

# HoloNet: Towards Robust Emotion Recognition in the Wild

Anbang Yao  
Intel Labs China,  
Beijing, 100190, China  
anbang.yao@intel.com

Dongqi Cai  
Intel Labs China,  
Beijing, 100190, China  
dongqi.cai@intel.com

Ping Hu  
Intel Labs China,  
Beijing, 100190, China  
ping1.hu@intel.com

Shandong Wang  
Intel Labs China  
Beijing, 100190, China  
shandong.wang@intel.com

Liang Sha  
673 Lab, Beihang University  
Beijing, 100191, China  
liang.sha@buaa.edu.cn

Yurong Chen  
Intel Labs China  
Beijing, 100190, China  
yurong.chen@intel.com

## ABSTRACT

In this paper, we present HoloNet, a well-designed Convolutional Neural Network (CNN) architecture regarding our submissions to the video based sub-challenge of the Emotion Recognition in the Wild (EmotiW) 2016 challenge. In contrast to previous related methods that usually adopt relatively simple and shallow neural network architectures to address emotion recognition task, our HoloNet has three critical considerations in network design. (1) To reduce redundant filters and enhance the non-saturated non-linearity in the lower convolutional layers, we use a modified Concatenated Rectified Linear Unit (CReLU) instead of ReLU. (2) To enjoy the accuracy gain from considerably increased network depth and maintain efficiency, we combine residual structure and CReLU to construct the middle layers. (3) To broaden network width and introduce multi-scale feature extraction property, the top layer is designed as a variant of inception-residual structure. The main benefit of grouping these modules into the HoloNet is that both negative and positive phase information implicitly contained in the input data can flow over it in multiple paths, thus deep multi-scale features explicitly capturing emotion variation can be well extracted from multi-path sibling layers, and then can be further concatenated for robust recognition. We obtain competitive results in this year's video based emotion recognition sub-challenge using an ensemble of two HoloNet models trained with given data only. Specifically, we obtain a mean recognition rate of 57.84%, outperforming the baseline accuracy with an absolute margin of 17.37%, and yielding 4.04% absolute accuracy gain compared to the result of last year's winner team. Meanwhile, our method runs with a speed of several thousands of frames per second on a GPU, thus it is well applicable to real-time scenarios.

**CCS Concepts:** •Computing methodologies ~ Appearance and texture representations •Computing methodologies ~ Neural networks

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [Permissions@acm.org](mailto:Permissions@acm.org).

ICMI'16, November 12–16, 2016, Tokyo, Japan  
© 2016 ACM. 978-1-4503-4556-9/16/11...\$15.00  
<http://dx.doi.org/10.1145/2993148.2997639>

## Keywords

Emotion Recognition, EmotiW 2016 Challenge, Deep Learning, Convolutional Neural Networks

## 1. INTRODUCTION

<sup>1</sup>The video based Emotion Recognition in the Wild (EmotiW) challenge has been continuously held for three years [5–7]. In the task, competitors are asked to use their methods to automatically assign one label from seven candidate emotion categories (i.e., Angry, Disgust, Fear, Happy, Neutral, Sad and Surprise) to each of test video clips collected from real Hollywood movies. In contrast to the previous challenges, this year's challenge further introduces video clips extracted from real TV programs for testing only [8]. Generally, this challenge provides a comprehensive evaluation dataset, covering a broad range of difficulties such as diverse illumination, viewpoint variation in face, partial occlusion, cluttered background and context. Since emotion recognition in unconstrained conditions has great potential for many applications, such as intelligent human-machine interaction, smart robotics and online advertising, the EmotiW challenge series has attracted ever growing interest both in academia and industry.

According to the report [7] summarized by the organizers, deep learning based methods, especially popular Convolutional Neural Networks (CNNs), have been adopted by most of the teams in the past challenges. This is mainly due to the fact that CNNs have made performance breakthroughs on many computer vision tasks such as image classification [13, 22, 23] and face recognition [26]. One important fact we noticed is that the network architectures of former winning teams are both simple and shallow. The EmoNet proposed by the 2013 winner team [12] is a four-stage architecture which has three convolutional layers and one fully connected layer with 7 softmax units. The winner team of 2014 [16] presents a multiple-kernel learning method to study the potential of combining hand-crafted features and CNN features. However, the architectures of two CNNs they adopted are also shallow and simple. 2015 winner team [31] also uses a five-layer shallow CNN model to enhance the performance of their proposed relation based features. In this paper, we explore the method on how to design a deep yet computationally efficient CNN architecture for advancing emotion recognition in the wild. Our method is inspired

<sup>1</sup> This work was done by Dongqi Cai, Ping Hu, Shandong Wang and Liang Sha with equal contribution, led by Anbang Yao.

