

One for All: An End-to-End Compact Solution for Hand Gesture Recognition

Monu Verma

Computer Science and Engineering
Malaviya National Institute of Technology
Jaipur, India, 302017
Email: monuverma.cv@gmail.com

Ayushi Gupta

olx people, Bangalore,
Karnataka, 560034
Email: ayushigup26@gmail.com

Santosh K. Vipparthi

Computer Science and Engineering
Malaviya National Institute of Technology
Jaipur, India, 302017
Email: skvipparthi@mmit.ac.in

Abstract—The HGR is a quite challenging task as its performance is influenced by various aspects such as illumination variations, cluttered backgrounds, spontaneous capture, etc. The conventional CNN networks for HGR are following two stage pipeline to deal with the various challenges: complex signs, illumination variations, complex and cluttered backgrounds. The existing approaches needs expert expertise as well as auxiliary computation at stage 1 to remove the complexities from the input images. Therefore, in this paper, we proposes an novel end-to-end compact CNN framework: fine grained feature attentive network for hand gesture recognition (Fit-Hand) to solve the challenges as discussed above. The pipeline of the proposed architecture consists of two main units: FineFeat module and dilated convolutional (Conv) layer. The FineFeat module extracts fine grained feature maps by employing attention mechanism over multi-scale receptive fields. The attention mechanism is introduced to capture effective features by enlarging the average behaviour of multi-scale responses. Moreover, dilated convolution provides global features of hand gestures through a larger receptive field. In addition, integrated layer is also utilized to combine the features of FineFeat module and dilated layer which enhances the discriminability of the network by capturing complementary context information of hand postures. The effectiveness of Fit-Hand is evaluated by using subject dependent (SD) and subject independent (SI) validation setup over seven benchmark datasets: MUGD-I, MUGD-II, MUGD-III, MUGD-IV, MUGD-V, Finger Spelling and OUHANDS, respectively. Furthermore, to investigate the deep insights of the proposed Fit-Hand framework, we performed ten ablation study

I. INTRODUCTION

Hand gestures represent specific finger and hand movements that depict a particular message in non-verbal communication. Gestures also reinforce verbal communication by conveying human’s intentions in certain conversations. Hand gesture recognition is perceptual computing that allows machines to identify hand gestures and execute the relevant action. The current situation of coronavirus (COVID19) pandemic outbreak has caused sudden need of HGR in various domains such as: consumer electronics market, transit sector, gaming, touch-less smartphones, defence, home automation, robotics, automated sign language translation etc. Thus, there is a need to amalgamate the HGR with AI to design and develop custom-made touch-less interface to carry out daily activities while maintaining physical distance. Thus, a robust HGR system is needed that can work efficiently on memory-limited devices for real-life applications.

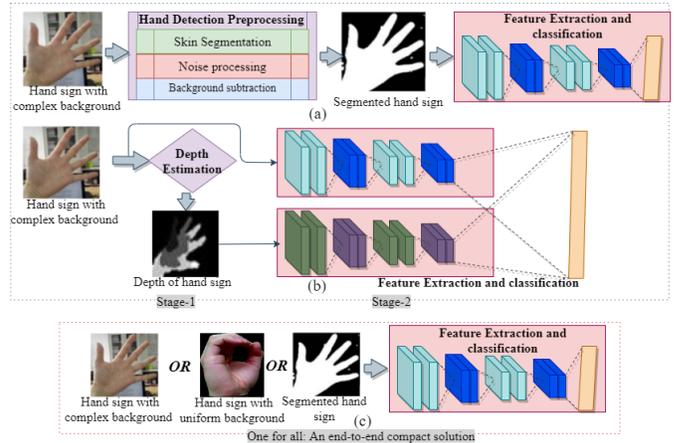


Fig. 1. Architectural comparison between existing (a) [1], [2] -(b) [3], [4] two stage network and proposed one for all: an end-to-end solution for HGR. The visual representation implies the proposed frameworks efficacy to handle the all kind of challenges: complex signs, illumination variations, complex and cluttered backgrounds. While existing HGR frameworks aids extra computation to deal with the complex backgrounds.

The HGR methods can be divided into two broad categories: sensor and vision-based techniques. Sensor based techniques [5], [6] used gloves and other electronic devices to measure the joint angles position of the fingers, position of the hands to extract the features of hands. Although, glove-based techniques have sufficient cues to identify hand gestures, but gloves with wires and sensors are too expensive and makes people uncomfortable to wear. However, vision-based techniques can analyze hand gestures in a non-intrusive manner without any involvement of gloves and electromagnetic devices. Moreover, vision-based techniques are divided into two categories: 3D hand and appearance-based model. 3D hand-based models [7], [8] represent hand structure by defining geometrical shapes of hands as joint angles of wrist, joints of fingers, space between fingers etc. In appearance-based methods, texture features are extracted from visual appearance of hands. Appearance-based methods further can be split into two groups: pre-designed and learning based. Pre-designed based models encoded hand structure by imposing handcrafted feature descriptors. Pre-designed feature descriptor [9], [10] has got promising results

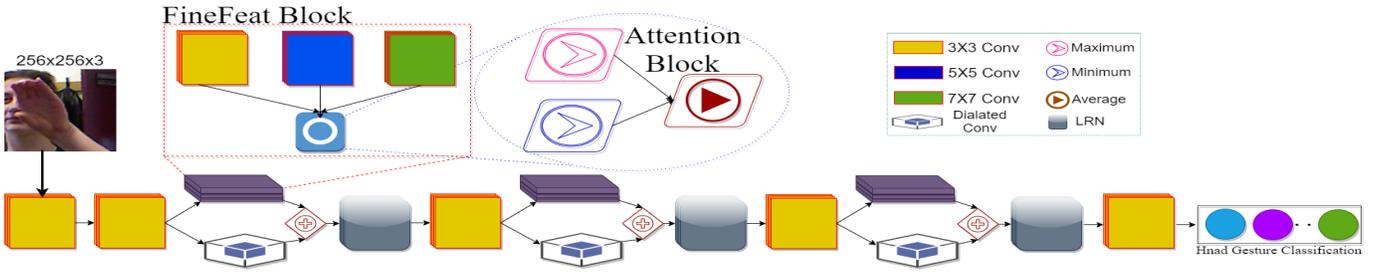


Fig. 2. Architecture of the proposed Fit-Hand.

in field of computer vision. However, these approaches fail to derive efficient feature in real time scenarios with variant challenges as noise, complex background conditions, low resolution and initial contour sequences. Whereas, learning based approaches capture specific features of hand gestures by updating filter weights gradually.

