

# ROBUST LATENT REPRESENTATIONS VIA CROSS-MODAL TRANSLATION AND ALIGNMENT

Vandana Rajan<sup>1</sup>, Alessio Brutti<sup>2</sup>, Andrea Cavallaro<sup>1</sup>

<sup>1</sup>Centre for Intelligent Sensing, Queen Mary University of London, UK

<sup>2</sup>Fondazione Bruno Kessler, Trento, Italy

## ABSTRACT

Multi-modal learning relates information across observation modalities of the same physical phenomenon to leverage complementary information. Most multi-modal machine learning methods require that all the modalities used for training are also available for testing. This is a limitation when signals from some modalities are unavailable or severely degraded. To address this limitation, we aim to improve the testing performance of uni-modal systems using multiple modalities *during training only*. The proposed multi-modal training framework uses cross-modal translation and correlation-based latent space alignment to improve the representations of a worse performing (or weaker) modality. The translation from the weaker to the better performing (or stronger) modality generates a multi-modal intermediate encoding that is representative of both modalities. This encoding is then correlated with the stronger modality representation in a shared latent space. We validate the proposed framework on the AVEC 2016 dataset (RECOLA) for continuous emotion recognition and show the effectiveness of the framework that achieves state-of-the-art (uni-modal) performance for weaker modalities.

**Index Terms**— Cross-modal knowledge transfer, multi-modal training uni-modal testing, emotion recognition

## 1. INTRODUCTION

The term *modality* refers to the particular form in which something exists or is experienced or expressed [1]. Most physical phenomena we experience consist of multiple modalities; for example, we can *see, hear and touch* the rain; objects around us may have their own characteristic shape, sound and smell. The information to explain an event is often unevenly spread across the modalities that capture this event. *Multi-modal* machine learning uses multiple modalities to model or explain an event, whereas *uni-modal* or *mono-modal* machine learning uses only one of these modalities [2].

The uni-modal performance of individual modalities on any task may be significantly different [3]; modalities whose individual performance is comparatively better (worse) are referred to as *stronger* (*weaker*) modalities [4]. For a given task, we can rank the available modalities according to their uni-modal performance. *Multi-modal fusion* methods combine the supplementary and complementary information provided by these modalities to improve performance compared to uni-modal methods [5][6][7]. However, in general, most multi-modal fusion techniques require for the testing phase the simultaneous presence of all the modalities that were used during the model training phase [1]. This requirement becomes a severe limitation in case one or more sensors are missing or their signals are severely corrupted by noise during testing, unless such situations are explicitly handled by the modelling framework [8]. Thus, it would

be desirable to improve the testing performance of individual modalities using other modalities during training [3][9][10]. In particular, since the individual modalities corresponding to the same physical phenomenon might not perform equally well on the downstream task, our aim is to improve the uni-modal testing performance of a weaker modality by exploiting a stronger modality during training.

To this end, we propose Stronger Enhancing Weaker (SEW), a framework for improving the testing performance of a weaker modality by exploiting a stronger modality *during the training phase only*. SEW is a supervised neural network framework for knowledge transfer across modalities. During training, the stronger modality serves as an auxiliary or guiding modality that helps to create weaker-modality representations that are more discriminative than the representations obtained using uni-modal training for a classification or regression task. We achieve this by combining weaker-to-stronger modality translation and feature alignment with the stronger modality representations. This solution is based on the intuition that inter-modal translation can create intermediate representations that capture joint information between both modalities. Explicitly aligning the intermediate and the stronger modality representations further encourages the framework to discover components of the weaker modality that are maximally correlated with the stronger modality. Note that, after using SEW for training, the stronger modality is no longer required at testing. We show the effectiveness of our framework on the AVEC 2016 audio-visual continuous emotion recognition tasks and show that SEW improves the uni-modal performance of weaker modalities.

