

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

İlke Çuğu, Eren Şener, Emre Akbaş

Department of Computer Engineering  
Middle East Technical University

IPTA 2019

## 1 Introduction

## 2 Search for a Compact FER Model

- Architecture: Max pooling vs. No pooling
- Dataset: Random split vs. Subject-independent split
- Performance: Model size and speed

## 3 Knowledge Distillation for FER

- Regularization: Model size vs. Teacher's Supervision
- Hyperparameters: Temperature Analysis

## 4 Future Research Directions

# Facial Expression Recognition

- Automatic recognition of basic emotions
- Anger, contempt, disgust, fear, happy, sadness, surprise
- Datasets:
  - CK+ <sup>1</sup>
  - Oulu-CASIA <sup>2</sup>

---

<sup>1</sup>Lucey, Patrick, et al. "The extended cohn-kanade dataset (ck+): A complete dataset for action unit and emotion-specified expression." 2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition-Workshops. IEEE, 2010.

<sup>2</sup>Zhao, Guoying, et al. "Facial expression recognition from near-infrared videos." Image and Vision Computing 29.9 (2011): 607-619.

# Model Family

Model	# of neurons in <i>fc1</i>
M	256
S	64
XS	32
XXS	16

- 2 conv layers (conv1, conv2)
- 2 fully-connected layers (fc1, fc2)
- Rectified linear units (ReLU)<sup>3</sup> as activation functions.
- Most of the parameters are at fully-connected layers
- Grid search to test size/performance trade-off

---

<sup>3</sup>Nair, Vinod, and Geoffrey E. Hinton. "Rectified linear units improve restricted boltzmann machines." Proceedings of the 27th international conference on machine learning (ICML-10). 2010.

# Max pooling vs. No pooling

- Facial expressions are located mostly on eyes and mouth <sup>4</sup>
- Hypothesis: max-pooling layers hurt the performance of a FER model as the expressions are sensitive to small, pixel-wise changes around the eye and the mouth.
  - **FAILED**
  - However, testing environment shapes the problem definition (memorization vs. learning)
  - Specifically, for train/val/test set separation:
    - If random split, **YES**
    - If subject-independent split, **NO**

---

<sup>4</sup>Ekman, Rosenberg. What the face reveals: Basic and applied studies of spontaneous expression using the Facial Action Coding System (FACS). Oxford University Press, USA, 1997.

# Notation

- $v$ : no pooling layer
- $p_1$ : only one max pooling layer after conv1
- $p_2$ : only one max pooling layer after conv2
- $p_{12}$ : each conv layer is followed by a max pooling layer

# Train/Val/Test Set Separation

	Model	CK+	Oulu-CASIA	Model	CK+	Oulu-CASIA
Random	$v_M$	97.93%	97.68%	$v_{XS}$	<b>93.41%</b>	<b>88.73%</b>
	$p1_M$	<b>97.99%</b>	<b>97.79%</b>	$p1_{XS}$	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%	$v_{XXS}$	<b>81.91%</b>	<b>73.64%</b>
	$p1_S$	<b>96.73%</b>	<b>93.22%</b>	$p1_{XXS}$	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%

Random split

- Model may see images of the same subject both in training & testing
- Images are numerically different, but visually very similar
- **No pooling gives the best result for XS and XXS models**
- **Q: Does preserving the pixel information ease memorization?**

