

Facial Key Points Detection using Deep Convolutional Neural Network - NaimishNet

Naimish Agarwal, IIIT-Allahabad (*irm2013013@iiita.ac.in*)
 Artus Krohn-Grimberghe, University of Paderborn (*artus@aisbi.de*)
 Ranjana Vyas, IIIT-Allahabad (*ranjana@iiita.ac.in*)

Abstract—Facial Key Points (FKPs) Detection is an important and challenging problem in the fields of computer vision and machine learning. It involves predicting the co-ordinates of the FKPs, e.g. nose tip, center of eyes, etc, for a given face. In this paper, we propose a LeNet adapted Deep CNN model - NaimishNet, to operate on facial key points data and compare our model's performance against existing state of the art approaches.

Index Terms—Facial Key Points Detection, Deep Convolutional Neural Network, NaimishNet

I. INTRODUCTION

Facial Key Points (FKPs) detection is an important and challenging problem in the field of computer vision, which involves detecting FKPs like centers and corners of eyes, nose tip, etc. The problem is to predict the (x, y) real-valued co-ordinates in the space of image pixels of the FKPs for a given face image. It finds its application in tracking faces in images and videos, analysis of facial expressions, detection of dysmorphic facial signs for medical diagnosis, face recognition, etc.

Facial features vary greatly from one individual to another, and even for a single individual there is a large amount of variation due to pose, size, position, etc. The problem becomes even more challenging when the face images are taken under different illumination conditions, viewing angles, etc.

In the past few years, advancements in FKPs detection are made by the application of deep convolutional neural network (DCNN), which is a special type of feed-forward neural network with shared weights and local connectivity. DCNNs have helped build state-of-the-art models for image recognition, recommender systems, natural language processing, etc. Krizhevsky et al. [1] applied DCNN in ImageNet image classification challenge and outperformed the previous state-of-the-art model for image classification.

Wang et al. [2] addressed FKPs detection by first applying histogram stretching for image contrast enhancement, followed by principal component analysis for noise reduction and mean patch search algorithm with correlation scoring and mutual information scoring for predicting left and right eye centers. Sun et al. [3] estimated FKPs by using a three level convolutional neural network, where at each level, outputs of multiple networks were fused for robust and accurate estimation. Longpre et al. [4] predicted FKPs by first applying data augmentation to expand the number of training examples, followed by testing different architectures of convolutional

| Layer Number | Layer Name | Layer Shape |
|--------------|----------------------------|---------------|
| 1 | Input ₁ | (1, 96, 96) |
| 2 | Convolution2d ₁ | (32, 93, 93) |
| 3 | Activation ₁ | (32, 93, 93) |
| 4 | Maxpooling2d ₁ | (32, 46, 46) |
| 5 | Dropout ₁ | (32, 46, 46) |
| 6 | Convolution2d ₂ | (64, 44, 44) |
| 7 | Activation ₂ | (64, 44, 44) |
| 8 | Maxpooling2d ₂ | (64, 22, 22) |
| 9 | Dropout ₂ | (64, 22, 22) |
| 10 | Convolution2d ₃ | (128, 21, 21) |
| 11 | Activation ₃ | (128, 21, 21) |
| 12 | Maxpooling2d ₃ | (128, 10, 10) |
| 13 | Dropout ₃ | (128, 10, 10) |
| 14 | Convolution2d ₄ | (256, 10, 10) |
| 15 | Activation ₄ | (256, 10, 10) |
| 16 | Maxpooling2d ₄ | (256, 5, 5) |
| 17 | Dropout ₄ | (256, 5, 5) |
| 18 | Flatten ₁ | (6400) |
| 19 | Dense ₁ | (1000) |
| 20 | Activation ₅ | (1000) |
| 21 | Dropout ₅ | (1000) |
| 22 | Dense ₂ | (1000) |
| 23 | Activation ₆ | (1000) |
| 24 | Dropout ₆ | (1000) |
| 25 | Dense ₃ | (2) |

Table I
NAIMISHNET LAYER-WISE ARCHITECTURE

neural networks like LeNet [5] and VGGNet [6], and finally used a weighted ensemble of models. Nouri et al. [7] used six specialist DCNNs trained over pre-trained weights. Oneto et al. [8] applied a variety of data pre-processing techniques like histogram stretching, Gaussian blurring, followed by image flipping, key point grouping, and then finally applied LeNet.

Taigman et al. [9] provided a new deep network architecture, DeepFace, for state-of-the-art face recognition. Li et al. [10] provided a new DCNN architecture for state-of-the-art face alignment.

