# Do Deep Neural Networks Learn Facial Action Units
# When Doing Expression Recognition?

Pooya Khorrami          Tom Le Paine          Thomas S. Huang

Beckman Institute for Advanced Science and Technology
University of Illinois at Urbana-Champaign
{pkhorra2,paine1,t-huang1}@illinois.edu

## Abstract

*Despite being the appearance-based classifier of choice in recent years, relatively few works have examined how much convolutional neural networks (CNNs) can improve performance on accepted expression recognition benchmarks and, more importantly, examine what it is they actually learn. In this work, not only do we show that CNNs can achieve strong performance, but we also introduce an approach to decipher which portions of the face influence the CNN's predictions. First, we train a zero-bias CNN on facial expression data and achieve, to our knowledge, state-of-the-art performance on two expression recognition benchmarks: the extended Cohn-Kanade (CK+) dataset and the Toronto Face Dataset (TFD). We then qualitatively analyze the network by visualizing the spatial patterns that maximally excite different neurons in the convolutional layers and show how they resemble Facial Action Units (FAUs). Finally, we use the FAU labels provided in the CK+ dataset to verify that the FAUs observed in our filter visualizations indeed align with the subject's facial movements.*

## 1. Introduction

Facial expressions provide a natural and compact way for humans to convey their emotional state to another party. Therefore, designing accurate facial expression recognition algorithms is crucial to the development of interactive computer systems in artificial intelligence. Extensive work in this area has found that only a small number of regions change as a human changes their expression and are located around the subject's eyes, nose and mouth. In [7], Paul Ekman proposed the Facial Action Coding System (FACS) which enumerated these regions and described how every facial expression can be described as the combination of multiple action units (AUs), each corresponding to a particular muscle group in the face. However, having a computer
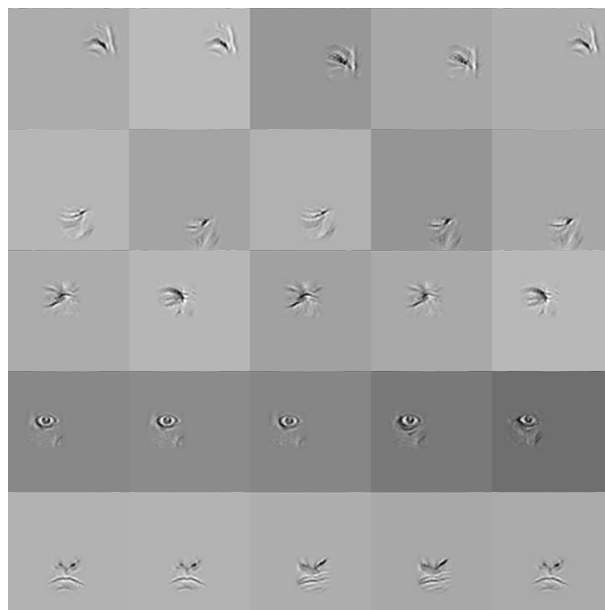


Figure 1. Visualization of facial regions that activate five selected filters in the 3rd convolutional layer of a network trained on the Extended Cohn-Kanade (CK+) dataset. Each row corresponds to one filter in the conv3 layer and we display the spatial patterns from the top 5 images.

accurately learn the parts of the face that convey emotion has proven to be a non-trivial task.

Previous work in facial expression recognition can be split into two broad categories: AU-based/rule-based methods and appearance-based methods. AU-based methods [29, 30] would detect the presence of individual AUs explicitly and then classify a person's emotion based on the combinations originally proposed by Friesen and Ekman in [8]. Unfortunately, each AU detector required careful hand-engineering to ensure good performance. On the other hand, appearance-based methods [1, 2, 31, 33] modeled a person's expression from their general facial shape and texture.

In the last few years, many well-established problems in computer vision have greatly benefited from the rise of convolutional neural networks (CNNs) as an appearance-based classifier. Tasks such as object recognition [14], object detection [9], and face recognition [28] have seen huge boosts in performance on several accepted benchmarks. Unfortunately, other tasks such as facial expression recognition have not experienced performance gains of the same magnitude. Little work has been done to see how much deep CNNs can help on accepted expression recognition benchmarks.

In this paper, we seek the answer to the following questions: Can CNNs improve performance on emotion recognition datasets/baselines and what do they learn? We propose to do this by training a CNN on established facial expression datasets and then analyzing what they learn by visualizing the individual filters in the network. In this work, we apply the visualization techniques proposed by Zeiler and Fergus [32] and Springenberg et al. [25] where individual neurons in the network are excited and their corresponding spatial patterns are displayed in pixel space using a deconvolutional network. When visualizing these discriminative spatial patterns, we find that many of the filters are excited by regions in the face that corresponded to Facial Action Units (FAUs). A subset of these spatial patterns is shown in Figure 1.

Thus, the main contributions of this paper are as follows:

1. We show that CNNs trained for the emotion recognition task learn features that correspond strongly with the FAUs proposed by Ekman [7]. We demonstrate this result by first visualizing the spatial patterns that maximally excite different filters in the convolutional layers of our networks, and then using the ground truth FAU labels to verify that the FAUs observed in the filter visualizations align with the subject's facial movements.