Recently, convolutional neural networks (CNN) have shown tremendous performance in different computer vision fields like object detection, man pose estimation, anomaly detection, face verification, emotion recognition and many more. There are many deep CNN models postulated in literature such as Inception V3 [11], ResNet [12], ResNet Inc [13], DenseNet [14], Mobilenet [15], Mobilenet V2 [16] and NasMobileNet [17] etc. Even in HGR field, CNN based approaches [18], [19], [20], [21] also have been shown impressive performance. These CNN based HGR frameworks are followed two-stage pipeline. Where, in first stage they have utilized handcrafted techniques to computes optical flow or heat maps. Whereas in second stage, CNN architectures are used for feature extraction and classification. Two-stage framework based approaches have gain impressive attention and successfully resolved the problem of multi-view, noise, low resolution etc. However, the performance of the these two stage models is limited by the pre-design techniques that are strongly dependent on prior expertise. To overcome the limitation of two-stage frameworks, advanced CNN based approaches [18], [19], [20] have been introduced for hand gesture recognition without including any pre-designed feature extraction method. However, most of the advanced CNN models are designed to solve specific problems like [18] works for black and white hand posture images, [19], [20] color images. Thus above mentioned work is not proving generic solution for all types of hand gesture images. Moreover, existing approaches were evaluated over subject dependent setup and have been gained high accuracy. However, they perform poorly when evaluated for unseen subjects' hand gestures (subject independent setup). The existing CNN networks also need huge computation with large parameters and incapable to work on handheld and portable devices.

Inspired with the above challenges, in this paper, one for all: an end-to-end compact network: fine grained feature attentive network for hand gesture recognition is introduced. The Fit-Hand model consist of two units: FineFeat module and dilated Conv layer to extract the effective fine features

and abstract features, respectively. Furthermore, to learn the complementary information, resultant feature maps of FineFeat and dilated layer are combined by utilizing integrated layer. Complementary features allows network to learn both micro and high level features, which makes Fit-Hand a robust method to deal with black and white (segmented) as well as color full complex background hand gesture images. The proposed Fit-Hand also reduce the complexity of the HGR model by eliminating the need of hand segmentation. The visual demonstration of the comparison between state-of-the-art two-stage approaches and proposed one for all: FitHand framework is presented in Fig. 1. The main contribution of the proposed network is summarized as follows.

- 1) We proposed a light weighted end-to-end fine grained feature attentive network for hand gesture recognition as one solution for different challenges.
- 2) A FineFeat module is proposed to extract abstract and detailed variation of the hand postures by utilizing both global and local receptive field information.
- 3) A novel attention mechanism is introduced to preserve effective edge information by utilizing averaging behavior of multi-scale receptive responses.
- 4) The dilated Conv layer is used to extract high-level feature of the hand-posture and improves the discriminability of the Fit-Hand.
- 5) Effectiveness of proposed Fit-Hand is validated on seven benchmark datasets: MUGD-I, MUGD-II, MUGD-III, MUGD-IV, MUGD-V, Finger Spelling (ASL), OUHANDS, with subject dependent and subject independent evaluation strategies.

II. RELATED WORK

With recent advent of technologies CNN based approaches has gained good achievements in hand gesture recognition. Jose et al. [22] designed two different CNN frameworks based on LeNET architecture [23] to extract the prominent features of hand gesture structures. Further, Oyebade et al. [2] applied CNN network with auto-encoder to represents the features of hand gestures. Sérgio et al. [24] designed a feature fusion-based convolutional neural network (FFCNN) that incorporated auxiliary features extracted by gabor with CNN network. Moreover, Dadashzadeh et al. [25] proposed a two stage fusion network named as HGR Net. Where, in

first stage, hand regions are detected by applying pixel-level semantic segmentation and second-stage network comprises two-stream CNN to determine the label of hand gesture. Furthermore, multi-task information sharing based approaches [26], [27] has been introduced for hand pose estimation. They extract the features of hands by decomposing them into sub-task through 2D and 3D- heat maps. Mohanty et al. [20] proposed a deep learning model named as DeepGestures for static hand gesture recognition. The DeepGestures model have been designed to handle the various challenges like variation in hand sizes, spatial location variations in the image and complex clutter background. Neethu et al. [28] introduced a CNN based hand gesture detection and recognition framework. Where, first they utilized mask images for hand region extraction and then segment fingers of images from the image though CNN. Further, the adaptive histogram equalization technique is used for image enhancement. Finally, the segmented fingers are fed to CNN model for hand gesture classification. Zhan et al. [18] introduced a CNN model to solve the black and white hand gestures. Furthermore Islam et al. [29] utilized augmentation techniques and increase the hand gesture data sample to enhance the performance of CNN network. Adithya et al. [19] introduce a CNN model for static hand gesture recognition without including any segmentation technique.

III. PROPOSED METHOD

Various researchers have exploit the learning capabilities of the pre-trained models [30], [31] for HGR. Some of the existing approaches [32], [33] have taken advantages of pre-designed feature descriptors and aid in CNN models to boost their performance. While, some of the CNN networks have been deigned to learn the features for specific hand postures. Furthermore, some other HGR approaches have achieve good results, but require huge computation cost. All above explained aspects limit the performance of HGR in practical scenarios. This motivated us to design a generic and portable end-to-end CNN model for HGR which does not have dependency on neither pre-designed descriptors nor pre-trained weights. The detailed architecture of the proposed network is demonstrating in Fig. 2.

Primarily network employ two consecutive Conv layers with 3×3 sized filters to extract variation patterns of hand poses. Let $I(l, m)$ be an input image and $\eta_S^{u,v,d}(o)$ represents Conv function, where S implies for stride, d is depth, u and v represent the size of filter. Then response features R_f of first two layers are calculated by Eq 1.

$$R_f = \eta_2^{3,3,32} \left\{ \eta_2^{3,3,32} \{I(l, m)\} \right\} \quad (1)$$

Further resultant feature maps are simultaneously forwarded to fine feature extraction (FineFeat) module and dilated Conv layer to preserve contextual information of hand postures.

