

Bounded Residual Gradient Networks (BReG-Net) for Facial Affect Computing

Behzad Hasani, Pooran Singh Negi, and Mohammad H. Mahoor
Department of Electrical & Computer Engineering, University of Denver, USA

Abstract—Residual-based neural networks have shown remarkable results in various visual recognition tasks including Facial Expression Recognition (FER). Despite the tremendous efforts have been made to improve the performance of FER systems using DNNs, existing methods are not generalizable enough for practical applications. This paper introduces Bounded Residual Gradient Networks (BReG-Net) for facial expression recognition, in which the shortcut connection between the input and the output of the ResNet module is replaced with a differentiable function with a bounded gradient. This configuration prevents the network from facing the vanishing or exploding gradient problem. We show that utilizing such non-linear units will result in shallower networks with better performance. Further, by using a weighted loss function which gives a higher priority to less represented categories, we can achieve an overall better recognition rate. The results of our experiments show that BReG-Nets outperform state-of-the-art methods on three publicly available facial databases in the wild, on both the categorical and dimensional models of affect.

I. INTRODUCTION

Facial expressions are one of the most important nonverbal channels for expressing internal emotions during face-to-face communication. Six expressions of anger, disgust, fear, happiness, sadness, and surprise are defined as the basic emotional expressions by Ekman *et al.* [5]. Automated Facial Expression Recognition (FER) has been a topic of study for decades. Although there have been many achievements in developing automated FER systems, the majority of existing methods lack the required generalization due to a use of controlled data in developing methods [25]. This is predominant because there are significant variations in facial images owing to variable scene lighting, background variation, camera view, and subjects' head pose, gender, and ethnicity [21]. A comprehensive way of studying facial expressions is to approach the task through the concept of *affective computing*. Affect is a psychological term for describing the external exhibition of internal emotions and feelings. Affective computing attempts to develop systems that can interpret and estimate human affects through different channels (e.g. visual, auditory, biological signals, etc.) [29].

The dimensional modeling of affect can distinguish between subtle differences in exhibiting affect and encode small changes in the intensity of each emotion on a continuous scale, such as *valence* and *arousal* where valence shows how positive or negative an emotion is, and arousal indicates how much an event is intriguing/agitating or calming/soothing [24]. This paper focuses on developing automated algorithms for computation of the categorical and

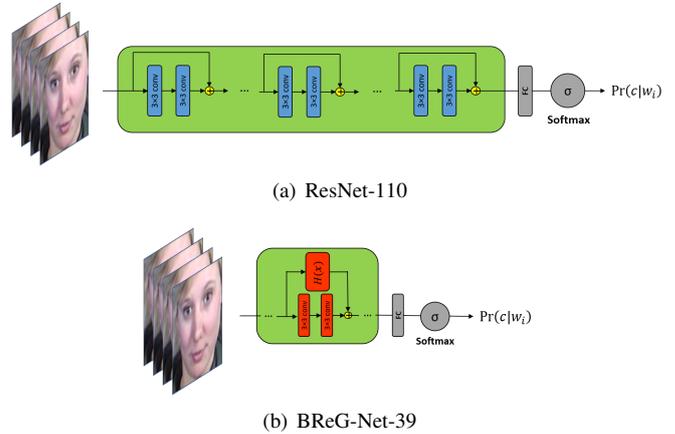


Fig. 1. Comparison of a) ResNet-110 and b) BReG-Net-39. ResNet-110 has more layers, while is slower and less accurate. BReG-Net-39 is shallower, faster, and more accurate.

dimensional models of affect.

In the field of machine learning, one of the main tasks is to optimize a function or distribution estimation with respect to a defined measure. Based on the connectionist principle [23], deep neural networks allow us to build very complex classes of functions. A wide variety of network topologies and activation functions have been proposed in the recent years and they seem to play a crucial role in design and improving the underlined class of reproducible functions available to DNNs. To pave the way of training very deep DNNs, current methods focus on improving neuron saturation or the efficiency of the gradient flow across various networks layers. Such approaches are evident in the ReLU class of non-linear functions, and the use of identity mappings in Deep Residual Networks [11]. While having deeper architectures has shown to improve the result of recognition, one possibility is to design more complex neurons to extract more useful information at each layer of the network which results in shallower networks and less parameters but more comprehensive information and a higher recognition rate.

This paper proposes and evaluates BReG-Net (Figure 1), in which the aforementioned identity mapping is replaced with a differentiable function with a bounded gradient that results in a shallower network with a considerably better recognition rate. We evaluate our proposed method using three in the wild facial expression databases (AffectNet [20], Affect-in-the-wild [32], and FER2013 [1]) in computation of both the categorical and dimensional models of affect.