# Train/Val/Test Set Separation

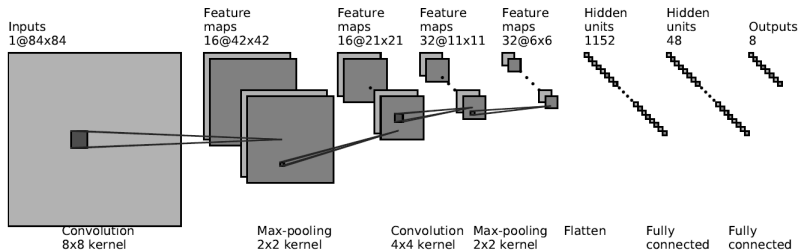
	Model	CK+	Oulu-CASIA	Model	CK+	Oulu-CASIA
Subject-independent	$v_M$	81.23%	60.87%	$v_{XS}$	77.14%	53.73%
	$p1_M$	<b>81.57%</b>	<b>62.46%</b>	$p1_{XS}$	77.14%	53.41%
	$p2_M$	78.77%	60.21%	$p2_{XS}$	78.42%	57.51%
	$p12_M$	79.95%	60.53%	$p12_{XS}$	<b>79.78%</b>	<b>57.54%</b>
Subject-independent	$v_S$	79.73%	58.18%	$v_{XXS}$	71.36%	44.33%
	$p1_S$	<b>81.25%</b>	<b>59.49%</b>	$p1_{XXS}$	67.04%	34.04%
	$p2_S$	78.75%	57.37%	$p2_{XXS}$	76.91%	54.62%
	$p12_S$	79.71%	57.25%	$p12_{XXS}$	<b>78.44%</b>	<b>55.03%</b>

## Subject-independent split

- Model trains with a set of subjects  $s_1$
- It is tested with another set of subjects  $s_2$
- where  $s_1 \cap s_2 = \emptyset$
- Having 2 max-pooling gives the best result for XS and XXS models
- Q: Does information loss improve generalization?



# MicroExpNet Architecture



The final architecture with max-pooling layers

# Model size and speed

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 <sub>M</sub>	900K	10.88	0.89 ms	1.13 ms	1.74 ms
StudentExpNet <sub>S</sub>	232K	2.91	0.78 ms	1.08 ms	1.69 ms
StudentExpNet <sub>X<sub>S</sub></sub>	121K	1.52	0.63 ms	0.97 ms	1.63 ms
<b>MicroExpNet</b>	<b>65K</b>	<b>0.88</b>	<b>0.53 ms</b>	<b>0.97 ms</b>	<b>1.52 ms</b>

Memory requirements and average per-image running times

## Recap: Knowledge Distillation <sup>5</sup>

Let  $p_t$  and  $p'_s$  be the softened softmax of the student and teacher respectively whereas  $p_s$  is the vanilla softmax of the student:

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

Then the cost function becomes:

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

---

<sup>5</sup>Hinton, Geoffrey, Oriol Vinyals, and Jeff Dean. "Distilling the knowledge in a neural network." arXiv preprint arXiv:1503.02531 (2015).

# Teacher and Student Models

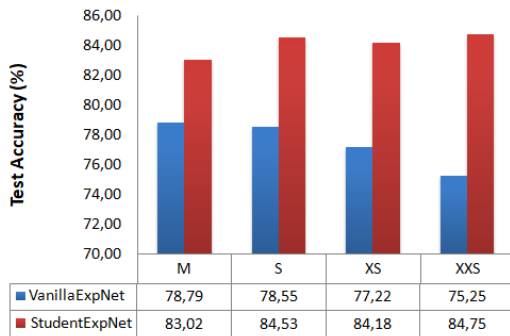
- TeacherExpNet: Inception\_v3<sup>6</sup> network trained on ImageNet<sup>7</sup>
- StudentExpNet:
  - $p_{12}$ : each conv layer is followed by a max pooling layer
  - M, S, XS, XXS

---

<sup>6</sup>Szegedy, Christian, et al. "Rethinking the inception architecture for computer vision." Proceedings of the IEEE conference on computer vision and pattern recognition. 2016.

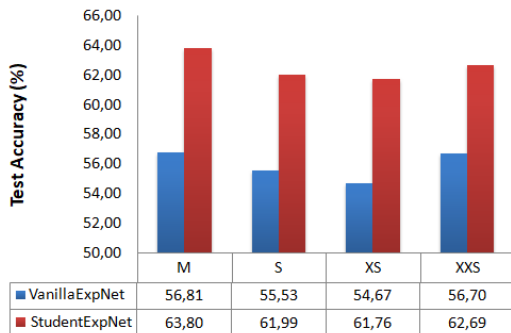
<sup>7</sup>Russakovsky, Olga, et al. "Imagenet large scale visual recognition challenge." International journal of computer vision 115.3 (2015): 211-252.