1) *Fine Feature Extraction Module*: The aim of designing FineFeat module is to preserve fine grained edge information for discriminative feature representation of hand postures. The FineFeat module mainly comprises of three laterally connected multi-scale Conv of size 3×3 , 5×5 and 7×7 with minimal

parameters as show in Fig. 2. The multi-scale filters are liable to capture scale invariant features with multi-scale receptive fields. Further, attention block is employed to fetch only effective edges and neglects others by establishing averaging concept over response multi-receptive fields. Response of FineFeat (Im_f^d) module is calculated by using Eq. 2.

$$Im_f^d = \delta \{ \eta_1^{7 \times 7 \times d}(R_f), \eta_1^{5 \times 5 \times d}(R_f), \eta_1^{3 \times 3 \times d}(R_f) \} \quad (2)$$

where, d represents depth of the Conv. Filters. attention block $\delta(o)$ is calculated by using Eq. (3 – 5).

$$\delta(f_1, f_2, f_3) = \varphi(f_1, f_2, f_3) + \min(\gamma(f_1, f_2, f_3)) \quad (3)$$

$$\gamma(f_1, f_2, f_3) = |\varphi(f_1, f_2, f_3) - (f_1, f_2, f_3)| \quad (4)$$

$$\varphi(f_1, f_2, f_3) = \frac{1}{2} (\max(f_1, f_2, f_3) + \min(f_1, f_2, f_3)) \quad (5)$$

where, f_1, f_2, f_3 are implies Conv layer holding multi-scale filters of size 3×3 , 5×5 and 7×7 respectively.

2) *Dilated Convolution layer*: The dilated Conv layer [34] is embedded in FitHand network to extract global spatial features of hand gestures by refining inputs in high resolution. Dilated Conv layer allows to conserve more comprehensive context knowledge from input with reducing trainable parameters. Kernel size of dilated Conv is calculated by using Eq. 6.

$$R_i = i + (i - 1)(D - 1) \quad (6)$$

Where, i is the kernel size and D represent the dilation rate. For Fit-Hand, we have used the 2 dilation. Moreover, Fit-Hand utilized the integrated layer [12] to accumulate preserved feature maps of FineFeat module and dilated Conv layer. Integrated layer captures distinctive edge features and enhances robustness of Fit-Hand to define the disparities between different types of hand gestures problems. Final outcome of Fit-Hand can be computed by using Eq. (7 – 9).

$$H_f = FC \left[\eta_2^{3 \times 3 \times 128} \left(LRN \{ Im_f^{96}(\chi_1) + Dil_2^{3 \times 3 \times 96}(\chi_1) \} \right) \right] \quad (7)$$

$$\chi_1 = LRN \{ Im_f^{64}(\chi_2) + Dil_2^{3 \times 3 \times 64}(\chi_2) \} \quad (8)$$

$$\chi_2 = LRN \{ Im_f^{32}(R_f) + Dil_2^{3 \times 3 \times 32}(R_f) \} \quad (9)$$

where, $Dil_S^{u \times v \times d}$ dilated convolution function, where S implies for stride, d is depth, u and v represent the size of filter. LRN and FC implies for local response normalization and fully connected layer.

Since resultant responses are carrying different scale information, local response normalization (LRN) and L2 normalization is used to normalize them. Therefore, these normalization techniques help to reduce the over-fitting and improve the prediction of the network.

TABLE I
 RECOGNITION ACCURACY ON MUGD, FINGER SPELLING, OUHANDS IN SD AND SI SETUPS. *Here, Fing. Spell, IncV3 and ResNet Inc, stands for finger spelling, inception V3 and resnet inception, respectively.*

Method	SD						SI				
	MUGD					Fing. Spell.	MUGD			Fing. Spell.	OUHANDS
	I	II	III	IV	V		I	II	V		
IncV3 [11] CVPR (2016)	18.0	23.5	11.5	15.0	12.0	51.5	8.05	8.80	3.38	25.5	34.8
ResNet50 [12] CVPR (2016)	82.8	89.2	78.6	83.2	70.2	95.0	66.2	75.8	37.7	49.8	63.4
DeepGestures [20] CVIP (2016)	69.3	83.0	81.8	80.6	56.4	87.2	54.0	75.0	37.2	46.5	46.3
ResNetInc [13] AAAI (2017)	45.0	50.6	39.2	40.0	36.2	69.6	56.7	33.6	3.33	22.8	35.2
Dense121 [14] CVPR(2017)	82.0	87.5	72.0	77.5	72.1	95.0	64.7	69.2	33.8	55.7	64.9
MobileNet [15] (2017)	72.6	81.4	66.0	69.8	68.4	94.8	56.2	67.7	39.4	51.0	59.0
DeepHand [20] Neu. Comp. (2017)	N/A	N/A	N/A	N/A	N/A	91.33	N/A	N/A	N/A	N/A	N/A
MobileV2 [16] CVPR (2018)	50.6	52.2	38.2	32.8	42.6	84.8	38.9	37.2	20.0	45.4	53.0
NASMob [17] CVPR (2018)	21.5	21.5	18.5	19.5	15.5	68.5	3.33	3.60	2.27	37.8	37.9
HandShape [4] ICCCV (2018)	N/A	N/A	N/A	N/A	N/A	90.60	N/A	N/A	N/A	N/A	N/A
HandGes [18] IRI (2019)	6.2	6.4	3.0	3.4	3.2	34.75	10.0	12.0	6.0	22.4	23.0
DeepConv [19] PCS-Elsevier (2020)	71.2	80.4	74.8	86.4	63.4	96.0	54.7	69.0	35.1	43.5	56.0
DDaNet [3] IEEE-Acc. (2020)	N/A	N/A	N/A	N/A	N/A	94.10	N/A	N/A	N/A	N/A	N/A
Fit-Hand	85.2	91.6	98.6	98.8	72.2	95.8	67.0	79.2	40.0	58.8	65.0