## 2. RELATED WORK

Most works on multi-modal training for uni-modal performance enhancement are designed for tasks where the different modalities are different types of images. For example, multi-modal training using RGB and depth images improves the uni-modal performance for hand gesture recognition [3]. This is achieved by forcing the modality-specific parts of the network to learn a common correlation matrix for their intermediate feature maps. Depth images also improve the test-time performance of RGB images for action recognition using an adversarial loss for feature alignment [10]. However, the modalities considered in these methods are images of equal size and the uni-modal networks have the same architecture, thus preventing their direct application to distinct modalities like audio, video and text, whose feature types and dimensionality differ. A few works have been proposed to address this problem [9][11][12]. For sentiment analysis, a sequence-to-sequence network with cyclic translation across modalities generates an intermediate representation that is robust to missing modalities during testing [11]. A multi-modal co-learning framework improves the uni-modal performance

of the text modality via training using audio, video and text modalities [12]. However, these methods primarily benefit the uni-modal performance of the text modality, which is the stronger modality for the task. In contrast, we aim to explicitly improve the weaker modality using the stronger modality during training. A joint audio-visual training and cross-modal triplet loss can be used to develop a face/speech emotion recognition system using multi-modal training [9]. However, in such a system the weaker modality may degrade the performance of the stronger modality [9].

### 3. PROPOSED METHOD

In this section, we describe our proposed *Stronger Enhancing Weaker* (SEW), a supervised neural network framework, which uses jointly the stronger and the weaker modality representations during training to improve the testing performance of the weaker modality (Figure 1). The key concepts of our framework are inter-modality translation and feature alignment. These concepts are implemented using four main modules: an inter-modal translator, an intra-modal auto-encoder, a feature alignment module and a task-specific regressor or classifier. These modules are described next.

The inter-modal translator contains an encoder,  $W_E$  and a decoder,  $S_{D1}$ . The translator takes the features of the weaker modality,  $M_W$ , as input and produces the features of the stronger modality,  $\hat{M}_S$ , as output. The encoder of the inter-modal translator, creates intermediate representations,  $m_{sw}$ , that capture joint information across modalities. This is achieved by using a translation loss,  $\mathcal{L}_{tr}$ , between the true,  $M_S$ , and the predicted,  $\hat{M}_S$ , features of the stronger modality:

$$\mathcal{L}_1 = \mathcal{L}_{tr}(M_S, \hat{M}_S). \quad (1)$$

$W_E$  is encouraged to discover components of the weaker modality that are inclined towards the stronger modality by increasing the alignment between  $m_{sw}$  and the representations of the stronger modality. For this purpose, we project the stronger modality features into the same latent space as  $m_{sw}$ . We use an intra-modal auto-encoder to create stronger modality representations,  $m_{ss}$ , of the same dimensionality as that of the inter-modal translator representations,  $m_{sw}$ . To this end, we employ an auto-encoding loss,  $\mathcal{L}_{ae}$ , between true,  $M_S$ , and predicted,  $\hat{M}_S$ , features:

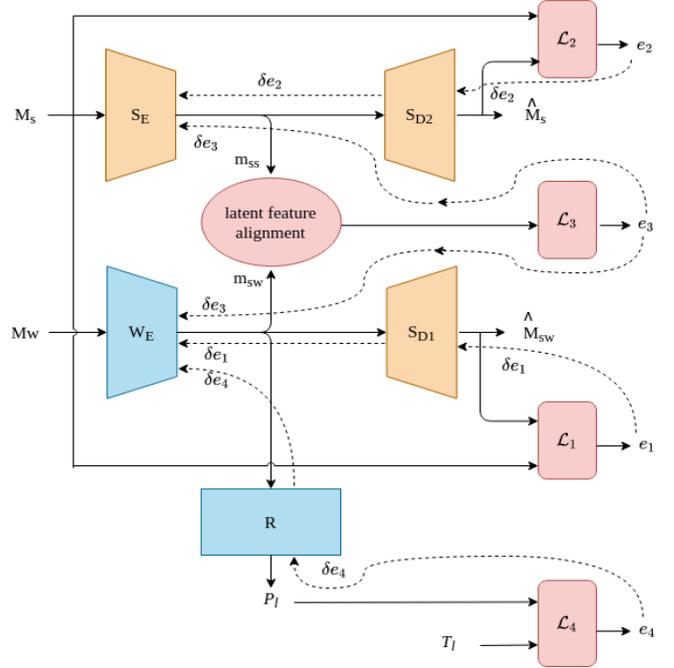
$$\mathcal{L}_2 = \mathcal{L}_{ae}(M_S, \hat{M}_S). \quad (2)$$

For modality reconstructions, we use Mean-Square-Error (MSE) as  $\mathcal{L}_{tr}$  and  $\mathcal{L}_{ae}$  [11].