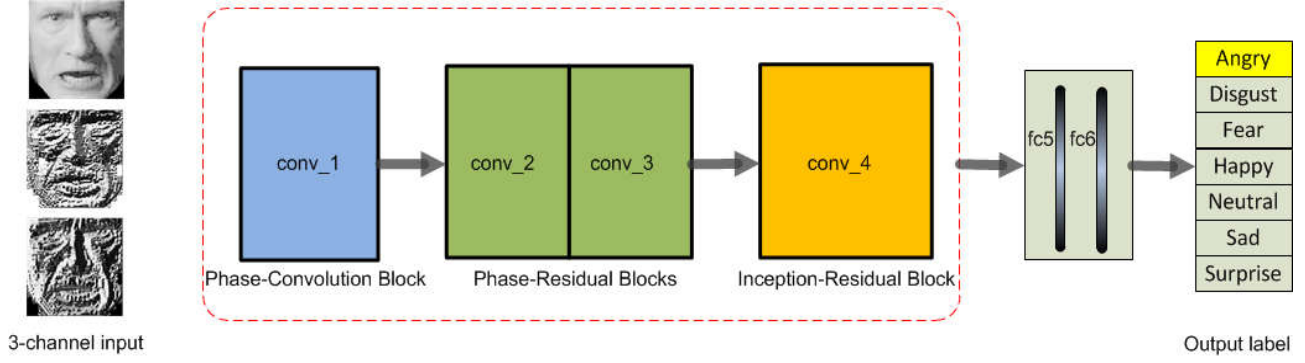


Figure 1. The flowchart of our HoloNet based emotion recognition framework.

by three important facts. First, the latest study [20] shows that the filters in the lower convolutional layers of a deep CNN form pairs in phase. This means the filters in the lower convolutional layers contain considerable redundancy, thus the number of filters can be largely reduced while the accuracy of the whole network will not decrease through a simple yet effective activation scheme. Second, recent progress in deep learning sufficiently demonstrates that residual structure [11] can resolve the dilemma between training considerably deeper network and gaining accuracy from largely increased depth. Finally, it has been well validated in the previous works [1, 15, 19] that multi-scale feature extraction scheme can bring performance improvement in many face related tasks. We conjecture that a more powerful CNN architecture could be introduced for emotion recognition in unconstrained conditions if above three facts can be jointly considered. To test our hypothesis, we present HoloNet, a deep yet computationally efficient CNN architecture, coupling all above facts naturally: (1) To reduce redundant filters and improve the non-saturated non-linearity in the lower convolutional layers, we apply a modified Concatenated Rectified Linear Unit (CReLU) [20] instead of basic ReLU [13]. (2) To reap the accuracy gain from considerably increased depth and maintain efficiency, we combine powerful residual structure and CReLU to construct the middle layers. (3) To broaden network width and introduce multi-scale feature extraction property, the topper layers are designed as a variant of inception-residual structure [24]. Besides, features from different modalities can be well combined for improved accuracy. Extensive experiments on EmotiW 2016 challenge well validate the efficacy of our method.

We will detail HoloNet architecture and show its performance on the video based emotion recognition sub-challenge of EmotiW 2016 in the following sections.

## 2. THE PROPOSED METHOD

Figure 1 shows the flowchart of our HoloNet framework. It can be seen that the core modules (parts within the rectangle with red dashed border) of HoloNet are four stacks of convolutional layers abbreviated as conv\_1, conv\_2, conv\_3 and conv\_4. As shown in Figure 1, we present three different building blocks namely Phase-Convolution Block, Phase-Residual Block and Inception-Residual Block, to construct these stacked convolutional layers. To make the whole paper self-contained, in what follows, we will start the description with data pre-processing and network input.

### 2.1 Data Pre-processing and Network Input

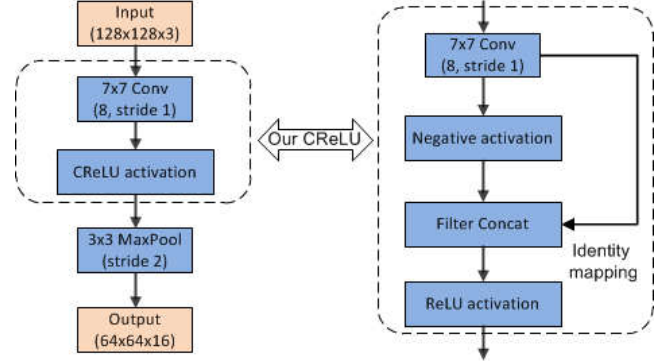


Figure 2. Structure of the Phase-Convolution Block.

The data pre-processing module includes four steps. We use a self-developed Adaboost based multi-view face detector [27] to locate target face in the first frame of each video clip. After obtaining a face region in the first frame, Supervised Descent Method (SDM) [30] is applied to track facial features over time. Based on the tracked facial features, related face image regions are cropped out and scaled to a standard size of  $156 \times 156$  pixels. Subsequently, face frontalization is performed by back-projecting each scaled face image to an estimated 3D face coordinate system [10]. The frontalized face images are further scaled to a size of  $136 \times 136$  pixels. The last pre-processing step is to compensate for lighting effect by a popular Discrete Cosine Transform (DCT) based method [2].