A. Comparative Study with Existing Approaches

The existing CNN networks: VGGNet [35] and ResNet [12] gain impressive results by using the sequential coupling behavior of Conv layers. However, in a deep dense network, linearly connected Conv layers may drop some salient features due to recurrence of cross-correlation, which has an important role to define a gesture class. Moreover, deep networks are failing to achieve good performance over smaller sized datasets [36]. To resolve this problem, we proposed a light-weighted end-to-end shallow network which is more appropriate in HGR systems. In addition, most of the challenging hand gesture datasets are captured with complex and cluttered backgrounds. Existing approaches needs hand segmentation to remove the complex background. Although, in literature various hand segmentation techniques like hand shape, skin, color segmentation etc. were proposed for hand segmentation. However, the same segmentation technique is not work with all types of backgrounds and limits the practical usability of the HGR. The proposed Fit-Hand utilized the integration layer to collect the complementary context features from FineFeat module and dilated Conv layer. FineFeat module provides fine grained features with the help of attention mechanism, which is able to capture effective edge information. While, dilated Conv layer generates global representation of hand gestures. Therefore, complementary features of FineFeat and dilated layer boosts the robustness of Fit-Hand to extract edges of hand postures and surpass the background information. Thus, the proposed Fit-Hand does not need hand segmentation. Also, Fit-Hand can easily learn the features from segmented or black and white hand gesture images. There, we conclude that Fit-Hand is a generic HGR framework that is capable to learn features in practical scenarios.

Fit-Hand incorporated a novel attention block, which has capability to extract only effective edges from the multi-scaled feature responses. Whereas, existing module inception layer [13] simply concatenates previously extracted scale variant feature maps and let the neural network to learn relevant weights at the time of training. Thus, inception layer increases

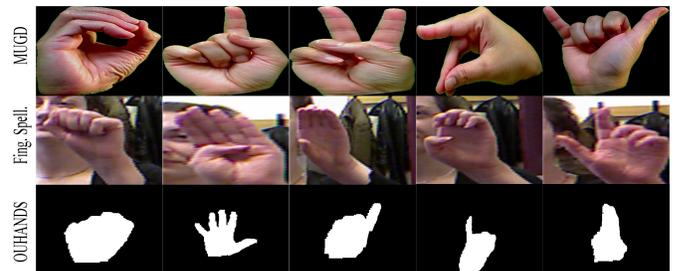


Fig. 3. Sample images of different challenges as complex finger gestures, cluttered backgrounds, segmented gestures etc. in datasets: (a) MUGD, (b) Finger Spelling (ASL) and (c) OUHANDS, respectively.

the complexity of network. Furthermore, Fit-Hand exploits the effectiveness of the dilated Conv layer and preserved the global context features of hand gestures. Moreover, Fit-Hand embedded Conv layer with stride 2, to down-sample input size instead max pooling to preserve minute variation information. Max pooling executed max function to scale down the input size, which eliminates the minute edge features. Sometimes, small variations in a gesture may change the interpretation of its class, thus micro level edge variation representation is also playing a significant role to define a gesture. In literature some studies [37], [38] have validated that Conv with stride instead of pooling adds inter-feature dependencies and improves the learnability of neurons.

IV. EXPERIMENTAL SETUP AND ANALYSIS

In this section, we examined the proposed network for HGR on seven benchmark databases/datasets: massey university gesture dataset part-I (MUGD-I), MUGD-II, MUGD-III, MUGD-IV, MUGD-V [39], Finger Spelling (American Sign Language) [40] and OUHANDs [41]. The quantitative and graphical results in terms of recognition accuracy and F1-Score is verified with the state-of-the-art methods for HGR. The qualitative results are demonstrated to visualise the effectiveness of Fit-Hand as compare to existing HGR approaches. Further, ten supplementary experiments are conducted for ablation study

