

# SS-VAERR: Self-Supervised Apparent Emotional Reaction Recognition from Video

Marija Jegorova<sup>1</sup>, Stavros Petridis<sup>1,2</sup>, and Maja Pantic<sup>1,2</sup>

<sup>1</sup> Meta Reality Labs, London, United Kingdom

<sup>2</sup> Department of Computing, Imperial College London, United Kingdom

**Abstract**— This work focuses on the apparent emotional reaction recognition (AERR) from the video-only input, conducted in a self-supervised fashion. The network is first pre-trained on different self-supervised pretext tasks and later fine-tuned on the downstream target task. Self-supervised learning facilitates the use of pre-trained architectures and larger datasets that might be deemed unfit for the target task and yet might be useful to learn informative representations and hence provide useful initializations for further fine-tuning on smaller more suitable data. Our presented contribution is two-fold: (1) an analysis of different state-of-the-art (SOTA) pretext tasks for the video-only apparent emotional reaction recognition architecture, and (2) an analysis of various combinations of the regression and classification losses that are likely to improve the performance further. Together these two contributions result in the current state-of-the-art performance for the video-only spontaneous apparent emotional reaction recognition with continuous annotations.

## I. INTRODUCTION

Apparent emotional reaction recognition (AERR) is a broadly applicable branch of computer vision. In this paper we are going to focus on specifically the video-only domain for AERR for several reasons. First, the audio stream is not always available, and not every apparent emotional reaction is accompanied by a sound. Second, in audio-visual domain active speaker detection is a whole new problem in case of multiple speakers in the video. Finally, generalising to noisy environments can represent certain challenges for audio. Hence it would be useful to explore the efficient AERR restricted solely to the video modality for the sake of prediction robustness and broader applicability.

Further, this work explores predicting the continuous emotion characteristics - arousal and valence (in this paper we call these *continuous emotions*) instead of more traditional AERR that is concerned with classifying the *categorical emotions* (sadness, fear, surprise, etc). The reason being that the categorical emotion theory is limited in its ability to express subtle and disparate emotions [1].

Current state-of-the-art for video-only AERR are [2] and [3]. First presents a model based on probabilistic modeling of the temporal context, presenting compelling results on SEWA dataset [4]. A somewhat comparable performance is achieved by [5], using spatio-temporal higher-order convolutional neural net. Secondly, for RECOLA dataset [6] the current SOTA is TS-SATCN [3], a two-stage spatio-temporal

This work has been supported by Meta Reality Labs. With the exception of training the pre-text architectures on LRS3 dataset, which has been conducted on the servers of Imperial College London.

attention temporal convolution network. The only additional video-only AERR method to be found is Visual ResNet-50 presented in [7], also evaluated on RECOLA.

Shortage of annotated data for specific tasks and domains often represents a challenge. This can be addressed from several angles, e.g. transfer or semi-supervised learning and self-supervised learning (SSL). We focus on the SSL approach, that can use labelled and unlabelled data within the same model. It relies on the pretext training to leverage the additional data, and then serve as an initialization to the downstream training, solving target tasks.

The works that contributed to the SSL paradigm for facial data in adjacent domains are [8] and [9]. One presented a SSL framework for a number of tasks, including the AERR from images, providing results on AffectNet, large-scale facial expression image database [10], and is the SOTA for self-supervised AERR on images [8]. The other, [9], describes contrastive-learning across the video-sequences, for specifically categorical emotions on acted dataset Oulu-CASIA [11], and is SOTA for acted AERR from video. Additionally, [12] offered a unified framework for multiple tasks, but it does not surpass [7] and [3] on RECOLA [6].

To the best of our knowledge video-only self-supervised framework for natural apparent emotional reaction recognition has not yet been explored, which is what we present in this paper. We compare 3 different SSL methods for the pretext training and investigate the impact of a variety of loss functions during downstream training. We evaluate our proposed method on two different natural emotional reactions datasets (SEWA and RECOLA, [4], [6]) and achieve an improvement by up to 10% on previously published models.

**Our main contributions** can be summarized as follows: (1) a review of several pretext tasks for apparent emotional reaction recognition from video for their downstream performance across several *spontaneous* emotion datasets; (2) analysis of the impact of the combined regression and classification losses, data augmentations, and downstream learning parameters; (3) adding up to the first to our knowledge Self-Supervised Visual Apparent Emotional Reaction Recognition method for spontaneous emotions with continuous annotations, SS-VAERR. Please check Tab. I for the results.

