

PROBLEM

Deep convolutional neural network has achieved great success on large-scale image classification task. However, how to effectively train deep network on small dataset is still a challenging problem.

Conventional method is to fine-tune a deep network trained on face recognition dataset to adapt to the facial expression recognition task. This simple strategy has two notable problems:

1. The fine-tuned face net may still contain information useful for subject identification.
2. The network designed for the face recognition domain often has a large capacity, thus the overfitting issue is still severe.

METHOD

The distribution function of the high-level neurons can be formulated as follows:

$$f(X^l) = C_p \cdot e^{-||X^l||_p^p} \tag{1}$$

To incorporate the knowledge of a face net, we propose to extend (1) to have the following form, i.e.,:

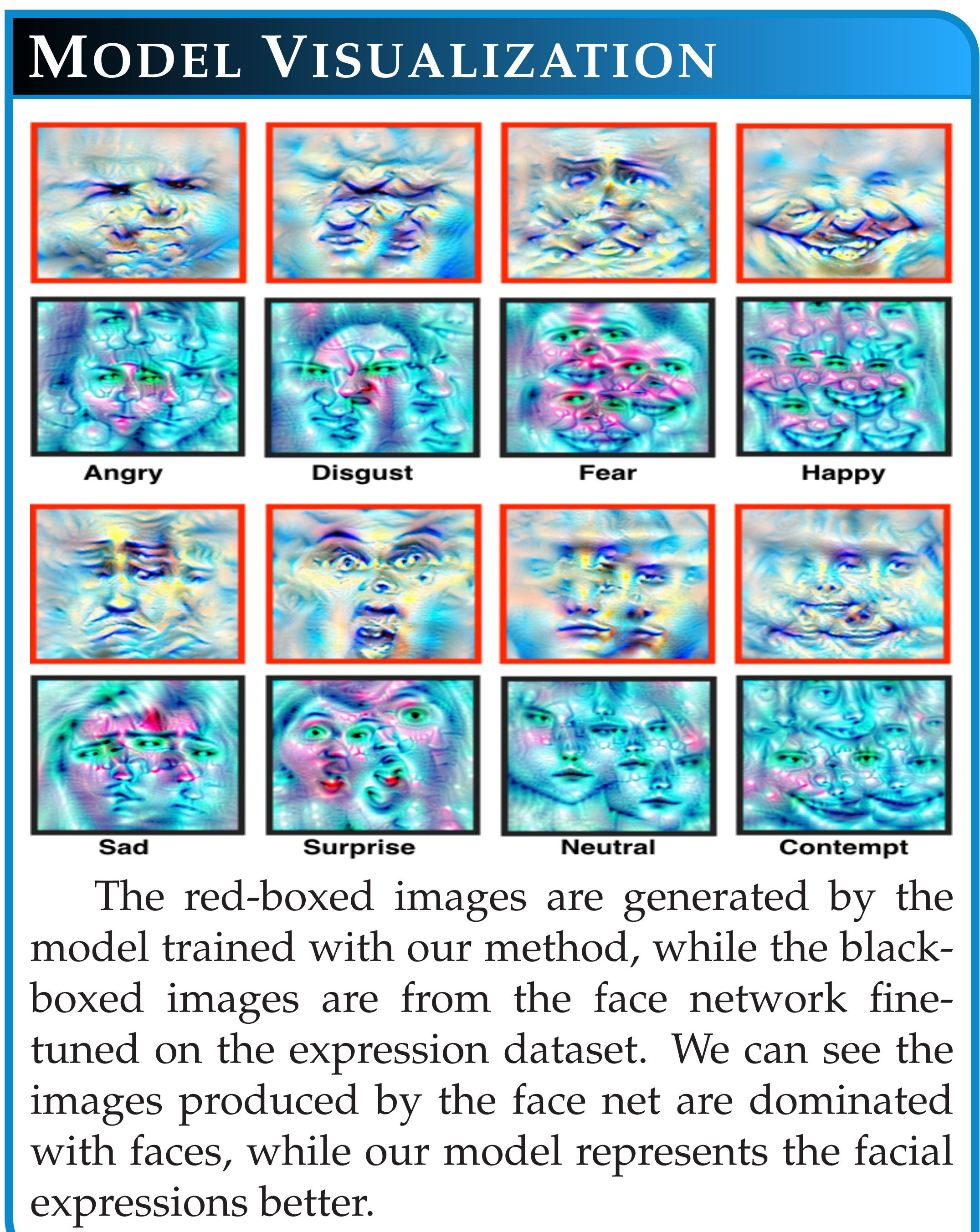
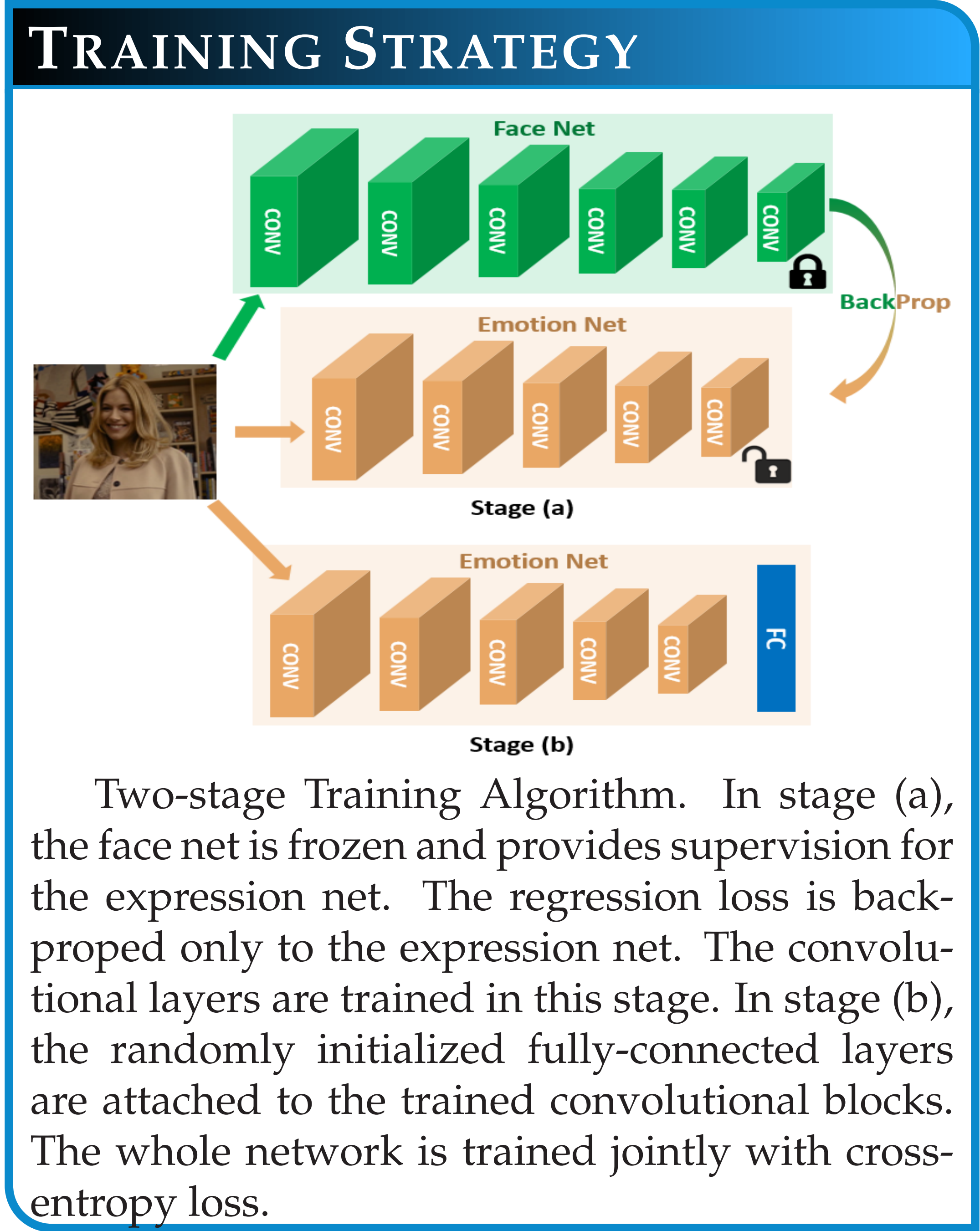
$$f(X^l) = C_p \cdot e^{-||X^l - \mu||_p^p} \tag{2}$$

The mean is modeled by the face net. This is motivated by the observation that the fine-tuned face net already achieves competitive performance on the expression dataset, so it should provide a good initialization point for the expression net.