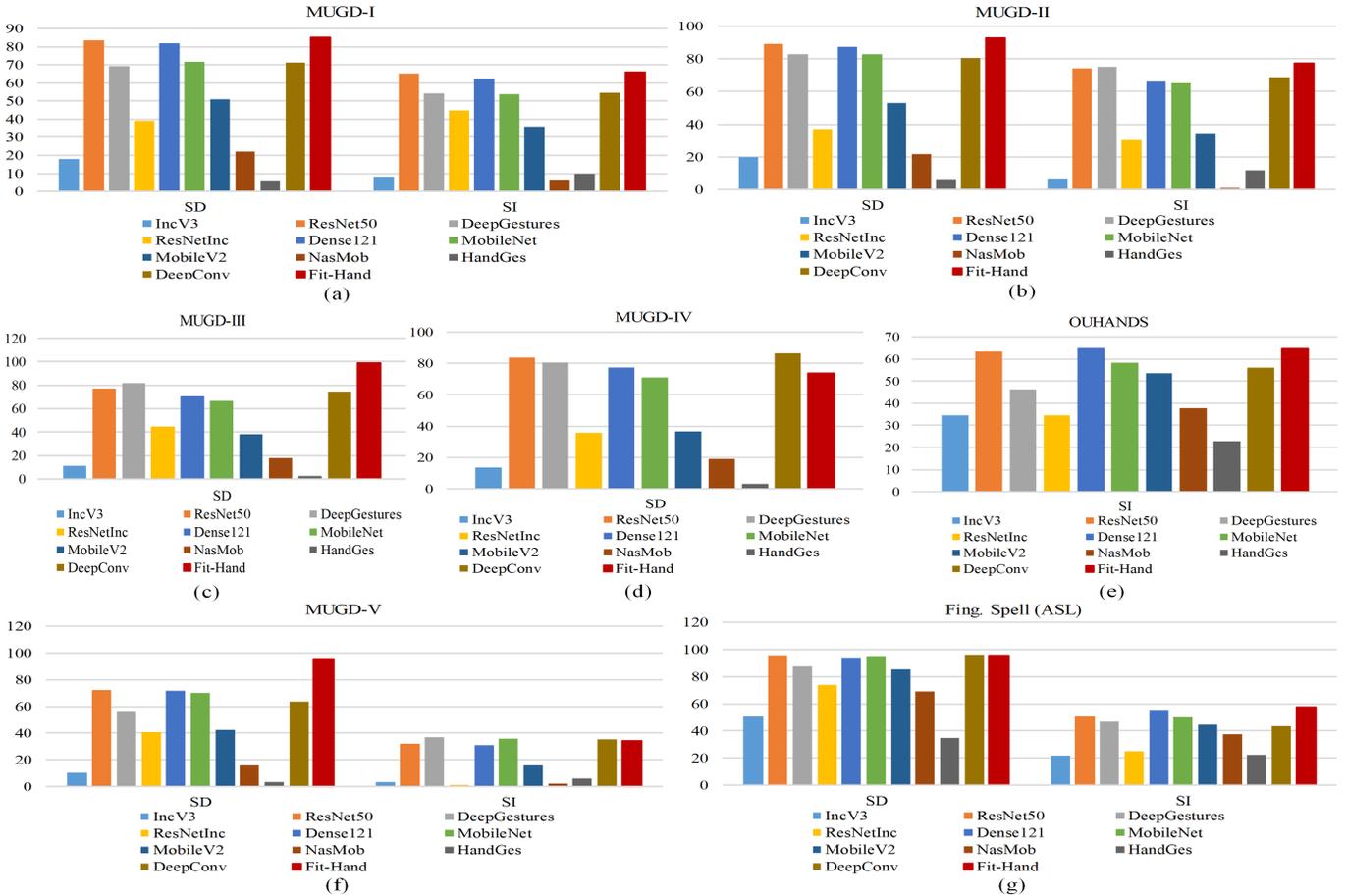


Fig. 4. The graphical plot representing the comparative analysis between existing: IncV3, ResNet50, DeepGestures, ResNetInc, Dense121, MobileNet, MobileV2, NASMob, HandGes, DeepConv and proposed: Fit-Hand in-terms of F1-Score over (a) MUGD-I, (b) MUGD-II, (c) MUGD-III, (d) MUGD-IV, (e) OUHANDS, (f) MUGD-V and (g) Finger Spelling (ASL) datasets for SD and SI experimental settings, respectively.

to validate the effectiveness of each module in the proposed method on OUHANDS. Furthermore, complexity analysis between proposed and state-of-the-art models is represented to validate the portability of the Fit-Hand.

A. Implementation details

All experiments of Fit-Hand are conducted using Keras open-source deep-learning library with the Tensor flow in the backend. The cross-entropy is used as the cost function and the SGD optimizer is used for optimization. The Fit-Hand is trained with learning rate 0.0001 for all experiments. The input image size has been fixed with 256×256 for training and inference of the model. The Nvidia GeForce RTX 2080 GPU with Xeon processor, 16-core CPU, and 11 GB RAM under Cuda 10.0. on Tensorflow-GPU 2.0.0 is used for the experiments.

Moreover, to examine the effectiveness of the proposed Framework, we have compared our results with other state-of-the-art approaches. The researchers have taken up various dataset selection procedures and experimental settings. Therefore, it is hard to make valid comparison between the various published results. To ensure fair comparison of HGR networks:

DeepGestures, DeepConv and HandGes, we have implemented all of them according to our experimental setups. In addition to general networks: Inception V3, ResNet50, ResNet-Inception, DenseNet, MobileNdet, MobileNetV2, NASMobileNet, we have fine-tuned pre-trained weights with our experimental hyper-parameters over 10 epochs. As all versions of MUGD datasets are limited in number, and deep Convolutional neural networks require a large database to learn the most significant features, the datasets are augmented offline to enhance the generalization of the model and prevent overfitting. The following transformations are applied for data augmentation: rotation in between $[-45^\circ, 45^\circ]$ with increment of 15° , horizontal flip and histogram equalization. Finally, one image instance is converted into 10 images.

B. Experimental Setup

In literature most of the researchers utilized N-fold cross validation scheme for evaluation. In N-fold cross validation, datasets are divided into random N folds, where N-1 folds are used for training purpose and one fold is used for inference. The same procedure is followed for each fold and average of all folds are considered as final results. However, N-fold

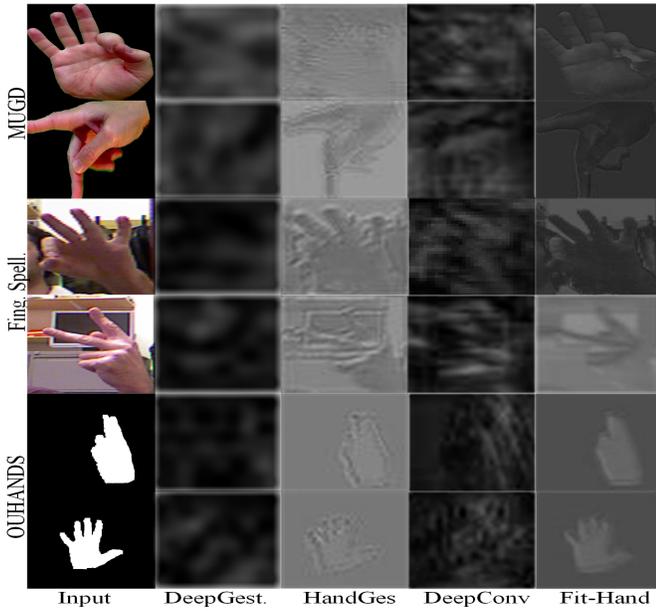


Fig. 5. The qualitative comparison between the feature maps generated by state-of-the-art HGR network: DeepGestures, HandGes, DeepConv and proposed Fit-Hand, over six different gestures.

cross validation strategy is a subject dependent evaluation due to random division of folds and not ensure the performance of models for unseen data. Therefore, these approaches have been gained high accuracy for seen data samples. However, they perform poorly when evaluated for unseen subjects' hand gestures (subject independent setup) and not suitable for real time data validation. Thus, for fair performance of the Fit-Hand we have adopted two validation schemes: subject- dependent and -independent. In subject dependent (SD), datasets are randomly partitioned into 80:20 ration such that 80% dataset is processed for training and 20% for testing set. While, in subject independent (SI), three subjects' hand gestures are used in training and remaining hand gestures are used for inference for MUGD-I, MUGD-II and Finger Spelling datasets. For MUGD-V hand gesture, one subject is included in training set and other one used for testing purpose. The data division for SI is done purely in mutually exclusive manner. Moreover, OUHANDS dataset contains two parts: training and testing set with 2000 and 1000 images, respectively.