2. We also show that our CNN model, based on works originally proposed by [20, 21], can achieve, to our knowledge, state-of-the-art performance on the extended Cohn-Kanade (CK+) dataset and the Toronto Face Dataset (TFD).

## 2. Related Work

In most facial expression recognition systems, the main machinery matches quite nicely with the traditional machine learning pipeline. More specifically, a face image is passed to a classifier that tries to categorize it as one of several (typically 7) expression classes: 1. anger, 2. disgust, 3. fear, 4. neutral, 5. happy, 6. sad, and 7. surprise. In most cases, prior to being passed to the classifier, the face image is pre-processed and given to a feature extractor. Up until rather recently, most appearance-based expression recognition techniques relied on hand-crafted features, specifically

Gabor wavelets [1, 2], Haar features [31] and LBP features [33], in order to make representations of different expression classes more discriminative.

For some time, systems based on hand-crafted features were able to achieve impressive results on accepted expression recognition benchmarks such as the Japanese Female Facial Expression (JAFFE) database [19], the extended Cohn-Kanade (CK+) dataset [18], and the Multi-PIE dataset [10]. However, the recent success of deep neural networks has caused many researchers to explore feature representations that are learned from data. Not surprisingly, almost all of the methods used some form of unsupervised pre-training/learning to initialize their models. We hypothesize this may be because the scarcity of labeled data prevented the authors from training a completely supervised model that did not experience heavy overfitting.

In [17], the authors trained a multi-layer boosted deep belief network (BDBN) and achieved state-of-the-art accuracy on the CK+ and JAFFE datasets. Meanwhile in [23], the authors used a convolutional contractive auto-encoder (CAE) as their underlying unsupervised model. They then performed a semi-supervised encoding function called Contractive Discriminant Analysis (CDA) to separate discriminative expression features from the unsupervised representation.

A few works based on unsupervised deep learning have also tried to analyze the relationship between FAUs and the learned feature representations. In [15, 16], the authors learned a patch-based filter bank using K-means as their low-level feature. These features were then used to select receptive fields corresponding to specific FAU receptive fields which were subsequently passed to multi-layer restricted Boltzmann machines (RBMs) for classification. The FAU receptive fields were selected using a mutual information criterion between the image feature and the expression label. An earlier work by Susskind et al. [27], showed that the first layer features a deep belief network trained to generate facial expression images appeared to learn filters that were sensitive to face parts. We conduct a similar analysis except we use a CNN as our underlying model and we visualize the spatial patterns that excite higher-level neurons in the network.

To the authors' knowledge, the only works that previously applied CNNs to expression data were that of Kahou et al. [13, 12] and Jung et al. [11]. In [13, 12], the authors developed a system for doing audio/visual emotion recognition for the Emotion Recognition in the Wild Challenge (EmotiW) [6, 5] while in [11], the authors trained a network that incorporated both appearance and geometric features when doing recognition. However, one key point is that these works dealt with emotion recognition of video / image sequence data and therefore, actively incorporated temporal data when computing their predictions.
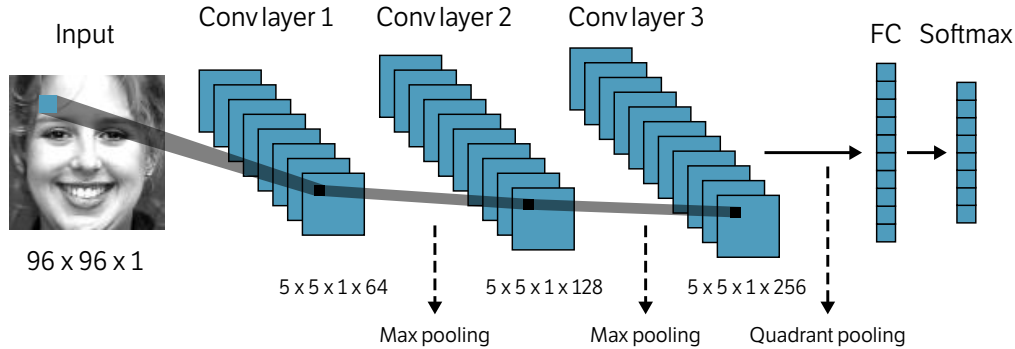
Figure 2. Network Architecture - Our network consists of three convolutional layers containing 64, 128, and 256 filters, respectively, each of size 5x5 followed by ReLU (Rectified Linear Unit) activation functions. We add 2x2 max pooling layers after the first two convolutional layers and quadrant pooling after the third. The three convolutional layers are followed by a fully-connected layer containing 300 hidden units and a softmax layer.

In contrast, our work deals with emotion recognition from a single image, and will focus on analyzing the features learned by the network. Thus, not only will we demonstrate the effectiveness of CNNs on existing emotion classification baselines but we will also qualitatively show that the network is able to learn patterns in the face images that correspond to Facial Action Units (FAUs).

## 3. Our Approach

### 3.1. Network Architecture

For all of the experiments we present in this paper, we use a classic feed-forward convolutional neural network. The networks we use, shown visually in Figure 2 consist of three convolutional layers with 64, 128, and 256 filters, respectively, and with filter sizes of 5x5 followed by ReLU (Rectified Linear Unit) activation functions. Max pooling layers are placed after the first two convolutional layers while quadrant pooling [3] is applied after the third. The quadrant pooling layer is then followed by a full-connected layer with 300 hidden units and, finally, a softmax layer for classification. The softmax layer contains anywhere between 6-8 outputs corresponding to the number of expressions present in the training set.