A feature alignment loss,  $\mathcal{L}_{al}$ , ensures that the intermediate representations of the inter-modal translator are maximally aligned to the stronger modality representations:

$$\mathcal{L}_3 = \mathcal{L}_{al}(m_{ss}, m_{sw}). \quad (3)$$

Following [13][14], we use Canonical Correlation Analysis (CCA) for feature alignment, such that  $\mathcal{L}_{al} = -\text{CCA}$ . CCA for deep neural networks, also known as Deep CCA or DCCA, is a method to learn complex nonlinear transformations of data from two different modalities, such that the resulting representations are highly linearly correlated [15]. For a training set of size  $p$ ,  $M_s \in \mathbb{R}^{d_1 \times p}$  and  $M_w \in \mathbb{R}^{d_2 \times p}$  are the input matrices corresponding to the stronger and the weaker modalities, respectively.  $m_{ss} \in \mathbb{R}^{d \times p}$  and  $m_{sw} \in \mathbb{R}^{d \times p}$  are the representations obtained by nonlinear transformations introduced by the layers in the encoders  $S_E$  and  $W_E$ , respectively. Note



**Figure 1:** The proposed SEW training framework.  $(M_S, M_W)$  denotes a pair of stronger and weaker modality instances,  $S_E$  and  $S_{D2}$  represent intra-modal autoencoder,  $W_E$  and  $S_{D1}$  represent inter-modal translator,  $\hat{M}_S$  and  $\hat{M}_{SW}$  are the reconstructions of stronger modality from the encodings of stronger and weaker modalities respectively,  $R$  denotes the regressor/classifier connected to the inter-modal encoder,  $T_l$  and  $P_l$  stands for true and predicted labels, respectively,  $m_{ss}$  and  $m_{sw}$  represent the two latent representations,  $\mathcal{L}_1$ - $\mathcal{L}_4$  represent the 4 components of the total loss and  $e_1$ - $e_4$  are their respective error values. Dotted arrows represent the back-propagation of component error gradients. Only the blocks in cyan are retained during the deployment/inference phase.

that  $S_E$  and  $W_E$  bring the individual modalities with dimensions  $d_1$  and  $d_2$  into a common latent dimension  $d$ . If  $\theta_{es}$  and  $\theta_{ew}$  denote the vectors of all parameters of  $S_E$  and  $W_E$ , respectively, then the goal of DCCA is to jointly learn parameters for both the views such that correlation,  $(\rho)$ , between  $m_{ss}$  and  $m_{sw}$  is as high as possible, i.e.,

$$\begin{aligned} (\theta_{es}^*, \theta_{ew}^*) &= \arg \max_{\theta_{es}, \theta_{ew}} \rho(m_{ss}, m_{sw}) \\ &= \arg \max_{\theta_{es}, \theta_{ew}} \rho(S_E(M_S; \theta_{es}), W_E(M_W; \theta_{ew})). \end{aligned} \quad (4)$$

If  $\bar{m}_{ss}$  and  $\bar{m}_{sw}$  are the mean-centred versions of  $m_{ss}$  and  $m_{sw}$ , respectively, then the total correlation of the top-K components of  $m_{ss}$  and  $m_{sw}$  is the sum of the top-K singular values of the matrix,  $T = \Sigma_s^{-1/2} \Sigma_{sw} \Sigma_w^{-1/2}$ , in which the self ( $\Sigma_s, \Sigma_w$ ) and cross covariance ( $\Sigma_{sw}$ ) matrices are given by

$$\Sigma_{sw} = \frac{1}{p-1} \bar{m}_{ss} \bar{m}_{sw}^T. \quad (5)$$

$$\Sigma_s = \frac{1}{p-1} \bar{m}_{ss} \bar{m}_{ss}^T + r_1 I. \quad (6)$$

$$\Sigma_w = \frac{1}{p-1} \bar{m}_{sw} \bar{m}_{sw}^T + r_2 I. \quad (7)$$

**Table 1:** Unimodal results on the RECOLA development set. KEY - CCC: Concordance Correlation Coefficient, Acc: Binary Classification Accuracy (%), geo: geometric, app: appearance. The best uni-modal results in arousal and valence are highlighted in bold.