C. Quantitative Analysis

This section demonstrates the effectiveness of proposed network over all datasets: MUGD, Finger Spelling and OUHANDS, in terms of recognition accuracy over two experimental setups: SD and SI respectively. Comparative analysis of existing and Fit-Hand is tabulated in Table I. Specifically, Fit-Hand achieved 15.9%, 8.6%, 16.8%, 18.2%, 15.8% and 14%, 11.2%, 23.8%, 12.4%, 8.8% more accuracy as compared to DeepGestures and DeepConv HGR models over MUGD dataset for part I-V in SD setup, respectively. Similarly, in SI setup, Fit-Hand gained 13%, 4.2%, 2.8%, 12.3%, 18.7% and 12.3%, 10.2%, 4.9%, 15.3%, 9% more accuracy as compare to

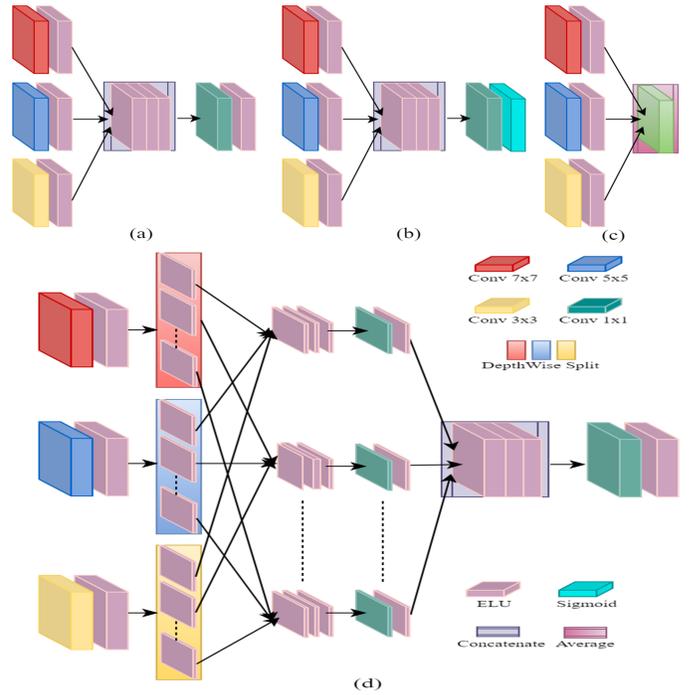


Fig. 6. Different structures of FineFeat module for ablation study a) FineFeat with concatenation (FineFeat_Cat), b) FineFeat with concatenation and Sigmoid (fineFeat_CatSig), c) FineFeat with average (FineFeat_Avg) and d) FineFeat with deep median (FineFeat_DpMed).

DeepGestures and DeepConv models for MUGD-I, MUGD-II, MUGD-V, Finger Spelling and OUHANDS datasets. Moreover, performance of Fit-Hand in terms of F1-Score are graphically demonstrated in Fig. 4. From the Table I and Fig. 4 results it is clear that all HGR methods: proposed as well as state-of-the-art generates high results in subject dependent setup as compare to subject independent. Moreover, the proposed FitHand outperformed the two-stage networks; DeepHand, HandShape and DDaNet with 4.45%, 5.2% and 1.7% high accuracy over ASL finger spelling dataset. From the results, it is validated that proposed framework is robust to all kind of challenges presents in the HGR, which reflect the efficacy of the model to real-life applications. Also from the results, it is proven that subject independent validation strategy is more significant as compare to subject dependent to examined the performance of any CNN model. In addition some methods like Inception V3, NasMobile and HandGes are under-fitted and not suitable for small size datasets.

D. Qualitative Analysis

This section elaborates the effectiveness of Fit-Hand through the visual representation of neurons. Fig. 5 depicted response maps of two different gestures captured at first Conv layer of ResNet-50, ResNet-Inc, MobileNet, MobileNet V2 and Fit-Hand. To represent the significance of all response maps, we have calculated mean response for each network. From the figure it is clear that Fit-Hand extracts more fine edge variations and highlighted regions like figure lines, palm lines,

TABLE II

ABLATION RESULTS IN TERMS OF ACCURACY AND COMPLEXITY. *Here, Param, Acc, M, S and MB stands for parameters, accuracy, millions, seconds and megabytes.*

Method	Acc.	#Param (M)	#Mem. (MB)	#Time (S)
Fit-Hand_WImp	61.7	0.5	3.3	1.08
Fit-Hand_WDil	60.9	1.4	11.9	2.32
Fit-Hand_WL	47.3	1.8	12.9	3.26
Fit-Hand_2Stack	51.7	1.0	6.5	.82
Fit-Hand_4Stack	62.0	3.4	27.1	5.07
Fit-Hand_Kul	61.7	1.8	12.9	2.75
FineFeat_Cat	57.9	1.6	13.3	2.58
FineFeat_CatSig	60.1	1.6	13.3	2.66
FineFeat_Avg	58.6	1.5	12.9	2.63
FineFeat_DpMed	58.3	1.6	16.0	30.8
Fit-Hand	65.0	1.8	12.9	4.00

thumb articulates etc, which plays a significant role to define distinctiveness between hand postures. From above, it is clear that Fit-Hand has capability to preserve prominent features of hand gestures. Therefore, we can conclude that Fit-Hand has preserved more relevant feature responses to outperform the existing CNN based networks ResNet50, Res-Net-Inc, MobileNet and MobileNet V2 for different hand postures.

E. Ablation Study

In order to investigate the deep insights of Fit-HandeNet, we have conducted ten more ablation experiments for detail study as represented in Table II over OUHANDS Dataset. This section fully explores the contribution of each module (FineFeat module, dilation layer, L2-Normalization, loss function and attention block) of the network in terms of performance and network complexity. Specifically, to validate the importance of FineFeat module and dilated layer in Fit-Hand, we have evaluated results for Fit-Hand without FineFeat module (Fit-Hand_WImp) and Fit-Hand without dilated layer (Fit-Hand_WDil). From the Table II, it is clear that both FineFeat module and dilated layer play a significant role in Fit-Hand and improve the performance of the network. To analyze the role of L2 normalization, results are evaluated by dropping L2 normalization (Fit-Hand_WL). Evaluated results in Table II, validated the effect of L2 normalization with high performance.

