# MicroExpNet: An Extremely Small and Fast Model For Expression Recognition From Face Images

Ilke Cugu, Eren Sener, Emre Akbas

*Department of Computer Engineering*, *Middle East Technical University*, Ankara, Turkey

{cugu.ilke,sener.eren,eakbas}@metu.edu.tr

*Abstract*—This paper is aimed at creating extremely small and fast convolutional neural networks (CNN) for the problem of facial expression recognition (FER) from frontal face images. To this end, we employed the popular knowledge distillation (KD) method and identified two major shortcomings with its use: 1) a fine-grained grid search is needed for tuning the temperature hyperparameter and 2) to find the optimal size-accuracy balance, one needs to search for the final network size (or the compression rate). On the other hand, KD is proved to be useful for model compression for the FER problem, and we discovered that its effects get more and more significant with decreasing model size. In addition, we hypothesized that translation invariance achieved using max-pooling layers would not be useful for the FER problem as the expressions are sensitive to small, pixel-wise changes around the eye and the mouth. However, we have found an intriguing improvement in generalization when max-pooling is used. We conducted experiments on two widely-used FER datasets, CK+ and Oulu-CASIA. Our smallest model (MicroExpNet), obtained using knowledge distillation, is less than 1MB in size and works at $1851$ frames per second on an Intel i7 CPU. Despite being less accurate than the state-of-the-art, MicroExpNet still provides significant insights for designing a microarchitecture for the FER problem.

## I. INTRODUCTION

Expression recognition from frontal face images is an important aspect of human-computer interaction and has many potential applications, especially in mobile devices. Face detection models have long been deployed in mobile devices, and relatively recently, face recognition models are also being used, e.g. for face-based authentication. Arguably, one of the next steps is the mobile deployment of facial expression recognition models. Therefore, creating small and fast models is an important goal. In order to asses the current situation, we looked at the size and runtime speeds of two representatives, current state-of-the-art models, namely PPDN [1] and FN2EN [2]. In terms of the number of total parameters in the network, both models are in the order of millions (PPDN has 6M and FN2EN has 11M). In terms of speed, both models run at $9-11$ms per image on a GTX 1050 GPU, however, on an Intel i7 CPU, while PPDN takes 57.18 ms, FN2EN takes 96.08 ms (further details in Table V).

The central question that motivated the present work was how much we could push the size and speed limits so that we end up with a compact expression recognition model that still works reasonably well. To this end, we focused only on frontal face images and first explored training a large model on two widely used benchmarks FER datasets, CK+ [3] and Oulu-Casia [4], by using the Inception_v3 [5] model. Then, using the "knowledge distillation" (KD) method [6], we created a family of small and fast models. In the KD method, there is a large, cumbersome model called the *teacher* (Inception_v3 in our case) and a relatively much smaller model called the *student*. The student is trained to "mimic" the softmax values of the teacher via a *temperature* hyperparameter (see Eq. 2). We have experimented on four student networks with different sizes. The smallest one, called **MicroExpNet**, is 100x smaller in size and has 335x fewer parameters compared to the teacher.

We found two major shortcomings of the KD method. First, the temperature hyperparameter does not seem to have any meaningful relationship with the accuracy of the student model. We found that the accuracy fluctuates between low and high values as the temperature is swept across a wide range. In order to find a high-accuracy temperature, one needs to do a fine-grained grid search. Second, in the KD method, the final student model size (i.e. the compression rate) is given as input. Therefore, to find the optimal size-accuracy balance, one needs to search for the size, too.

We also hypothesized that invariance to translation achieved using max-pooling layers would not be useful for the FER problem as the expressions are sensitive to small, pixel-wise changes around the eye and the mouth. However, we empirically found that this is **not** the case. The best results are obtained when there is a max-pooling layer after each convolutional layer. Another important and related finding of our work is that the effect of max-pooling is reversed depending on how the dataset is split into training and testing sets. When the dataset is split into train and test sets in such a way that a human subject can appear only in one of them (this is called **subject-independent** split in the literature), max-pooling has a positive effect on performance. However, when the dataset is split purely randomly; that is, images from the same subject may appear in both train and test sets, max-pooling has a negative effect. Although each image in the dataset is numerically different, **random** split transform the FER problem to a memorization problem. We validate this proposition with our empirical analysis (Table II). Our findings raise three important questions:

**(1) Is information loss essential for generalization?** We show that, considering it as a tool for information loss where numerically dominant values suppress the others, max-pooling

improves the classification performance. However, when the problem is transformed into a memorization challenge, having no max-pooling layers yields the best results.