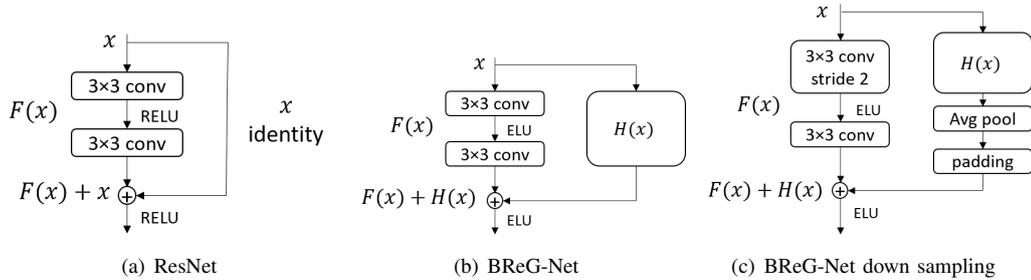


Fig. 2. Block diagram of a) ResNet b) BReG-Net and c) BReG-Net with Down-Sampling building blocks

II. RELATED WORK

A. Facial expression recognition

In recent years, “Convolutional Neural Networks” (CNNs) have become the most popular approach in the field of computer vision and pattern recognition. AlexNet and GoogLeNet are among the first successful architectures proposed based on CNNs. AlexNet consists of several convolution layers followed by max-pooling layers and Rectified Linear Units (ReLU)s. Szegedy *et al.* [27] introduced GoogLeNet which is composed of multiple “Inception” layers. Inception applies several convolutions on the feature map in different scales which extends the model both in depth and width. Mollahosseini *et al.* [19], [21] have used the Inception layer for the task of facial expression recognition and achieved state-of-the-art results. Following the success of Inception layers, several variations of them have been proposed [13]. Moreover, Inception layer is combined with a residual unit introduced by He *et al.* [10] and shows that the resulting architecture accelerates the training of Inception networks significantly [26]. Hasani *et al.* proposed a modification of ResNets for the task of facial expression recognition [7] and valence/arousal prediction of emotions [6]. While these methods use very deep architectures, the question of whether having a more complex building block results in a shallower and more efficient network remains unanswered. In the following, we will review some of the works that have looked into this concept.

B. Dimensional model of affect

A few studies have been conducted on the dimensional model of affect in the literature. Nicolaou *et al.* [22] trained bidirectional Long Short Term Memory (LSTM) architecture on multiple engineered features extracted from audio, facial geometry, and shoulders. They achieved Root Mean Square Error (RMSE) of 0.15 and Correlation Coefficient (CC) of 0.79 for valence as well as RMSE of 0.21 and CC of 0.64 for arousal. He *et al.* [12] won the AVEC 2015 challenge by training multiple stacks of bidirectional LSTMs (DBLSTM-RNN) on engineered features extracted from audio (LLDs features), video (LPQ-TOP features), 52 ECG features, and 22 EDA features. They achieved RMSE of 0.104 and CC of 0.616 for valence as well as RMSE of 0.121 and CC of 0.753 for arousal. Koelstra *et al.* [15] trained Gaussian naive Bayes classifiers on EEG, physiological signals, and

multimedia features by binary classification of low/high categories for arousal, valence, and liking on their proposed database DEAP. They achieved F1-score of 0.39, 0.37, and 0.40 on arousal, valence, and liking categories respectively.

III. PROPOSED METHOD

In this paper, we propose a residual-based network in which the shortcut connection between the input and the output of the module is replaced with a differentiable function with bounded gradient. In the following, we explain each of the aforementioned concepts in detail.

A. BReG-Net

The shortcut path in the ResNet module, which connects the input and output of the residual unit proposed, results in accelerating the convergence of the loss and simultaneously prevents the problem of vanishing/exploding gradient. The residual unit can be expressed as:

$$\begin{aligned} y_l &= H(x_l) + F(x_l, W_l) \\ x_{l+1} &= f(y_l) \end{aligned} \quad (1)$$

where x_l and x_{l+1} are the input and the output of the l -th unit and F is a residual function. In [9], $H(x_l) = x_l$ is a shortcut path, and f is an ReLU function. Later on in [11], different combination of components both on F and the shortcut was investigated. Hasani *et al.* [7] proposed a 3D ResNet based model for the task of facial expression recognition in which the shortcut was replaced with element-wise multiplication of the weight function ω and the input layer x_l as follows:

$$\begin{aligned} y_l &= \omega(L, P) \circ x_l + F(x_l, W_l) \\ x_{l+1} &= f(y_l) \end{aligned} \quad (2)$$

in which \circ denotes the Hadamard product symbol and the weight values gradually decrease when pixels P get farther away from the facial landmark points L . This shows that having a more complex function than a simple shortcut (identity mapping) can help the network to extract more effective features in less number of layers which results in a shallower network and less number of parameters to be trained.