One modification that we introduce to the classical configuration is that we ignore the biases of the convolutional layers. This idea was introduced first by Memisevic et al. in [20] for fully-connected networks and later extended by Paine et al. in [21] to convolutional layers. In our experiments, we found that ignoring the bias allowed our network to train very quickly while simultaneously reducing the number of parameters to learn.

### 3.2. Network Training

When training our network, we train from scratch using stochastic gradient descent with a batch size of 64, mo-

mentum set to 0.9, and a weight decay parameter of 1e-5. We use a constant learning rate of 0.01 and do not use any form of annealing. The parameters of each layer are randomly initialized by drawing from a Gaussian distribution with zero mean and standard deviation $\sigma = \frac{k}{N_{\text{FAN\_IN}}}$ where $N_{\text{FAN\_IN}}$ is the number of input connections to each layer and k is drawn uniformly from the range: $[0.2, 1.2]$.

We also use dropout and various forms of data augmentation to regularize our network and combat overfitting. We apply dropout to the fully-connected layer with a probability of 0.5 (i.e. each neuron's output is set to zero with probability 0.5). For data augmentation, we apply a random transformation to each input image consisting of: translations, horizontal flips, rotations, scaling, and pixel intensity augmentation. All of our models were trained using the anna software library [1].

## 4. Experiments and Analysis

We use two facial expression datasets in our experiments: the extended Cohn-Kanade database (CK+) [18] and the Toronto Face Dataset (TFD) [26]. The CK+ database contains 327 image sequences, each of which is assigned one of 7 expression labels: anger, contempt, disgust, fear, happy, sad, and surprise. For fair comparison, we follow the protocol used by previous works [15, 17], and use the first frame of each sequence as a neutral frame in addition to the last three expressive frames to form our dataset. This leads to a total of 1308 images and 8 classes total. We then split the frames into 10 subject independent subsets in the manner presented by [15] and perform 10-fold cross-validation.

TFD is an amalgamation of several facial expression datasets. It contains 4178 images annotated with one of 7 expression labels: anger, disgust, fear, happy, neutral, sad, and surprise. The labeled samples are divided into 5 folds,

---

[1] https://github.com/ifp-uiuc/anna

Table 1. Recognition accuracy on the Toronto Face Dataset (TFD) - 7 classes - A: Data Augmentation, D: Dropout

| Method | Accuracy |
|---|---|
| Gabor+PCA [4] | 80.2% |
| Deep mPoT [22] | 82.4% |
| CDA [23] | 85.0% |
| Zero-bias CNN | 79.0% $\pm$ 1.1% |
| Zero-bias CNN+D | 81.8% $\pm$ 2.1% |
| Zero-bias CNN+A | 88.4% $\pm$ 1.7% |
| **Zero-bias CNN+AD** | **88.6% $\pm$ 1.5%** |

Table 2. Recognition accuracy on the Extended Cohn-Kanade (CK+) Dataset - 8 classes - A: Data Augmentation, D: Dropout

| Method | Accuracy |
|---|---|
| AURF [15] | 92.22% |
| AUDN [16] | 93.70% |
| Zero-bias CNN | 78.2% $\pm$ 5.7% |
| Zero-bias CNN+D | 82.3% $\pm$ 4.0% |
| Zero-bias CNN+A | 94.6% $\pm$ 3.3% |
| **Zero-bias CNN+AD** | **95.1% $\pm$ 3.1%** |

Table 3. Recognition accuracy on the Extended Cohn-Kanade (CK+) Dataset - 6 classes - A: Data Augmentation, D: Dropout

| Method | Accuracy |
|---|---|
| CSPL [34] | 89.89% |
| LBPSVM [24] | 95.10% |
| **BDBN [17]** | **96.70%** |
| Zero-bias CNN+AD | 95.7% $\pm$ 2.5% |

each containing a train, validation, and test set. We train all of our models using just the training set of each fold, pick the best performing model using each split's validation set, then we evaluate on each split's test set and average the results over all 5 folds.

In both datasets, the images are grayscale and are of size 96x96 pixels. In the case of TFD, the faces have already been detected and normalized such that all of the subjects' eyes are the same distance apart and have the same vertical coordinates. Meanwhile for the CK+ dataset, we simply detect the face in the 640x480 image and resize it to 96x96. The only other pre-processing we employ is patch-wise mean subtraction and scaling to unit variance.

## 4.1. Performance on Toronto Face Database (TFD)

First, we analyze the discriminative ability of the CNN by assessing its performance on the TFD dataset. Table 1 shows the recognition accuracy obtained when training a zero-bias CNN from a random initialization with no other regularization as well as CNNs that have dropout (D), data augmentation (A) or both (AD). We also include recognition accuracies from previous methods. From the results in Table 1, there are two main observations: (i) not surprisingly, regularization significantly boosts performance (ii) data augmentation improves performance over the regular CNN more than dropout (9.4% vs. 2.8%). Furthermore, when both dropout and data augmentation are used, our model is able to exceed the previous state-of-the-art performance on TFD by 3.6%.

