

Emotion Recognition in the Wild via Convolutional Neural Networks and Mapped Binary Patterns

Gil Levi
The Open University of Israel
gil.levi100@gmail.com

Tal Hassner
USC / Information Sciences Institute
hassner@isi.edu
and
The Open University of Israel

ABSTRACT

We present a novel method for classifying emotions from static facial images. Our approach leverages on the recent success of Convolutional Neural Networks (CNN) on face recognition problems. Unlike the settings often assumed there, far less labeled data is typically available for training emotion classification systems. Our method is therefore designed with the goal of simplifying the problem domain by removing confounding factors from the input images, with an emphasis on image illumination variations. This, in an effort to reduce the amount of data required to effectively train deep CNN models. To this end, we propose novel transformations of image intensities to 3D spaces, designed to be invariant to monotonic photometric transformations. These are applied to CASIA Webface images which are then used to train an ensemble of multiple architecture CNNs on multiple representations. Each model is then fine-tuned with limited emotion labeled training data to obtain final classification models. Our method was tested on the Emotion Recognition in the Wild Challenge (EmotiW 2015), Static Facial Expression Recognition sub-challenge (SFEW) and shown to provide a substantial, 15.36% improvement over baseline results (40% gain in performance).

Categories and Subject Descriptors

I.5.4 [Pattern Recognition]: Applications—*Computer Vision*

Keywords

Emotion Recognition; Deep Learning; Local Binary Patterns; EmotiW 2015 Challenge

1. INTRODUCTION

Facial expressions play a vital role in social communications between humans. It is therefore unsurprising that automatic facial emotion recognition has become a subject of much recent research. Additional motivation comes from

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.

EmotiW'15, November 9, 2015, Seattle, WA, USA.

© 2015 ACM. ISBN 978-1-4503-3983-4/15/11 ...\$15.00.

DOI: <http://dx.doi.org/10.1145/2823327.2823333>.

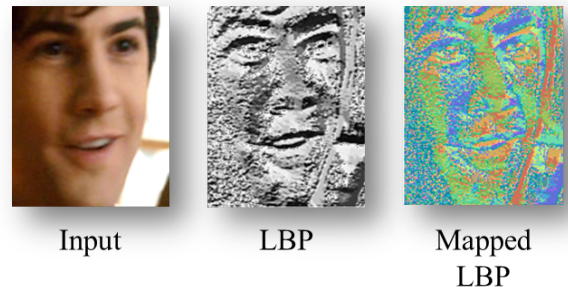


Figure 1: Image intensities (left) are converted to Local Binary Pattern (LBP) codes (middle), shown here as grayscale values. We propose to map these values to a 3D metric space (right) in order to use them as input for Convolutional Neural Network (CNN) models. 3D codes in the right image are visualized as RGB colors.

potential applications of effective systems for emotion recognition, with some examples including affect-sensitive Man-Machine-Interaction systems and auto tutors [12, 11, 29]. Yet despite this work, present automatic capabilities are still far from meeting the needs of commercial systems. This gap in performance is especially noteworthy when considering the substantial leaps in performance recently reported for the related task of face recognition (e.g., [33, 38, 37]).

Possible reasons for this performance gap between expression and face recognition systems include the significant difference in the data available for developing, training and testing automatic systems: Much of the recent achievements in machine face recognition have been due to deep Convolutional Neural Networks (CNN) which require massive amounts of labeled training data [33], and these are not yet available for emotion recognition. Exacerbating this is the particular nature of the emotion recognition problem, which involves large intra-class and small inter-class appearance variations. These are added to the many confounding factors of face image analysis under unconstrained *in-the-wild* conditions, such as those considered here.

We propose a method designed to jointly address two of the key challenges of automatic emotion recognition from still images. The first is the small amount of labeled data available for training deep CNN models and the second is appearance variation, specifically, illumination variability. We

define a novel mapping from intensities to illumination invariant, 3D spaces, based on the well-known Local Binary Patterns (LBP) feature transform [25, 27, 26]. Unlike LBP codes, our mapping produces values in a metric space which can be processed by CNN models (Fig. 1). Transformed images from the CASIA webface collection are used to train an ensemble of CNN models using different network architectures and applied to different representations. These are then fine-tuned on a substantially smaller set of emotion-labeled face images. We test our method on the Emotion Recognition in the Wild Challenge (EmotiW 2015), Static Facial Expression Recognition sub-challenge (SFEW) [10]. Our results demonstrate a remarkable 15.36% improvement over baseline scores (40% gain in performance).

2. RELATED WORK

Automatic facial emotion recognition has received increasing interest in the last two decades. The nature of the problems considered by previous work is reflected in the benchmarks used to measure and report performances. For example, early relevant data-sets such as those of [23, 18, 30] contained only constrained images taken in laboratory controlled conditions. More recently, unconstrained, in-the-wild photos have been considered with the release of the EmotiW challenge [9].