In Equation (3), it can be seen that the identity bypass mapping (x) is a simple choice and is not contributing to feature learning. In fact, the original motivation for using x in the residual connection was to have bounded feedbacks from

the loss layer to every other layers of the network. Building on this observation, we studied developing more complex residual connections with bounded gradient which enrich feature learning through the residual parts of the network. This results in richer feature maps and therefore shallower networks. We investigated several functions and replaced the shortcut path in the network with those functions. There are few limitations on choosing the suitable function and not all the functions can be used, as the network will not converge otherwise. The reason behind this is that in the training phase, we need to calculate the gradient. An improper choice of the function will cause facing with either vanishing or exploding gradient. To have a better understanding of this concept we start with the ResNet’s residual unit formulation. In this case, since we have an identical mapping of the inputs for the function $H(x_l)$, Equation (1) and its derivative will be re-written as follows:

$$\begin{aligned} y_l &= x + F(x_l, W_l) \\ y'_l &= 1 + F'(x_l, W_l) \end{aligned} \quad (3)$$

It is obvious that $H(x_l) = x$ is differentiable and its derivative is constant which means that it is also bounded. This allows the ResNet to converge and prevents the vanishing/exploding gradient problem. Therefore, any other function that is the replacement of x needs to have the same properties.

We observed several functions that have the aforementioned properties. Our experiments show that by incorporating any of these functions, the network will still converge and this is not surprising, based on the aforementioned argument. Hence, it is a matter of choosing the right function to have the best results for the facial expression task and valence/arousal prediction. Among the functions we investigated, the followings showed the most promising results:

$$\begin{aligned} H_1(x) &= x - \log(e^x + 1), H'_1(x) = \frac{1}{1 + e^x} \\ H_2(x) &= x \tan^{-1}(x) - \frac{1}{2} \log(x^2 + 1), H'_2(x) = \tan^{-1}(x) \\ H_3(x) &= \tan^{-1}(x), H'_3(x) = \frac{1}{1 + x^2} \end{aligned} \quad (4)$$

and the corresponding derivative of these functions is as follows:

Figure 3 shows the plots of these three functions and their derivatives. As shown, all of these functions are differentiable at any point and their derivatives are also bounded which shows that previously mentioned conditions hold for all of these functions. We call our network **Bounded Residual Gradient Network** (BReG-Net). Figure 2(b) shows the resulting building block of BReG-Net module. In our proposed network, similar to ResNet, we have dimension reductions of the tensor, achieved by down sampling (stride 2) on the first convolution layer of $F(x)$ (Figure 2(c)). As explained in the experiments section, we stack up 39 layers of these blocks in all of our experiments and compare the results on different databases.

B. Weighted loss

Facial expression databases are usually highly skewed. This form of imbalance is commonly referred to as *intrinsic* variation, i.e., it is a direct result of the nature of expressions in the real world. This phenomenon exists in both the categorical and dimensional models of affect. For instance, in AffectNet database well represented categories like happiness have almost 30 times more number of samples than less represented categories like contempt. The problem of learning from imbalanced data has two downsides. First, training data with an imbalanced distribution often causes learning algorithms to perform poorly on the less represented category [8]. Second, imbalanced test/validation data can affect the performance metrics drastically resulting in an unrealistic image of method’s performance. Jeni *et al.* [14] studied the influence of skew on imbalanced datasets. This study shows that except for of area under the ROC curve (AUC), many other evaluation metrics such as accuracy, F1-score, Cohens kappa [3], Krippendorfs alpha [16], and area under Precision-Recall curve (AUC-PR) are affected by skewed distributions dramatically. In order to minimize skew-biased estimates of performance, the study suggests reporting both skew-normalized metrics as well as the original evaluation.