To investigate that how three stacks of FineFeat modules help to learn the adequate information in Fit-Hand, we have performed two supplementary experiments with 2 FineFeat module (2 stacked) and Fit-Hand with 4 FineFeat module (4 Stacked). From the results tabulated in Table II, it is concluded that, proposed Fit-Hand outperforms other dept combinations of FineFeat module. Moreover, to validate the effectiveness of cross-entropy loss function, we have computed results for Fit-Hand by replacing cross-entropy by Kullback Leibler Divergence Loss (Fit-Hand_Kul). Computed results confirmed that cross-entropy loss function is most suitable for hand gestures classification in Fit-Hand.

Furthermore, to examine the performance of attention block, we have implemented four different FineFeat modules by replacing pivot with a) concatenation (FineFeat_Cat), b)

TABLE III

COMPLEXITY ANALYSIS COMPARISON BETWEEN EXISTING AND PROPOSED FIT-HAND NETWORK. *Here, M, K, MB, KB and S represents millions, thousands, megabytes, kilobytes and seconds, respectively.*

Method	#Param	#Mem.	#Time (S)
IncV3 [11]	22M	179.3MB	17.64
ResNet50 [12]	31M	208.4MB	18.17
DeepGestures [20]	10K	128KB	0.48
ResNetInc [13]	4M	40.5MB	60.93
Dense121 [14]	7.5M	61.3MB	24.34
Mobile [15]	5M	36.5MB	8.26
MobileV2 [16]	3.5M	24.4MB	15.26
NASMob [17]	4.7M	41.7MB	45.65
HandGes [18]	16K	252KB	0.67
DeepConv [19]	1M	1.32MB	1.00
Fit-Hand	1.8M	12.9MB	4.00

concatenation with sigmoid (FineFeat_CatSig) , c) average (FineFeat_Avg) and d) deep median (FineFeat_DpMed), as shown in Fig. 6. From Table II, it is evident that proposed FineFeat module with attention block outperforms all other combinations of the FineFeat module.

F. Complexity Analysis

This section provides a comparative analysis of the computational complexity between the existing and proposed network. The total number of parameters, memory space and testing time involved in each network are tabulated in Table III. The proposed Fit-Hand has very lesser number of parameters 1.8M as compared to other state-of-the-art models like: Inception V3: 22M, ResNet50: 31M, ResNetInception: 4M, DenseNet 121: 7.5M, MobileNet: 5M, MobileNet V2: 3.5M and NASNet Mobile: 4.7M. Moreover, Fit-Hand trainable model captures smaller memory storage as compared to others. Fit-Hand query response time is also very less as compare to existing approaches. However, from the Table III, it is also observable that complexity of some existing HGR models [[19][20][18]] is less as compared to proposed Fit-Hand. However, it is clear from Table I that, these approaches are not providing generic solution as they failed to maintain good performance for all datasets. Therefore, on the basis of experimental and computational complexity results, we conclude that the proposed FitHand framework is a portable and generic solution for the different hand gestures.

V. CONCLUSION

We proposed an one for all: end-to-end compact solution named as Fit-Hand: fine grained feature attentive network for HGR, which is responsible to identify distinct classes of hand gestures. Fit-Hand contains two main blocks: FineFeat module and dilated Conv layer. FineFeat module conserves features of minute as well as major edge variation regions and further employ a attention block that has ability to fetch only pertinent features. Similarly, dilated layer is incorporated to capture global features. Further, Integrated layer added feature map of both blocks and enhance the learnability of Fit-Hand. Cohesively both layers allow Fit-Hand to learn the imperative features of hand postures and define disparities between them.

Furthermore, variants of Fit-Hand were evaluated to verify the effectiveness of proposed Fit-Hand.