In [36] a Boosted-LBP descriptor was proposed, designed by learning the most discriminative LBP features using AdaBoost. Images represented using these features were then classified for different emotions using Support Vector Machines (SVM) [8]. In [5] a novel classification tree, based on sparse coding [24] was presented. The authors of [22] proposed a deep architecture which models face expressions by utilizing a set of local Action Unit features. Finally, the authors of [16] suggest using high-dimensional image features produced by dense census transformed vectors [44] based on locations defined by fiducial landmark detections.

Similar to some of the methods listed above, our own also uses deep CNNs. Unlike them, we propose applying these networks to pre-processed images, transformed in a manner designed to reduce variability due to illumination changes. Furthermore, we use network ensembles, rather than just single networks, where each one is trained using different image representations (different pre-processing applied to the input images) as well as different network architectures. Rigorous testing shows that these different networks provide complementary information which, when combined, provides a substantial boost to recognition performance.

Deep CNNs. Though CNNs have been introduced more than three decades ago [20], it is only recently that they become a predominant method in image classification tasks. With the emergence of very large classification data-sets, the increase in computation power and algorithmic improvements in training those models, the huge number of model parameters is no longer a limiting factor in applying CNNs in practical settings. Thus, in recent years, deep CNNs have been applied in various image classification problems, including, e.g., object recognition [19], scene recognition [45], face verification [33, 38, 37], age and gender classification [21], and more.

Local Binary Patterns. Local Binary Patterns (LBP) were originally developed as a means of describing texture images [25, 27, 26]. They have since been successfully applied to a wide range of other image recognition tasks, most notably face recognition [1, 41], age estimation [6] and gender classification [35]. To our knowledge, we are the first to propose the use of LBP features as input to CNN models. We show that doing so boosts performance well beyond that obtained with CNN models trained on RGB alone.

3. METHOD

We next describe the various components of our approach. We assume that images have been preprocessed by converting them to gray scale and cropping them to the region of the face. We further assume that faces are in-plane aligned. Though frontalization [14] may presumably be used here to reduce appearance variability further, we have not tested this approach in our pipeline. In practice, we use the images provided by the EmotiW 2015 challenge [10], aligned using the Zhu and Ramanan facial feature detector [46].

Each face image is processed as follows:

1. We begin by applying LBP encoding to the pixels of each image using different values for the LBP radius parameter (Section 3.1). Each encoding converts image intensity values to one of 256 LBP code values.
2. LBP codes are mapped to a 3D space using the mapping obtained by applying Multi Dimensional Scaling (MDS) using code-to-code dissimilarity scores based on an approximation to the Earth Mover’s Distance (Section 3.2).
3. The original RGB image, along with the mapped code images, are then used to train multiple, separate CNN models to predict one of seven emotion classes. A final classification is obtained by a weighted average over the outputs of the network ensemble taking the predicted emotion class to be the one with the maximum average probability (Section 3.3).

Each of these steps is described in detail next.

3.1 LBP code extraction

LBP codes have been widely used for nearly two decades; we refer to previous work for detailed description of how these are produced and of their various properties [25, 27, 26]. Here, we provide a cursory overview related to their use in our work.

LBP codes capture local image micro-textures. They are produced by applying thresholds on the intensity values of the pixels in small neighborhoods using the intensity of each neighborhood’s central pixel as the threshold. The resulting pattern of 0s (lower than the threshold) and 1s (higher than the threshold) is then treated as the pixel’s representation or *code*. When the neighborhood contains eight other pixels, this binary string is treated as an eight-bit number between 0 and 255. These codes are typically pooled over image regions using a histogram of code frequencies. Histograms are then used to represent the image region (see, e.g., [42]). In our work we use these codes in an entirely different manner.

Before continuing, we stop to consider specific advantage of LBP codes and the reason for their use here. By basing each pixel’s encoding on a threshold value applied to its

neighbors, the resulting representation is inherently invariant to monotonic photometric transformations; that is, any photometric transformation which does not change the order of image intensities. This includes, but is not limited to, additive and multiplicative intensity transformations, gamma correction and contrast manipulations.

Used with Support Vector Machines (SVM) classifiers, LBP code histograms have been key to the success of face recognition systems [41, 42, 14]. In contrast to these earlier methods, we wish to process LBP codes directly, without pooling, using CNN models. LBP codes, however, are, by their nature, not well suited as CNN inputs.

To understand why, note that the basic operation performed by a CNN on its input values is a convolution, which is equivalent to a weighted average of these values. When these values are codes from an unordered set, the outcome of a convolution is meaningless. To illustrate this, consider the following three LBP codes: $a = (0, 0, 0, 0, 0, 0, 0, 0)$, $b = (1, 0, 0, 0, 0, 0, 0, 0)$ and $c = (0, 0, 0, 0, 0, 0, 0, 1)$. Both b and c differ from a by just one bit, meaning the distance in the binary LBP space between a and b is the same as the distance between a and c . This can be taken to imply that they represent very similar local arrangements of pixel values. Simply treating these codes as eight-bit integer values, however, we get that $a = 0$, $b = 128$ and $c = 1$. Hence, the euclidean distance between b and a is much larger than the euclidean distance between c and a . Our mapping, described next, addresses this very issue.