**(2) Is a smaller model more open to teacher's supervision?** We show that, compared to training from scratch, KD becomes more effective as the network gets smaller. To the best of our knowledge, this is a novel finding Whether this effect is specific to the FER problem is yet to be seen.

**(3) Why does the classification accuracy fluctuate as the *temperature* hyperparameter is changed?** We show that, regardless of how the dataset is split (random or subject-independent), the accuracy fluctuates between low and high values as the temperature is swept across a wide range.

## II. RELATED WORK

### A. Facial expression recognition (FER)

We categorize the previous work as image (or frame) based and sequence-based. While image-based methods analyze individual images independently, sequence-based methods exploit the spatio-temporal information between frames.

*a) Image based.:* There are three groups of work. Models that use 1) hand-crafted features (HCFs), 2) deep representations, and 3) both. Our work falls into the second group.

We do not focus on HCF models [7]–[10] here because they are obsolete (with the emergence of deep models) and in general, they do not achieve competitive results.

Deep representations learned from face images are the main ingredients of [1], [2], [11]–[13]. Liu et al. [11] proposed a loopy boosted deep belief network framework for feature learning, then used them in an AdaBoost classifier. Mollahosseini et al. [12] introduced an inception network for FER. Their model is much larger compared to ours considering the two large fully connected layers at the end of their network. Zhao et al. [1] proposed a peak-piloted GoogLeNet [14] model which uses both peak and non-peak expression images during training. Training peak and non-peak images in pairs naturally requires their proposed back-propagation algorithm which adds complexity to implementation compared to our work. FN2EN [2] employs a multi-staged model production for FER. First, they train convolutional layers by mimicking [15] a pre-trained FaceNet [16]. Then, they append a $fc$ layer to the model for retraining. Recently, Kim et al. [13] introduced a deep generative contrastive model for FER. They combined encoder-decoder networks and CNNs into a unified network that simultaneously learns to generate, compare, and classify samples on a dataset.

Finally, [17] form a hybrid approach. They train CNNs with both the original input images and 3D mappings of local binary patterns [18], then finalize via fine-tuning.

*b) Sequence-based.:* We can categorize sequence-based facial expression classifiers in the same three groups as in the case of image-based classifiers. We do not focus on HCF based sequence models [19]–[22] for the same reasons with the image-based case.

Liu et al. [23] proposed manifold modeling of videos based on representations gathered via learned spatio-temporal filters.

Kahou et al. [24] fused CNNs with recurrent neural networks (RNNs). CNN is used on static images to gather high-level representations which are then used by the RNN training. Jung et al. [25] proposed a hybrid approach via two deep models. First, a 3D-CNN to extract the temporal appearance features from image sequences. Second, a fully connected model which captures geometrical information about the motion of the facial landmark points.

### B. Model size reduction

The knowledge distillation (KD) method, which is the first and the most popular teacher-student style compression method, [6] is described in Section III. Ba and Caruana [15] proposed using the $L_2$ loss on the logits to mimic the teacher (without any temperature hyperparameter). Both ideas are combined in FitNets [26]: they use the $L_2$ loss on logits in the pre-training phase for better initialization, then, they train the whole student network using the KD method. In contrast to FitNets, we choose a model that is much shallower than the teacher and avoid any pre-training of the student to prevent increasing the complexity of the overall training procedure. Iandola et al. [27] (SqueezeNet) proposed a CNN with no fully connected layers to reduce the model size, and preserved the classification performance via their fire modules.