## 4.2. Performance on the Extended Cohn-Kanade Dataset (CK+)

We now present our results on the CK+ dataset. The CK+ dataset usually contains eight labels (anger, contempt, disgust, fear, happy, neutral, sad, and surprise). However, many works [34, 24, 17] ignore the samples labeled as neutral or contempt, and only evaluate on the six basic emotions. Therefore, to ensure fair comparison, we trained two separate models. We present the eight class model results in Table 2 and the six class model results in Table 3. For

the eight class model, we conduct the same study we did on the TFD and we observe rather similar results. Once again, regularization appears to play a significant role in obtaining good performance. Data augmentation gives a significant boost in performance (16.4%) and when combined with dropout, leads to a 16.9% increase. For the eight class and six class models, we achieve state-of-the-art and near state-of-the-art accuracy respectively on the CK+ dataset.

## 4.3. Visualization of higher-level neurons

Now, with a strong discriminative model in hand, we will analyze which facial regions the neural network identifies as the most discriminative when performing classification. To do this, we employ the visualization technique presented by Zeiler and Fergus in [32].

For each dataset, we consider the third convolutional layer and for each filter, we find the N images in the chosen split's training set that generated the strongest magnitude response. We then leave the strongest neuron high and set all other activations to zero and use the deconvolutional network to reconstruct the region in pixel space. For our experiments, we chose N=10 training images.

We further refine our reconstructions by employing a technique called "Guided Backpropagation" proposed by Springenberg et al. in [25]. "Guided Backpropagation" aims to improve the reconstructed spatial patterns by not solely relying on the masked activations given by the top-level signal during deconvolution but by also incorporating knowledge of which activations were suppressed during the forward pass. Therefore, each layer's output during the deconvolution stage is masked twice: (i) once by the ReLU of

Table 4. Correspondences between CK+ visualization plots shown in Figure 4 and the FAU whose activation distribution had the highest KL divergence value. The KL divergence values of all the FAUs computed for each filter are shown in Figure 5.

| Filter Number | FAU with the Largest KL Divergence Value |
|:---:|:---:|
| 1 | AU25: Lips Part |
| 2 | AU12: Lip Corner Puller |
| 3 | AU9: Nose Wrinkler |
| 4 | AU5: Upper Lid Raiser |
| 5 | AU17: Chin Raiser |
| 6 | AU12: Lip Corner Puller |
| 7 | AU24: Lip Pressor |
| 8 | AU27: Mouth Stretch |
| 9 | AU12: Lip Corner Puller |
| 10 | AU1: Inner Brow Raiser |

the deconvotional layer and (ii) again by the mask generated by the ReLU of the layer's matching convolutional layer in the forward pass.

First, we will analyze patterns discovered in the Toronto Face Dataset (TFD). In Figure 3, we select 10 of the 256 filters in the third convolutional layer and for each filter, we present the spatial patterns of the top-10 images in the training set. From these images, the reader can see that several of the filters appear to be sensitive to regions that align with several of the Facial Actions Units such as: AU12: Lip Corner Puller (row 1), AU4: Brow Lowerer (row 4), and AU15: Lip Corner Depressor (row 9).

Next, we display the patterns discovered in the CK+ dataset. In Figure 4, we, once again, select 10 of the 256 filters in the third convolutional layer and for each filter, we present the spatial patterns of the top-10 images in the training set. The reader will notice that the CK+ discriminative spatial patterns are very clearly defined and correspond nicely with Facial Action Units such as: AU12: Lip Corner Puller (rows 2, 6, and 9), AU9: Nose Wrinkler (row 3) and AU27: Mouth Stretch (row 8).

## 4.4. Finding Correspondences Between Filter Activations and the Ground Truth Facial Action Units (FAUs)

In addition to categorical labels (anger, disgust, etc.), the CK+ dataset also contains labels that denote which FAUs are present in each image sequence. Using these labels, we now present a preliminary experiment to verify that the filter activations/spatial patterns learned by the CNN indeed match with the actual FAUs shown by the subject in the image. Our experiment aims to answer the following question: For a particular filter i, which FAU j has samples whose activation values most strongly differ from the activations of

samples that do not contain FAU j, and does that FAU accurately correspond with the visual spatial patterns that maximally excite filter i?

Given a training set of M images ($X$) and their corresponding FAU labels ($Y$), let $F_{\ell i}(x)$ be the activations of sample x at layer $\ell$ for filter $i$. Since we are examining the 3rd convolutional layer in the network, we set $\ell = 3$. Then, for each of the 10 filters visualized in Figure 4, we do the following:

(i) We consider a particular FAU j and place the samples $X$ that contain j in set S where:
$$S = \{x_m \mid j \in y_m\}, \ \forall m \in \{1, ..., M\}$$

(ii) We then build a histogram of the maximum activations of the samples that contained FAU j:
$$Q_{ij}(x) = P(F_{3i}(x) \mid S), \ \forall (x, y) \in (X, Y)$$

(iii) We then, similarly, build a distribution over maximum activations of the samples that do not contain FAU j:
$$R_{ij}(x) = P(F_{3i}(x) \mid S^c), \ \forall (x, y) \in (X, Y)$$

(iv) We compute the KL divergence between $Q_{ij}(x)$ and $R_{ij}(x)$, $D_{KL}(Q_{ij} \parallel R_{ij})$, and repeat the process for all of the other FAUs.