After data pre-processing, the gray-scale face image together with its corresponding basic Local Binary Patterns (LBP) and mean LBP feature images [17, 32] are feed-forwarded through our HoloNet, as a three-channel input image. It is worth noting that other types of feature images may also be fed to HoloNet.

### 2.2 Phase-Convolution Block

During training or feed-forwarding, the input of our HoloNet is a fixed-size image ( $128 \times 128 \times 3$  pixels). The image is first passed through conv\_1 which is a stack of composite layers formed by one Phase-Convolution Block. The structure of our Phase-Convolution Block is shown in Figure 2. Unlike basic ReLU [13], Phase-Convolution Block utilizes a different feature activation scheme. Basic ReLU popularly used in deep CNNs retains the phase information but discards the modulus information when the phase of a filter response is negative. As validated in [20], the filters in the lower convolutional layers of a deep CNN form pairs

in phase. This means the filters in the lower convolutional layers contain considerable redundancy. To reduce redundant filters and enhance the non-saturated non-linearity in the lower convolutional layers, we apply a modified Concatenated Rectified Linear Unit (CReLU) instead of basic ReLU. As shown in Figure 2, our CReLU makes an identical copy of the linear responses right after convolution, first negates responses, then concatenates copied responses and negated responses along channels, and finally applies ReLU altogether to obtain activation maps. In this way, both the positive and the negative phase information are well preserved with no additional hyper-parameters. This is the reason why such kind of building blocks is named as Phase-Convolution Block. Compared with the basic ReLU, the other merit of CReLU is that the number of convolutional filters can be reduced to half under the same number of response maps. In this paper, we apply batch normalization right after each convolution and before activation, following [11].

### 2.3 Phase-Residual Block

Following conv\_1, there are two deep stacks of convolutional layers in HoloNet, namely conv\_2 and conv\_3 which are built by two Phase-Residual Blocks.

Recently, residual networks with largely increased depth have shown leading performance in the famous ImageNet Large Scale Visual Recognition Challenge (ILSVRC) 2015 [11]. However, the network used in the ILSVRC 2015 is substantially deep (hundreds of layers), which means heavy computational cost both in training and feed-forwarding. Furthermore, such kind of deep residual networks is prone to incur over-fitting problem, thus increasing difficulties on training significantly. To overcome these problems, especially in emotion recognition task, our Phase-Residual Block combines aforementioned CReLU with residual structure.

Figure 3 shows the structures of our two Phase-Residual Blocks used for constructing conv\_2 and conv\_3. It can be seen that each Phase-Residual Block is composed of two basic “bottleneck” building blocks sharing similar layer structures.

#### 2.3.1 Bottleneck Building Blocks

In [11], the authors have shown that bottleneck building block is more economical and powerful than non-bottleneck ones when constructing deep residual networks. For HoloNet, a “bottleneck” building block variant is designed, in which we replace originally basic ReLU with our CReLU and present a new residual variant, just as the blue dashed rectangle shown in Figure 3.

Specifically, our “bottleneck” building block is composed of three stacked convolutional layers with  $1 \times 1$ ,  $3 \times 3$  and  $1 \times 1$  filters, where the first  $1 \times 1$  convolution is responsible for dimensionality reduction and the other is used to increase dimensionality, leaving the middle  $3 \times 3$  convolutional layer as a “bottleneck”. On top of the  $3 \times 3$  convolution, CReLU operation is added. Our “bottleneck” building block is ended with a residual operation, performing element-wise addition channel by channel. Mathematically, the operations in a “bottleneck” building block can be formulated as

$$y = F(x, \{W_i\}) + W_s x, \quad (1)$$

where  $x$  is the input,  $y$  is the output, and

$$F(x, \{W_i\}) = W_3 \sigma(W_2(W_1 x)) \quad (2)$$

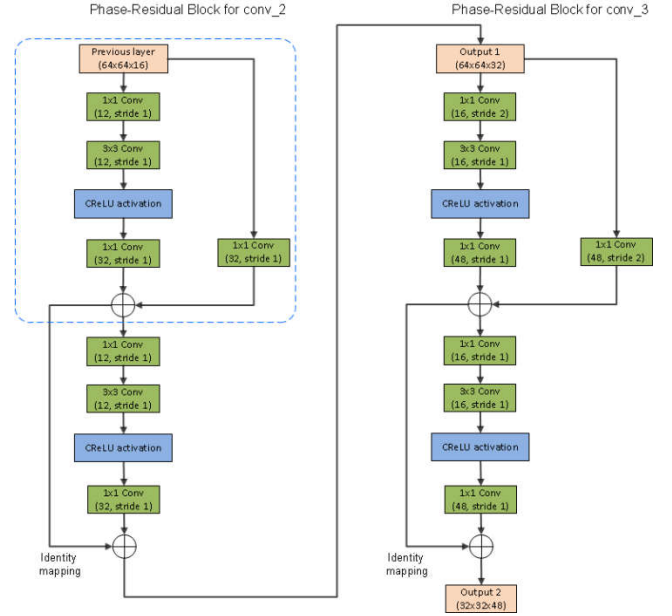


Figure 3. Structure of two paired Phase-Residual Blocks.