## III. METHODOLOGY

### A. Knowledge distillation

Formally, let $p_t$ be the softened output of the teacher's softmax, $z_i$ be the logits of the teacher, $p_s$ be the hard and $p'_s$ be the soft output of the student's softmax, $v_i$ be the logits of the student, $\lambda$ be the weight of distillation, $y$ be the ground truth labels, $N$ be the batch size, $T$ temperature and function $\mathcal{H}$ refers to the cross-entropy. Then:

$$p_t = \frac{e^{z_i/T}}{\sum_j e^{z_j/T}}, \quad p'_s = \frac{e^{v_i/T}}{\sum_j e^{v_j/T}}, \quad p_s = \frac{e^{v_i}}{\sum_j e^{v_j}} \quad (1)$$

and the loss becomes

$$\mathcal{L} = \lambda\left(\frac{1}{N}\sum_{n=1}^{N}\mathcal{H}(p_t, p'_s)\right) + (1-\lambda)\left(\frac{1}{N}\sum_{n=1}^{N}\mathcal{H}(y, p_s)\right). \quad (2)$$

### B. Network architectures

*a) Teacher network:* We use the Inception_v3 [5] network as the teacher for its proven record of success on classification tasks [28].

*b) Student network:* Our student network has a very simple architecture: two convolutional layers (conv1, conv2) and two fully-connected layers (fc1, fc2). We used rectified linear units (ReLU) [29] as activation functions. There are max-pooling layers after each conv layer (analysis on pooling vs. no-pooling is given in Section IV-A). We created four different versions of the student network (namely M, S, XS, XXS) by varying the number of neurons at fc1. Table I presents the sizes of these networks. We trained them using the KD [6] method. We compare their classification performances in sections IV-B and IV-C, speeds & memory requirements in Section IV-E.

| Model | # of neurons in $fc1$ | total # of parameters |
|---|---|---|
| M | 256 | 900920 |
| S | 64 | 232184 |
| XS | 32 | 120728 |
| XXS | 16 | 65000 |

| | Model | CK+ | Oulu-CASIA | Model | CK+ | Oulu-CASIA |
|---|---|---|---|---|---|---|
| **Random** | $v_M$ | 97.93% | 97.68% | $v_{XS}$ | **93.41%** | **88.73%** |
| | $p1_M$ | **97.99%** | **97.79%** | $p1_{XS}$ | 91.85% | 80.16% |
| | $p2_M$ | 97.41% | 96.64% | $p2_{XS}$ | 86.84% | 77.88% |
| | $p12_M$ | 97.39% | 97.47% | $p12_{XS}$ | 88.07% | 77.04% |
| | $v_S$ | 96.65% | 92.95% | $v_{XXS}$ | **81.91%** | **73.64%** |
| | $p1_S$ | **96.73%** | **93.22%** | $p1_{XXS}$ | 69.05% | 52.99% |
| | $p2_S$ | 94.09% | 88.61% | $p2_{XXS}$ | 77.74% | 66.84% |
| | $p12_S$ | 94.39% | 88.72% | $p12_{XXS}$ | 78.52% | 61.71% |
| **Subject-independent** | $v_M$ | 81.23% | 60.87% | $v_{XS}$ | 77.14% | 53.73% |
| | $p1_M$ | **81.57%** | **62.46%** | $p1_{XS}$ | 77.14% | 53.41% |
| | $p2_M$ | 78.77% | 60.21% | $p2_{XS}$ | 78.42% | 57.51% |
| | $p12_M$ | 79.95% | 60.53% | $p12_{XS}$ | **79.78%** | **57.54%** |
| | $v_S$ | 79.73% | 58.18% | $v_{XXS}$ | 71.36% | 44.33% |
| | $p1_S$ | **81.25%** | **59.49%** | $p1_{XXS}$ | 67.04% | 34.04% |
| | $p2_S$ | 78.75% | 57.37% | $p2_{XXS}$ | 76.91% | 54.62% |
| | $p12_S$ | 79.71% | 57.25% | $p12_{XXS}$ | **78.44%** | **55.03%** |

### C. Implementation

*a) CK+ & Oulu-CASIA:* For each image in CK+, we apply the Viola Jones [30] face detector, and for each image in Oulu-CASIA we use the already cropped versions. All images are converted to grayscale. Then, in order to augment the data, we extract 8 crops (4 from each corner and 4 from each side) from an image with dimensions of 84x84 for students and 256x256 for the teacher. There is no difference on hyperparameter selections for the trainings on CK+ and Oulu-CASIA. As done in previous work, we report the average 10-fold cross-validation (CV) performance. For both the teacher and students, trainings are finalized after 3000 epochs.

*b) Teacher Network:* We employ a Inception_v3 network trained on ImageNet [28], and fine-tune it on FER datasets. The base learning rate is set as $10^{-4}$ and remained constant through iterations, mini-batch size is 64, and the learning method is Adam [31].