# Regularization on CK+ Performance



The effect of supervision on CK+ for 3000 epochs of training

# Regularization on Oulu-CASIA Performance

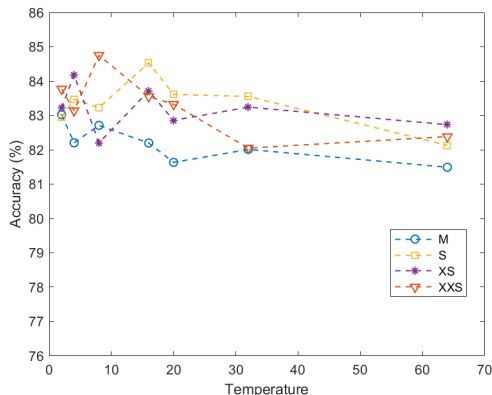


The effect of supervision on Oulu-CASIA for 3000 epochs of training

# Temperature Analysis

- Grid search for temperatures:  $T \in [2, 4, 8, 16, 20, 32, 64]$
- Random split vs. subject-independent split

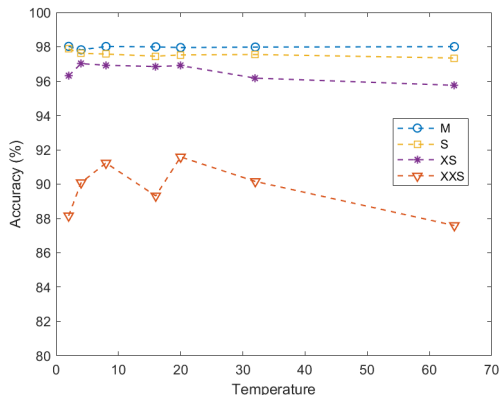
# Temperature Analysis using CK+



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

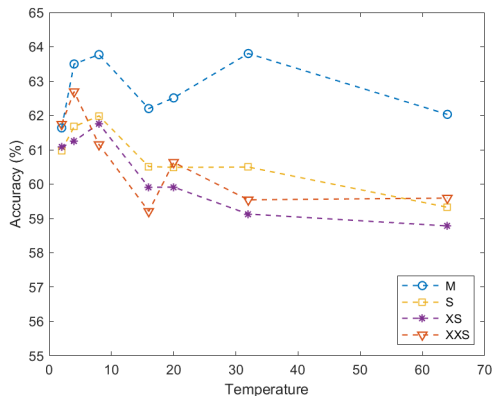


# Temperature Analysis using CK+



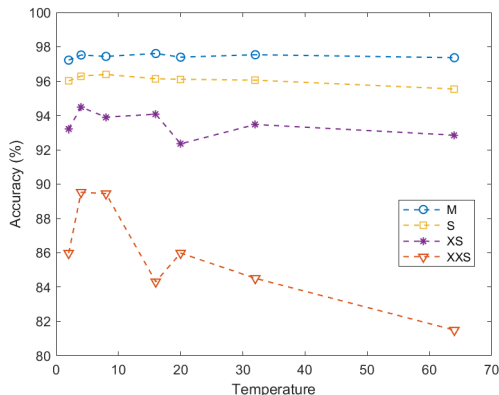
Classification performances of the student networks across different temperatures on the CK+ dataset using **random splits**

# Temperature Analysis using Oulu-CASIA



Classification performances of the student networks across different temperatures on the Oulu-CASIA dataset using **subject-independent splits**

# Temperature Analysis using Oulu-CASIA



Classification performances of the student networks across different temperatures on the Oulu-CASIA dataset using **random splits**

# Future Research Directions

- Is information loss essential for generalization?
- Is a smaller model more open to teacher's supervision?
- Why does the classification accuracy fluctuate as the temperature  $T$  is changed?