Using the maximum likelihood estimation (MLE) procedure, we can derive the loss function as:

$$\begin{aligned} \max_{\theta_1} L_1 &= \max_{\theta_1} \log f(X^l) \\ &= \max_{\theta_1} \log C_p \cdot e^{-||X^l - \mu||_p^p} \\ &= \min_{\theta_1} ||g_{\theta_1}(I) - G(I)||_p^p \end{aligned} \tag{3}$$

- Which layer to transfer?
- Our experiment results suggest that late middle layer, such as pool5, is a good tradeoff between supervision richness and representation discriminativeness.



EXPERIMENTS

- Classification Results

|                          |  |                  |                          |  |                  |
|--------------------------|--|------------------|--------------------------|--|------------------|
| CK+                      |  |                  | OULU-CAS                 |  |                  |
| Method                   |  | Average Accuracy | Method                   |  | Average Accuracy |
| CSPL [16]                |  | 89.9%            | HOG 3D [38]              |  | 70.63%           |
| AdaGabor [35]            |  | 93.3%            | AdaLBP [28]              |  | 73.54%           |
| LBPSVM [36]              |  | 95.1%            | Atlases [39]             |  | 75.52%           |
| 3DCNN-DAP [21]           |  | 92.4%            | STM-ExpLet [20]          |  | 74.59%           |
| BDBN [19]                |  | 96.7%            | DTAGN [22]               |  | 81.46%           |
| STM-ExpLet [20]          |  | 94.2%            | LOMo [37]                |  | 82.10%           |
| DTAGN [22]               |  | 97.3%            | PPDN [9]                 |  | 84.59%           |
| Inception [23]           |  | 93.2%            | Train From Scratch (BN)  |  | 76.87%           |
| LOMo [37]                |  | 95.1%            | VGG Fine-Tune (baseline) |  | 83.26%           |
| PPDN [9]                 |  | 97.3%            | FN2EN                    |  | <b>87.71%</b>    |
| FN2EN                    |  | <b>98.6%</b>     |                          |  |                  |
| AUDN [18]                |  | 92.1%            |                          |  |                  |
| Train From Scratch (BN)  |  | 88.7%            |                          |  |                  |
| VGG Fine-Tune (baseline) |  | 89.9%            |                          |  |                  |
| FN2EN                    |  | <b>96.8%</b>     |                          |  |                  |

|                          |  |                  |                            |  |                  |
|--------------------------|--|------------------|----------------------------|--|------------------|
| TFD                      |  |                  | SFEW                       |  |                  |
| Method                   |  | Average Accuracy | Method                     |  | Average Accuracy |
| Gabor + PCA [40]         |  | 80.2%            | AUDN [18]                  |  | 26.14%           |
| Deep mPoT [41]           |  | 82.4%            | STM-ExpLet [20]            |  | 31.73%           |
| CDA+CCA [42]             |  | 85.0%            | Inception [23]             |  | 47.70%           |
| disRBM [43]              |  | 85.4%            | Mapped LBP [8]             |  | 41.92%           |
| bootstrap-recon [44]     |  | 86.8%            | Train From Scratch (BN)    |  | 39.55%           |
| Train From Scratch (BN)  |  | 82.5%            | VGG Fine-Tune (baseline)   |  | 41.23%           |
| VGG Fine-Tune (baseline) |  | 86.7%            | FN2EN                      |  | <b>48.19%</b>    |
| FN2EN                    |  | <b>88.9%</b>     | Transfer Learning [25]     |  | 48.50%           |
|                          |  |                  | Multiple Deep Network [24] |  | 52.29%           |
|                          |  |                  | FN2EN                      |  | <b>55.15%</b>    |

- Top hidden layer neuron visualization

|             |  |             |             |  |             |
|-------------|--|-------------|-------------|--|-------------|
| CK+         |  |             | Oulu-CASIA  |  |             |
| Neuron #11  |  | Neuron #16  | Neuron #17  |  | Neuron #74  |
| Neuron #53  |  | Neuron #84  | Neuron #53  |  | Neuron #88  |
| Neuron #68  |  | Neuron #29  | Neuron #5   |  | Neuron #138 |
| Neuron #29  |  | Neuron #5   | Neuron #138 |  |             |
| Neuron #5   |  | Neuron #138 |             |  |             |
| Neuron #138 |  |             |             |  |             |

SUMMARY

We propose a probabilistic distribution function to model the high level neuron response based on already fine-tuned face net, thereby leading to feature level regularization that exploits the rich face information in the face net. In the second stage, we perform label supervision to boost the final discriminative capability.

As a result, FaceNet2ExpNet improves visual feature representation and outperforms various state-of-the-art methods on four public datasets. In future, we plan to apply this training method to other domains with small datasets.