*c) Vanilla & Student Networks:* We have the same hyperparameters across all of the different model sizes for both vanilla and student trainings. "Vanilla" training means that the network is trained from scratch without any teacher guidance. Weights and biases are initialized using Xavier initialization [32]. Network architectures are implemented using Tensorflow [33]. We used Adam [31] with a learning rate of $10^{-4}$. The dropout [34] rate is 0.5, mini-batch is 64 and the weight of the distillation $\lambda$ is 0.5 (see Section III-A) for all student models. Selected model sizes are 900K, 232K, 121K and 65K parameters respectively, which are produced by decreasing the size of the $fc1$ layer (see Table I). Training operations are finalized after 3000 epochs for all models and the XXS student model is denoted as **MicroExpNet**. Empirical results are given in Table IV; note that for student networks we only list the best performers across different temperatures (selected using cross-validation). Furthermore, student models are used in temperature selection tests (for detailed explanation see Section IV-D). The results we report for these models are obtained by averaging the 10-fold cross-validation performances.

## IV. EXPERIMENTS

### A. Max. Pooling vs. No Pooling Analysis

Facial expressions are located mostly on eyes and mouth [35], and they form only a small fraction of a frontal face image. The idea is to capture these subtle indicators of emotion by preserving the pixel information across layers. Therefore, our starting point was a CNN with no pooling layers. However, in order to validate our intuition, we build three variations containing max-pooling layers for each student. All pooling layers have 2x2 filters with stride 2. All hyperparameters mentioned at Section III-C apply to these variations as well. We call them candidate expression networks. These candidates are explained in Table II.

From the results in Table II, we draw the following conclusions. When models are large enough, the network capacity for learning dominates pooling effects. For instance, for the size M, classification performances of candidates are very close to each other. For size S, poolings in later layers drop the performance but early pooling is still the most profitable. After this point (XS and XXS), we begin to see an interesting difference between the results of random and subject-independent split experiments. **For random split, we see the advantage of not having any pooling layers with significant gains in performance.** Since the trained candidates see the same subjects in both training and test (for $\approx 80\%$ of the subjects), although the images are numerically different, we think that the resemblance transforms the FER problem to a memorization challenge. Hence, the information loss caused by the pooling layers drops the performance. **On the contrary, for subject-independent split where test subjects are not seen during the training, we observe the advantage of having pooling layers.** Thus, the intuition mentioned above does not seem to hold. It is also interesting to observe that the second pooling layer seems to be a much critical point of improvement than the first pooling layer. Nevertheless, combining these observations with our intention to reduce the

model size, we decided to employ the architecture with two pooling layers as the foundation of our student networks.

Note that adding a pooling layer drops the number of parameters; thus, prevents a proper performance comparison. Therefore, we did two modifications to increase the model size, in order to make it a fair comparison. First, when we add a pooling layer after the first convolutional layer, we decrease the stride of the first conv layer from 4 to 2. This directly recovers all parameters that has been lost. Second, when we add a pooling layers after the second convolutional layer, we increase the number of outputs of the first fully connected layer by 3-fold. This results in having slightly less parameters than the original one ($v$).

### B. The CK+ dataset

CK+ is a widely used benchmark database for facial expression recognition. This database is composed of 327 image sequences with eight emotion labels: anger, contempt, disgust, fear, happiness, sadness, surprise and neutral. There are 123 subjects. As done in previous work, we extract the last three and the first frames of each expression sequence when images are labeled. When unlabeled, we only extracted first frames as neutral. The total number of images is 1574 (see at Table III), which is split into 10 folds.

*a) Training in Isolation:* We evaluate the pre-trained Inception_v3 via fine-tuning on CK+. Then, we train four models, namely VanillaExpNet$_M$, VanillaExpNet$_S$, VanillaExpNet$_{XS}$, and VanillaExpNet$_{XXS}$, from scratch. At this stage, we did not employ knowledge distillation. For all models, we used 3000 epochs for training, and the classification performances are shown in Table IV. Although Zhao et al. [1] seem to achieve better performance than Inception_v3 (in Table IV), they use only 6 emotion categories, whereas we use all of the 8 emotion categories. In the light of these results, we chose Inception_v3 as the teacher for the knowledge distillation stage.