In the result section, we report the skew-normalized metrics of our methods in addition to Matthews Correlation Coefficient (MCC) [18] and Positive Predictive Value (PPV) which is often called precision. Moreover, in order to improve the recognition rate of different categories of emotions in our methods, we assign higher priority to the less represented categories of the databases in the loss calculation layer of our networks. We weigh the loss function for each of the classes by their relative proportion in the training dataset. In other words, the loss function highly penalizes the networks for misclassifying examples from under-represented categories, while it penalizes the networks less for misclassifying examples from well-represented categories. The entropy loss formulation for a training example (x_i, l) is defined as:

$$E = - \sum_{i=1}^K H_{l,i} \log(\hat{p}_i) \quad (5)$$

where $H_{l,i}$ denotes row l penalization factor of class i . K is the number of classes and \hat{p}_i is the predictive softmax with values in interval $[0, 1]$ indicating the predicted probability of each class as:

$$\hat{p}_i = \frac{\exp(z_i)}{\sum_{j=1}^K \exp(z_j)} \quad (6)$$

When $H = I$ (I is the identity matrix), the proposed weighted-loss approach will turn to the traditional cross-entropy loss function. In other words, if the training data is completely balanced, the weighted-loss method is equal to the conventional cross-entropy loss function. We implemented this loss function in our TensorFlow model and we

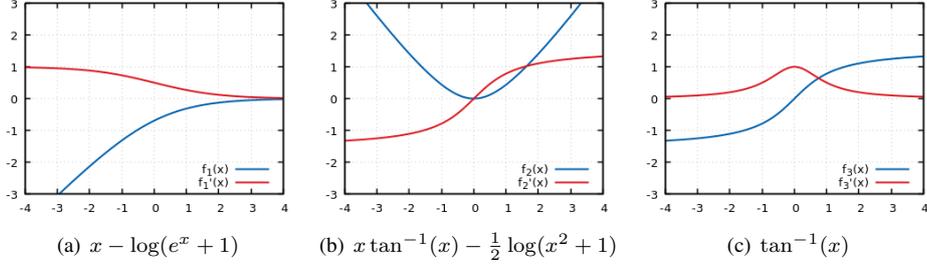


Fig. 3. Plots of proposed functions and their derivatives for the BReG-Net (best in color)

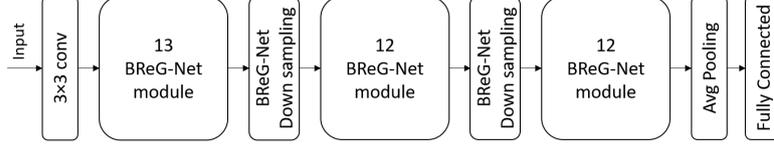


Fig. 4. General architecture of the proposed method

define the diagonal matrix H_{ij} as:

$$H_{ij} = \begin{cases} \frac{f_{min}}{f_i}, & \text{if } i = j \\ 0, & \text{otherwise} \end{cases} \quad (7)$$

where f_i is the number of samples in the i^{th} category and f_{min} is the number of samples in the least-represented category. As mentioned earlier, this will cause the loss function to highly penalize the network for misclassifying examples from under-represented categories. In the results section we show that this improves the network recognition of under represented categories and has an overall better recognition rate.

IV. EXPERIMENTS AND RESULTS

In this section, we briefly review the face databases used for evaluating our proposed method. We then report the results of our experiments using these databases evaluated on different metrics on both categorical and dimensional model of affect.

A. Face databases

As noted earlier, many of the traditional facial expression databases are assembled in a controlled environment while for developing a practical methods, these databases do not yield satisfying results. Therefore, we chose databases that are captured in the wild setting which contain a variety of backgrounds, lighting, pose, subject ethnicity, etc. These databases are AffectNet [20], Affect-in-Wild [32], and FER2013 [1] of which AffectNet contains labels of both categorical and dimensional models. Affect-in-Wild contains only labels of dimensional model, and FER2013 contains only labels of categorical model. AffectNet contains more than one million facial images collected from the Internet by querying three major search engines using 1250 emotion related keywords in six different languages. Affect-in-Wild contains 300 videos of different subjects watching videos of various TV shows and movies. FER2013 was created using

the Google image search API. Faces are labeled with any of the six basic expressions, along with neutral. The resulting database contains 35,887 images in the wild settings.

B. Evaluation metrics of dimensional model

In order to evaluate our methods, we calculate and report Root Mean Square Error (RMSE), Correlation Coefficient (CC), Concordance Correlation Coefficient (CCC), and Sign AGREement (SAGR) metrics for our methods. In the following, we briefly describe the definitions of these metrics.

Root Mean Square Error (RMSE) is the most common evaluation metric in a continuous domain which is defined as:

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{\theta}_i - \theta_i)^2} \quad (8)$$

where $\hat{\theta}_i$ and θ_i are the prediction and the ground-truth of i^{th} sample, and n is the number of samples. RMSE-based evaluation metrics can heavily weigh the outliers [2], and they do not consider covariance of the data.