is the residual mapping in which  $\{W_i\}$  denotes weight parameters to be learnt,  $\sigma$  denotes CReLU, and  $W_s$  denotes the identity mapping or a linear projection matrix. If the dimensions of  $x$  and  $F$  are the same, identity mapping is used to form short connections, otherwise a linear projection by the shortcut connections is used to match dimensions.

#### 2.3.2 Phase-Residual Block

In HoloNet, conv\_2 and conv\_3 are constructed by two similar Phase-Residual Blocks, as shown in Figure 3. It is clear that each Phase-Residual Block includes two “bottleneck” building blocks that have the same filter banks in respective stacked convolutional layers ( $1 \times 1$ ,  $3 \times 3$  and  $1 \times 1$  filters). In each Phase-Residual Block, the first “bottleneck” building block makes short connections by the identity mapping and the other makes short connections by a linear projection. For each Phase-Residual Block, two “bottleneck” building blocks compose a couple, increasing network depth while enhancing feature representation. In the first Phase-Residual Block (i.e. conv\_2), a couple of two “bottleneck” building blocks extracts feature map with the same spatial size to the input. And in the second Phase-Residual Block (i.e. conv\_3), the couple starts with a  $1 \times 1$  convolution with stride 2 to down-sample feature map, and then performs subsequent operations. In this way, conv\_3 can mine face characteristic through a larger scale compared with conv\_2. With above Phase-Residual Blocks, our HoloNet can well enjoy accuracy gain from considerably increased depth and maintain efficiency both in training and feed-forwarding.

### 2.4 Inception-Residual Block

The last stack of convolutional layers of HoloNet is conv\_4 which is formed as a variant of Inception-Residual Block. Our Inception-Residual Block is inspired by two crucial factors. First, it has been well proved that high dimensional feature extracted in a multi-scale manner is necessary to achieve competitive results in many face related tasks. As well validated in the work of [1], with multi-scale facial feature extraction strategy, increasing feature dimensionality from 1K to over 100K can bring over 6% accuracy improvement in face recognition for five popular local descriptors.

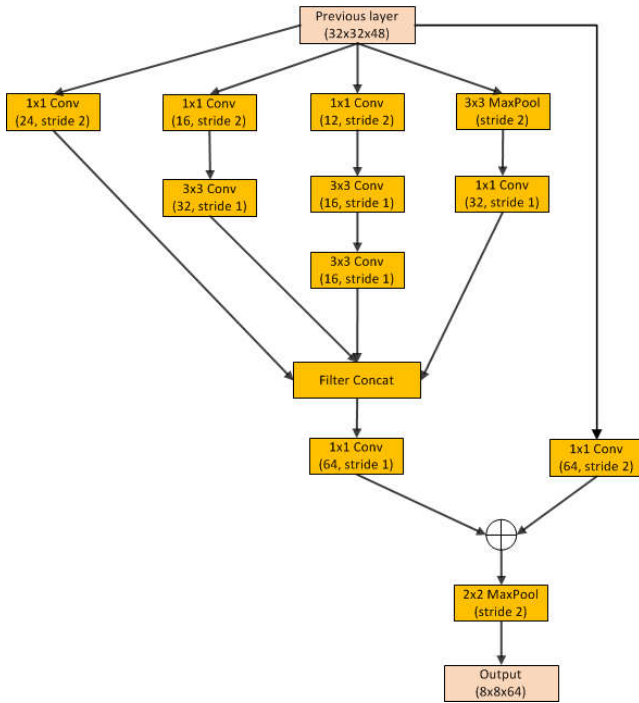


Figure 4. Structure of the Inception-Residual Block.

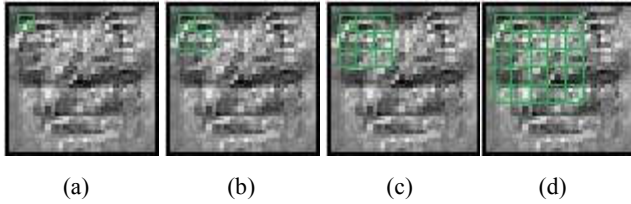


Figure 5. An illustration of multi-scale feature extraction over multi-path sibling branches of an Inception-Residual Block. Green grids denote convolutions with different filter sizes.