*b) Training with Supervision:* We evaluate four students, namely StudentExpNet$_M$, StudentExpNet$_S$, StudentExpNet$_{XS}$, and StudentExpNet$_{XXS}$, via knowledge distillation on CK+. At this stage, we use the teacher's supervision to improve the learning. As explained in Section III-C, we need to tune the *temperature* for each student since it is considered to be correlated with model size. Therefore, we conducted an extensive experiment on classification performances for a wide range of *temperatures*. The results are reported in Figure 1. According to these results, fluctuations between performances are increased while models are getting smaller. Consequently, it suggests that large networks are more tolerant to the changes in the temperature. This observation also holds for the random split case as shown in Figure 2.

Best performers, based on their average classification performances for 10-fold cross-validation, across different temperatures are then used for performance comparison in Table IV. Our findings show that KD can be used to gain back some of the performance lost by decreasing the model size.

### C. The Oulu-CASIA dataset

Oulu-CASIA has 480 image sequences taken under *dark, strong, weak* illumination conditions. In this experiment, as also done in previous work, we used only videos with *strong* condition captured by a VIS camera. In total, there are 80 subjects and six expressions: anger, disgust, fear, happiness, sadness, and surprise. Similar to CK+, the first frame is always neutral while the last frame has the peak expression. All studies we have encountered on Oulu-CASIA database use only the last three frames of the sequences, so we also use the same frames. Therefore, the total number of images is 1440. As in the earlier studies, a 10 fold CV is performed, and the split is subject independent.

*a) Training in isolation:* The same approach taken for CK+ is employed for Oulu-CASIA. The classification performances are shown in Table IV. According to the table, Inception_v3 performs on par with the state-of-the-art solutions whereas our vanilla models failed to achieve competitive results.

*b) Training with supervision:* The same explanations on students for CK+ also apply to Oulu-CASIA experiments. The results are reported in Figure 3 from which, we can observe a similar fluctuating behavior as seen in the CK+ experiments. Once again, we can see that large networks are more tolerant to the changes in the temperature than the smaller ones. In addition, as in CK+ experiments, this observation also holds for the random split case as shown in Figure 4.

Best performers across different temperatures are then used for performance comparison in Table IV. We can still observe that the student models perform better than vanilla models (which are trained from scratch without any teacher supervision) for facial expression recognition.

### D. Temperature analysis

Temperature is a hyperparameter to control the uncertainty in teacher's output. This uncertainty may be used as similarity information between different classes to enhance the training. However, there is no formulation for selecting the most effective temperature; it is set empirically. Hence, we did a grid search for temperatures of [2, 4, 8, 16, 20, 32, 64] with 10-fold cross-validation across all of our student networks using both CK+ (see Figure 1) and Oulu-CASIA (see Figure 3) datasets using a subject-independent train & validation split. Moreover, we did a grid search for a random train & validation split as well (see Figures 2 and 4).

According to the results, smaller models are more prone to temperature changes in general, and performances for a given temperature seem rather stochastic. However, large models show different characteristics for the random split case and subject-independent split case. When the subject-independent split is used, we observe fluctuations in performance for all models regardless of their size. Whereas for the random split case, large models have relatively stable performances. Nevertheless, when calibrated adequately, KD improves the overall FER performance for the subject-independent split case.

TABLE III
THE NUMBER OF IMAGES PER EXPRESSION CLASS IN CK+ AND OULU-CASIA.

| | Anger | Contempt | Disgust | Fear | Happy | Sad | Surprise | Neutral | All |
|---|---|---|---|---|---|---|---|---|---|
| CK+ | 135 | 54 | 177 | 75 | 207 | 84 | 249 | 593 | 1574 |
| Oulu-CASIA | 240 | - | 240 | 240 | 240 | 240 | 240 | - | 1440 |

TABLE IV
AVERAGE CLASSIFICATION ACCURICIES (OVER 10-FOLDS) OF DIFFERENT
METHODS ON CK+ AND OULU-CASIA DATASETS.