Pearson's Correlation Coefficient (CC) overcomes this problem [22] and it is defined as:

$$CC = \frac{COV\{\hat{\theta}, \theta\}}{\sigma_{\hat{\theta}}\sigma_{\theta}} = \frac{E[(\hat{\theta} - \mu_{\hat{\theta}})(\theta - \mu_{\theta})]}{\sigma_{\hat{\theta}}\sigma_{\theta}} \quad (9)$$

where COV is covariance function.

Concordance Correlation Coefficient (CCC) is another metric [30] and combines CC with the square difference between the means of two compared time series:

$$\rho_c = \frac{2\rho\sigma_{\hat{\theta}}\sigma_{\theta}}{\sigma_{\hat{\theta}}^2 + \sigma_{\theta}^2 + (\mu_{\hat{\theta}} - \mu_{\theta})^2} \quad (10)$$

where ρ is the Pearson correlation coefficient (CC) between two time-series (e.g., prediction and ground-truth), $\sigma_{\hat{\theta}}^2$ and σ_{θ}^2 are the variance of each time series, $\sigma_{\hat{\theta}}$ and σ_{θ} are the standard deviation of each, and $\mu_{\hat{\theta}}$ and μ_{θ} are the mean value of each. Unlike CC, the predictions that are well correlated

with the ground-truth but shifted in value are penalized in proportion to the deviation in the CCC.

The value of valence and arousal fall within the interval of $[-1,+1]$ and correctly predicting their signs are essential in many emotion-prediction applications. Therefore, we use Sign AGREement (SAGR) metric as proposed in [22] to evaluate the performance of a valence and arousal prediction system with respect to the sign agreement. SAGR is defined as:

$$SAGR = \frac{1}{n} \sum_{i=1}^n \delta(\text{sign}(\hat{\theta}_i), \text{sign}(\theta_i)) \quad (11)$$

where δ is the Kronecker delta function, defined as:

$$\delta(a, b) = \begin{cases} 1, & a = b \\ 0, & a \neq b \end{cases} \quad (12)$$

C. Results

Figure 4 shows the general structure of the network. Our experiments show that $H_3(x)$ yields better results in terms of both prediction rate and convergence speed. We also investigated a variety of BReG-Net architectures with shallower and deeper depths. Our experiments indicated that when the network is too shallow, the number of parameters is not enough to distinguish the subtle facial muscle changes. Figure 5 shows the results of different depths in both categorical and dimensional models of affect while using $H_3(x) = \tan^{-1}(x)$ as residual function in our proposed method. Thus, we propose the architecture in Figure 4 for two tasks of prediction of categorical and dimensional model of affect. We provide the results of our experiment for each of these tasks separately. All of the proposed methods are implemented using a combination of TensorFlow [17] and TfLearn [4] toolboxes. We used Momentum optimization method with a weight decay of 0.0001, and learning rate of 0.01. Mean square error is used for the loss function of the dimensional model experiments.

1) *Categorical model*: Table I shows the results of our experiments with the three functions in Equation (4) as the residual function. We can see that $H_3(x) = \tan^{-1}(x)$ has the best result compared to the other functions. This was true throughout all of the experiments. Therefore, due to space limitation, all of the reported results from this point are the result of $H_3(x)$ function. Table II shows the result of our experiments in the categorical model of affect on AffectNet and FER2013 databases. It can be seen that weighted loss further improves the recognition rates in both databases. However, weighted-loss is data dependent while our proposed method improves the recognition rate regardless of the distribution of the data. All of the reported numbers, are the result of our experiments only on the validation set of these databases as their test sets are not publicly available for any of the databases. As it can be seen, our proposed modification of the ResNet module achieves better recognition rates compared to ResNet-110 and it also outperforms the existing methods on both AffectNet and FER2013 databases. We need to mention that [20] uses AlexNet, Wiles *et al.* [31] achieved 74.4 for AUC, and [19] uses an Inception-based method to classify

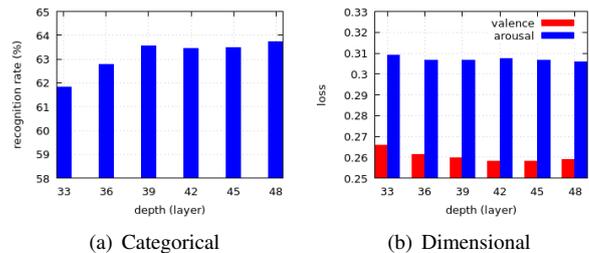


Fig. 5. Result of experimenting different depth on categorical and dimensional model (best in color)

the expressions, and [28] trained deep learning methods combined with SVMs. Our proposed method is considerably shallower than many of the methods proposed in the field.