We present a DCNN architecture – NaimishNet, based on LeNet, which addresses the problem of facial key points detection by providing a learning model for a single facial key point.

II. NAIMISHNET ARCHITECTURE

NaimishNet consists of 4 convolution2d layers, 4 maxpooling2d layers and 3 dense layers, with sandwiched dropout and activation layers, as shown in table I.

A. Layer-wise Details

The following points give the details of every layer in the architecture:

| Layer Name | Number of Filters | Filter Shape |
|----------------------------|-------------------|--------------|
| Convolution2d ₁ | 32 | (4, 4) |
| Convolution2d ₂ | 64 | (3, 3) |
| Convolution2d ₃ | 128 | (2, 2) |
| Convolution2d ₄ | 256 | (1, 1) |

Table II
FILTER DETAILS OF CONVOLUTION2D LAYERS

- Input₁ is the input layer.
- Activation₁ to Activation₅ use Exponential Linear Units (ELUs) [11] as activation functions, whereas Activation₆ uses Linear Activation Function.
- Dropout [12] probability is increased from 0.1 to 0.6 from Dropout₁ to Dropout₆, with a step size of 0.1.
- Maxpooling2d₁ to Maxpooling2d₄ use a pool shape of (2, 2), with non-overlapping strides and no zero padding.
- Flatten₁ flattens 3d input to 1d output.
- Convolution2d₁ to Convolution2d₄ do not use zero padding, have their weights initialized with random numbers drawn from uniform distribution, and the specifics of number of filters and filter shape are shown in table II.
- Dense₁ to Dense₃ are regular fully connected layers with weights initialized using Glorot uniform initialization [13].
- Adam [14] optimizer, with learning rate of 0.001, β_1 of 0.9, β_2 of 0.999 and ϵ of $1e-08$, is used for minimizing Mean Squared Error (MSE).
- There are 7,488,962 model parameters.

B. Design Decisions

We have tried various different deep network architectures before finally settling down with NaimishNet architecture. Most of the design decisions were governed primarily by the hardware restrictions and the time taken per epoch for different runs of the network. Although we wanted one architecture to fit well to all the different FKPs as target, we have mainly taken decisions based on the performance of the model on left eye center FKP as target. We have run the different architectures for at most 40 epochs and forecasted its performance and generalizability based on our experience we gathered while training multiple deep network architectures.

We experimented with 5 convolution2d and 5 maxpooling2d layers and 4 dense layers with rectified linear units as activation functions, but that over fitted. We experimented with 3 convolution2d and 3 maxpooling2d layers and 2 dense layers with rectified linear units as activation functions, but that in turn required pre-training the weights to achieve good performance. We wanted to avoid pre-training to reduce the training time. We did not experiment with zero padding and strides.

We sandwiched batch normalization layers [15] after every maxpooling2d layer but that increased the per epoch training time by fivefold. We used dropout layers with dropout probability of 0.5 for every dropout layer but that showed a poor fit. We removed the dropout layers but that resulted in overfitting. We tried different initialization schemes like He normal initialization [15], but that did not show nice performance in the first 20 epochs, in fact uniform initialization performed better

than that for us. We experimented with different patience levels (PLs) for early stopping callback function and concluded that PL of 10% of total number of epochs, in our case 300, would do well in most of the cases. For the aforementioned cases, we used RMSProp optimizer [16].

The main aim was to reduce the training time along with achieving maximum performance and generalizability. To achieve that we experimented with exponential linear units as activation function, Adam optimizer, and linearly increasing dropout probabilities along the depth of the network. The specifics of the filter were chosen based on our past experience with exponentially increasing number of filters and linearly decreasing filter size. The batch size was kept 128 because vector operations are optimized for sizes which are in powers of 2. Also, given over 7 million model parameters, and limited GPU memory, batch size of 128 was an ideal size for us. The number of main layers i.e. 4 convolutional layers (convolution2d and maxpooling2d) and 3 dense layers was ideal since exponential linear units perform well in cases where the number of layers are greater than 5. The parameters of Adam optimizer were used as is as recommended by its authors and we did not experiment with them.

NaimishNet, is inspired by LeNet [5], since LeNet's architecture follows the pattern $input \rightarrow conv2d \rightarrow maxpool2d \rightarrow conv2d \rightarrow maxpool2d \rightarrow conv2d \rightarrow maxpool2d \rightarrow dense \rightarrow output$. Works of Longpre et al. [4], Oneto et al. [8], etc. have shown that LeNet styled architectures have been successfully applied for Facial Key Points Detection problem.