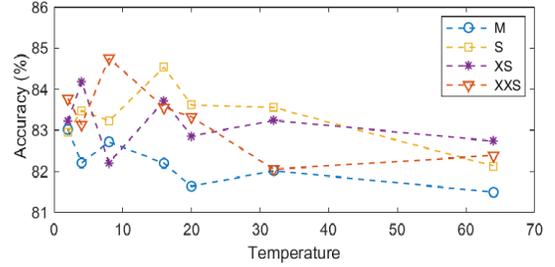| Method | CK+ | | Oulu-CASIA |
|---|---|---|---|
| | 6 cls | 8 cls | |
| CSPL [10] | 89.9% | - | - |
| 3DCNN-DAP [36] | 92.4% | - | - |
| Inception [12] | 93.2% | - | - |
| AdaGabor [7] | 93.3% | - | - |
| AdaLBP [4] | - | - | 73.54% |
| STM-ExpLet [23] | 94.2% | - | 74.59% |
| Atlases [19] | - | - | 75.52% |
| LOMo [22] | 95.1% | - | 82.10% |
| LBPSVM [8] | 95.1% | - | - |
| BDBN [11] | 96.7% | - | - |
| DTAGN [25] | 97.3% | - | 81.46% |
| FN2EN [2] | 98.6% | 96.8% | 87.71% |
| DCN [37] | 98.9% | - | |
| PPDN [1] | 99.3% | - | 84.59% |
| AU-Aware [38] | - | 92.1% | |
| GCNet [13] | - | 97.3% | 86.39% |
| **TeacherExpNet** | - | **97.6%** | **85.83%** |
| VanillaExpNet$_M$ | - | 78.8% | 56.81% |
| VanillaExpNet$_S$ | - | 78.6% | 55.53% |
| VanillaExpNet$_{XS}$ | - | 77.2% | 54.67% |
| VanillaExpNet$_{XXS}$ | - | 75.3% | 56.71% |
| StudentExpNet$_M$ | - | 83.1% | 63.81% |
| StudentExpNet$_S$ | - | 83.6% | 62.01% |
| StudentExpNet$_{XS}$ | - | 83.7% | 61.76% |
| **MicroExpNet** | - | **84.8%** | **62.69%** |



Fig. 1. Classification performances of the student networks across different temperatures on the CK+ dataset using **subject-independent splits**.
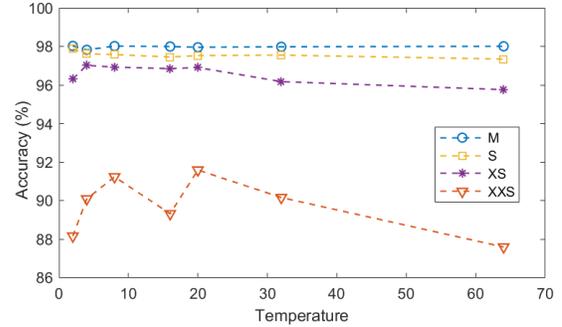


Fig. 2. Classification performances of the student networks across different temperatures on the CK+ dataset using **random splits**.

### E. Model size and speed analysis

We show the comparison of the model sizes in megabytes in Table V. Our smallest FER model MicroExpNet takes less than 1 MB to store which is 100x smaller than our teacher network (Inception_v3), and, it has 335x fewer parameters than the teacher. We also show the running times to process one image in milliseconds (average of 1000 runs) in Table V. According to the table, MicroExpNet achieves the best performance by classifying the facial expression in an image in less than 1 ms on an Intel i7 CPU. Ultimately, when compared to its teacher, MicroExpNet is 234x faster on Intel i7, and 85x faster on GTX1050.

### V. CONCLUSION

We presented an extensive analysis of the creation of a microarchitecture, called the MicroExpNet, for facial expression recognition (FER) from frontal face images.

From our experimental work, we have drawn the following conclusions. (1) Information loss achieved via max-pooling and the KD method both improve the performance especially when the network is small, (2) we showed that a simple change in the approach taken for the separation of train & validation sets results in drastic changes in the problem definition, and

thus in the performance observations, (3) KD's effect gets more prominent as the network size decreases. Whether this effect is generalizable to other problems/datasets is yet to be seen in future work. (4) The temperature hyperparameter (in KD) should be tuned carefully for optimal performance. Especially when the network is small, the final performance fluctuates with temperature.
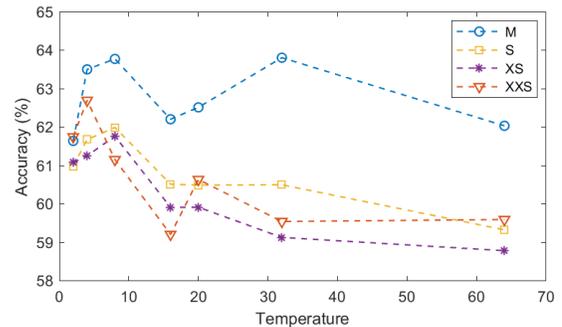