In order to further investigate the effect of the weighted-loss method, we calculated F1-score, alpha, kappa, MCC, and PPV metrics in both cases of regular loss and weighted-loss. Tables III and IV show the results for these losses, respectively. The skew normalization is performed by random under-sampling of the classes in the test set. This process is repeated 200 times, and the skew-normalized score is the average of the score on multiple trials. It can be seen that in most cases there is an improvement of correlation in the weighted loss case which shows that our weighted loss addition to the network has a positive impact in recognition of different categories. It is important to note that the FER2013 database is an almost balanced database. Therefore, the reported results for original and skew-normalized cases have almost the same value.

2) *Dimensional model*: Table V shows the results of our experiments in the dimensional model of affect on the validation set of the AffectNet and Affect-in-Wild databases (test set was not released for either of the databases). It is important to point out that [20] uses AlexNet, and [6] uses an Inception-ResNet-based method to classify the expressions. The reported results are RMSE values, as other methods have only provided this metric in their work. Table V shows that our proposed method outperforms the state-of-the-art methods in terms of RMSE for both databases. Our results show significant improvement compared to methods reported in the AffectNet paper [20]. Also, as shown in the categorical model experiments, we can see significant improvement using the BReG-Net comparing to ResNet-110. Figure 6 shows that our proposed method has a higher reduction rate compared to ResNet-110 and eventually reaches a lower loss value on both training and validation sets during training.

In order to further investigate the effect of BReG-Net in the dimensional model of affect, we report the results by using the metrics of CC, CCC, and SAGR. Tables VI and VII show the values of these metrics on BReG-Net and ResNet-110, respectively. It can be seen that the sign agreement is significantly improved when using BReG-Net, and also correlation of the predicted values is higher than the ones for ResNet. Also, we can see that predicted valence values have lower RMSE while have higher correlation with ground-truth

TABLE I
RECOGNITION RATE (%) OF PROPOSED FUNCTIONS IN EQUATION (4) IN CATEGORICAL MODEL

	$H_1(x)$		$H_2(x)$		$H_3(x)$	
	regular loss	weighted loss	regular loss	weighted loss	regular loss	weighted loss
AffectNet	57.37	58.83	59.43	64.02	60.03	63.54
FER2013	65.80	66.21	65.16	67.66	68.74	69.49

TABLE II
RECOGNITION RATES (%) IN CATEGORICAL MODEL OF AFFECT

	ResNet-110	proposed method		state-of-the-art methods
		regular loss	weighted loss	
AffectNet	58.20	60.03	63.54	58.0 [20], 57.31 [33]
FER2013	66.48	68.74	69.49	69.3 [28], 66.4 [19]

TABLE III
RESULTS OF WEIGHTED-LOSS EXPERIMENTS ON CATEGORICAL MODEL OF AFFECT

	F1-score		kappa		alpha		MCC		PPV	
	Orig*	Norm*	Orig	Norm	Orig	Norm	Orig	Norm	Orig	Norm
AffectNet	0.63	0.68	0.58	0.63	0.58	0.64	0.59	0.64	0.62	0.71
FER2013	0.67	0.67	0.62	0.62	0.62	0.62	0.62	0.62	0.69	0.68

*Orig and Norm stand for **O**riginal and skew-**N**ormalized, respectively.

TABLE IV
RESULTS OF REGULAR-LOSS EXPERIMENTS ON CATEGORICAL MODEL OF AFFECT

	F1-score		kappa		alpha		MCC		PPV	
	Orig*	Norm*	Orig	Norm	Orig	Norm	Orig	Norm	Orig	Norm
AffectNet	0.58	0.60	0.52	0.54	0.52	0.54	0.52	0.54	0.58	0.62
FER2013	0.67	0.68	0.61	0.62	0.61	0.62	0.62	0.62	0.69	0.69

*Orig and Norm stand for **O**riginal and skew-**N**ormalized, respectively.