## II. RELATED WORK

**Apparent Emotional Reaction Recognition** is a vast research field spread across different methods and domains.

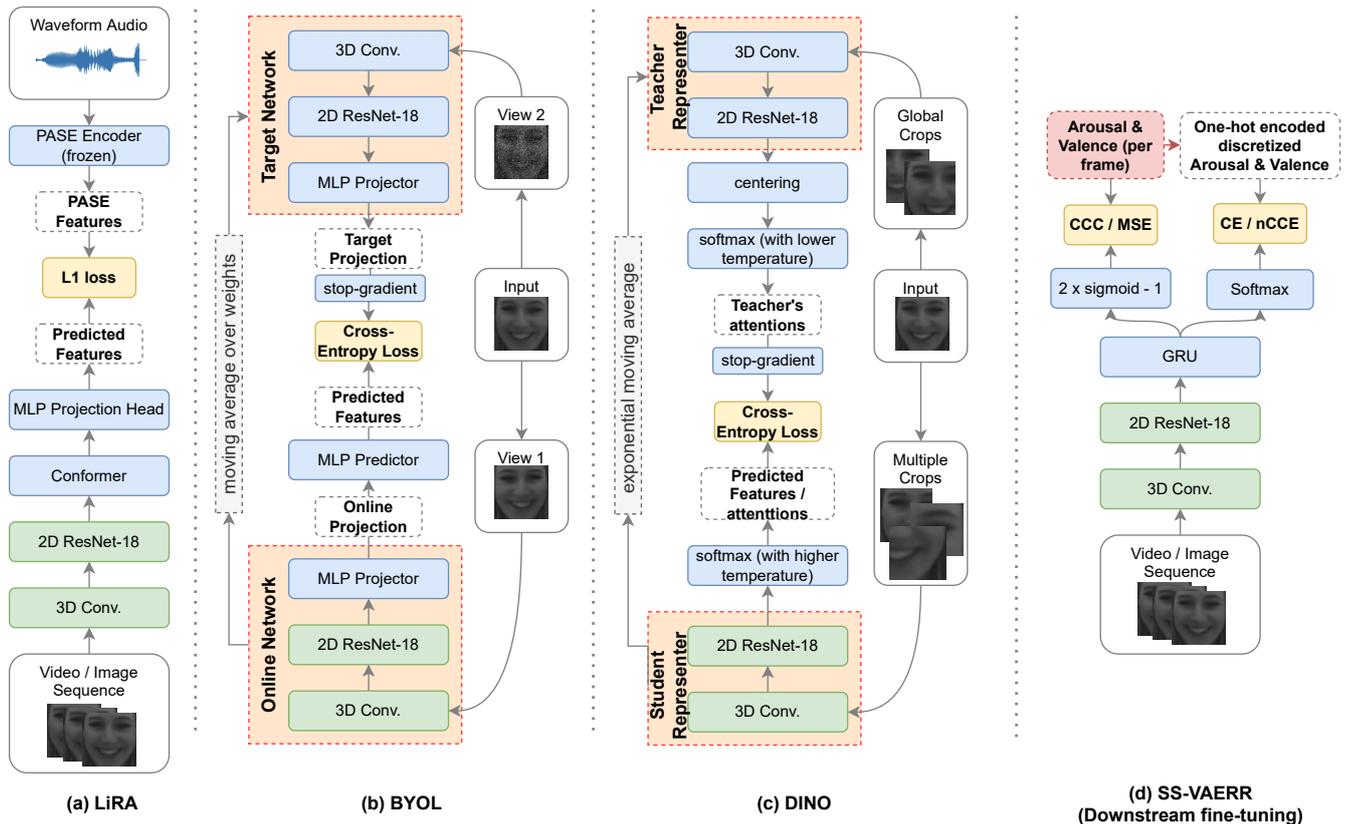


Fig. 1. A comparative overview of all the reviewed pretext architectures. Please note ResNet18 architecture is used instead of the original ResNet50 and transformers for BYOL [13] and DINO [14] correspondingly for comparability. Blue and green boxes represent the network blocks (layers, standard structures, activations/normalization). Green are the blocks used to initialize the downstream architecture. Yellow boxes represent training losses, and orange are networks with interconnected weights.

Domain-wise there is audio-based, image-based, video-based, and audio-visual AERR, additionally separated into acted and spontaneous/natural AER. We mostly focus on spontaneous and visual AER here. The results across the field are reported on different datasets, complicating the comparison, that is why SOTA is reported per dataset. There are also multiple datasets for AERR of different modalities, the ones discussed in this paper (Sec. IV-B), and others, such as AffNet [15], Oulu-CASIA [11], and AffectNet [10].

*Audio-visual or multi-modal AERR* tends to yield better results than video-only. Specifically audio is known to provide better signal for arousal [16]–[18]. Multi-modal AERR works include [19], a BLSTM-based method, using joint discrete and continuous emotion representation for AERR that holds the current SOTA for multi-modal AERR on RECOLA [6] development set. However, a ResNet50-based method presented by [7] achieves SOTA results for RECOLA on the test set, also presenting video- and audio-only results.

The most prominent examples of the *image-based AERR* include EmoFan [20] - an approach for direct estimation of facial landmarks, discrete and continuous emotions with a single neural net from facial images, and [8] - an SSL framework, purposed for a variety of downstream face-related applications, including the SOTA results for AERR on images presented on AffectNet [10].

*Video-only AERR* is a less explored field, the current SOTA being Affective Processes [2] and [3]. First is a neural processes model with a global stochastic contextual representation, task-aware temporal context modelling, and temporal context selection. Second is a two-stage spatio-temporal attention temporal convolution network.

**Self-supervised learning (SSL)** focuses on minimizing the use of human-generated annotations at training time. It is often used to leverage the large amounts of unlabelled data to aid learning on significantly smaller annotated datasets. The SSL is rooted in the assumption that solving a seemingly unrelated self-supervised *pretext task* can help to learn useful visual representations. These would serve as a good initialization point for a task of interest, *downstream task*, given that the model is well generalized and the tasks are similar enough in kind [21]–[24]. If these assumptions are violated a *negative transfer* might occur [25] - the performance would be worse than that of a model trained entirely from scratch.

There is plenty of research showing the benefits of SSL for general image datasets [26], [27]. SSL techniques vary by both the downstream and pretext tasks. Traditional pretext tasks include transformation classification [28], image inpainting [29], image colorization (from grayscale) [30], [31], and solving jigsaw puzzles [32]. A more recent field of SSL is *contrastive learning*, that relies on minimizing the distance

between the learned representations of the *positive pairs* - differently augmented versions of the same image, and maximising it for the *negative pairs* - different augmentations of different images, [33], [34]. On the plus side contrastive learning does not require labels, on the downside it is sensitive to the choice of the negative pairs as well as to the choice of the augmentations used.

An ultimate improvement on the contrastive techniques involved getting rid of the negative pairs, accomplished by BYOL[13] and DINO [14]. Instead those rely on different teacher-student-like architectures trained on a variety of image crops and augmentations, that can be interpreted as positive pairs. Not relying on negative examples accounts for a potentially better generalization because of their lower vulnerability to the systematic biases in training data.

**For the face imagery** some SSL techniques were developed for tasks such as action unit detection [35] and lip-reading [36], [37]. For AERR, SSL research can be sparse depending on the modality and target (i.e. categorical emotion vs valence and arousal). The prominent works include [8] - universal facial representation learning for images, [38] - contrastive learning method for recognition of categorical emotions from multi-view images of emoted faces.