3.2 Mapping LBP codes

The key to our LBP code mapping is the use of *Multi Dimensional Scaling* (MDS) [2, 34] to transform the unordered LBP code values to points in a metric space. In this way, transformed points may be averaged together using convolution operations yet their distances approximate the original code-to-code distances.

To this end, we first define a distance (dissimilarity) $\delta_{i,j}$ between LBP codes, $C_i, C_j \in 2_2^8$. This distance should reflect the underlying similarity of the image intensity patterns used to produce each LBP code string. A full dissimilarity matrix, representing the distances between all possible code values, can then be defined as:

$$\Delta := \begin{bmatrix} \delta_{11} & \delta_{12} & \delta_{13} & \dots & \delta_{1n} \\ \delta_{21} & \delta_{22} & \delta_{23} & \dots & \delta_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \delta_{d1} & \delta_{d2} & \delta_{d3} & \dots & \delta_{dn} \end{bmatrix},$$

For a given dissimilarity matrix Δ , MDS seeks a mapping of the codes to a low dimensional metric space, such that:

$$\delta_{i,j} \approx \|V_i - V_j\| = \|MDS(C_i) - MDS(C_j)\|. \quad (1)$$

Here, LBP codes C_i , and C_j are mapped to V_i, V_j resp.

Defining a binary LBP pattern dissimilarity. Ostensibly, the difference between two LBP codes can be estimated by the number of different bit values they have between them; that is by their Hamming distance. This, however, may not accurately reflect the differences in the intensity patterns which produced these codes: The *locations* of differing bit values, not just their number, are also important when considering code similarity.

To illustrate this, again consider the following three binary LBP vectors: $a = (1, 0, 0, 0, 0, 0, 0, 0)$, $b = (0, 1, 0, 0, 0, 0, 0, 0)$ and $c = (0, 0, 0, 0, 1, 0, 0, 0)$ (see also, Eq. 2). The number of different bits (and hence the Hamming distance) for the pair a and b is the same number as for a and c : two. The original texture patterns which produced these patterns, however, are very different: the pair a and b are related by a slight, one-pixel rotation of the pattern around the central pixel, whereas code c represents a mirror of the intensity pattern represented by a and the two are hence less similar.

$$a = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \times & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad b = \begin{bmatrix} 0 & 1 & 0 \\ 0 & \times & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad c = \begin{bmatrix} 0 & 0 & 0 \\ 0 & \times & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (2)$$

In order to account for differences in spatial locations of pixel codes rather than Hamming distance we use their *Earth Mover's Distance* (EMD) [31].

Generally speaking, EMD is defined to reflect the smallest effort required to convert one distribution into another. It is used here as a measure of the difference between two LBP codes. Formally, the EMD between two codes, P and Q is defined as follows:

$$EMD(P, Q) = \frac{\min_{\{f_{kl}\}} \sum_{kl} f_{kl} d_{kl}}{\sum_{kl} f_{kl}}, \text{ s.t.}, \quad (3)$$

$$0 \leq f_{kl}, \quad \sum_l f_{kl} \leq P_k, \quad \sum_k f_{kl} \leq Q_l, \quad \text{and}$$

$$\sum_{kl} f_{kl} = \min \left(\sum_k P_k, \sum_l Q_l \right).$$

Intuitively, this reflects the effort required to shift values from one code to another, where $f_{kl} \in \{0, 1\}$ is the flow of the value from bit k in P (P_k) to bit l in Q (Q_l), $\{f_{kl}\}$ is the entire flow from P to Q and

$$d_{kl} = \|k - l\|_2, \quad (4)$$

is taken to be the standard definition of a *ground distance* between two bit positions, and represents the effort required to move the value between bit positions.

EMD approximation. In practice, rather than compute the true EMD between the two code strings (Eq. 3), we approximate it by making the (here, often untrue) assumption that both codes have the same number of bits set to 1. This allows the use of the closed form solution to EMD [40, 7]:

$$EMD(P, Q) = \|CDF(P) - CDF(Q)\|_1, \quad (5)$$

with CDF being the cumulative distribution function (cumulative sum) of bit values. This not only allows faster computation than general EMD, but, more importantly, it better reflects code distances whenever the numbers of bits set to 1 is different. To illustrate this, the true EMD distance between the codes $a = (0, 0, 0, 0, 0, 0, 0, 0)$ and $b = (1, 1, 1, 1, 1, 1, 1, 1)$ would be 0, as there are no bit values of 1 to move from any position in a to b . This would wrongly imply that they encode similar local appearance. The value computed using Eq. 5, however, would be 36, correctly reflecting the difference in the local appearance represented by the two codes. The entire code-to-code distance matrix