III. EXPERIMENTS

The details of the experiment are given below.

A. Dataset Description

We have used the dataset from Kaggle Competition – Facial Key Points Detection [17]. The dataset was chosen to benchmark our solution against the existing approaches which address FKPs detection problem.

There are 15 FKPs per face image like left eye center, right eye center, left eye inner corner, left eye outer corner, right eye inner corner, right eye outer corner, left eyebrow inner end, left eyebrow outer end, right eyebrow inner end, right eyebrow outer end, nose tip, mouth left corner, mouth right corner, mouth center top lip and mouth center bottom lip. Here, left and right are from the point of view of the subject.

The greyscale input images, with pixel values in range of [0, 255], have size of 96×96 pixels. The train dataset consists of 7049 images, with 30 targets, i.e. (x, y) for each of 15 FKPs. It has missing target values for many FKPs for many face images. The test dataset consists of 1783 images with no target information. The Kaggle submission file consists of 27124 FKPs co-ordinates, which are to be predicted.

The Kaggle submissions are scored based on the Root Mean Squared Error (RMSE), which is given by

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}}$$

where y_i is the original value, \hat{y}_i is the predicted value and n is the number of targets to be predicted.

| Target Facial Key Points | Target Facial Key Points |
|-----------------------------|------------------------------|
| Left Eye Center X | Right Eye Center X |
| Left Eye Center Y | Right Eye Center Y |
| Left Eye Inner Corner X | Right Eye Inner Corner X |
| Left Eye Inner Corner Y | Right Eye Inner Corner Y |
| Left Eye Outer Corner X | Right Eye Outer Corner X |
| Left Eye Outer Corner Y | Right Eye Outer Corner Y |
| Left Eyebrow Inner Corner X | Right Eyebrow Inner Corner X |
| Left Eyebrow Inner Corner Y | Right Eyebrow Inner Corner Y |
| Left Eyebrow Outer Corner X | Right Eyebrow Outer Corner X |
| Left Eyebrow Outer Corner Y | Right Eyebrow Outer Corner Y |
| Mouth Left Corner X | Mouth Right Corner X |
| Mouth Left Corner Y | Mouth Right Corner Y |

Table III

TARGET FACIAL KEY POINTS WHICH NEED TO BE SWAPPED WHILE HORIZONTALLY FLIPPING THE DATA

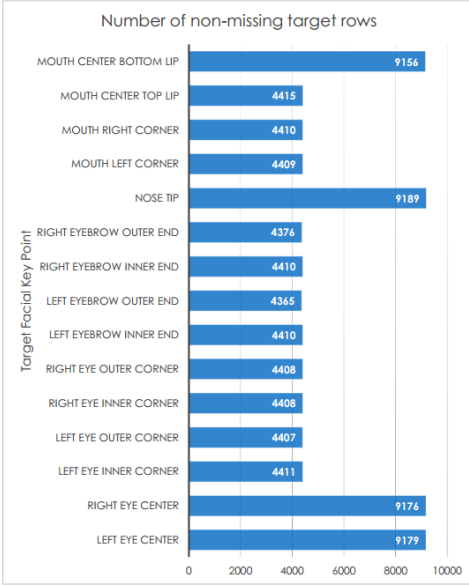


Figure 1. Number of non-missing target rows in the augmented and pre-processed train dataset

B. Proposed Approach

We propose the following approach for Facial Key Points Detection.

1) *Data Augmentation*: Data Augmentation helps boost the performance of a deep learning model when the training data is limited by generating more training data. We have horizontally flipped [7] the images for which target information for all the 15 FPKs are available, and also swapped the target values according to table III. Then, we vertically stacked the new horizontally flipped data under the original train data to create the augmented train dataset.

2) *Data Pre-processing*: The image pixels are normalized to the range $[0, 1]$ by dividing by 255.0, and the train targets are zero-centered to the range $[-1, 1]$ by first dividing by 48.0, since the images are 96×96 , and then subtracting 48.0.

3) *Pre-training Analysis*: We have created 15 NaimishNet models, since figure 1 shows that there are different number of non-missing target rows for different FPKs.

4) *Training*: For each of the 15 NaimishNet models, filter (keep) the rows with non-missing target values and split the



Figure 2. Annotated Faces with 15 Facial Key Points, marked in blue for original and in red for predicted.

filtered train data into 80% train dataset (T) and 20% validation dataset (V). Compute feature-wise mean M on T, and zero-center T and V by subtracting M , and also store M . Reshape features of T and V to $(1, 96, 96)$, and start the training. Finally, validate on V, and store the model's loss history.