However, making a high dimensional feature will incur high cost on training, testing and storage. Second, the inception architecture proposed by Szegedy et al. [25] provide an easy and safe way of incorporating multi-scale feature extraction into the network by increasing both the depth and the width, yielding better accuracy. However, increasing network depth and width means a larger number of learnable parameters, which will result in dramatically increased use of computational resources, and the network may also be more prone to overfitting. Recently, Szegedy et al. [24] further give clear empirical evidence that the combinations of inception and residual connections can well accelerate the training. In their paper, there is also experimental evidence to show that several variants of such inception residual networks outperform similarly expensive pure inception networks by a thin accuracy margin. Although inception residual networks have been proved to be useful in image classification task, how to design such a hybrid network which is suitable for emotion recognition is still an open problem. Here, we present an Inception-Residual Block variant, a novel and efficient module which elaborately embeds multi-scale hierarchical sparse feature extraction into the last stack of convolutional layers of HoloNet.

The detailed structure of our Inception-Residual Block is shown in Figure 4. Note that there are four sibling branches whose inputs are the output of conv\_3, acting as a multi-scale sparse feature re-

Table 1. Summary of FLOPs in our HoloNet model.

Layer name	Output size	FLOPs (million)
conv_1	128×128	19.27
conv_2	64×64	22.22
conv_3	32×32	10.85
conv_4	16×16	5.86
fc5	1024	16.79
fc6	7	0.01
Total		75.00

presentation. From left to right, the first sibling branch in Figure 4 is a 1×1 convolution, akin to a 1×1 width filter illustrated in Figure 5 (a). The second sibling branch is a stack of 1×1 and 3×3 convolutions, akin to a 3×3 width filter illustrated in Figure 5 (c). The third and fourth sibling branches have their stacked layer configurations as shown in Figure 4, akin to a 5×5 width filter and a 2×2 width filter illustrated in Figure 5 (d) and (b), respectively. It can be seen that all sibling branches except the first have filter expansion layers (1×1 or 3×3 convolution) which scale up the dimensions of the filter banks to compensate for reduction in spatial feature channels. The sibling branches are followed by a grouping layer in which the output feature vectors are sequentially concatenated into a single output feature vector. At the next stage, the concatenated feature vector is further compressed by a 1×1 convolution. Subsequently, there is a residual operation layer where the final multi-scale feature vector is added to the conv\_3's output with a linear projection of shortcut connections. Finally, we perform a max pooling operation for down-sampling directly after residual layer. HoloNet is ended with two fully connected layers: the first has 1024 channels, and the second performs 7-way (one for each emotion category) classification with a softmax.

It should be emphasized here that our Inception-Residual Block is uniquely designed for emotion recognition, in contrast to any existing inception-residual variant [24]. Our design has two main merits. On the one hand, informative facial features, which are useful for discriminating target emotion categories, can be well extracted from micro scale to macro scale through inception unit. On the other hand, residual unit strengthens fast convergence and computational efficiency.

## 2.5 Model Efficiency

Besides the novel architecture of our HoloNet, it is also worth emphasizing its computational merit. Compared with popular classification CNN models such as AlexNet [13], VGG16 [22] and 18-layer ResNet [11], our HoloNet model has much fewer convolutional filters and thus has much lower computational cost, although they have similar depths. Generally, our HoloNet model has 75 million FLOPs (multiply-add), which is only 6.6% of AlexNet (1.14 billion FLOPs), 4.17% of 18-layer ResNet (1.8 billion FLOPs) and 0.48% of VGG16 (15.5 billion FLOPs). Considering that NVIDIA GPU such as K40 can process with a speed of over 250 image crops per second when using AlexNet model, thus our HoloNet model can be well run on any mobile platform with real-time processing requirement. The summary of FLOPs of our HoloNet model is given in Table 1.

## 3. PERFORMANCE EVALUATION

### 3.1 Parameter Settings

**Challenge Data.** We evaluate the performance of our HoloNet based method on the video based task of 4<sup>th</sup> Emotion Recognition in the Wild (EmotiW) 2016 challenge. The video base emotion recognition task includes audio-video clips containing seven basic emotion categories, namely Angry, Disgust, Fear, Happy, Neutral, Sad and Surprise. In the previous challenges, all video clips are collected from Hollywood real movie records, thus numerous variations in viewpoint, lighting, background, occlusion, context and etc. are introduced. The task is to assign one unique emotion label from seven candidate categories to each test video clip. The major change in this year's challenge as compared to the earlier challenges is the introduction of reality TV data into the test set while not into the training and validation sets. Another minor change is that the volume of training video clips is also increased. Similarly, this year's data, i.e., AFEW 6.0 is also split into three parts: training set (773 video clips), validation set (383 video clips) and test set (593 video clips, and 54 out of them are real TV clips). It should be highlighted that all our models described below are trained on the given data only. That is, in the challenge, we do not apply any outside data for augmenting the performance of our models.