	Audio		Video-geo		Video-app	
	CCC	Acc	CCC	Acc	CCC	Acc
Arousal	<b>0.761</b>	<b>81.3</b>	0.482	66.4	0.492	69.1
Valence	0.543	74.2	<b>0.643</b>	<b>81.9</b>	0.489	68.4

where  $r_1 > 0$  and  $r_2 > 0$  are regularisation constants. We use the gradient of correlation obtained on the training data to determine  $(\theta_{es}^*, \theta_{ew}^*)$ .

Finally, the task-specific regressor or classification module, which takes the inter-modal translator representations as input, ensures the discriminative ability of the resulting latent space. We use a prediction loss,  $\mathcal{L}_{pr}$ , that operates on the true,  $T_l$ , and predicted task labels,  $P_l$ , as:

$$\mathcal{L}_4 = \mathcal{L}_{pr}(T_l, P_l). \quad (8)$$

The total training loss,  $\mathcal{L}$  combines the four components:

$$\mathcal{L} = \alpha\mathcal{L}_1 + \beta\mathcal{L}_2 + \gamma\mathcal{L}_3 + \mathcal{L}_4, \quad (9)$$

where  $\alpha$ ,  $\beta$  and  $\gamma$  are hyper-parameters. After training, all the components except the encoder,  $W_E$ , and the regressor,  $R$ , are removed and the stronger modality is not required at the testing (deployment) phase.

## 4. VALIDATION

In this section, we compare the performance of SEW with other uni-modal methods [16][17][18][19] and a state-of-the-art cross-modal knowledge transfer method [9]. We describe the dataset, the evaluation metrics, the details about the architecture and training, and present an ablation study, which quantifies the contributions of different parts of SEW.

### 4.1. Dataset and Evaluation Measures

We use RECOLA, the AVEC 2016 emotion recognition dataset [16], which contains audiovisual recordings of spontaneous and natural interactions from 27 French-speaking participants. Continuous dimensional emotion annotations (in the range [-1,1]) in terms of both *arousal* (level of activation or intensity) and *valence* (level of positiveness or negativness) are provided with a constant frame rate of 40 ms for the first five minutes of each recording, by averaging the annotations from all annotators and also taking the inter-evaluator agreement into consideration [16]. The dataset is equally divided into three sets, by balancing gender, age, and mother-tongue of the participants with each set consisting of nine unique recordings, resulting in 67.5k segments in total for each part (training, development and test). Since, the test labels are not publicly available, we report the results on the development set. We have used the same audio and video features as the AVEC 2015 and 2016 baselines [16] for a fair comparison with the previous literature. These are 88-D extended Geneva Minimalistic Acoustic Parameter Set (eGeMAPS) features extracted using openSMILE, LGBP-TOP based 168-D video-appearance features and 49 facial landmarks based 632-D video-geometric features. It is to be noted that the dataset provides

separate features for arousal and valence. As in [9][16], to compensate for the delay in annotation, we shift the ground-truth labels back in time by 2.4 s. This dataset is ideal for our objective, since the uni-modal performance of audio and video features varies considerably for arousal and valence, as reported in [9] and confirmed by our experiments (see Table 1). As in the AVEC 2016 challenge, we use the Concordance Correlation Coefficient (CCC) (eq. 10) as the primary evaluation metric:

$$CCC = \frac{2\sigma_{xy}^2}{\sigma_x^2 + \sigma_y^2 + (\mu_x - \mu_y)^2}, \quad (10)$$

where  $x$  and  $y$  are the true and the predicted labels, respectively, and  $\mu_x$ ,  $\mu_y$ ,  $\sigma_x$ ,  $\sigma_y$  and  $\sigma_{xy}$  refer to their means, variances and covariance, respectively. We also evaluate the binary classification results by separating the true and predicted annotations into negative [-1,0] and positive (0,1] classes.