The only current example of the self-supervised AERR from video is based on spatio-temporal contrastive learning and delivers SOTA results for synthetic categorical emotion recognition [9]. Unfortunately, it is not directly comparable to ours because 1) it is designed to perform on a lab recorded and most importantly acted dataset Oulu-CASIA [11]; 2) it aims to predict discrete emotions rather than arousal and valence. Meaning not only that it uses non-realistic posed emotion depictions, but also that every discrete acted emotion is present for every individual in the dataset. Whereas we aim for the spontaneous apparent emotional reactions with their natural distribution.

Although some work has been done in multi-modal self-supervised in-the-wild AERR [19], to our knowledge video-only self-supervised continuous spontaneous facial AERR has not been explored yet.

### III. METHODOLOGY

#### A. Shared Architecture and Pretext Tasks

As mentioned before, the similarity of the pretext and downstream tasks can significantly contribute to the usefulness of the visual representations learned during pretext training to the downstream task, hence in this study we review and evaluate the self-supervised pretext methods designed to learn visual representations.

We examine three suitable pretext methods: LiRA [36], BYOL[13], and DINO [14]. Please refer to Fig. 1 for the flowcharts complementing the explanations below. For comparability sake we have incorporated the ResNet18, [39], as part of the pretext architectures to facilitate the transfer for the downstream task. LiRA architecture, Fig. 1(a), remains unchanged. For BYOL it is a change from ResNet50 to ResNet18, Fig. 1(b). For DINO it is a change from trans-

former, [40], to ResNet18, Fig. 1(c); in principle DINO supports any architecture for both of its networks.

**Learning visual speech Representations from Audio (LiRA) [36]** is a self-supervised method for predicting visual representations of acoustic features from unlabelled speech videos. Its architecture features a ResNet18 followed by a conformer. Once it is trained, we only use the ResNet18 weights for the downstream initialization. LiRA uses the random flip ( $Prob = 50\%$ ) and the random crop of  $80 \times 80$  (out of  $96 \times 96$ ) augmentations during training.

**Bootstrap Your Own Latent (BYOL) [13]** is an approach to self-supervised visual representation learning that does not rely on negative pairs typical for previous SSL methods. Its architecture consists of online and target networks, interacting and learning from one another. The online network is trained to predict the invariant visual representation of the same image under different augmentations, while the target network learns via a slow-moving average from the online network. At inference time only the 3D convolutional layer and the ResNet18 of the online network are preserved, the visual representations are extracted from the final layer of ResNet for the downstream task. BYOL uses the following training time data augmentations: random cropping, random flip, color jittering (brightness, contrast, saturation and hue of an image, shifted by a uniformly random offset), Gaussian blurring, and solarization. For details see Appendix B in [13].

**Self-Distillation with NO labels (DINO) [14]** is yet another two-network SSL architecture, an attention-based self-distillation method using no labels, and reportedly an improvement on BYOL. The principle is also rather similar to BYOL: the two networks are student and teacher networks, they have the exact same architecture and teacher is updated as an exponential moving average of the student. There are certain additional tricks, such as centering and sharpening for the teacher network. Sharpening is a technique introducing the temperature parameter into the softmax of both networks. Temperature is lower for the teacher than for the student, it reduces noise but encourages a potential *mode collapse* (a phenomenon where network systematically produces same outputs for different inputs), whereas centering (a type of normalization specific to DINO technique with respect to teacher's previous logits) is meant to compensate for that and prevent the mode collapse. In principle DINO allows for both student and teacher to draw from a broad range of potential network architectures. We have adapted it to rely on the ResNet18 instead of the typical transformer networks for comparability with the other reviewed pretext models. DINO uses the same augmentations as BYOL. However, one of the key differences of DINO is that the teacher network only sees the global crops (covering most of the image) while the student gets to see both global and local crops and should derive similar or ideally the same representations from both.

Please note that the data augmentations during the pretext training are as in the corresponding, whereas the ones at the downstream training time are discussed in Sec. V-D.

## B. Downstream VAERR Architecture

As can be seen on Figure 1(d), the downstream task architecture uses 3D convolutional layer, followed by ResNet18, and finally by GRU, which then returns regression or classification result or both depending on what loss is chosen. For all of the pretext tasks the weights of the 3D convolutional layer and ResNet18 are passed from the pretext to downstream architecture, as the initialisation of the latter. The intuition behind the choice of the architecture is follows: ResNet18 is a reliable choice for initial image processing resulting in meaningful latent features, powered by pre-text tasks, whereas GRU is meant to capture the temporal component within the videos.

*Fine-tuned* version of each pretext method only initializes the downstream task with the shared weights from the pretext. It is later free to update these weights in accordance to the downstream training under the chosen loss function. *Frozen* version of the pretext freezes some of the layers, so they cannot be updated during the downstream training of the model. We have compared freezing the entire set of shared weights and found that freezing just the first 3D-convolutional layer of the network generally shows better results than freezing 3D-convolutional layer and ResNet18. Therefore, "Frozen" in this paper means "frozen first layer".

## C. AERR Discrete and Continuous Labels

Humans tend to generalize and discretize the facial emotions into 7-8 categorical classes such as happiness, sadness, surprise, anger, rage, disgust, boredom (and neutral). From the computational point of view, there are multiple perks to using the continuous emotion labels - valence (positivity/negativity of the emotion) and arousal (the magnitude of it) [20]. Therefore, a number of sentiment analysis datasets are annotated with continuous emotion labels only. In this paper we focus specifically on correctly predicting apparent arousal and valence.

Furthermore, there are several publications in the field suggesting that using combined classification and regression losses, called *composite losses* in this paper, on both continuous and discretized versions of the same labels improves the prediction quality drastically [2], [19]. Hence, we study composite loss functions as well as various pretext tasks.