Each model is trained with a batch size of 128 and maximum number of epochs are 300. Two callback functions are used which execute at the end of every epoch – Early Stopping Callback (ESC) and Model Checkpoint Callback (MCC).

ESC stops the training when the number of contiguous epochs without improvement in validation loss are 30 (Patience Level). MCC stores the weights of the model with the lowest validation loss.

5) *Evaluation*: For each of the 15 NaimishNet models, reload M , zero center test features by subtracting M and reshape the samples to $(1, 96, 96)$. Reload the best checkpointed model weights, predict on test features, and store the predictions.

IV. RESULTS

The annotated faces with 15 FPKs of 6 individuals is shown in figure 2. It is visible that the predicted locations are very close to the actual locations of FPKs.

In figure 3, high $RMSE_{left}$ is because of the presence of less number of images, in the training dataset, with the side face pose. Most of the images in the training dataset are that of the frontal face with nearly closed lips. So, NaimishNet delivers lower $RMSE_{center}$ and $RMSE_{right}$.

Total number of epochs are 3507, as shown in figure 4. Total training time of 15 NaimishNet models is approximately 20 hours.

In the following points, we present the analysis of train and validation error for each of the 15 NaimishNet models taken 2 at a time:

- Both figure 5 and figure 6 show that both train and validation RMSE decrease over time with small wiggles along the curve showcasing that the batch size was rightly chosen and end up being close enough, thus showing that the two NaimishNet models generalized well. The best check-pointed models were achieved before 300 epochs, thus certifying the decision of fixing 300 epochs as the upper limit.
- Both figure 7 and figure 8 show that both train and validation RMSE decrease over time with small wiggles and end up being close enough. The best check-pointed models were achieved before 300 epochs. Also, for figure 8, since the train RMSE does not become less than validation RMSE implies that the Patience Level (PL) for ESC could have been increased to 50. Thus, the training could have been continued for some more epochs.
- Both figure 9 and figure 10 show that both train and validation RMSE decrease over time with very small wiggles and end up being close enough. The best check-pointed models were achieved close to 300 epochs. For figure 10, PL of ESC could have been increased to 50 so that the train RMSE becomes less than validation RMSE.
- Both figure 11 and figure 12, show that the train and validation RMSE decreases over time with small wiggles, and end up being close enough. The best check-pointed models were achieved far before 300 epochs. No further improvement could have been seen by increasing PL of ESC.
- Both figure 13 and figure 14 show that both train and validation RMSE decrease over time with small wiggles and end up being close in the check-pointed model, which can be seen 30 epochs before the last epoch on x-axis.
- Both figure 15 and figure 16 show that both train and validation RMSE decrease over time and they end up being close enough in the final check-pointed model. For figure 16, PL of ESC could have been increased to 50 but not necessarily there would have been any visible improvement.
- Both figure 17 and figure 18 show that the train and validation RMSE decrease over time, with small wiggles, and the upper limit of 300 epochs was just right enough for the models to show convergence.
- As per figure 19, the train and validation RMSE decreases over time and the wiggly nature of the validation curve increases towards the end. We could have tried to increase the batch size to 256 if it could be handled by our hardware. The curves are close together and show convergence before 300 epochs.

As per figure 20, the NaimishNet generalizes well, since the average train RMSE / validation RMSE is 1.03, where average loss is calculated by



Figure 3. The worst, medium and best predictions, where the $RMSE_{left} = 6.87$, $RMSE_{center} = 1.04$, and $RMSE_{right} = 0.37$, the original FKPs are shown in blue and the predicted FKPs are shown in red.

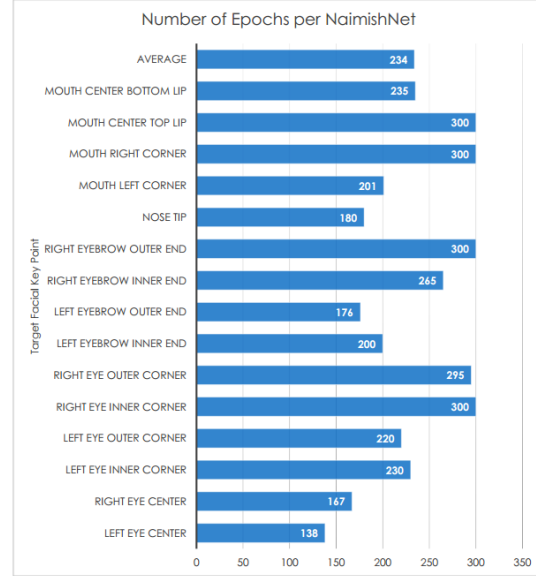


Figure 4. Number of Epochs per NaimishNet

$$\text{Average RMSE} = \sqrt{\frac{\sum_{i=1}^{15} RMSE_i^2}{15}}$$

As per figure 21, the NaimishNet outperforms the approaches used by Longpre et al. [4], Oneto et al. [8], Wang et al. [2], since lower the score, better the model.