**Implementation.** Our visual models are mainly based on HoloNet. We train two HoloNet models (denoted as A and B for simplicity), whose network architectures are exactly the same, for evaluation. The popular Caffe tool [21] is used to train our HoloNet models. In the implementation, we use the pre-processing steps described in Section 2.1 to normalize face image. The face image is finally scaled to a standard size of 136×136 pixels. A crop of 128×128 pixels is randomly sampled from the normalized face image or its horizontal flip for data augmentation. We initialize the weights as in [11] and train two HoloNet models from scratch. The standard SGD with a batch size of 128 is set. The learning rate starts from 0.01 and is divided by 10 per 20000/40000 iterations or the error plateaus. Totally, the models are trained for up to 70000/120000 iterations. We use a weight decay of 0.0002 and a momentum of 0.9. Considering that spatio-temporal hand-crafted features have shown to be useful for improving the accuracy of CNN models when handling video based emotion recognition task [12, 16, 31], we also extract improved Dense Trajectories (iDT) [29] over each video clip to characterize the dynamics of visual facial motions. Given a face video clip, three popular hand-crafted descriptors, namely Histogram of Oriented Gradients (HoG) [4], Histogram of Optical Flow (HoF) [3] and Motion Boundary Histogram (MBH) [28], are extracted. Fisher vectors for each of three descriptors are calculated first, and then are concatenated together as the final iDT descriptor. We use a linear Support Vector Machine (SVM) classifier to train our iDT model. Besides above visual models, similar to all previous participants, we also train an audio model to describe acoustic context cues. Our audio model is an SVM classifier with polynomial kernel (we set  $c=2.87$  and  $g=0.0025$ ). We use popular Opensmile tool [9] to extract 1582-D acoustic feature over every labeled video clip for training an audio model.

### 3.2 Results on AFEW 6.0

**Results on the Validation Set.** We first use validation set to explore the performance of different fusion strategies. Note we have four basic models: two HoloNet model A and B, an audio model and an iDT model. For a test video clip, both the audio model and the iDT model can predict the emotion scores of seven categories directly. However, the emotion scores of the HoloNet model are obtained in two steps. It sequentially operates on every

frame first, and then the average of the predicted emotion scores over all frames is calculated and used as the final emotion score for the given test video clip. Following this evaluation procedure, on the validation set, the recognition rate of each of above four models is 42.82%, 44.57%, 35.77% and 32.11%, respectively. Since our visual method is mainly based on HoloNet, so we first test the fusion strategy using each HoloNet model and the audio model. As a result, we get the recognition rate of 47.78% and 48.83% when fusing HoloNet model A and B separately with the audio model, obtaining 4.96 and 4.26 percent accuracy gain to that of respective individual HoloNet model. Such improvements are on par with those reported from the previous winner teams [12, 16, 31]. Another common fact is that the ensemble of several same typed CNN models usually performs better than single one [11, 13, 22]. Inspired by this, we further test the performance of fusing our two HoloNet models with the audio model. Specifically, we get the recognition rate of 50.03% and 50.91% in another two fusions, yielding higher accuracy gains against former two fusions. The main difference of two new fusions is only in the contribution portions of different models. As demonstrated in [12, 16], deep hierarchical CNN features and conventional hand-crafted features are complementary in nature. Therefore, in the last fusion, we use all our four models and achieve an impressive recognition rate of 51.96% which is the best among our 5 fusion strategies. Detailed fusion results are summarized in Table 2.

**Results on the Test Set.** Our 5 submissions to the challenge are based on the aforementioned 5 fusion strategies evaluated on the validation set, respectively. The summary of the mean recognition rate obtained from each of our 5 submissions in the challenge is shown in Table 2. On the test set, our best recognition rate is 57.84%, outperforming the baseline of 40.47% with an absolute margin of 17.37%, and yielding 4.04% absolute accuracy gain compared to the result of last year's winner team [31]. For each of 5 fusion strategies, it can be concluded from Table 2 that larger accuracy improvement is consistently obtained on the test set than on the validation set. This may be partially attributed to training data augmentation by introducing the validation set. Generally, these results well demonstrate the efficacy of our HoloNet method. The confusion matrices regarding the fusion method of our best submission on the validation and the test sets are shown in Figure 6 (a) and (b), respectively.

### 3.3 Discussions

**Result Analysis.** Although the proposed HoloNet based method gets competitive results in the challenge, we find its capability in classifying predefined seven basic emotion categories is not well balanced. Table 3 gives the per-class recognition rate concerning the method of our best submission. It can be seen that our method can well recognize Angry, Happy and Neutral, usually with above 70% recognition rate both on the test and the validation sets, yet it shows much lower accuracy in classifying the other four emotions. We can also see that our method gets an accuracy rate of 15% in recognizing Disgust on the validation set, but on the test set it produces a zero recognition rate. Similar phenomena can be also found in previous top-performing methods [12, 16, 31]. We think this may be attributed to two main factors regardless of common difficulties like different viewpoints, poses, lighting conditions, etc. First, in the dataset, there exists non-negligible ambiguity among different emotion categories (as can be seen from Figure 6 to some extent). Second, unbalanced training data distribution is another critical factor that may make the training procedure bias towards emotion categories with more training data.