### 4.2. Experiments

In order to identify the stronger and weaker modalities, we first assess the unimodal performances of audio, video-geometric and video-appearance features for arousal and valence using a regressor similar to [9]. The regressor consists of 4 single time-step GRU-RNN layers, each made up of 120 neurons, followed by a linear layer and trained using the MSE loss. The unimodal results thus obtained are shown in Table 1. For arousal, the performance of audio surpasses both video-geometric and video-appearance features. For valence, the video-geometric features outperform audio and video-appearance features. Thus, we have 5 cases for cross-modal knowledge transfer from stronger to weaker modalities, namely video-geo(+audio) and video-app(+audio) for arousal and audio(+video-geo), video-app(+audio) and video-app(+video-geo) for valence, where the modality in parenthesis indicates the stronger modality.

### 4.3. Architecture and Training

Because the proposed method combines multi-modal data with different characteristics, it was necessary to adjust various architectural parameters according to the characteristics of the given modalities rather than solving the problem using a generic model. Hence, the encoders and decoders of both inter-modal translator and intra-modal auto-encoder of the 5 multi-modal combinations vary from each other. Specifically, the encoder and decoder for each modality differ in terms of the number of linear layers and the number of neurons in each layer. Since the provided video-appearance features were already refined using PCA, we did not reduce the dimensionality further and used a single linear layer of size 168 for both its encoder and decoder. Thus, for all combinations that contain video-appearance features, the size of the latent dimension was 168. For all the rest, it was 128. The encoder and decoder for video-geometric features use linear layers of size [512, 256, 128] and [256, 512, 632], respectively with tanh activation between layers. For audio features, these were [108, 128] and [108, 88]. Note that 632 and 88 were chosen to match the dimensionality of the video-geometric and audio features, respectively. All the models were developed, trained and tested using PyTorch. We used the SGD optimiser with learning rate 0.001, momentum 0.7 and weight decay regularisation. The batch size was 32. The number of CCA components,  $K$ , was 10 in all the experiments. The contribution of each loss component was equally important:  $\alpha = \beta = \gamma = 1$ .

**Table 2:** Ablation results for SEW in terms of CCC and binary classification accuracy, Acc (%). KEY - geo: geometric, app: appearance.

	Arousal				Valence					
	video-geo(+audio)		video-app(+audio)		audio(+video-geo)		video-app(+audio)		video-app(+video-geo)	
	CCC	Acc	CCC	Acc	CCC	Acc	CCC	Acc	CCC	Acc
SEW	<b>0.565</b>	<b>73.6</b>	<b>0.544</b>	<b>73.6</b>	0.552	76.3	<b>0.554</b>	<b>72.2</b>	<b>0.549</b>	<b>74.1</b>
- $S_{D2}$	0.532	71.1	0.519	71.5	0.486	72.4	0.539	68.7	0.540	74.0
-CCA	0.512	70.7	0.508	69.5	0.496	73.8	0.532	67.9	0.546	74.1
- $S_{D1}$	0.514	71.0	0.523	71.0	<b>0.556</b>	<b>76.3</b>	0.514	67.8	0.505	69.3
-(CCA & $S_{D1}$ )	0.484	69.1	0.497	68.4	0.545	75.9	0.497	67.3	0.491	68.9
uni-modal	0.482	66.4	0.492	66.1	0.543	74.2	0.489	68.4	0.489	68.4

**Table 3:** Performance comparison of SEW with other methods in terms of CCC. KEY - geo: geometric, app: appearance. Best and second best results are shown in bold and italics, respectively.