Figure 5 shows the bar charts of the KL divergences computed for all of the FAUs for each of the 10 filters displayed in Figure 4. The FAU with the largest KL divergence value is denoted in red and its corresponding name is documented in Table 4 for each filter. From these results, we see that in the majority of the cases, the FAUs listed in Table 4 match the facial regions visualized in Figure 4. This means that the samples that appear to strongly influence the activations of these particular filters are indeed those that possess the AU shown in the corresponding filter visualizations. Thus, we show that certain neurons in the neural network implicitly learn to detect specific FAUs in face images when given a relatively "loose" supervisory signal (i.e. emotion type: anger, happy, sad, etc.).

What is most encouraging is that these results appear to confirm our intuitions about how CNNs work as appearance-based classifiers. For instance, filter 2, 6, and 9 appear to be very sensitive to patterns that correspond to AU 12. This is not surprising as AU 12 (Lip Corner Puller) is almost always associated with smiles and from the visualizations in Figure 4, a subject often shows their teeth when smiling, a highly distinctive appearance cue. Similarly, for filter 8, it is not surprising that FAU 25 (Lips Part) and FAU 27 (Mouth Stretch) had the most different activation distributions given that the filter's spatial patterns corresponded to the "O" shape made by the mouth region in surprised faces, another visually salient cue.
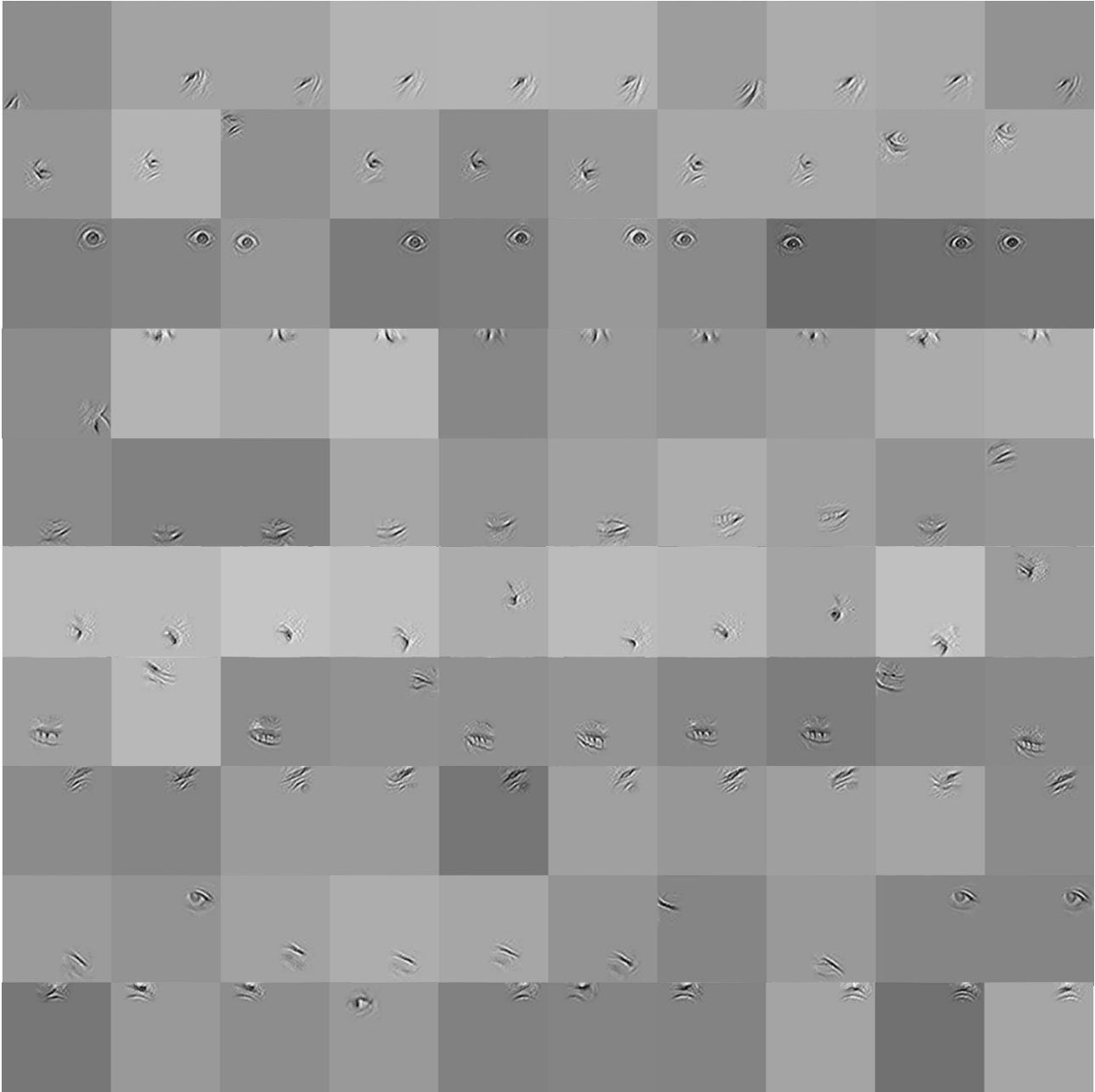
Figure 3. Visualization of spatial patterns that activate 10 selected filters in the conv3 layer of our network trained on the Toronto Face Dataset (TFD). Each row corresponds to one filter in the conv3 layer. We display the top 10 images that elicited the maximum magnitude response. Notice that the spatial patterns appear to correspond with some of the Facial Action Units.