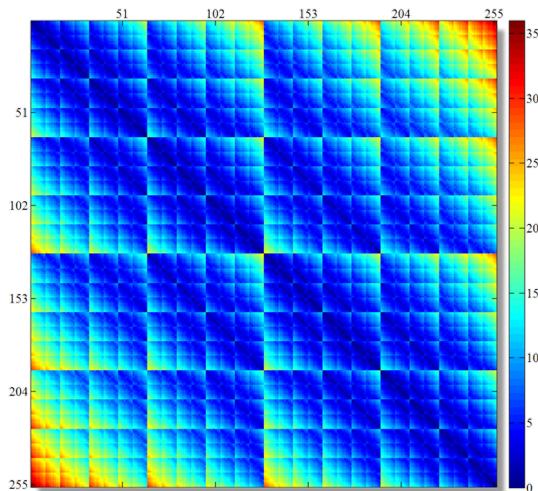


Figure 2: Visualization of code-to-code distance matrix Δ . Left is the 256×256 matrix Δ of code distances (color coded), computed using the EMD approximation of Eq. 5. On the right is the distances color bar. Evidently, this distance penalizes two codes for different positions and different numbers of of bits.

Δ produced using this approximation is visualized in Fig. 2. Clearly, the distances between codes reflect both the difference in positions and numbers of on bits.

Distances for cyclic codes. LBP codes are cyclic by design: the least significant bit represents a pixel adjacent to the pixel represented by the most significant bit (see, e.g., Eq. 2). When using the standard EMD of Eq. 3 this can be expressed by modifying the ground distance (Eq. 4). In order to employ the approximation of Eq. 5, however, we chose the following method of accounting for the cyclic nature of LBP codes.

Given two LBP codes $P, Q \in 2_8^2$, we append a single 0-valued bit as the new least significant bit of each code (increasing code sizes to nine bits) and denote the modified codes as P' and Q' , respectively. The modified, *cyclic distance*, $\delta'(P, Q)$ is defined by

$$\delta'(P, Q) = \min(\delta(P', Q'), \delta(\text{rev}(P'), Q'), \delta(P', \text{rev}(Q'))), \quad (6)$$

where δ is the the distance computed by Eq. 5 and $\text{rev}()$ rearranges code values in reverse order. The distances computed following this modification are illustrated in Fig. 3. A visual comparison of the distances in Fig. 3 to those in Fig. 2 demonstrates that indeed smaller distances (more similar codes) are assigned whenever bit values can be moved from one end of the code to another.

3.3 Ensemble of Deep CNN

Previous work has shown the benefits of employing multiple image representations and multiple similarity measures in face recognition tasks [41, 42]. Here, to our knowledge for the first time, we propose doing the same using multiple CNN architectures, multiple representations and multiple similarity measures.

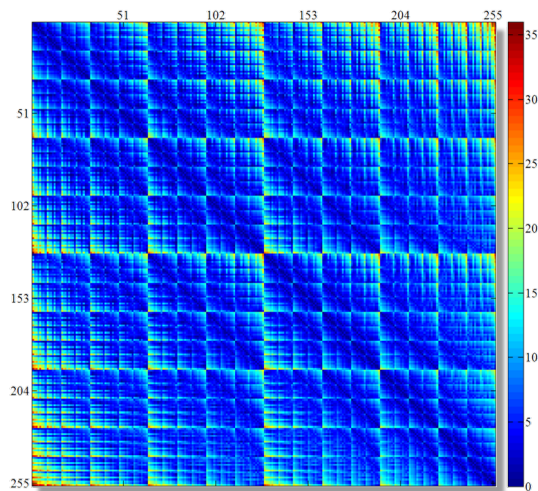


Figure 3: Visualization of code-to-code distance matrix Δ . Left is the 256×256 matrix Δ of code distances (color coded), computed using Eq. 6 to account for the cyclic nature of LBP codes. On the right is the distances color bar. Compared with the distances visualized in Fig. 2, distances are smaller (codes are more similar) whenever codes have bits set to 1 across code rotations.

Specifically, we employ four different, existing network architectures: the *VGG_S*, *VGG_M-2048* and *VGG_M-4096* networks presented in [4] and the BVLC *GoogLeNet* network presented in [39]. Please see relevant papers for the details of each network design and architecture. In all cases, CNNs were trained to predict 7D vectors of emotion class probabilities using the labeled training data (seven different emotion classes and an additional “neutral” class).

Images were represented using both RGB values as well as by extracting LBP codes with three different radius parameter values: the default of 1, as well as 5 and 10. All three LBP variants were processed using the encoding described in Sec. 3.2 with the cyclic distance of Eq. 6. In order to compare the influence of the cyclic adaptation to the original EMD approximation (Eq. 5) LBP codes with radius set to 1 were additionally encoded using the EMD approximation directly. All told, we use four CNN architectures and five representations for an ensemble of 20 networks.