	Arousal		Valence	
	video-geo	video-app	audio	video-app
SVR + offset [16]	0.379	0.483	0.455	0.474
MTL (RE) [17]	0.502	0.512	0.519	0.529
MTL (PU) [18]	0.508	0.502	0.506	0.468
DDAT (RE) [19]	<i>0.544</i>	0.539	0.508	0.528
DDAT (PU) [19]	0.513	0.518	0.498	0.514
EmoBed [9]	0.527	<b>0.549</b>	<i>0.521</i>	<b>0.564</b>
SEW	<b>0.565</b>	<i>0.544</i>	<b>0.552</b>	<i>0.554</i>

#### 4.4. Results

Table 2 reports the results using the full SEW framework as well as after ablating individual components. The bottom row provides the uni-modal results for the weaker modalities for ease of comparison with the SEW results. Comparing the last 2 rows, we can see that the SEW-(CCA& $S_{D1}$ ) results are close to the uni-modal results of the weaker modality. This is as expected since SEW-(CCA& $S_{D1}$ ) contains only the  $W_E$  and regressor with no interaction with the stronger modality. In all the 5 cases, SEW was able to improve the results from the uni-modal and SEW-(CCA& $S_{D1}$ ) models both in terms of CCC and binary accuracy. For arousal video-geo(+audio) and video-app(+audio), removing the CCA based alignment causes a drop of 0.053 and 0.036, respectively in CCC and 2.9% and 4.1%, respectively in binary accuracy. The corresponding numbers for valence audio(+video-geo) are 0.056 and 2.5%, respectively. These observations support the significance of the CCA based distribution alignment in the SEW framework. For valence video-app(+audio) and video-app(+video-geo), removing the decoder of the inter-modal translator causes a drop of 0.040 and 0.044, respectively in CCC and 4.4% and 4.8%, respectively in binary accuracy, which indicates the effectiveness of the weaker-to-stronger modality translation.

In Table 3, we compare the best uni-modal results of SEW with the 4 most relevant uni-modal models [16][17][18][19] and a cross-modal training method [9] in terms of CCC. [16] provides the baseline results on the RECOLA dataset for the AVEC 2016 challenge. The uni-modal baseline used an SVM based classifier on the individual features. SEW significantly outperforms the baseline uni-modal results for all the weaker modalities considered. Our method is also able to improve the uni-modal results for all the cases from [19], which uses difficulty awareness based training and [17][18] which uses multi-task learning. SEW outperforms EmoBed [9] for arousal video-geo(+audio) and valence

audio(+video-geo) by a margin of 0.038 and 0.031, respectively, in CCC. For arousal video-app(+audio), the performance of SEW and EmoBed are very close (0.544 and 0.549, respectively). However, for the valence video-app features, EmoBed outperforms SEW. The top and bottom rows of Table 2 show that SEW improves the uni-modal performance of the weaker modalities. Specifically, to the best of our knowledge, the best results to date on the uni-modal performance of arousal video-geometric features and valence audio features have been achieved by SEW.

## 5. CONCLUSION

We exploited the gap between the uni-modal performance of different modalities on a task to improve, using a stronger modality in a new training framework, the performance of the weaker modality. Our proposed framework, *Stronger Enhancing Weaker* (SEW), enables cross-modal knowledge transfer from the stronger to the weaker modality. The results of SEW on the RECOLA dataset for the task of continuous emotion recognition show its ability to improve the uni-modal performance of a weaker modality using a stronger modality during training.

Future work includes applying the SEW framework to other tasks involving different features and modalities as well as extending SEW to cope with multi-modal sequential data.

**Acknowledgement.** This research used QMUL’s Apocrita HPC facility, supported by QMUL Research-IT.

## 6. REFERENCES

- [1] Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe Morency, “Multimodal machine learning: A survey and taxonomy,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 41, no. 2, pp. 423–443, 2018.
- [2] Louis-Philippe Morency, Tadas Baltrušaitis, “Tutorial on Multimodal Machine Learning,” <https://www.cs.cmu.edu/~morency/MMML-Tutorial-ACL2017.pdf>, Online; accessed 01 September 2020.
- [3] Mahdi Abavisani, Hamid Reza Vaezi Joze, and Vishal M Patel, “Improving the performance of unimodal dynamic hand-gesture recognition with multimodal training,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1165–1174.
- [4] Nawid Sayed, Biagio Brattoli, and Björn Ommer, “Cross and learn: Cross-modal self-supervision,” in *Proceedings of the German Conference on Pattern Recognition*. Springer, 2018, pp. 228–243.