V. CONCLUSION

NaimishNet – a learning model for a single facial key point, is based on LeNet [5]. As per figure 21, LeNet based

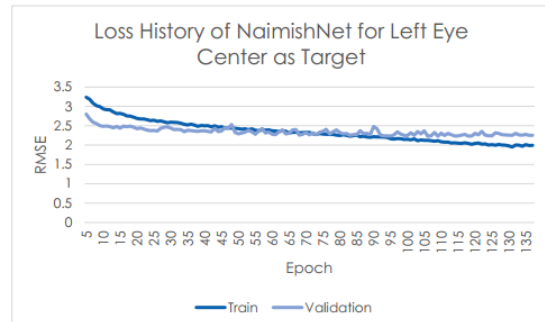


Figure 5. Loss History of NaimishNet for Left Eye Center as Target

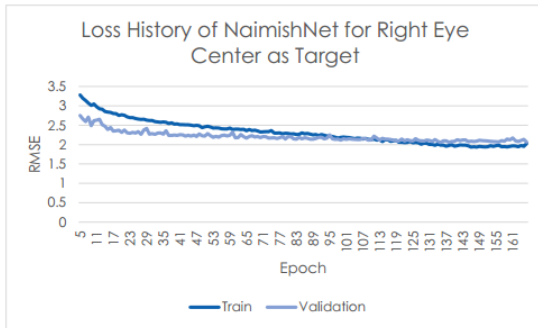


Figure 6. Loss History of NaimishNet for Right Eye Center as Target

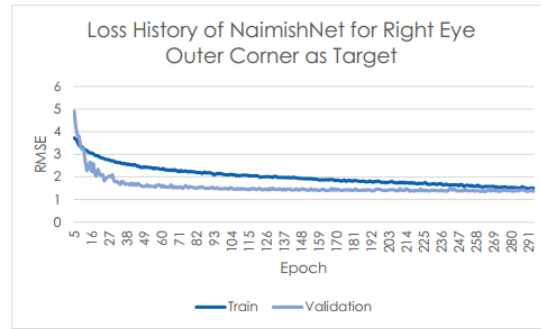


Figure 10. Loss History of NaimishNet for Right Eye Outer Corner as Target

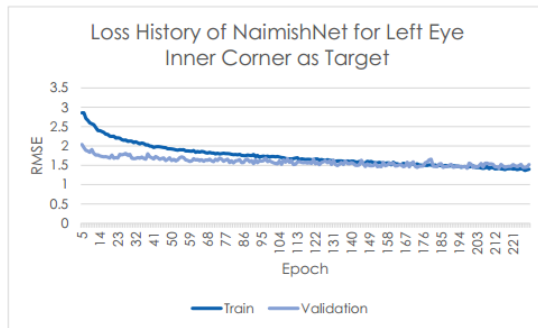


Figure 7. Loss History of NaimishNet for Left Eye Inner Corner as Target

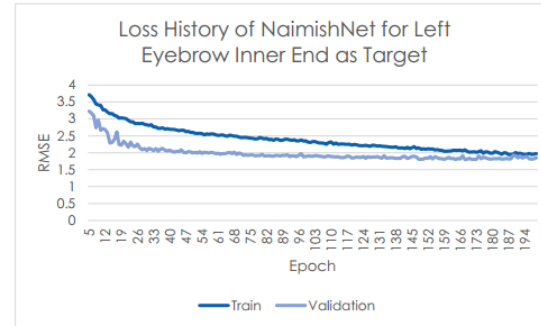


Figure 11. Loss History of NaimishNet for Left Eyebrow Inner End as Target

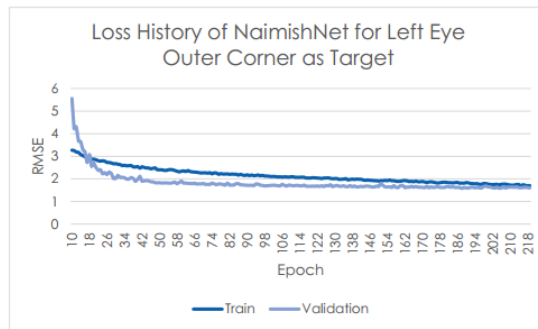


Figure 8. Loss History of NaimishNet for Left Eye Outer Corner as Target

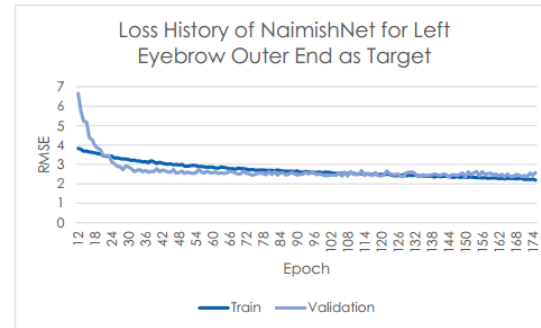


Figure 12. Loss History of NaimishNet for Left Eyebrow Outer End as Target

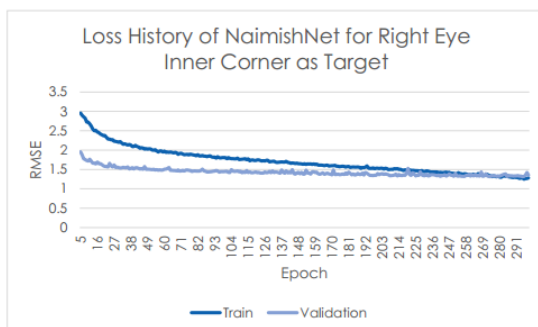


Figure 9. Loss History of NaimishNet for Right Eye Inner Corner as Target

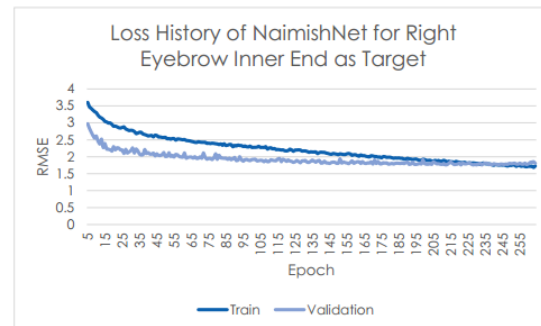