In order to predict emotion labels, we take a weighted average of the 7D output vectors produced by our 20 ensemble models. The selected class is the one with the highest probability in the resulting 7D average prediction vector. Weights reflect the relative importance of each ensemble component. These were determined empirically by random searching through different weight combinations using the training data to evaluate the quality of each combination. The best performing weights were selected for our tests and are visualized in Fig. 4. Curiously, of the top ten models, only one uses the original RGB values. Despite being the common practice for CNN based approaches, an RGB representation seems to provide inferior results than those obtained using more robust features as input.

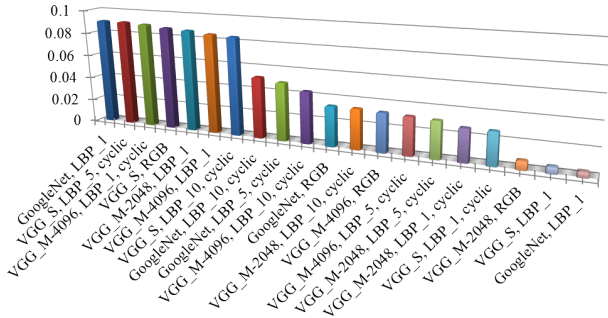


Figure 4: Relative weights of CNN ensemble components. Sorted left to right from the most important to the least. Interestingly, of the top ten most important networks, only one uses original RGB values, despite being the common practice in CNN based approaches.

Training vs. fine-tuning. Due to the huge number of model parameters, deep CNN are prone to overfitting when they are not trained with a sufficiently large training set. The EmotiW challenge contains only 891 training samples, making it dramatically smaller than other image classification data-sets commonly used for training deep networks (e.g, the Imagenet dataset [32]).

To alleviate this problem, we train our models in two steps: First, we fine-tune pre-trained object classification networks on a large face recognition data-set, namely the CASIA WebFace data-set [43]. This allows the network to learn general features relevant for face classification related problems. Then, we fine-tune the resulting networks for the problem of emotion recognition using the smaller training set given in the challenge.

Oversampling the input representations. *Oversampling* [15] is the process of providing a network with several, slightly translated versions of the same input image. In the related work of [21] oversampling was shown to provide superior age and gender classification performance on the Adience benchmark [13]. Here, we employ the same process with all representations and all CNN architectures. Specifically, our results, reported in Sec. 4 provide the performances obtained using the following techniques:

- **Center crop:** The CNN was applied to image regions of size 224×224 pixels, cropped from the center of the input representation.
- **Oversampling:** We produce five 227×227 pixel regions cropped from the input representation as follows: four regions aligned with the four corners of each input representation and one from its center. These five regions, along with their mirrored versions were presented to the CNN. The prediction values of the CNN were then averaged over all ten predictions.

As we later show, similarly to [21], oversampling generally provided better prediction accuracy and we therefore employed it with all our networks in our ensemble.

Table 1: The EmotiW SFEW 2.0 Challenge. Break-down of the SFEW 2.0 benchmark into the different emotions classes.

	Anger	Disgust	Fear	Happy	Neutral	Sad	Surprise	Total
Train	178	52	78	184	144	161	94	891
Val	77	23	46	72	84	73	56	431

4. EXPERIMENTS

LBP encoding and mapping, as described in Sections 3.1 and 3.2, was implemented in Matlab. Training and testing the networks was done using the Caffe open source framework for Deep Convolutional Neural Networks [15]. We have used two Amazon Web Services g2.8xlarge GPU machines, each with 4 NVIDIA GRID GPUs, where each GPU has 1,536 CUDA cores and 4 GB of video memory.

Fine-tuning each of the networks on the CASIA WebFace Database [43] required approximately four days. Additional fine-tuning of the resulting networks the training set for the EmotiW 2015 data set [10] required one more day. Emotion prediction for a single image requires ≈ 500 ms. This time, however, may potentially be substantially reduced if prediction was performed on image batches rather than individual images. Please see our project webpage¹ for code and additional details.

Our tests were performed on the EmotiW 2015 benchmark [10] which includes data from version 2.0 of the **Static Facial Expression in the Wild** benchmark [9]. It was assembled by selecting frames from different videos of the **Acted Facial Expressions in the Wild (AFEW)**, and then assigning them one of the following seven emotion labels: *angry*, *disgust*, *happy*, *sad*, *surprise*, *fear* and *neutral*. Images from this data set are unconstrained and cover a wide range of head poses and ages, both genders, different occlusions and resolution qualities. Table 1 additionally provides a break-down of the number of images in each emotion category for the training and validation subsets.