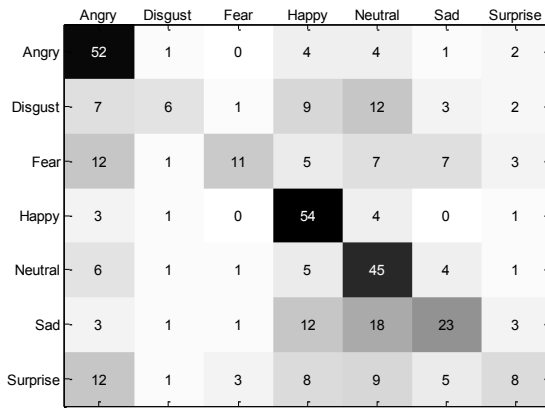


**Table 2. Total recognition accuracy of our 5 submissions to AFEW 6.0, both on the validation and the test sets.**

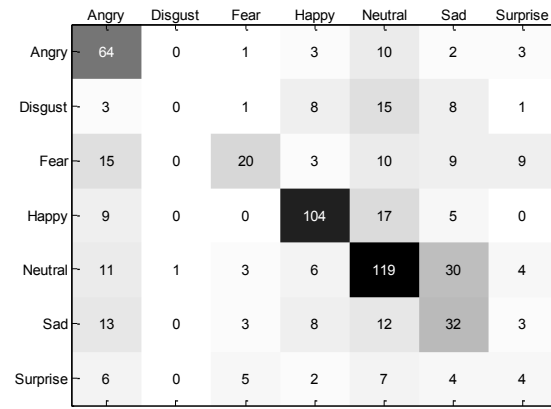
Submission#	Validation (%)	Test (%)	Method
1	47.78	54.30	Fusion of HoloNet model A + 1 audio model
2	48.83	55.14	Fusion of HoloNet model B + 1 audio model
3	50.13	56.83	1 <sup>st</sup> Fusion of HoloNet model A&B + 1 audio model
4	50.91	55.14	2 <sup>nd</sup> Fusion of HoloNet model A&B + 1 audio model
5	<b>51.96</b>	<b>57.84</b>	Fusion of HoloNet model A&B + 1 audio model + 1 iDT model

**Table 3. Per-class recognition rate of our best submission to AFEW 6.0, both on the validation and the test sets.**

Accuracy (%)	Angry	Disgust	Fear	Happy	Neutral	Sad	Surprise
<b>Validation</b>	81.25	15.00	23.91	85.71	71.43	37.70	17.36
<b>Test</b>	77.11	0.00	30.30	77.04	68.39	45.07	14.29



(a)



(b)

**Figure 6. Confusion matrices of our best submission to AFEW 6.0. For the results on the validation set, the models are trained with given training data only. For the results on the test set, the models are trained with the union of given training and validation sets. In the figures, the darker the grid, the higher the recognition rate.**

**Future Improvement.** According to above analysis, we believe that the performance can be further improved in two ways. One is to develop a tree based hierarchical method. In the method, seven emotion categories can be divided into several predefined non-overlapping sub-groups. Emotion categories sharing sufficient ambiguity will be merged into a sub-group. In this way, a tree based model can be trained to recognize target emotions in a coarse-to-fine manner [14]. The other one is to develop a robust learning strategy that can well handle unbalanced training data. Through bridging the performance gap from unbalanced training data, better recognition accuracy can be expected. We will explore them in the future.

#### 4. CONCLUSIONS

This paper presents HoloNet, a deep yet computational efficient convolutional neural network for emotion recognition in the wild. Our framework elaborately combines several state-of-the-art CNN backbones, including CReLU, residual structure and inception-residual structure, which are by design naturally complementary

to each other when handling tasks like emotion recognition in unconstrained conditions. It achieves considerably better accuracy in the video based task of EmotiW 2016 challenge compared with the baseline and previous winner counterparts. Besides, it is also computationally efficient during feed-forward inference, enabling the possible use in real-time applications. We believe that our framework naturally enjoys the benefits from the latest progress in the CNN field.