## D. Composite Losses

The success of the models is assessed via the Concordance Correlation Coefficient (*CCC*) metric calculated on the combined videos of the test set. Rather than optimizing the performance solely for CCC, we examine the combinations of the following losses.

**Regression loss** for continuous predictions, is calculated as  $1 - CCC$ , where CCC is the Concordance Correlation Coefficient of the validation set. The idea behind this metric is to assess the correlation between the predictions and the targets, while also penalising the signals with the different means more. It can be interpreted as a version of Pearson

Coefficient weighted towards predictions with higher errors.

$$CCC(Y, \hat{Y}) = 2 \frac{\mathbb{E}(Y - \mu_Y)(\hat{Y} - \mu_{\hat{Y}})}{\sigma_{\hat{Y}} + \sigma_Y + (\mu_{\hat{Y}} + \mu_Y)^2} \quad (1)$$

where  $Y$  are ground truth labels and  $\hat{Y}$  are the predicted values, and  $\mu$  and  $\sigma$  are their mean and variance.

**Mean Squared Error (MSE)** (for continuous labels) shows how close the predicted values are to the target values. We were inspired to use it by the EmoFAN paper [20], which suggests that optimising with respect to the MSE tends to improve the performance with respect to the CCC as well.

$$MSE(\hat{Y}, Y) = \mathbb{E}((Y - \hat{Y})^2) \quad (2)$$

**Cross-Entropy Loss (CE)** is a classification loss for the discretized labels, penalising the divergence of the predicted probability from the actual label. Discretization was conducted as following: valence and arousal per frame labels have been split into 20 bins/classes each with uniformly distributed bin boundaries (as in [2]). These classes then have been presented as one-hot vector labels of length 20 (per frame). So the cross-entropy presented as

$$CE(Y, \hat{Y}) = - \sum_{i=1}^L Y_i \log(\hat{Y}_i) \quad (3)$$

for  $L$  classes, here 20 by number of the discrete bins, and  $\hat{Y}_i$  is computed as a Softmax probability of each class per label per frame.

**Cost-sensitive cross-entropy loss (nCCE)** function [41] with the cost norm loss for discretized labels, similar to [19]:

$$nCCE(Y, \hat{Y}) = \frac{1}{F} \sum_{f=1}^F C_{norm}(Y_f, \hat{Y}_f) \sum_{l=1}^L Y_f^{(l)} \cdot \log \hat{Y}_f^{(l)} \quad (4)$$

for  $f = 1, \dots, F$  being the number of frames and a cost norm function  $C_{norm}$ , inspired by [19], that takes into consideration the spatial relation which helps the stability of the training / fine-tuning:

$$C_{norm}(Y_f, \hat{Y}_f) = 1 + \left\| \sum_{l=1}^L K^{(l)}(Y_f^{(l)} - \hat{Y}_f^{(l)}) \right\|_2 \quad (5)$$

where  $K^{(l)}$  is the centroid of the label  $l$  in k-means classification ( $k = 20$  in our case). We analyse the effect of these losses and some of their combinations in Sec. V-C.

The total loss function can be described as

$$L = w_{ccc} CCC + w_{mse} MSE + w_{ce} CE + w_{nCCE} nCCE \quad (6)$$

## IV. EXPERIMENTAL SETUP

### A. Datasets and Preprocessing

**Pretext Dataset** In order to maintain the comparability across different pre-training methods in this paper all of the pretext tasks are trained on Lip Reading Sentences 3 dataset (LRS3) [42], containing thousands of spoken sentences from TED and TEDx videos.

**Downstream Datasets** The downstream task results are presented for the following facial video datasets: SEWA [4], and RECOLA [6], the most popular academic datasets for AERR. SEWA database consists of the videos of volunteers watching adverts chosen to elicit apparent emotional reactions, and later discussing what they have seen. SEWA has annotations of valence and arousal per frame. SEWA dataset has been collected across the residents of 6 countries:

TABLE I

OUR PROPOSED MODEL VS STATE-OF-THE-ART. RECOLA: RESULTS ARE PRESENTED ON DEVELOPMENT SET AS THE TEST SET IS NOT PUBLIC. FOR [2] (AP+DET.+ATT.) STANDS FOR AFFECTIVE PROCESSES WITH COMBINED LATENT AND DETERMINISTIC LAYERS WITH SELF-ATTENTION.

Methods	SEWA		RECOLA	
	Arous.	Val.	Arous.	Val.
HO-CPCov [5]	0.520	0.750		
Affective Processes (AP+Det.+Att.) [2]	0.662	0.672		
Affective Processes Best [2]	0.640	0.750		
End-to-End Visual ResNet-50 [7]			0.371	0.637
TS-SATCN [3]			0.659	<b>0.690</b>
Baseline: 3Dconv+ResNet18+GRU From Scratch	0.588	0.609	0.344	0.538
Our SS-VAERR backbone	0.678	0.737	0.630	0.607
Our SS-VAERR (+ augmentations + composite loss)	<b>0.713</b>	<b>0.771</b>	<b>0.675</b>	0.626

the UK, Germany, Hungary, Serbia, Greece, and China. *RECOLA* is a database of multi-domain data recordings of native French-speaking participants completing a collaborative task in pairs during a video conference call, collected in France. Although *RECOLA* in the wild possesses a rich choice of modalities, we only use the pre-processed video data and continuous arousal and valence labels recorded per frame, averaged across the annotators.

**Pre-processing** All of the above datasets are converted into gray-scale videos and cropped around the face to  $96 \times 96$  based on the landmark detection. More specifically we use RetinaFace face detector [43] and the Face Alignment Network (FAN) [44] to detect 68 facial landmarks and crop the face based on these. Annotations include arousal and valence, labelled per frame, averaged across multiple annotators.

### B. Training details

Breakdown into training, validation, and test set is conducted in the same manner as in [2] for SEWA (train./val./test sets containing 435/53/53 instances), and as in [19] for RECOLA (train./val. containing 197/152 instances, with results reported on the validation set, as the test set for RECOLA is not publicly available). Both of the datasets are normalized with their respective means and standard deviations throughout.