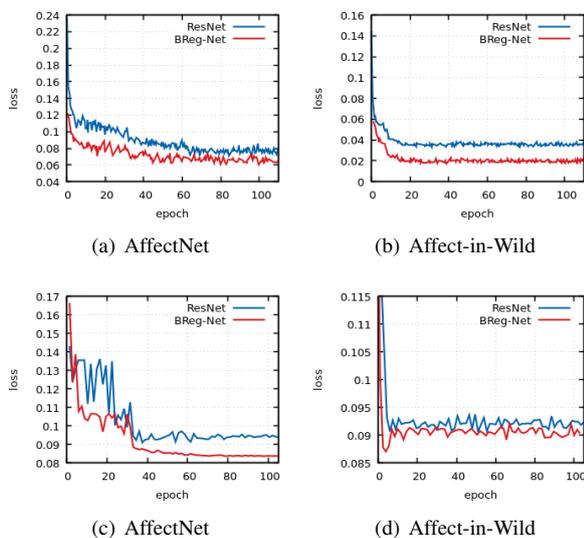


Fig. 6. Mean square loss of **training** (a and b) and **validation** (c and d) for ResNet-110 and BReG-Net (best in color)

compared to their corresponding arousal values. This is not surprising as RMSE and correlation coefficient measure two different aspects of distribution of the data. These tables also show that the Affect-in-Wild database is a more challenging database as the predicted values have less correlation with the ground-truth ones.

In order to compare the computational cost of BReG-Net and ResNet, we recorded the computation time of training the model for one epoch on AffectNet database in categorical model. The average processing time of an epoch on AffectNet for BReG-Net with 4.9M parameters is 750.21 seconds and for ResNet-110 with 7.2M parameters is 836.04 seconds on a GeForce GTX 1080 Ti GPU. Therefore, our proposed method is trained considerably faster than ResNet-110 as it has less number of parameters to train.

V. CONCLUSION

This paper introduces BReG-Net, a new residual-based network architecture using a differentiable and bounded gradient function instead of a shortcut path between the input and the output of the residual block for the task of affect estimation in both categorical and dimensional models of affect. Our experiments showed that recruiting more complex units will result in shallower networks with better performance. We also used weighted loss function in the categorical model, where our method gives higher priority to the under represented categories, resulting in a better recognition rate. We evaluated our proposed method on three databases of facial images captured in wild settings. Our experiments showed that the proposed method outperforms state-of-the-art methods in both tasks.

TABLE V
RMSE VALUES OF EXPERIMENTS ON DIMENSIONAL MODEL OF AFFECT

	ResNet-110		proposed method		state-of-the-art methods	
	valence	arousal	valence	arousal	valence	arousal
AffectNet	0.2693	0.3082	0.2597	0.3067	0.37 [20]	0.41 [20]
Affect-in-Wild	0.2733	0.3309	0.2661	0.3265	0.27 [6]	0.36 [6]

TABLE VI
RESULTS OF BREG-NET ON DIMENSIONAL MODEL

	CC		CCC		SAGR	
	valence	arousal	valence	arousal	valence	arousal
AffectNet	0.66	0.84	0.66	0.82	0.73	0.84
Affect-in-Wild	0.45	0.40	0.43	0.34	0.63	0.77

TABLE VII
RESULTS OF RESNET-110 ON DIMENSIONAL MODEL

	CC		CCC		SAGR	
	valence	arousal	valence	arousal	valence	arousal
AffectNet	0.66	0.84	0.63	0.82	0.66	0.84
Affect-in-Wild	0.41	0.41	0.38	0.35	0.61	0.75

VI. ACKNOWLEDGEMENT

This paper is based upon work partially supported by the National Science Foundation under Grant No. CNS-1427872. We also thank NVIDIA for donation of a GPU to the University of Denver.