4.1 Results

Table 2 summarizes the results on the validation set for all of the different network architectures and image presentations considered. Subscripts used for the LBP representations denote the value of the radius parameter values used (default of 1, as well as 5 and 10). The use of the approximate EMD without modifying it for cyclic codes – that is, direct application of Eq. 5 to compute code-to-code distances – is referred to as “w.o. cyclic”. We use “cyclic” to refer to the modified distance of Eq. 6.

Several interesting observations can be made from these results. First, in line with the ensemble weights reported in Fig. 4, using RGB values as input does not necessarily provide the best performance. This is at odds with the common practice of applying CNNs directly to input images. Presumably, pre-processing images using robust feature alleviates some of the challenges CNNs face when adopting to a particular domain; learning can better focus on information important to the recognition task, rather than filtering out confounding factors such as illumination variations.

Also important are the differences between ensembles and single CNNs. In all cases, whether by combining differ-

¹Available: http://www.openu.ac.il/home/hassner/projects/cnn_emotions

Table 2: Emotion classification results. The accuracy over all emotion classes is listed. Subscripts for the LBP representations denote the values used for their radius parameters; *w.o. cyclic* refers to the use of the approximate EMD of Eq. 5; and finally, *cyclic* refers to the modified distance of Eq. 6.

	RGB	LBP ₁ , w.o. cyclic	Validation			Test
			LBP ₁ , cyclic	LBP ₅ , cyclic	LBP ₁₀ , cyclic	
Baseline (provided by the Challenge authors)			35.33%			39.13%
GoogleNet - single crop	41.68%	39.34%	41.45%	41.68%	40.28%	—
GoogleNet - Oversampling	41.21%	39.57%	40.98%	40.98%	41.45%	—
VGG_S - single crop	41.45%	39.34%	41.92%	41.92%	39.34%	—
VGG_S - Oversampling	40.04%	43.79%	42.38%	43.09%	41.21%	—
VGG_M-2048 - single crop	37.93%	38.17%	36.06%	40.98%	32.31%	—
VGG_M-2048 - Oversampling	40.74%	42.62%	36.76%	42.62%	33.72%	—
VGG_M-4096 - single crop	37.00%	41.92%	40.28%	40.28%	37.47%	—
VGG_M-4096 - Oversampling	37.93%	42.85%	44.73%	42.85%	40.98%	—
Ensemble	46.13	48.94	47.3	47.54	47.3	—
Ensemble of all methods			51.75%			54.56%

ent networks applied to the same representation or different representations used with the same architectures, ensembles seem superior in performance than single networks. Though this should come as no surprise, considering earlier related work (e.g., [42]), evidence of the advantages provided by CNN ensembles are still scarce.

Finally, modifying the approximate EMD distance in order to address rotations of the LBP codes slightly degrades results. On the other hand, combining all distances appears to provide complementary information and contributes to improving overall accuracy.

Our results should be compared against the baseline performance for the benchmark. These were obtained using features produced from pyramids of Histogram of Gradients [3] and Local Phase Quantization [28] extracted from the aligned faces and classified using a fusion of separate support vector machines. In the majority of cases, our individual models outperformed the baseline, though not by large margins. Ensemble results, however, boosted performance substantially, by a remarkable 40% gain in performance (15.36% improvement).

Tables 3 and 4 provide confusion matrices for the validation and test sets respectively. Evidently, the *disgust* emotion was never classified correctly. This is consistent with results previously reported for the 2013 Emotion Recognition in the Wild Challenge (e.g., [17]). This performance may be due to the class being inherently more challenging to classify or simply due to the relatively small number of examples available for the class (see Table 1).

Finally, we provide a selection of correct and wrong classification results for all of the remaining six classes is provided in Fig. 5. These show that at least in some cases, poor performance may be traced to failures in the face alignment step, rather than the recognition pipeline.

5. CONCLUSIONS

We present a substantial improvement over existing baseline results on the Emotion Recognition in the Wild Challenge (EmotiW 2015), Static Facial Expression Recognition sub-challenge (SFEW). To achieve this performance boost, we make a number of novel contributions: We propose to apply CNNs to pre-processed images rather than RGB, in order to eliminate confounding factors and focus the network’s efforts on variations due to class labels. To this end, we convert images to LBP codes, designed to be robust to illumination changes. In order to apply CNNs to these bi-

nary codes, we further describe a unique mapping of codes to metric space by applying an approximation of the EMD. Finally, multiple CNN architectures and representations are combined in an ensemble by a weighted average of their predictions. Our results clearly demonstrate the advantage of looking beyond RGB as the input space for CNNs, as well as the complementary information offered by multiple representations and network architectures.

Acknowledgments

This research is based upon work supported in part by the Office of the Director of National Intelligence (ODNI), Intelligence Advanced Research Projects Activity (IARPA), via IARPA 2014-14071600010. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of ODNI, IARPA, or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purpose notwithstanding any copyright annotation thereon.