Videos are kept at original lengths and sampled as fixed length segments at training time. Empirically, SEWA experiments yield better results with segments of 200 frames, whereas RECOLA experiments deliver better results when sampled as 500-frame-long video segments.

During training we have used AdamW optimizer with the weight decay of 0.0001 and initial learning rate ranging from 0.0003 to 0.00007, depending on the downstream dataset and other parameters, and batch-size ranging from 3 to 20. All models in this paper have been trained for 10 epochs. The augmentations are discussed in sections and IV.

## V. RESULTS

### A. Comparison with state-of-the-art

A slight complication natural to the field is that the results are being presented on a variety of datasets that do not coincide between the papers. Hence, we are restricting our benchmarks by the modality and type of the prediction - video-only natural AERR for valence and arousal. This

leaves us with only a few recent benchmarks. First, [2] and [5], demonstrated their results on SEWA dataset, one of our downstream task datasets. For RECOLA there is TS-SATCN [3], a two-stage spatio-temporal attention temporal convolution network, and a visual ResNet50 [7]. That is it for the spontaneous video-only AERR for arousal and valence.

Since these benchmarks were not evaluated on the same dataset, we compare our results in Table I to their reported results on the respective datasets: SEWA for [2], [5], and RECOLA for [3], [7]. Our final model is pre-trained in a self-supervised manner, using augmentations and composite losses (a detailed analysis for each of them can be found in sections V-B to V-D). We see that our model compares favourably to the reported results for most of these, confidently outrunning [2], [5]. It also outperforms [3] and the visual network from [7] for arousal, and only behind these a little (but still at a comparable level) for valence.

**Baseline** For completeness we present the results for our model stripped of the self-supervised component, trained from scratch, called 3Dconv+ResNet18+GRU in Table I.

### B. Empirical comparison of pretext tasks

First, we compare the performance of different pretext methods - LiRA, BYOL, and DINO used to pre-train 3D convolutional layer + ResNet18 on LRS3 and then fine-tune and assess the performance of the downstream architecture (3D convolutional layer + ResNet18 + GRU) in terms of the CCC across video-only facial datasets. Please see Table II for the results. Please note that the pretext techniques only have the basic CCC-based regression loss at this point.

It appears that LiRA pretext initialization fine-tuned on the downstream task seem to perform generally better than the other pretext methods. It achieves either the best or the second-best results across all the datasets for both arousal and valence. It even beats some of the benchmarks, despite not yet benefiting from the composite losses used by these methods. The DINO-ResNet18 also delivers results comparable to the benchmarks on most datasets.

There are several potential reasons why LiRA performs better than DINO and BYOL. First of all LiRA uses a temporal model, while others are trained per frame, which might affect the quality of learnt representations for video. It also uses the audio input for guided learning of the visual representations, the other two do not, which might

TABLE II  
COMPARISON OF THE PRETEXT TECHNIQUES ACROSS VARIOUS DATASETS FOR VIDEO-ONLY AERR.

		SEWA		RECOLA	
		Arous.	Val.	Arous.	Val.
<b>PRETEXT TECHNIQUES</b>	+ LIRA frozen	0.652	0.722	0.602	0.532
	+ LIRA fine-tuned	<b>0.678</b>	<b>0.737</b>	0.630	<b>0.607</b>
	+ Video-BYOL frozen	0.593	0.726	0.224	0.344
	+ Video-BYOL fine-tuned	0.604	<b>0.757</b>	0.307	0.446
	+ DINO-ResNet frozen	0.607	0.638	0.269	0.545
	+ DINO-ResNet fine-tuned	0.648	0.667	0.420	0.520

TABLE III  
COMPARISON OF THE VARIOUS LOSSES FOR THE DOWNSTREAM TASKS WITH LiRA PRE-TRAINING. ONLY NON-ZERO LOSS-WEIGHTS ARE PRESENTED. ‘AROUS.’ AND ‘VAL.’ SUPERSCRIPTS SPECIFY THE LOSS APPLIED SPECIFICALLY TO EITHER AROUSAL OR VALENCE PREDICTIONS.

		SEWA				RECOLA			
		Fine-Tuned		Frozen		Fine-Tuned		Frozen	
		Arous.	Val.	Arous.	Val.	Arous.	Val.	Arous.	Val.
<b>REGRESSION LOSSES</b>	$w_{ccc} = 1$	0.678	<b>0.737</b>	0.652	0.722	0.630	0.607	0.560	0.603
	$w_{mse} = 1$	0.664	0.726	0.648	0.710	0.399	0.596	0.394	0.596
<b>COMPOSITE LOSSES</b>	$w_{ccc} = 0.5, w_{ce} = 0.5$	0.671	<b>0.735</b>	0.650	<b>0.747</b>	0.454	0.625	0.513	0.606
	$w_{ccc} = 0.5, w_{ce} = 0.25, w_{mse} = 0.25$	<b>0.716</b>	0.731	0.699	<b>0.747</b>	0.473	0.611	0.469	0.610
	$w_{ccc}^{Val.} = 1, w_{ccc}^{Arous.} = 0.66, w_{ce}^{Arous.} = 0.34$	0.631	0.663	0.659	0.709	<b>0.675</b>	0.626	0.640	<b>0.668</b>
	$w_{ccc}^{Val.} = 1, w_{ccc}^{Arous.} = 0.66, w_{ce}^{Arous.}, w_{mse}^{Arous.} = 0.17$	0.638	0.716	0.658	0.691	0.664	0.644	<b>0.655</b>	0.605
	$w_{ccc} = 0.5, w_{nccc} = 0.5$	0.633	0.667	<b>0.701</b>	0.741	0.669	0.655	0.614	0.661
	$w_{ccc} = 0.5, w_{nccc} = 0.25, w_{mse} = 0.25$	0.669	0.716	0.633	0.733	0.606	<b>0.669</b>	0.626	0.623

help learning more emotion-relevant representations. Finally, DINO and BYOL usually rely on larger networks trained the datasets of scale that simply not available in AERR field.