- [5] George Caridakis, Ginevra Castellano, Loic Kessous, Amaryllis Raouzaiou, Lori Malatesta, Stelios Asteriadis, and Kostas Karpouzis, "Multimodal emotion recognition from expressive faces, body gestures and speech," in *Proceedings of the International Conference on Artificial Intelligence Applications and Innovations*. Springer, 2007, pp. 375–388.
- [6] Panagiotis Tzirakis, George Trigeorgis, Mihalis A Nicolaou, Björn Schuller, and Stefanos Zafeiriou, "End-to-end multimodal emotion recognition using deep neural networks," *IEEE Journal of Selected Topics in Signal Processing*, vol. 11, no. 8, pp. 1301–1309, 2017.
- [7] Di Hu, Chengze Wang, Feiping Nie, and Xuelong Li, "Dense multimodal fusion for hierarchically joint representation," in *Proceedings of the International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2019, pp. 3941–3945.
- [8] Trisha Mittal, Uttaran Bhattacharya, Rohan Chandra, Aniket Bera, and Dinesh Manocha, "M3er: Multiplicative multimodal emotion recognition using facial, textual, and speech cues.," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020, pp. 1359–1367.
- [9] Jing Han, Zixing Zhang, Zhao Ren, and Björn Schuller, "Emobed: Strengthening monomodal emotion recognition via training with crossmodal emotion embeddings," *IEEE Transactions on Affective Computing*, pp. 1–1, 2019.
- [10] Nuno C Garcia, Pietro Morerio, and Vittorio Murino, "Learning with privileged information via adversarial discriminative modality distillation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 42, no. 10, pp. 2581–2593, 2020.
- [11] Hai Pham, Paul Pu Liang, Thomas Manzini, Louis-Philippe Morency, and Barnabás Póczos, "Found in translation: Learning robust joint representations by cyclic translations between modalities," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2019, vol. 33, pp. 6892–6899.
- [12] Amir Zadeh, Paul Pu Liang, and Louis-Philippe Morency, "Foundations of multimodal co-learning," *Information Fusion*, vol. 64, pp. 188–193, 2020.
- [13] Zhongkai Sun, Prathusha K Sarma, William Sethares, and Erik P Bucy, "Multi-modal sentiment analysis using deep canonical correlation analysis," in *Proceedings of Interspeech*, 2019, pp. 1323–1327.
- [14] Sri Harsha Dumpala, Rupayan Chakraborty, and Sunil Kumar Kopparapu, "Audio-visual fusion for sentiment classification using cross-modal autoencoder," in *Proceedings of the 32nd Conference on Neural Information Processing Systems*, 2019, pp. 1–4.
- [15] Galen Andrew, Raman Arora, Jeff Bilmes, and Karen Livescu, "Deep canonical correlation analysis," in *Proceedings of the International Conference on Machine Learning*. PMLR, 2013, pp. 1247–1255.
- [16] Michel Valstar, Jonathan Gratch, Björn Schuller, Fabien Ringeval, Denis Lalanne, Mercedes Torres Torres, Stefan Scherer, Giota Stratou, Roddy Cowie, and Maja Pantic, "AVEC 2016: Depression, mood, and emotion recognition workshop and challenge," in *Proceedings of the 6th International Workshop on Audio/Visual Emotion Challenge*, 2016, pp. 3–10.
- [17] Jing Han, Zixing Zhang, Fabien Ringeval, and Björn Schuller, "Reconstruction-error-based learning for continuous emotion recognition in speech," in *Proceedings of the International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2017, pp. 2367–2371.
- [18] Jing Han, Zixing Zhang, Maximilian Schmitt, Maja Pantic, and Björn Schuller, "From hard to soft: Towards more human-like emotion recognition by modelling the perception uncertainty," in *Proceedings of the 25th ACM International Conference on Multimedia*, 2017, pp. 890–897.
- [19] Zixing Zhang, Jing Han, Eduardo Coutinho, and Björn Schuller, "Dynamic difficulty awareness training for continuous emotion prediction," *IEEE Transactions on Multimedia*, vol. 21, no. 5, pp. 1289–1301, 2018.