Fig. 3. Classification performances of the student networks across different temperatures on the Oulu-CASIA dataset using **subject-independent splits**.
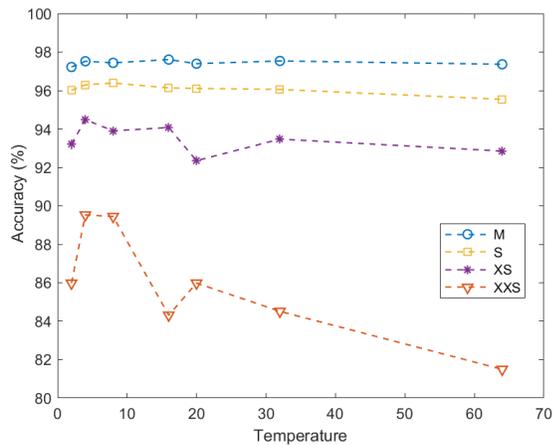
Fig. 4. Classification performances of the student networks across different temperatures on the Oulu-CASIA dataset using **random splits**.

TABLE V
MEMORY REQUIREMENTS AND AVERAGE PER-IMAGE RUNNING TIMES OF DIFFERENT FER MODELS.

| Model | # of params | Size (MB) | i7-7700HQ | GTX1050 | Tesla K40 |
|---|---|---|---|---|---|
| TeacherExpNet | 21.8M | 88.13 | 124.22 ms | 83.25 ms | - |
| FN2EN [2] | 11M | 42.42 | 96.08 ms | 23.81 ms | 13.09 ms |
| PPDN [1] | 6M | 23.93 | 57.18 ms | 9.12 ms | 13.11 ms |
| StudentExpNet$_M$ | 900K | 10.88 | 0.89 ms | 1.13 ms | 1.74 ms |
| StudentExpNet$_S$ | 232K | 2.91 | 0.78 ms | 1.08 ms | 1.69 ms |
| StudentExpNet$_{XS}$ | 121K | 1.52 | 0.63 ms | 0.97 ms | 1.63 ms |
| **MicroExpNet** | **65K** | **0.88** | **0.53 ms** | **0.97 ms** | **1.52 ms** |

## REFERENCES

[1] X. Zhao, X. Liang, L. Liu, T. Li, Y. Han, N. Vasconcelos, and S. Yan, "Peak-piloted deep network for facial expression recognition," in *ECCV*, 2016, pp. 425–442.

[2] H. Ding, S. K. Zhou, and R. Chellappa, "Facenet2expnet: Regularizing a deep face recognition net for expression recognition," in *FG*, 2017, pp. 118–126.

[3] 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 *CVPR Workshops*. IEEE, 2010, pp. 94–101.

[4] G. Zhao, X. Huang, M. Taini, S. Z. Li, and M. Pietikäinen, "Facial expression recognition from near-infrared videos," *Image and Vision Computing*, vol. 29, no. 9, pp. 607–619, 2011.

[5] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *CVPR*, 2016.

[6] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *arXiv:1503.02531*, 2015.

[7] 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*, 2005.

[8] X. Feng, M. Pietikäinen, and A. Hadid, "Facial expression recognition based on local binary patterns," *Pattern Recognition and Image Analysis*, vol. 17, no. 4, pp. 592–598, 2007.

[9] K. Sikka, T. Wu, J. Susskind, and M. Bartlett, "Exploring bag of words architectures in the facial expression domain," in *ECCV*, 2012.

[10] L. Zhong, Q. Liu, P. Yang, B. Liu, J. Huang, and D. N. Metaxas, "Learning active facial patches for expression analysis," in *CVPR*, 2012.

[11] P. Liu, S. Han, Z. Meng, and Y. Tong, "Facial expression recognition via a boosted deep belief network," in *CVPR*, 2014.

[12] A. Mollahosseini, D. Chan, and M. H. Mahoor, "Going deeper in facial expression recognition using deep neural networks," in *WACV*, 2016.

[13] Y. Kim, B. Yoo, Y. Kwak, C. Choi, and J. Kim, "Deep generative-contrastive networks for facial expression recognition," *arXiv:1703.07140*, 2017.