The results of this section certainly support the hypothesis that the self-supervised pre-training can be highly beneficial in the video-only spontaneous AERR scenarios.

### C. Empirical Comparison of Training Losses

In this section we explore the impact of the various auxiliary loss functions on the performance of the pretext method LiRA during the downstream fine-tuning. We focus on LiRA pretext from now on since it has been identified as the best pretext for our purposes in the previous section.

Please refer to the Table III for the results on the set of experiments concerning the downstream loss functions. The first line corresponds to the loss used during the pretext task analysis in Table II. The benchmarks are the same as in the Table II as well, except now the comparison to them is more fair as the presented results include the loss design. The benchmarks are relevant for both fine-tuned and frozen weights versions, however we enter them in the fine-tuned section since they themselves lack this distinction in their design, so closer to fine-tuned in fashion.

Evidently, adding the CE classification loss (Eq. 3) on SEWA improves the performance for apparent arousal, whereas valence seem to either benefit minimally or show worse results. Further adding MSE (Eq. 2) achieves state-of-the-art results on SEWA. For RECOLA neither CE classification loss nor MSE loss on their own seem to have a positive impact, however together they create a minor performance boost for valence. Further adjusting to only penalising the apparent arousal loss with CE and MSE leads to a considerable boost in arousal as well. These final results outperform

the current video-only SOTA for estimating arousal, while lacking a little for estimating valence [3].

The reason why estimating the apparent arousal might require some more careful loss design in this case is not entirely clear. However, it can be in part explained by the fact that apparent arousal tends to be more present and easier detected in audio, rather than in video, while valence tends to exhibit the opposite trend [16]–[18]. Given that we are restricted to the video-only modality it is only natural that achieving good results on arousal requires more parametrization than apparent valence.

Using nCCE (Eq. 4) instead of CE often yields close second or occasional best results, it can be viewed as a solid alternative to CE.

The hypothesis postulated earlier, as well as in [19] and [2], holds. Confirming that the better performing models tend to use combinations of various regression and classification losses, i.e. Equation 6 with various weight parameters, resulting in improvements over the classic CCC-based loss function (Eq. 1).

### D. The Impact of the Data Augmentation

Augmentations for the pretext tasks are preserved as in their corresponding publications [13], [14], [36]. The experimental results presented up until this point are performed without any augmentations during the downstream training. In this section we present the ablation for different types of augmentations applied during the downstream training (e.g. Fig. 2). We conduct these experiments on the best performing models for their respective datasets:

- SEWA: LiRA pretext with composite loss of  $w_{ccc} = 0.5, w_{ce} = 0.25, w_{mse} = 0.25$  fine-tuned;

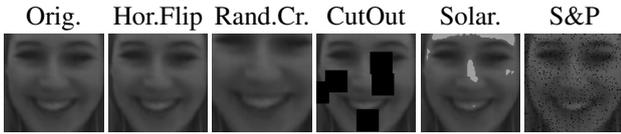


Fig. 2. Examples of augmentations used in Sec. V-D.

- RECOLA: LiRA pretext with composite loss of  $w_{cc}^{Val.} = 1, w_{cc}^{Arous.} = 0.66, w_{ce}^{Arous.} = 0.34$  fine-tuned.

TABLE IV

DOWNSTREAM DATA AUGMENTATION ABLATION RESULTS

LiRA with Composite Loss	SEWA		RECOLA	
	Arous.	Val.	Arous.	Val.
No Augmentations	0.746	0.747	0.675	<b>0.626</b>
Horizontal Flip ( $Prob = 0.5$ )	0.713	<b>0.771</b>	0.636	0.561
Random Crop	0.629	0.738	0.611	0.490
Crop-Out 25%	0.688	0.734	0.612	0.563
Missing Frames 20%	0.681	0.750	0.545	0.587
Solarization 20%	0.680	0.741	0.149	0.028
Salt & Pepper Noise	0.678	0.727	0.346	0.553
Rand.Crop + Horiz. Flip	0.637	0.653	0.601	0.395
Rand.Crop + Missing Frames	0.617	0.702	0.506	0.373
Horiz. Flip + Missing Frames	0.655	0.761	0.649	0.004
Horiz. Flip + Solarization	0.679	0.721	0.024	0.037
Horiz. Flip + Crop-Out	0.694	0.685	<b>0.701</b>	0.451
Horiz. Flip + S&P	<b>0.716</b>	0.761	0.688	0.445

For the detailed results please refer to Table IV. This list of the augmentations is a compilation of the LiRA and BYOL-recommended sets of augmentations used at the pretext training time, with some minor exceptions such as the colour-related augmentations, since we use the grayscale versions of the downstream datasets. We also explore promising combinations of the individual augmentations.

*Horizontal Flip* is a mirror reflection of an image, applied with 50% probability to training images. *Random Crop* is cropping an image to the smaller size at random (size being  $110 \times 110$  for SEWA and  $80 \times 80$  for RECOLA) for all training images. *Crop-Out* occludes several patches in the image (in our case 5 patches of square shape) with black boxes. *Missing Frames* means 20% of frames at random replaced by black frames at training time. *Solarization* is a phenomenon in photo-imaging where the image is wholly or partially reversed in tone. In this case solarization refers to an unnatural lighting effect, like in figure 2. Is not as popular in data augmentation techniques, however [13] found it beneficial for their model. We solarize above the average lightness, the effect is applied to 20% of the training images. *Salt & Pepper Noise* - classic noise with 50% salt vs pepper split, applied to all training images.

The best results on SEWA datasets are provided by the horizontal flip and its combination with the salt & pepper noise. The rest of the augmentations do not seem to bring any significant improvement and, in fact, often worsen the performance. For RECOLA augmentations seem to almost always have a negative effect on the performance. Nevertheless, there is a specific instance for the horizontal flip and

crop out combination where the performance for arousal gets close to even some of the multi-modal results [7].