#### 5. REFERENCES

- [1] Chen, D. et al. 2013. Blessing of dimensionality: High-dimensional feature and its efficient compression for face verification. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (2013)*, 3025–3032.
- [2] Chen, W. et al. 2006. Illumination compensation and normalization for robust face recognition using discrete cosine transform in logarithm domain. *IEEE Transactions on*

- Systems, Man, and Cybernetics, Part B (Cybernetics)*. 36, 2 (2006), 458–466.
- [3] Dalal, N. et al. 2006. Human detection using oriented histograms of flow and appearance. *Proceedings of the 9th European Conference on Computer Vision* (2006), 428–441.
  - [4] Dalal, N. and Triggs, B. 2005. Histograms of oriented gradients for human detection. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2005), 886–893.
  - [5] Dhall, A. et al. 2013. Emotion recognition in the wild challenge 2013. *Proceedings of the 15th ACM International Conference on Multimodal Interaction* (2013), 509–516.
  - [6] Dhall, A. et al. 2014. Emotion recognition in the wild challenge 2014: Baseline, data and protocol. *Proceedings of the 16th ACM International Conference on Multimodal Interaction* (2014), 461–466.
  - [7] Dhall, A. et al. 2015. Video and image based emotion recognition challenges in the wild: EmotiW 2015. *Proceedings of the 17th ACM International Conference on Multimodal Interaction* (2015), 423–426.
  - [8] Dhall, A. et al. 2016. EmotiW 2016: Video and group-level emotion recognition challenges. *Proceedings of the 18th ACM International Conference on Multimodal Interaction* (2016).
  - [9] Eyben, F. et al. 2010. Opensmile: the munich versatile and fast open-source audio feature extractor. *Proceedings of the 18th ACM International Conference on Multimedia* (2010), 1459–1462.
  - [10] Hassner, T. et al. 2015. Effective face frontalization in unconstrained images. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2015), 4295–4304.
  - [11] He, K. et al. 2016. Deep residual learning for image recognition. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2016).
  - [12] Kahou, S.E. et al. 2013. Combining modality specific deep neural networks for emotion recognition in video. *Proceedings of the 15th ACM International Conference on Multimodal Interaction* (2013), 543–550.
  - [13] Krizhevsky, A. et al. 2012. ImageNet classification with deep convolutional neural networks. *Advances in Neural Information Processing Systems* (2012), 1097–1105.
  - [14] Lee C.-C., et al. 2011. Emotion recognition using a hierarchical binary decision tree approach. *Speech Communication*. 53, 9 (2011), 1162–1171.
  - [15] Liao, S. et al. 2007. Learning multi-scale block local binary patterns for face recognition. *Proceedings of the International Conference on Biometrics* (2007), 828–837.
  - [16] Liu, M. et al. 2014. Combining multiple kernel methods on riemannian manifold for emotion recognition in the wild. *Proceedings of the 16th ACM International Conference on Multimodal Interaction* (2014), 494–501.
  - [17] Ojala, T. et al. 2002. Multiresolution gray-scale and rotation invariant texture classification with local binary patterns. *IEEE Transactions on Pattern Analysis and Machine Intelligence*. 24, 7 (2002), 971–987.
  - [18] Peng, K.-C. et al. 2015. A mixed bag of emotions: Model, predict, and transfer emotion distributions. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2015), 860–868.
  - [19] Shan, C. et al. 2009. Facial expression recognition based on local binary patterns: A comprehensive study. *Image and Vision Computing*. 27, 6 (2009), 803–816.
  - [20] Shang, W. et al. 2016. Understanding and improving convolutional neural networks via concatenated rectified linear units. *Proceedings of the 33rd International Conference on Machine Learning* (2016).
  - [21] Shelhamer, E. et al. 2014. DIY deep learning for vision: A hands-on tutorial with caffe. *Proceedings of the 13th European Conference on Computer Vision* (2014).
  - [22] Simonyan, K. and Zisserman, A. 2015. Very deep convolutional networks for large-scale image recognition. *Proceedings of the 3rd International Conference on Learning Representations* (2015).
  - [23] Szegedy, C. et al. 2015. Going deeper with convolutions. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2015).
  - [24] Szegedy, C. et al. 2016. Inception-v4, inception-resnet and the impact of residual connections on learning. *arXiv preprint arXiv:1602.07261*. (2016).
  - [25] Szegedy, C. et al. 2015. Rethinking the inception architecture for computer vision. *arXiv preprint arXiv:1512.00567*. (2015).
  - [26] Taigman, Y. et al. 2014. Deepface: Closing the gap to human-level performance in face verification. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2014), 1701–1708.
  - [27] Viola, P. and Jones, M. 2001. Rapid object detection using a boosted cascade of simple features. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2001), 511–518.
  - [28] Wang, H. et al. 2013. Dense trajectories and motion boundary descriptors for action recognition. *International Journal of Computer Vision*. 103, 1 (2013), 60–79.
  - [29] Wang, H. and Schmid, C. 2013. Action recognition with improved trajectories. *Proceedings of the IEEE International Conference on Computer Vision* (2013), 3551–3558.
  - [30] Xiong, X. and De la Torre, F. 2013. Supervised descent method and its applications to face alignment. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2013), 532–539.
  - [31] Yao, A. et al. 2015. Capturing au-aware facial features and their latent relations for emotion recognition in the wild. *Proceedings of the 17th ACM International Conference on Multimodal Interaction* (2015), 451–458.
  - [32] Zeng, Z. et al. 2009. A survey of affect recognition methods: Audio, visual, and spontaneous expressions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*. 31, 1 (2009), 39–58.