REFERENCES

- [1] Challenges in representation learning: Facial expression recognition challenge. <http://www.kaggle.com/c/challenges-in-representation-learning-facial-expression-recognition-challenge>.
- [2] S. Bermejo and J. Cabestany. Oriented principal component analysis for large margin classifiers. *Neural Networks*, 14(10):1447–1461, 2001.
- [3] J. Cohen. A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20(1):37, 1960.
- [4] A. Damien et al. Tlearn. <https://github.com/tlearn/tlearn>, 2016.
- [5] P. Ekman and W. V. Friesen. Constants across cultures in the face and emotion. *Journal of personality and social psychology*, 17(2):124, 1971.
- [6] B. Hasani and M. H. Mahoor. Facial affect estimation in the wild using deep residual and convolutional networks. In *CVPR Workshops*, pages 1955–1962. IEEE, 2017.
- [7] B. Hasani and M. H. Mahoor. Facial expression recognition using enhanced deep 3d convolutional neural networks. In *CVPR Workshops*, pages 2278–2288. IEEE, 2017.
- [8] H. He and E. A. Garcia. Learning from imbalanced data. *IEEE Transactions on Knowledge & Data Engineering*, 21(9):1263–1284, 2009.
- [9] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *arXiv preprint arXiv:1512.03385*, 2015.
- [10] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, June 2016.
- [11] K. He, X. Zhang, S. Ren, and J. Sun. Identity mappings in deep residual networks. *arXiv preprint arXiv:1603.05027*, 2016.
- [12] L. He, D. Jiang, L. Yang, E. Pei, P. Wu, and H. Sahli. Multimodal affective dimension prediction using deep bidirectional long short-term memory recurrent neural networks. In *Workshop on Audio/Visual Emotion Challenge*, pages 73–80. ACM, 2015.
- [13] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- [14] L. A. Jeni, J. F. Cohn, and F. De La Torre. Facing imbalanced data—recommendations for the use of performance metrics. In *ACII*, pages 245–251. IEEE, 2013.
- [15] S. Koelstra, C. Muhl, M. Soleymani, J.-S. Lee, A. Yazdani, T. Ebrahimi, T. Pun, A. Nijholt, and I. Patras. Deap: A database for emotion analysis; using physiological signals. *IEEE Transactions on Affective Computing*, 3(1):18–31, 2012.
- [16] K. Krippendorff. Estimating the reliability, systematic error and random error of interval data. *Educational and Psychological Measurement*, 30(1):61–70, 1970.
- [17] A. A. M. Abadi et al. Tensorflow: Large-scale machine learning on heterogeneous. *Software available from tensorflow.org*, 1, 6, 2015.
- [18] B. W. Matthews. Comparison of the predicted and observed secondary structure of t4 phage lysozyme. *Biochimica et Biophysica Acta (BBA)-Protein Structure*, 405(2):442–451, 1975.
- [19] A. Mollahosseini, D. Chan, and M. H. Mahoor. Going deeper in facial expression recognition using deep neural networks. In *WACV*, pages 1–10. IEEE, 2016.
- [20] A. Mollahosseini, B. Hasani, and M. H. Mahoor. Affectnet: A database for facial expression, valence, and arousal computing in the wild. *IEEE Transactions on Affective Computing*, 2017.
- [21] A. Mollahosseini, B. Hasani, M. J. Salvador, H. Abdollahi, D. Chan, and M. H. Mahoor. Facial expression recognition from world wild web. In *CVPR Workshops*, June 2016.
- [22] M. A. Nicolau, H. Gunes, and M. Pantic. Continuous prediction of spontaneous affect from multiple cues and modalities in valence-arousal space. *IEEE Transactions on Affective Computing*, 2(2):92–105, 2011.
- [23] D. E. Rumelhart, J. L. McClelland, P. R. Group, et al. Parallel distributed processing: Explorations in the microstructures of cognition. volume 1: Foundations, 1986.
- [24] J. A. Russell. A circumplex model of affect. *Journal of Personality and Social Psychology*, 39(6):1161–1178, 1980.
- [25] C. Shan, S. Gong, and P. W. McOwan. Facial expression recognition based on local binary patterns: A comprehensive study. *Image and Vision Computing*, 27(6):803–816, 2009.
- [26] C. Szegedy, S. Ioffe, and V. Vanhoucke. Inception-v4, inception-resnet and the impact of residual connections on learning. *arXiv preprint arXiv:1602.07261*, 2016.
- [27] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *CVPR*, pages 1–9, 2015.
- [28] Y. Tang. Deep learning using linear support vector machines. *arXiv preprint arXiv:1306.0239*, 2013.
- [29] J. Tao and T. Tan. Affective computing: A review. In *ACII*, pages 981–995. Springer, 2005.
- [30] M. Valstar, J. Gratch, B. Schuller, F. Ringeval, D. Lalanne, M. T. Torres, S. Scherer, G. Stratou, R. Cowie, and M. Pantic. AVEC 2016—depression, mood, and emotion recognition workshop and challenge. *arXiv preprint arXiv:1605.01600*, 2016.
- [31] O. Wiles, A. Koepke, and A. Zisserman. Self-supervised learning of a facial attribute embedding from video. *arXiv preprint*

arXiv:1808.06882, 2018.

- [32] S. Zafeiriou, D. Kollias, M. A. Nicolaou, A. Papaioannou, G. Zhao, and I. Kotsia. Aff-wild: Valence and arousal in-the-wild challenge. In *CVPR Workshop*, 2017.
- [33] J. Zeng, S. Shan, and X. Chen. Facial expression recognition with inconsistently annotated datasets. In *ECCV*, pages 222–237, 2018.