## VI. DISCUSSION AND CONCLUSION

To conclude, in this paper we have presented the first to our knowledge a self-supervised technique for the video-only natural apparent emotional reactions recognition, yielding the current state-of-the-art (or closely comparable) results for video-only natural AERR. Complete with comparative empirical study of the potential pretext methods, auxiliary loss functions, and downstream-time data augmentation ablation. We also found that the optimal parameter search is somewhat unsurprisingly data-dependent, whereas the self-supervised setting is on average beneficial.

Additionally, we argue that the facial apparent emotional reactions recognition is highly data-specific. Factors that should be considered can include: the source and distributions of the pretext and downstream data (acted vs spontaneous, lab-recorded vs in-the-wild, outdoors vs indoors, speaking vs passive listening faces), as well as the specific data preprocessing procedures, format and configuration of the annotations (per frame vs per video, categorical vs continuous emotion annotation), etc.

Previous research suggests that different modalities tend to provide better cues for different apparent emotion metrics: video tends to be a better indicator for the video-aided recognition, and arousal tends to be better detected from audio modality [16]–[18]. This makes the video-only AERR particularly challenging in terms of identifying the correct levels of arousal, and explains valence-arousal discrepancy for several results in this paper.

We present the results confirming that using the self-supervised setting alone helps beating (or at least reaching comparable results with) the current state-of-the-art without even touching upon the loss function design. Next we present the evidence that the careful composite loss design can further improve the performance. And finally we provide an ablation on potentially beneficial data augmentation techniques which can lead to further improvements.

Future work could be extended to a comparative analysis of the impact of the different pretext datasets, along with the pretext training parameters and data augmentations. Additional study could be conducted on the specifics of the downstream architectures, as well as investigating the effect of sharing the learned feature representations including temporal component (for architectures with such a component) in order to fully exploit the potential of the video domain.

## VII. ACKNOWLEDGEMENTS

We would like to thank Pingchuan Ma for for providing the pre-processing and LiRa pre-training code.

## REFERENCES

- [1] Y. Yang and Y. Sun, “Facial expression recognition based on arousal-valence emotion model and deep learning method,” in *2017 International Conference on Computer Technology, Electronics and Communication (ICCTEC)*, 2017, pp. 59–62.