REFERENCES

- [1] H.-Y. Chung, Y.-L. Chung, and W.-F. Tsai, "An efficient hand gesture recognition system based on deep cnn," in *2019 IEEE International Conference on Industrial Technology (ICIT)*. IEEE, 2019, pp. 853–858.
- [2] O. K. Oyedotun and A. Khashman, "Deep learning in vision-based static hand gesture recognition," *Neural Computing and Applications*, vol. 28, no. 12, pp. 3941–3951, 2017.
- [3] S.-H. Yang, W.-R. Chen, W.-J. Huang, and Y.-P. Chen, "Ddanet: Dual-path depth-aware attention network for fingerspelling recognition using rgb-d images," *IEEE Access*, 2020.
- [4] A. Rakowski and L. Wandzik, "Hand shape recognition using very deep convolutional neural networks," in *Proceedings of the 2018 International Conference on Control and Computer Vision*, 2018, pp. 8–12.
- [5] D. J. Sturman and D. Zeltzer, "A survey of glove-based input," *IEEE Computer Graphics and Applications*, vol. 14, no. 1, pp. 30–39, 1994.
- [6] L. Santos, N. Carbonaro, A. Tognetti, J. L. González, E. De la Fuente, J. C. Fraile, and J. Pérez-Turiel, "Dynamic gesture recognition using a smart glove in hand-assisted laparoscopic surgery," *Technologies*, vol. 6, no. 1, p. 8, 2018.
- [7] J. M. Rehg and T. Kanade, "Visual tracking of high dof articulated structures: an application to human hand tracking," in *European conference on computer vision*. Springer, 1994, pp. 35–46.
- [8] T. Heap and D. Hogg, "Towards 3d hand tracking using a deformable model," in *Proceedings of the Second International Conference on Automatic Face and Gesture Recognition*. Ieee, 1996, pp. 140–145.
- [9] P. K. Pisharady, P. Vadakkepat, and A. P. Loh, "Attention based detection and recognition of hand postures against complex backgrounds," *International Journal of Computer Vision*, vol. 101, no. 3, pp. 403–419, 2013.
- [10] M. Jasim and M. Hasanuzzaman, "Sign language interpretation using linear discriminant analysis and local binary patterns," in *2014 International Conference on Informatics, Electronics & Vision (ICIEV)*. IEEE, 2014, pp. 1–5.
- [11] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2818–2826.
- [12] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [13] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. Alemi, "Inception-v4, inception-resnet and the impact of residual connections on learning," *arXiv preprint arXiv:1602.07261*, 2016.
- [14] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 4700–4708.
- [15] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," *arXiv preprint arXiv:1704.04861*, 2017.
- [16] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "Mobilenetv2: Inverted residuals and linear bottlenecks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 4510–4520.
- [17] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le, "Learning transferable architectures for scalable image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 8697–8710.
- [18] F. Zhan, "Hand gesture recognition with convolution neural networks," in *2019 IEEE 20th International Conference on Information Reuse and Integration for Data Science (IRI)*. IEEE, 2019, pp. 295–298.
- [19] V. Adithya and R. Rajesh, "A deep convolutional neural network approach for static hand gesture recognition," *Procedia Computer Science*, vol. 171, pp. 2353–2361, 2020.
- [20] A. Mohanty, S. S. Rambhatla, and R. R. Sahay, "Deep gesture: static hand gesture recognition using cnn," in *Proceedings of International Conference on Computer Vision and Image Processing*. Springer, 2017, pp. 449–461.
- [21] S. W. Yahaya, A. Lotfi, M. Mahmud, P. Machado, and N. Kubota, "Gesture recognition intermediary robot for abnormality detection in human activities," in *2019 IEEE Symposium Series on Computational Intelligence (SSCI)*. IEEE, 2019, pp. 1415–1421.
- [22] C. J. L. Flores, A. G. Cutipa, and R. L. Enciso, "Application of convolutional neural networks for static hand gestures recognition under different invariant features," in *2017 IEEE XXIV International Conference on Electronics, Electrical Engineering and Computing (INTERCON)*. IEEE, 2017, pp. 1–4.
- [23] Y. LeCun *et al.*, "Lenet-5, convolutional neural networks," URL: <http://yann.lecun.com/exdb/lenet>, vol. 20, no. 5, p. 14, 2015.
- [24] S. F. Chevchenko, R. F. Vale, V. Macario, and F. R. Cordeiro, "A convolutional neural network with feature fusion for real-time hand posture recognition," *Applied Soft Computing*, vol. 73, pp. 748–766, 2018.
- [25] A. Dadashzadeh, A. T. Targhi, M. Tahmasbi, and M. Mirmehdi, "Hgr-net: a fusion network for hand gesture segmentation and recognition," *IET Computer Vision*, vol. 13, no. 8, pp. 700–707, 2019.
- [26] X. Nie, J. Feng, and S. Yan, "Mutual learning to adapt for joint human parsing and pose estimation," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 502–517.
- [27] K. Du, X. Lin, Y. Sun, and X. Ma, "Crossfonet: Multi-task information sharing based hand pose estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 9896–9905.
- [28] P. Neethu, R. Suguna, and D. Sathish, "An efficient method for human hand gesture detection and recognition using deep learning convolutional neural networks," *Soft Computing*, pp. 1–10, 2020.
- [29] M. Z. Islam, M. S. Hossain, R. ul Islam, and K. Andersson, "Static hand gesture recognition using convolutional neural network with data augmentation," in *2019 Joint 8th International Conference on Informatics, Electronics & Vision (ICIEV) and 2019 3rd International Conference on Imaging, Vision & Pattern Recognition (icIVPR)*. IEEE, 2019, pp. 324–329.
- [30] U. Côté-Allard, C. L. Fall, A. Drouin, A. Campeau-Lecours, C. Gosselin, K. Glette, F. Laviolette, and B. Gosselin, "Deep learning for electromyographic hand gesture signal classification using transfer learning," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 27, no. 4, pp. 760–771, 2019.
- [31] U. Cote-Allard, C. L. Fall, A. Campeau-Lecours, C. Gosselin, F. Laviolette, and B. Gosselin, "Transfer learning for semg hand gestures recognition using convolutional neural networks," in *2017 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*. IEEE, 2017, pp. 1663–1668.
- [32] H.-I. Lin, M.-H. Hsu, and W.-K. Chen, "Human hand gesture recognition using a convolution neural network," in *2014 IEEE International Conference on Automation Science and Engineering (CASE)*. IEEE, 2014, pp. 1038–1043.
- [33] B. Liao, J. Li, Z. Ju, and G. Ouyang, "Hand gesture recognition with generalized hough transform and dc-cnn using realsense," in *2018 Eighth International Conference on Information Science and Technology (ICIST)*. IEEE, 2018, pp. 84–90.
- [34] F. Yu and V. Koltun, "Multi-scale context aggregation by dilated convolutions," *arXiv preprint arXiv:1511.07122*, 2015.
- [35] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [36] A. Schindler, T. Lidy, and A. Rauber, "Comparing shallow versus deep neural network architectures for automatic music genre classification," in *FMT*, 2016, pp. 17–21.
- [37] M. Verma, S. K. Vipparthi, G. Singh, and S. Murala, "Learnnet: Dynamic imaging network for micro expression recognition," *IEEE Transactions on Image Processing*, vol. 29, pp. 1618–1627, 2019.
- [38] M. Verma, S. K. Vipparthi, and G. Singh, "Hinet: Hybrid inherited feature learning network for facial expression recognition," *IEEE Letters of the Computer Society*, vol. 2, no. 4, pp. 36–39, 2019.
- [39] A. Barczak, N. Reyes, M. Abastillas, A. Piccio, and T. Susnjak, "A new 2d static hand gesture colour image dataset for asl gestures," 2011.
- [40] N. Pugeault and R. Bowden, "Spelling it out: Real-time asl fingerspelling recognition," in *2011 IEEE International conference on computer vision workshops (ICCV workshops)*. IEEE, 2011, pp. 1114–1119.
- [41] M. Matilainen, P. Sangi, J. Holappa, and O. Silvén, "Ouhands database for hand detection and pose recognition," in *2016 Sixth International Conference on Image Processing Theory, Tools and Applications (IPTA)*. IEEE, 2016, pp. 1–5.