[14] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *CVPR*, 2015, pp. 1–9.

[15] J. Ba and R. Caruana, "Do deep nets really need to be deep?" in *NeurIPS*, 2014.

[16] O. M. Parkhi, A. Vedaldi, A. Zisserman *et al.*, "Deep face recognition." in *BMVC*, vol. 1, no. 3, 2015, p. 6.

[17] G. Levi and T. Hassner, "Emotion recognition in the wild via convolutional neural networks and mapped binary patterns," in *ICMI*. ACM, 2015, pp. 503–510.

[18] T. Ojala, M. Pietikäinen, and D. Harwood, "A comparative study of texture measures with classification based on featured distributions," *Pattern recognition*, vol. 29, no. 1, pp. 51–59, 1996.

[19] Y. Guo, G. Zhao, and M. Pietikäinen, "Dynamic facial expression recognition using longitudinal facial expression atlases," in *ECCV*. Springer, 2012, pp. 631–644.

[20] D. Ghimire and J. Lee, "Geometric feature-based facial expression recognition in image sequences using multi-class adaboost and support vector machines," *Sensors*, vol. 13, no. 6, pp. 7714–7734, 2013.

[21] A. Sanin, C. Sanderson, M. T. Harandi, and B. C. Lovell, "Spatio-temporal covariance descriptors for action and gesture recognition," in *WACV*. IEEE, 2013, pp. 103–110.

[22] K. Sikka, G. Sharma, and M. Bartlett, "Lomo: Latent ordinal model for facial analysis in videos," in *CVPR*, 2016, pp. 5580–5589.

[23] M. Liu, S. Shan, R. Wang, and X. Chen, "Learning expressionlets on spatio-temporal manifold for dynamic facial expression recognition," in *CVPR*, 2014, pp. 1749–1756.

[24] S. Ebrahimi Kahou, V. Michalski, K. Konda, R. Memisevic, and C. Pal, "Recurrent neural networks for emotion recognition in video," in *ICMI*. ACM, 2015, pp. 467–474.

[25] H. Jung, S. Lee, J. Yim, S. Park, and J. Kim, "Joint fine-tuning in deep neural networks for facial expression recognition," in *ICCV*, 2015.

[26] A. Romero, N. Ballas, S. E. Kahou, A. Chassang, C. Gatta, and Y. Bengio, "Fitnets: Hints for thin deep nets," *arXiv:1412.6550*, 2014.

[27] F. N. Iandola, S. Han, M. W. Moskewicz, K. Ashraf, W. J. Dally, and K. Keutzer, "Squeezenet: Alexnet-level accuracy with 50x fewer parameters and¡ 0.5 mb model size," *arXiv:1602.07360*, 2016.

[28] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy *et al.*, "Imagenet large scale visual recognition challenge," *IJCV*, vol. 115, no. 3, pp. 211–252, 2015.

[29] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," in *ICML*, 2010, pp. 807–814.

[30] P. Viola and M. J. Jones, "Robust real-time face detection," *IJCV*, vol. 57, no. 2, pp. 137–154, 2004.

[31] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv:1412.6980*, 2014.

[32] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks," in *AiStats*, 2010, pp. 249–256.

[33] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin *et al.*, "Tensorflow: Large-scale machine learning on heterogeneous distributed systems," *arXiv:1603.04467*, 2016.

[34] N. Srivastava, G. E. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: a simple way to prevent neural networks from overfitting." *JMLR*, vol. 15, no. 1, pp. 1929–1958, 2014.

[35] P. Ekman and E. L. Rosenberg, *What the face reveals: Basic and applied studies of spontaneous expression using the Facial Action Coding System (FACS)*. Oxford University Press, USA, 1997.

[36] M. Liu, S. Li, S. Shan, R. Wang, and X. Chen, "Deeply learning deformable facial action parts model for dynamic expression analysis," in *ACCV*. Springer, 2014, pp. 143–157.

[37] Z. Zhang, P. Luo, C. C. Loy, and X. Tang, "From facial expression recognition to interpersonal relation prediction," *arXiv:1609.06426*, 2016.

[38] M. Liu, S. Li, S. Shan, and X. Chen, "Au-aware deep networks for facial expression recognition," in *FG*. IEEE, 2013, pp. 1–6.