- [2] E. Sanchez, M. K. Tellamekala *et al.*, "Affective processes: Stochastic modelling of temporal context for emotion and facial expression recognition," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021*. Computer Vision Foundation / IEEE, 2021, pp. 9074–9084.
- [3] M. Hu, Q. Chu *et al.*, "A two-stage spatiotemporal attention convolution network for continuous dimensional emotion recognition from facial video," *IEEE Signal Process. Lett.*, vol. 28, pp. 698–702, 2021.
- [4] J. Kossaifi, R. Walecki *et al.*, "SEWA DB: A rich database for audiovisual emotion and sentiment research in the wild," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 43, pp. 1022–1040, 2021.
- [5] J. Kossaifi, A. Toisoul *et al.*, "Factorized higher-order cnns with an application to spatio-temporal emotion estimation," *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [6] F. Ringeval, A. Sonderegger *et al.*, "Introducing the RECOLA multimodal corpus of remote collaborative and affective interactions," in *IEEE International Conference and Workshops on Automatic Face and Gesture Recognition, FG 2013*. IEEE Computer Society, 2013.
- [7] P. Tzirakis, G. Trigeorgis *et al.*, "End-to-end multimodal emotion recognition using deep neural networks," *IEEE J. Sel. Top. Signal Process.*, vol. 11, no. 8, pp. 1301–1309, 2017.
- [8] A. Bulat, S. Cheng *et al.*, "Pre-training strategies and datasets for facial representation learning," *CoRR*, vol. abs/2103.16554, 2021. [Online]. Available: <https://arxiv.org/abs/2103.16554>
- [9] S. Roy and A. Etemad, "Spatiotemporal contrastive learning of facial expressions in videos," in *International Conference on Affective Computing and Intelligent Interaction, ACII 2021*. IEEE, 2021.
- [10] A. Mollahosseini, B. Hassani, and M. H. Mahoor, "Affectnet: A database for facial expression, valence, and arousal computing in the wild," *arXiv*, 2017. [Online]. Available: <http://arxiv.org/abs/1708.03985>
- [11] G. Zhao, X. Huang *et al.*, "Facial expression recognition from near-infrared videos," *Image Vis. Comput.*, vol. 29, pp. 607–619, 2011.
- [12] D. Kollias, P. Tzirakis *et al.*, "Deep affect prediction in-the-wild: Aff-wild database and challenge, deep architectures, and beyond," *International Journal of Computer Vision*, vol. 127, 06 2019.
- [13] J. Grill, F. Strub *et al.*, "Bootstrap your own latent - A new approach to self-supervised learning," in *Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, H. Larochelle, M. Ranzato *et al.*, Eds., 2020.
- [14] M. Caron, H. Touvron *et al.*, "Emerging properties in self-supervised vision transformers," *CoRR*, vol. abs/2104.14294, 2021. [Online]. Available: <https://arxiv.org/abs/2104.14294>
- [15] D. Kollias and S. Zafeiriou, "Expression, affect, action unit recognition: Aff-wild2, multi-task learning and arcfac," in *30th British Machine Vision Conference 2019, BMVC 2019, Cardiff, UK, September 9-12, 2019*. BMVA Press, 2019, p. 297.
- [16] R. A. Calvo and S. K. D'Mello, "Affect detection: An interdisciplinary review of models, methods, and their applications," *IEEE Trans. Affect. Comput.*, vol. 1, no. 1, pp. 18–37, 2010.
- [17] F. Ringeval, B. W. Schuller *et al.*, "AV+EC 2015: The first affect recognition challenge bridging across audio, video, and physiological data," in *International Workshop on Audio/Visual Emotion Challenge, AVEC 2015*, F. Ringeval, B. W. Schuller *et al.*, Eds. ACM, 2015.
- [18] —, "AVEC 2018 workshop and challenge: Bipolar disorder and cross-cultural affect recognition," in *Audio/Visual Emotion Challenge and Workshop, AVEC@MM 2018*, F. Ringeval, B. W. Schuller *et al.*, Eds. ACM, 2018, pp. 3–13.
- [19] E. Albadawy and Y. Kim, "Joint discrete and continuous emotion prediction using ensemble and end-to-end approaches," in *International Conference on Multimodal Interaction, ICMI 2018*, S. K. D'Mello, P. G. Georgiou *et al.*, Eds. ACM, 2018, pp. 366–375.
- [20] A. Toisoul, J. Kossaifi *et al.*, "Estimation of continuous valence and arousal levels from faces in naturalistic conditions," *Nat. Mach. Intell.*, vol. 3, no. 1, pp. 42–50, 2021.
- [21] M. Oquab, L. Bottou *et al.*, "Learning and transferring mid-level image representations using convolutional neural networks," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2014*. IEEE Computer Society, 2014, pp. 1717–1724.
- [22] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *International Conference on Learning Representations, ICLR 2015, Conference Track Proceedings*, Y. Bengio and Y. LeCun, Eds., 2015.
- [23] R. B. Girshick, J. Donahue *et al.*, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2014*. IEEE Computer Society, 2014, pp. 580–587.
- [24] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015*, 2015, pp. 3431–3440.
- [25] S. J. Pan and Q. Yang, "A survey on transfer learning," *IEEE Trans. Knowl. Data Eng.*, vol. 22, no. 10, pp. 1345–1359, 2010.
- [26] J. E. van Engelen and H. H. Hoos, "A survey on semi-supervised learning," *Machine Learning*, vol. 109, no. 2, pp. 373–440, 2020.
- [27] L. Ericsson, H. Gouk *et al.*, "Self-supervised representation learning: Introduction, advances and challenges," *CoRR*, vol. abs/2110.09327, 2021. [Online]. Available: <https://arxiv.org/abs/2110.09327>
- [28] S. Gidaris, P. Singh, and N. Komodakis, "Unsupervised representation learning by predicting image rotations," in *International Conference on Learning Representations, ICLR 2018*. OpenReview.net, 2018.
- [29] D. Pathak, P. Krähenbühl *et al.*, "Context encoders: Feature learning by inpainting," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016*, 2016, pp. 2536–2544.
- [30] R. Zhang, P. Isola, and A. A. Efros, "Colorful image colorization," in *Computer Vision - ECCV 2016 - 14th European Conference, Proceedings, Part III*, ser. Lecture Notes in Computer Science, B. Leibe, J. Matas *et al.*, Eds., vol. 9907. Springer, 2016, pp. 649–666.
- [31] L. Fang, J. Wang *et al.*, "Hand-drawn grayscale image colorful colorization based on natural image," *Vis. Comput.*, vol. 35, no. 11, pp. 1667–1681, 2019.
- [32] M. Noroozi and P. Favaro, "Unsupervised learning of visual representations by solving jigsaw puzzles," in *Computer Vision - ECCV 2016 - 14th European Conference, Proceedings, Part VI*, ser. Lecture Notes in Computer Science, B. Leibe, J. Matas *et al.*, Eds., vol. 9910. Springer, 2016, pp. 69–84.
- [33] T. Chen, S. Kornblith *et al.*, "A simple framework for contrastive learning of visual representations," in *Proceedings of the 37th International Conference on Machine Learning, ICML 2020*, ser. Proceedings of Machine Learning Research, vol. 119. PMLR, 2020, pp. 1597–1607.
- [34] K. He, H. Fan *et al.*, "Momentum contrast for unsupervised visual representation learning," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020*. Computer Vision Foundation / IEEE, 2020, pp. 9726–9735.
- [35] Y. Li, J. Zeng *et al.*, "Self-supervised representation learning from videos for facial action unit detection," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019*. Computer Vision Foundation / IEEE, 2019, pp. 10924–10933.
- [36] P. Ma, R. Mira *et al.*, "Lira: Learning visual speech representations from audio through self-supervision," *CoRR*, vol. abs/2106.09171, 2021. [Online]. Available: <https://arxiv.org/abs/2106.09171>
- [37] W. Hsu, B. Bolte *et al.*, "Hubert: Self-supervised speech representation learning by masked prediction of hidden units," *IEEE ACM Trans. Audio Speech Lang. Process.*, vol. 29, pp. 3451–3460, 2021.
- [38] S. Roy and A. Etemad, "Self-supervised contrastive learning of multi-view facial expressions," in *ICMI '21: International Conference on Multimodal Interaction*, Z. Hammal, C. Busso *et al.*, Eds. ACM, 2021, pp. 253–257.
- [39] K. He, X. Zhang *et al.*, "Deep residual learning for image recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR*. IEEE Computer Society, 2016, pp. 770–778.
- [40] A. Dosovitskiy, L. Beyer *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," in *International Conference on Learning Representations, ICLR*, 2021.
- [41] D. Le, Z. Aldeneh, and E. M. Provost, "Discretized continuous speech emotion recognition with multi-task deep recurrent neural network," in *Interspeech 2017*, F. Lacerda, Ed. ISCA, 2017, pp. 1108–1112.
- [42] T. Afouras, J. S. Chung, and A. Zisserman, "LRS3-TED: a large-scale dataset for visual speech recognition," *CoRR*, vol. abs/1809.00496, 2018. [Online]. Available: <http://arxiv.org/abs/1809.00496>
- [43] J. Deng, J. Guo *et al.*, "Retinaface: Single-shot multi-level face localisation in the wild," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [44] A. Bulat and G. Tzimiropoulos, "How far are we from solving the 2d & 3d face alignment problem? (and a dataset of 230,000 3d facial landmarks)," in *2017 IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 1021–1030.