTM and multiple kernel learning, *Proc. Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA)*, Siem Reap, Cambodia, 2014, pp. 1–4.
- [38] F. Wieringer, J. Bergmann, B. Schuller, Introducing CURRENNT: the munich open-source CUDA recurrent neural network toolkit, *J. Mach. Learn. Res.* 16 (1) (Jan. 2015) 547–551.
- [39] F. Wieringer, F. Ringeval, E. Marchi, B. Schuller, Discriminatively trained recurrent neural networks for continuous dimensional emotion recognition from audio, *Proc. International Joint Conference on Artificial Intelligence (IJCAI)*, New York, NY, 2016, pp. 2196–2202.
- [40] M. Wöllmer, F. Eyben, A. Graves, B. Schuller, G. Rigoll, Bidirectional LSTM networks for context-sensitive keyword detection in a cognitive virtual agent framework, *Cogn. Comput.* 2 (3) (Apr. 2010) 180–190.
- [41] M. Wöllmer, F. Eyben, S. Reiter, B. Schuller, C. Cox, E. Douglas-Cowie, R. Cowie, Abandoning emotion classes-towards continuous emotion recognition with modelling of long-range dependencies, *Proc. INTERSPEECH*, Brisbane, Australia, 2008, pp. 597–600.
- [42] M. Wöllmer, F. Eyben, B. Schuller, Y. Sun, T. Moosmayr, N. Nguyen-Thien, Robust in-car spelling recognition-a tandem BLSTM-HMM approach, *Proc. INTERSPEECH*, Brighton, UK, 2009, pp. 2507–2510.
- [43] M. Wöllmer, M. Kaiser, F. Eyben, B. Schuller, G. Rigoll, LSTM-Modeling of continuous emotions in an audiovisual affect recognition framework, *Image Vis. Comput.* 31 (2) (Feb. 2013) 153–163.
- [44] D.H. Wolpert, Stacked generalization, *Neural Netw.* 5 (2) (Dec. 1992) 241–259.
- [45] C.-H. Wu, J.-C. Lin, W.-L. Wei, Survey on audiovisual emotion recognition: databases, features, and data fusion strategies, *APSIPA Transactions on Signal and Information Processing* 3, e12, 2014, 18 pages

- [46] Y.H. Yang, Y.C. Lin, Y.F. Su, H.H. Chen, A regression approach to music emotion recognition, *IEEE Trans. Audio Speech Lang. Process.* 16 (2) (Feb. 2008) 448–457.
- [47] Z. Zeng, M. Pantic, G.J. Roisman, T.S. Huang, A survey of affect recognition methods: audio, visual, and spontaneous expressions, *IEEE Trans. Pattern Anal. Mach. Intell.* 31 (1) (Jan. 2009) 39–58.
- [48] Z. Zhang, J. Pinto, C. Plahl, B. Schuller, D. Willett, Channel mapping using bidirectional long short-term memory for dereverberation in hand-free voice controlled devices, *IEEE Trans. Consum. Electron.* 60 (3) (Aug. 2014) 525–533.
- [49] Z. Zhang, F. Ringeval, B. Dong, E. Coutinho, E. Marchi, B. Schuller, Enhanced semi-supervised learning for multimodal emotion recognition, *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Shanghai, China, 2016, pp. 5185–5189.
- [50] Z. Zhang, F. Ringeval, J. Han, J. Deng, E. Marchi, B. Schuller, Facing realism in spontaneous emotion recognition from speech: feature enhancement by autoencoder with LSTM neural networks, *Proc. INTERSPEECH*, San Francisco, CA, 2016, pp. 3593–3597.