Figure 13. Loss History of NaimishNet for Right Eyebrow Inner End as Target

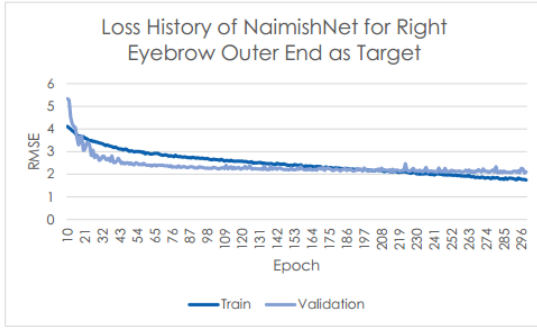


Figure 14. Loss History of NaimishNet for Right Eyebrow Outer End as Target

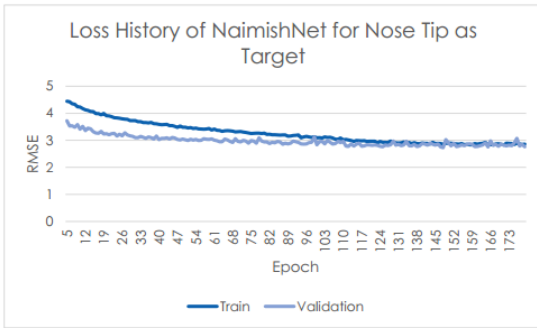


Figure 15. Loss History of NaimishNet for Nose Tip as Target

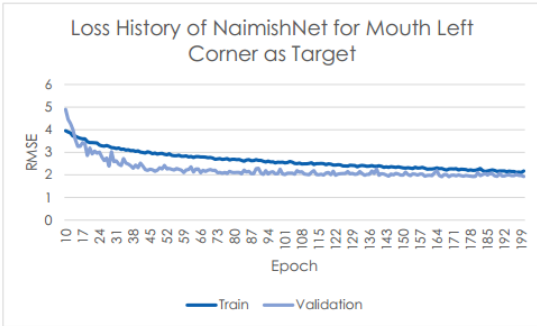


Figure 16. Loss History of NaimishNet for Mouth Left Corner as Target

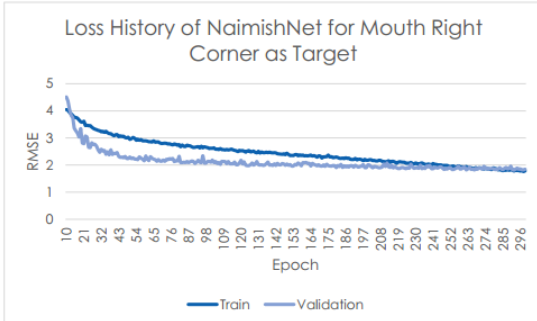


Figure 17. Loss History of NaimishNet for Mouth Right Corner as Target

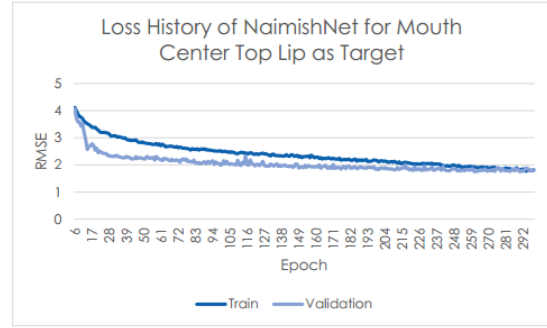


Figure 18. Loss History of NaimishNet for Mouth Center Top Lip as Target

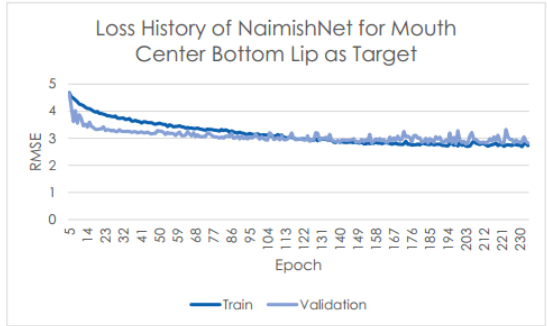


Figure 19. Loss History of NaimishNet for Mouth Center Bottom Lip as Target

architectures have proven to be successful for Facial Key Points Detection. Works of Longpre et al. [4], Oneto et al. [8], etc. have shown that LeNet styled architectures have been successfully applied for Facial Key Points Detection problem.

VI. FUTURE SCOPE

Given enough compute resources and time, NaimishNet can be further modified by experimenting with different initialization schemes, different activation functions, different number of filters and kernel size, different number of layers, switching from LeNet [5] to VGGNet [6], introducing Inception modules as in GoogleNet [18], or introducing Residual Networks [19], etc. Different image pre-processing approaches like histogram stretching, zero centering, etc. can be tried to check which approaches improve the model's accuracy. New approaches to Facial Key Points Detection can be based on the application of deep learning as a three stage process - face detection, face alignment, and facial key points detection, with the use of the state-of-the-art algorithms for each stage.

REFERENCES

- [1] A. Krizhevsky, I. Sutskever and G. E. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks," in Neural Information Processing Systems, 2012.
- [2] Y. Wang and Y. Song, "Facial Keypoints Detection," Stanford University, 2014.
- [3] Y. Sun, X. Wang and X. Tang, "Deep Convolutional Network Cascade for Facial Key Point Detection," in Conference on Computer Vision and Pattern Recognition, 2013.
- [4] S. Longpre and A. Sohmshetty, "Facial Keypoint Detection," Stanford University, 2016.