6. REFERENCES

- [1] T. Ahonen, A. Hadid, and M. Pietikainen. Face description with local binary patterns: Application to face recognition. *Trans. Pattern Anal. Mach. Intell.*, 28(12):2037–2041, 2006.
- [2] I. Borg and P. J. Groenen. *Modern multidimensional scaling: Theory and applications*. Springer Science & Business Media, 2005.
- [3] A. Bosch, A. Zisserman, and X. Munoz. Representing shape with a spatial pyramid kernel. In *Proc. Int. Conf. on Image and video retrieval*, pages 401–408. ACM, 2007.
- [4] K. Chatfield, K. Simonyan, A. Vedaldi, and A. Zisserman. Return of the devil in the details: Delving deep into convolutional nets. *arXiv preprint arXiv:1405.3531*, 2014.
- [5] H.-C. Chen, M. Z. Comiter, H. Kung, and B. McDanel. Sparse coding trees with application to emotion classification. In *Proc. Conf. Comput. Vision Pattern Recognition Workshops*. IEEE, 2015.
- [6] S. E. Choi, Y. J. Lee, S. J. Lee, K. R. Park, and J. Kim. Age estimation using a hierarchical classifier based on global and local facial features. *Pattern Recognition*, 44(6):1262–1281, 2011.

Table 3: Validation set confusion matrix

	Anger	Disgust	Fear	Happy	Neutral	Sad	Surprise
Anger	48.04	0.	16.67	12.22	3.3	10.42	6.25
Disgust	5.88	0.	0.	5.56	6.59	5.21	3.12
Fear	17.65	0.	41.67	3.33	6.59	8.33	15.62
Happy	2.94	0.	0.	64.44	2.2	8.33	3.12
Neutral	7.84	50.	8.33	2.22	58.24	13.54	15.62
Sad	5.88	0.	8.33	6.67	12.09	43.75	12.5
Surprise	11.76	50.	25.	5.56	10.99	10.42	43.75

Table 4: Test set confusion matrix

	Anger	Disgust	Fear	Happy	Neutral	Sad	Surprise
Anger	54.55	66.67	22.58	3.16	4.48	9.26	11.9
Disgust	2.6	0.	0.	8.42	4.48	5.56	2.38
Fear	16.88	16.67	19.35	3.16	4.48	12.96	19.05
Happy	10.39	16.67	9.68	75.79	4.48	14.81	0.
Neutral	2.6	0.	16.13	0.	62.69	14.81	2.38
Sad	5.19	0.	25.81	6.32	16.42	37.04	14.29
Surprise	7.79	0.	6.45	3.16	2.99	5.56	50.

**Figure 5: Example prediction results. For each emotion the left pair is a wrong classification and the right is a correct classification result.**

- [7] S. Cohen and L. Guibas. The earth mover’s distance: Lower bounds and invariance under translation. Technical report, DTIC Document, 1997.
- [8] C. Cortes and V. Vapnik. Support-vector networks. *Machine Learning*, 20(3):273–297, 1995.
- [9] A. Dhall, R. Goecke, S. Lucey, and T. Gedeon. Static facial expression analysis in tough conditions: Data, evaluation protocol and benchmark. In *Proc. Int. Conf. Comput. Vision Workshops*, pages 2106–2112. IEEE, 2011.
- [10] A. Dhall, O. R. Murthy, R. Goecke, J. Joshi, and T. Gedeon. Video and image based emotion recognition challenges in the wild: EmotiW 2015. In *Int. Conf. on Multimodal Interaction*. ACM, 2015.
- [11] S. D’ÁZMello, N. Blanchard, R. Baker, J. Ocumpaugh, and K. Brawner. Affect-sensitive instructional strategies. *Design Recommendations for Intelligent Tutoring Systems: Volume 2-Instructional Management*, 2:35, 2014.
- [12] S. D’Mello, R. W. Picard, and A. Graesser. Toward an affect-sensitive autotutor. *Intelligent Systems*, (4):53–61, 2007.
- [13] E. Eiding, R. Enbar, and T. Hassner. Age and gender estimation of unfiltered faces. *Trans. on Inform. Forensics and Security*, 9(12), 2014.
- [14] T. Hassner, S. Harel, E. Paz, and R. Enbar. Effective face frontalization in unconstrained images. In *Proc. Conf. Comput. Vision Pattern Recognition*, 2015.
- [15] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. *arXiv preprint arXiv:1408.5093*, 2014.
- [16] S. E. Kahou, P. Froumenty, and C. Pal. Facial expression analysis based on high dimensional binary features. In *European Conf. Comput. Vision*, pages 135–147. Springer, 2014.
- [17] S. E. Kahou, C. Pal, X. Bouthillier, P. Froumenty, Ç. Gülçehre, R. Memisevic, P. Vincent, A. Courville, Y. Bengio, R. C. Ferrari, et al. Combining modality specific deep neural networks for emotion recognition in video. In *Int. Conf. on Multimodal Interaction*, pages 543–550. ACM, 2013.
- [18] T. Kanade, J. F. Cohn, and Y. Tian. Comprehensive database for facial expression analysis. In *Automatic*