## 5. Conclusions

In this work, we showed both qualitatively and quantitatively that CNNs trained to do emotion recognition are indeed able to model high-level features that strongly correspond to FAUs. Qualitatively, we showed which portions of the face yielded the most discriminative information by visualizing the spatial patterns that maximally excited different filters in the convolutional layers of our learned networks. Meanwhile, quantitatively, we correlated the numerical activations of the visualized filters with the subject's actual facial movements using the FAU labels given in the CK+ dataset. Finally, we demonstrated how a zero-bias CNN can achieve state-of-the-art recognition accuracy on the extended Cohn-Kanade (CK+) dataset and the Toronto Face Dataset (TFD).
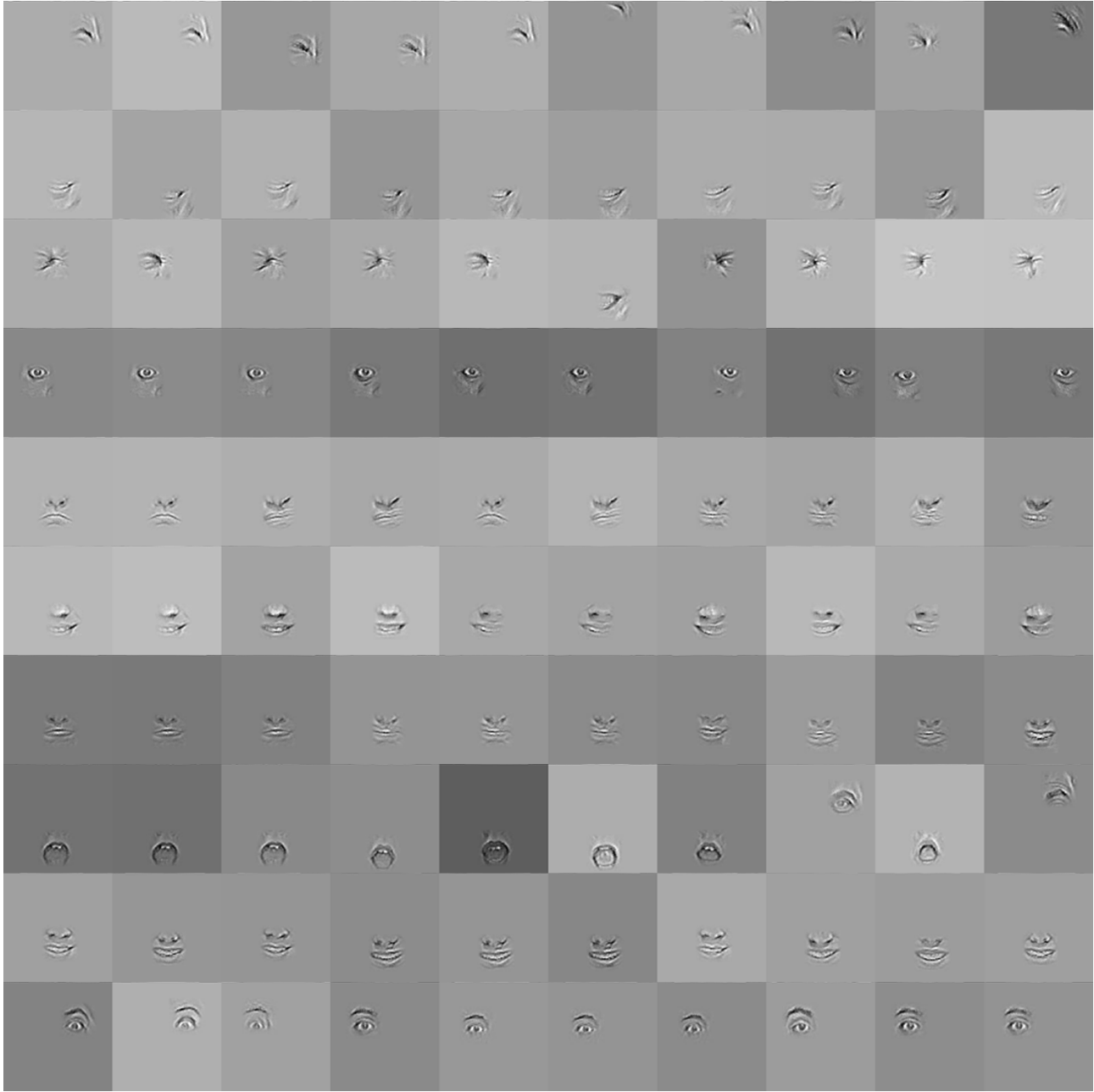
Figure 4. Visualization of spatial patterns that activate 10 selected filters in the conv3 layer of our network trained on the Cohn-Kanade (CK+) dataset. Each row corresponds to one filter in the conv3 layer. Once again, we display the top 10 images that elicited the maximum magnitude response. Notice that the spatial patterns appear to have very clear correspondences with some of the Facial Action Units.