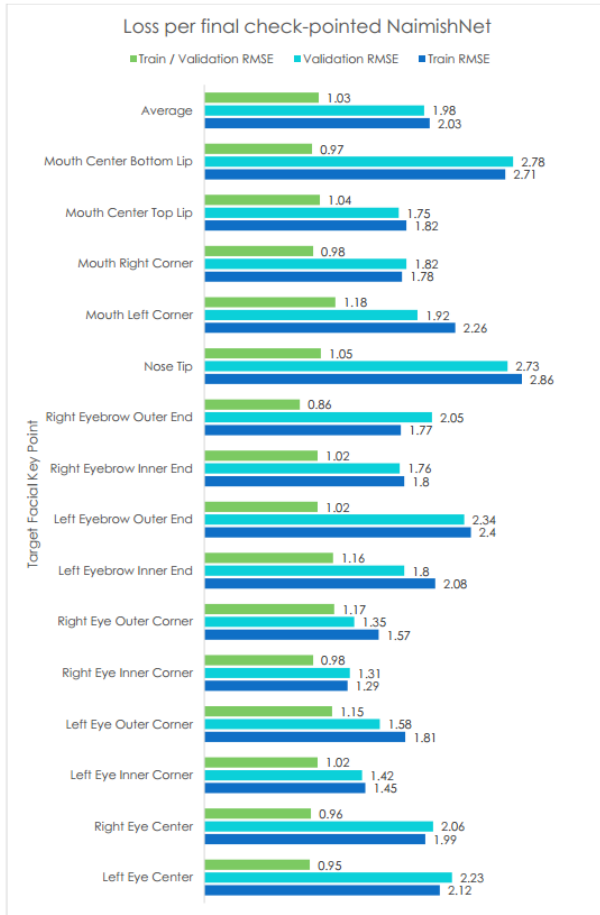


Figure 20. Loss per final check-pointed NaimishNet

- [5] Y. LeCun, L. Bottou, Y. Bengio and P. Haffner, "Gradient-Based Learning Applied to Document Recognition," in Proceedings of the IEEE, 1998.
- [6] K. Simoyan and A. Zisserman, "Very Deep Convolutional Networks for Large-scale Image Recognition," in International Conference on Learning Representation, 2015.
- [7] D. Nouri, "Using convolutional neural nets to detect facial keypoints tutorial," 17 12 2014. [Online]. Available: <http://danielnouri.org/notes/2014/12/17/using-convolutional-neural-nets-to-detect-facial-keypoints-tutorial/>.
- [8] M. Oneto, L. Yang, Y. Jang and C. Dailey, "Facial Keypoints Detection: An Effort to Top the Kaggle Leaderboard," Berkeley School of Information, 2015.
- [9] Y. Taigman, M. Yang, M. Ranzato and L. Wolf, "DeepFace: Closing the Gap to Human-Level Performance in Face Verification," in Computer Vision and Pattern Recognition, 2014.
- [10] P. Li, X. Liao, Y. Xu, H. Ling and C. Liao, "GPU Boosted Deep Learning in Real-time Face Alignment," in GPU Technology Conference, 2016.
- [11] D.-A. Clevert, T. Unterthiner and S. Hochreiter, "Fast and Accurate Deep Network Learning by Exponential Linear Units (ELUs)," in International Conference on Learning Representations, 2016.
- [12] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever and R. Salakhutdinov, "Dropout: A Simple Way to Prevent Neural Networks from Overfitting," Journal of Machine Learning Research, pp. 1929-1958, 2014.
- [13] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks," Journal of Machine Learning Research, pp. 249-256, 2010.
- [14] D. P. Kingma and J. L. Ba, "Adam: A Method for Stochastic Optimization," in International Conference on Learning Representation, 2015.
- [15] F.-F. Li, A. Karpathy and J. Johnson, "Neural Networks Part 2: Setting up the Data and Loss," 2016. [Online]. Available: <http://cs231n.github.io/neural-networks-2/>.
- [16] F.-F. Li, A. Karpathy and J. Johnson, "Neural Networks Part 3: Learning and Evaluation," 2016. [Online]. Available: <http://cs231n.github.io/neural-networks-3/>.
- [17] "Kaggle Competition - Facial Key Points Detection," 7 May 2013. [Online]. Available: <https://www.kaggle.com/c/facial-keypoints-detection>.
- [18] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke and A. Rabinovich, "Going Deeper with Convolutions," in Computer Vision and Pattern Recognition, 2014.
- [19] K. He, X. Zhang, S. Ren and J. Sun, "Deep Residual Learning for Image Recognition," in Computer Vision and Pattern Recognition, 2015.

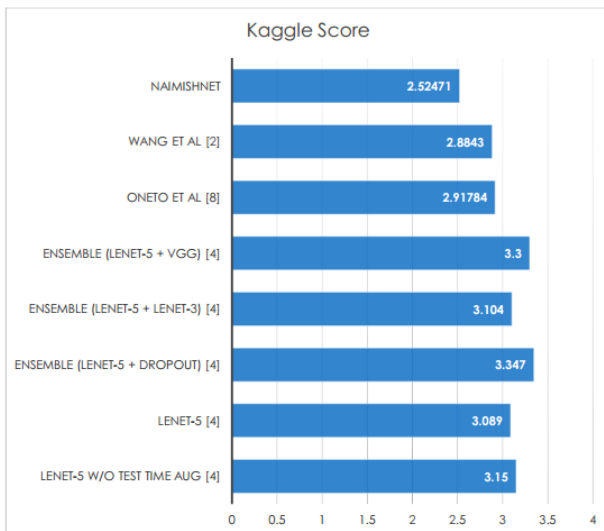


Figure 21. Comparison of Kaggle Scores