- Face and Gesture Recognition*, pages 46–53. IEEE, 2000.
- [19] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Neural Inform. Process. Syst.*, pages 1097–1105, 2012.
- [20] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel. Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1(4):541–551, 1989.
- [21] G. Levi and T. Hassner. Age and gender classification using convolutional neural networks. In *Proc. Conf. Comput. Vision Pattern Recognition Workshops*, June 2015.
- [22] M. Liu, S. Li, S. Shan, and X. Chen. AU-aware deep networks for facial expression recognition. In *Automatic Face and Gesture Recognition*. IEEE, 2013.
- [23] M. J. Lyons, S. Akamatsu, M. Kamachi, J. Gyoba, and J. Budynek. The japanese female facial expression (jaffe) database, 1998.
- [24] J. Mairal, F. Bach, J. Ponce, and G. Sapiro. Online dictionary learning for sparse coding. In *Int. Conf. Mach. Learning*, pages 689–696. ACM, 2009.
- [25] T. Ojala, M. Pietikäinen, and D. Harwood. A comparative study of texture measures with classification based on featured distributions. *Pattern recognition*, 29(1):51–59, 1996.
- [26] T. Ojala, M. Pietikäinen, and T. Mäenpää. A generalized local binary pattern operator for multiresolution gray scale and rotation invariant texture classification. In *Int. Conf. Pattern Recognition*, volume 1, pages 397–406. Springer, 2001.
- [27] T. Ojala, M. Pietikäinen, and T. Mäenpää. Multiresolution gray-scale and rotation invariant texture classification with local binary patterns. *Trans. Pattern Anal. Mach. Intell.*, 24(7):971–987, 2002.
- [28] V. Ojansivu and J. Heikkilä. Blur insensitive texture classification using local phase quantization. In *Image and signal processing*, pages 236–243. Springer, 2008.
- [29] M. Pantic and L. J. Rothkrantz. Toward an affect-sensitive multimodal human-computer interaction. *Proceedings of the IEEE*, 91(9):1370–1390, 2003.
- [30] M. Pantic, M. Valstar, R. Rademaker, and L. Maat. Web-based database for facial expression analysis. In *Int. Conf. on Multimedia and Expo*. IEEE, 2005.
- [31] Y. Rubner, C. Tomasi, and L. J. Guibas. The earth mover’s distance as a metric for image retrieval. *Int. J. Comput. Vision*, 40(2):99–121, 2000.
- [32] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *Int. J. Comput. Vision*, pages 1–42, 2014.
- [33] F. Schroff, D. Kalenichenko, and J. Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proc. Conf. Comput. Vision Pattern Recognition*, pages 815–823, 2015.
- [34] G. A. Seber. *Multivariate observations*, volume 252. John Wiley & Sons, 2009.
- [35] C. Shan. Learning local binary patterns for gender classification on real-world face images. *Pattern Recognition Letters*, 33(4):431–437, 2012.
- [36] 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.
- [37] Y. Sun, D. Liang, X. Wang, and X. Tang. Deepid3: Face recognition with very deep neural networks. *arXiv preprint arXiv:1502.00873*, 2015.
- [38] Y. Sun, X. Wang, and X. Tang. Deeply learned face representations are sparse, selective, and robust. *arXiv preprint arXiv:1412.1265*, 2014.
- [39] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. *arXiv preprint arXiv:1409.4842*, 2014.
- [40] M. Werman, S. Peleg, and A. Rosenfeld. A distance metric for multidimensional histograms. *Computer Vision, Graphics, and Image Processing*, 32(3):328–336, 1985.
- [41] L. Wolf, T. Hassner, and Y. Taigman. Descriptor based methods in the wild. In *European Conf. Comput. Vision Workshops*, 2008.
- [42] L. Wolf, T. Hassner, and Y. Taigman. Effective unconstrained face recognition by combining multiple descriptors and learned background statistics. *Trans. Pattern Anal. Mach. Intell.*, 33(10):1978–1990, 2011.
- [43] D. Yi, Z. Lei, S. Liao, and S. Z. Li. Learning face representation from scratch. *arXiv preprint arXiv:1411.7923*, 2014.
- [44] R. Zabih and J. Woodfill. Non-parametric local transforms for computing visual correspondence. In *European Conf. Comput. Vision*, pages 151–158. Springer, 1994.
- [45] B. Zhou, A. Lapedriza, J. Xiao, A. Torralba, and A. Oliva. Learning deep features for scene recognition using places database. In *Neural Inform. Process. Syst.*, pages 487–495, 2014.
- [46] X. Zhu and D. Ramanan. Face detection, pose estimation, and landmark localization in the wild. In *Proc. Conf. Comput. Vision Pattern Recognition*, pages 2879–2886. IEEE, 2012.