## Acknowledgments

## References

[1] M. S. Bartlett, G. Littlewort, M. Frank, C. Lainscsek, I. Fasel, and J. Movellan. Recognizing facial expression: machine learning and application to spontaneous behavior. In *CVPR*, pages 568–573, 2005. 1, 2

[2] M. S. Bartlett, G. Littlewort, M. Frank, C. Lainscsek, I. Fasel, and J. Movellan. Fully automatic facial action recognition in spontaneous behavior. In *FGR*, pages 223–230,
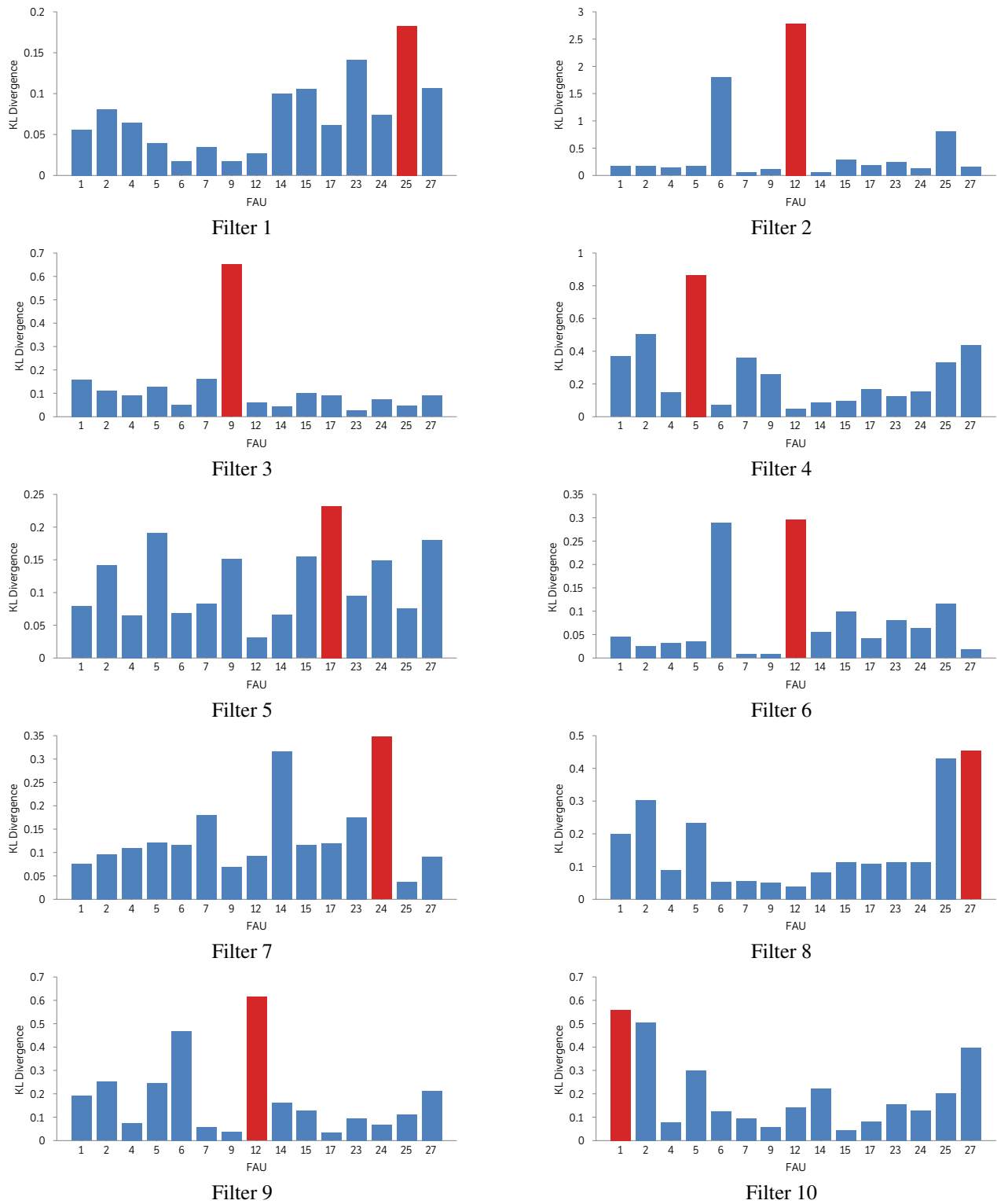
Figure 5. Bar charts showing which FAUs lead to the strongest shifts in the activation distributions of particular filters in the CNN. For each of the 10 filters visualized in Figure 4, we build histograms over the activations of training samples that contain a specific FAU j, and the activations of samples that do not contain FAU j. We then compute the KL divergence between the two distributions and plot them for each FAU above. The FAU with the largest KL divergence is displayed in red and its corresponding name is given in Table 4. (Best viewed in color).

2006. 1, 2

[3] A. Coates, A. Y. Ng, and H. Lee. An analysis of single-layer networks in unsupervised feature learning. In *International conference on artificial intelligence and statistics*, pages 215–223, 2011. 3

[4] M. N. Dailey, G. W. Cottrell, C. Padgett, and R. Adolphs. Empath: A neural network that categorizes facial expressions. *Journal of cognitive neuroscience*, 14(8):1158–1173, 2002. 4

[5] A. Dhall, R. Goecke, J. Joshi, K. Sikka, and T. Gedeon. Emotion recognition in the wild challenge 2014: Baseline, data and protocol. In *16th ACM International Conference on Multimodal Interaction*. ACM, 2014. 2

[6] A. Dhall, R. Goecke, J. Joshi, M. Wagner, and T. Gedeon. Emotion recognition in the wild challenge 2013. In *Proceedings of the 15th ACM on International conference on multimodal interaction*, pages 509–516. ACM, 2013. 2

[7] P. Ekman and W. V. Friesen. Facial action coding system. 1977. 1, 2

[8] W. V. Friesen and P. Ekman. Emfacs-7: Emotional facial action coding system. *Unpublished manuscript, University of California at San Francisco*, 2:36, 1983. 1

[9] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR*, pages 580–587, 2014. 2

[10] R. Gross, I. Matthews, J. Cohn, T. Kanade, and S. Baker. Multi-pie. *Image and Vision Computing*, 28(5):807–813, 2010. 2

[11] H. Jung, S. Lee, S. Park, I. Lee, C. Ahn, and J. Kim. Deep temporal appearance-geometry network for facial expression recognition. *arXiv preprint arXiv:1503.01532*, 2015. 2

[12] S. E. Kahou, X. Bouthillier, P. Lamblin, C. Gulcehre, V. Michalski, K. Konda, S. Jean, P. Froumenty, A. Courville, P. Vincent, et al. Emonets: Multimodal deep learning approaches for emotion recognition in video. *arXiv preprint arXiv:1503.01800*, 2015. 2

[13] 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 *ICMI*, pages 543–550, 2013. 2

[14] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, pages 1097–1105, 2012. 2

[15] M. Liu, S. Li, S. Shan, and X. Chen. Au-aware deep networks for facial expression recognition. In *FG*, pages 1–6, 2013. 2, 3, 4

[16] M. Liu, S. Li, S. Shan, and X. Chen. Au-inspired deep networks for facial expression feature learning. *Neurocomputing*, 159:126–136, 2015. 2, 4

[17] P. Liu, S. Han, Z. Meng, and Y. Tong. Facial expression recognition via a boosted deep belief network. In *CVPR*, pages 1805–1812, 2014. 2, 3, 4

[18] P. Lucey, J. F. Cohn, T. Kanade, J. Saragih, Z. Ambadar, and I. Matthews. The extended cohn-kanade dataset (ck+): A complete dataset for action unit and emotion-specified expression. In *CVPRW*, pages 94–101, 2010. 2, 3

[19] M. J. Lyons, J. Budynek, and S. Akamatsu. Automatic classification of single facial images. *PAMI*, 21(12):1357–1362, 1999. 2

[20] R. Memisevic, K. Konda, and D. Krueger. Zero-bias autoencoders and the benefits of co-adapting features. *stat*, 1050:10, 2014. 2, 3

[21] T. L. Paine, P. Khorrami, W. Han, and T. S. Huang. An analysis of unsupervised pre-training in light of recent advances. *arXiv preprint arXiv:1412.6597*, 2014. 2, 3

[22] M. Ranzato, J. Susskind, V. Mnih, and G. Hinton. On deep generative models with applications to recognition. In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pages 2857–2864. IEEE, 2011. 4

[23] S. Rifai, Y. Bengio, A. Courville, P. Vincent, and M. Mirza. Disentangling factors of variation for facial expression recognition. *ECCV 2012*, pages 808–822, 2012. 2, 4

[24] 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. 4

[25] J. T. Springenberg, A. Dosovitskiy, T. Brox, and M. Riedmiller. Striving for simplicity: The all convolutional net. *arXiv preprint arXiv:1412.6806*, 2014. 2, 4

[26] J. M. Susskind, A. K. Anderson, and G. E. Hinton. The toronto face database. *Department of Computer Science, University of Toronto, Toronto, ON, Canada, Tech. Rep*, 2010. 3

[27] J. M. Susskind, A. K. Anderson, G. E. Hinton, and J. R. Movellan. Generating facial expressions with deep belief nets. 2008. 2

[28] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf. Deepface: Closing the gap to human-level performance in face verification. In *CVPR*, pages 1701–1708, 2014. 2

[29] Y.-l. Tian, T. Kanada, and J. F. Cohn. Recognizing upper face action units for facial expression analysis. In *Computer Vision and Pattern Recognition, 2000. Proceedings. IEEE Conference on*, volume 1, pages 294–301. IEEE, 2000. 1

[30] Y. Tong, W. Liao, and Q. Ji. Facial action unit recognition by exploiting their dynamic and semantic relationships. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 29(10):1683–1699, 2007. 1

[31] J. Whitehill and C. W. Omlin. Haar features for facs au recognition. In *FGR*, 2006. 1, 2

[32] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *Computer Vision–ECCV 2014*, pages 818–833. Springer, 2014. 2, 4

[33] G. Zhao and M. Pietikainen. Dynamic texture recognition using local binary patterns with an application to facial expressions. *PAMI*, 29(6):915–928, 2007. 1, 2

[34] L. Zhong, Q. Liu, P. Yang, B. Liu, J. Huang, and D. N. Metaxas. Learning active facial patches for expression analysis. In *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pages 2562–2569. IEEE, 2012. 4