

Multi-Scale Region of Interest Pooling for Facial Action Unit Detection

BY

GUGLIELMO MENCHETTI

B.S, Università degli studi di Firenze, Firenze, Italy, 2017

THESIS

Submitted as partial fulfillment of the requirements
for the degree of Master of Science in Computer Science
in the Graduate College of the
University of Illinois at Chicago, 2020

Chicago, Illinois

Defense Committee:

Ahmet Enis Cetin, Chair and Advisor

Ugo Buy

Elena Zheleva

Mark Carman, Politecnico di Milano

ACKNOWLEDGMENTS

GM

TABLE OF CONTENTS

<u>CHAPTER</u>	<u>PAGE</u>
1 INTRODUCTION	1
1.1 Motivations	3
1.2 Document Structure	3
2 BACKGROUND	4
2.1 Convolutional Neural Networks	4
2.2 ResNet	6
2.3 Feature Pyramid Networks	7
2.4 Region of Interest Pooling	9
2.5 Face Alignment	12
2.6 Optimization	15
2.7 Transfer Learning	16
3 PREVIOUS WORK	17
3.1 Deep Region and Multi-label Learning for Facial Action Unit Detection	17
3.2 EAC-Net: Deep Nets with Enhancing and Cropping for Facial Action Unit Detection	18
3.3 Joint Facial Action Unit Detection and Face Alignment	20
3.4 Facial Action Unit Detection Using Attention and Relation Learning	22
4 MODEL STRUCTURE	24
4.1 Model Overview	24
4.2 Feature Pyramid Network Module	26
4.3 Regions of Interest Module	26
5 PRE-PROCESSING, TRAINING AND VALIDATION	29
5.1 Datasets	29
5.2 Pre-processing	30
5.2.1 Face Alignment	30
5.2.2 Extraction of RoIs	31
5.3 Training	32
5.4 Validation	34
6 RESULTS	37
6.1 Evaluation Metrics	37

TABLE OF CONTENTS (continued)

<u>CHAPTER</u>		<u>PAGE</u>
	6.2 Results on DISFA	38
	6.3 Ablation Studies	40
7	CONCLUSIONS AND FUTURE DIRECTIONS	43
	CITED LITERATURE	45
	VITA	51

LIST OF TABLES

<u>TABLE</u>		<u>PAGE</u>
I	ResNet-50. The fractions 1/4, 1/8, 1/16 and 1/32 refer to the height and width dimensions ratio with respect to the input image size	25
II	RoIs and face locations	31
III	Normalization values (mean and standard deviation) for each channel	33
IV	F1-score (in %) on DISFA dataset	40
V	Accuracy score (in %) on DISFA dataset	41
VI	F1-score (in %) on DISFA dataset for different model configurations	42
VII	F1-score (in %) on the DISFA dataset removing some training techniques	42

LIST OF FIGURES

<u>FIGURE</u>		<u>PAGE</u>
1	Examples of Action Units in the FACS system	2
2	An example of convolution with a filter of size 3x3, stride 1 and no padding	6
3	Residual Block introduced in ResNet	8
4	FPN	10
5	Panoptic FPN	10
6	RoI Pooling description	12
7	Example of the alignment process	14
8	High-level overview of the backbone model with enhancing layers .	19
21		
10	Architecture overview	25
11	FPN module used in our architecture	27
12	ROI module	28
13	Example of the extracted RoIs. In red, green and blue the RoIs of the upper, middle and lower face are showed, respectively.	32
14	Occurrences of AUs and Neutral Faces in DISFA	34
15	Occurrences of AUs and Neutral Faces in training data	35
16	Examples of the AUs in DISFA	39

LIST OF ABBREVIATIONS

AU	Action Unit
FACS	Facial Action Unit Coding System
DL	Deep Learning
SOTA	State-Of-The-Art
CNN	Convolutional Neural Network
CV	Computer Vision
MLP	Multi-Layer Perceptron
FPN	Feature Pyramid Network
RoI	Region of Interest
ML	Machine Learning
SGD	Stochastic Gradient Descent
LR	Learning Rate
DRML	Deep-Region and Multi-Label Learning
EAC	Enhancing And Cropping
FC	Fully Connected
ARL	Attention and Relation Learning
BB	Bounding Box

SUMMARY

Units in a single image or a video. In this document, we focus on single-frame action unit recognition, in which we aim at detecting the action units for a single frame.

In this work, we will examine the current state-of-the-art on Facial AU detection and we will introduce a novel Deep Learning architecture, that achieves competitive results with respect to the current state-of-the-art on a public benchmark.

We will compare the obtained results with other state-of-the-art methods and we investigate the effect of the different components of the model enclosing various ablation studies.

In conclusion, we will suggest possible future improvements for the current architecture.

CHAPTER 1

INTRODUCTION

Facial expression is the natural mean used by humans to express their intentions and emotions. The analysis of facial expressions captures the attention of researchers, due to its wide range of potential applications such as medical (pain detection [1][2], mental health diagnosis [3]), in education [4] and entertainment [5].

However, detecting facial expressions may be challenging due to the ambiguities of expressions. This ambiguity can be solved thanks to the representation of expressions towards the use of the Facial Action Unit Coding System (FACS) [6].

The FACS defines a taxonomy for 46 human facial movements by their change in appearance on the face. The system was first published in 1978 then updated in 2002 and since then has been used by human coders to characterize all the possible facial expressions, by deconstructing them in Facial Action Units (AUs) which are fine-grained facial movements defined for some local regions of the face. The process of AU annotation, however, requires a deep knowledge of the FACS system and a long time to manually annotate entire videos.

In recent years, many FACS coded databases [7][8][9][10] have been released by universities and laboratories, allowing researchers to develop new techniques for automatically annotate facial AUs. Moreover, the advances in Computer Vision related tasks brought by Deep Neural Networks, has made researchers apply neural network techniques to solve the task of AUs detection.

Upper Face Action Units					
AU 1	AU 2	AU 4	AU 5	AU 6	AU 7
					
Inner Brow Raiser	Outer Brow Raiser	Brow Lowerer	Upper Lid Raiser	Cheek Raiser	Lid Tightener
*AU 41	*AU 42	*AU 43	AU 44	AU 45	AU 46
					
Lid Droop	Slit	Eyes Closed	Squint	Blink	Wink
Lower Face Action Units					
AU 9	AU 10	AU 11	AU 12	AU 13	AU 14
					
Nose Wrinkler	Upper Lip Raiser	Nasolabial Deepener	Lip Corner Puller	Cheek Puffer	Dimpler
AU 15	AU 16	AU 17	AU 18	AU 20	AU 22
					
Lip Corner Depressor	Lower Lip Depressor	Chin Raiser	Lip Puckerer	Lip Stretcher	Lip Funneler
AU 23	AU 24	*AU 25	*AU 26	*AU 27	AU 28
					
Lip Tightener	Lip Pressor	Lips Part	Jaw Drop	Mouth Stretch	Lip Suck

Figure 1: Examples of Action Units in the FACS system

Nevertheless, the task of AU detection is not considered solved, since in a real-word environment AU recognition systems are sensible to illumination variation, head pose, and subject-dependence. Furthermore, the similarities between Facial Action Units make this task even more difficult and prone to errors in their recognition.

1.1 Motivations

The scope of this work is to examine the current state-of-the-art (SOTA) on Facial AU detection and to develop a novel Deep Learning (DL) architecture to solve this problem.

We designed a novel DL architecture taking inspiration from previous works on AU recognition. In particular, we enhanced the region learning technique, introducing a region pooling layer based on the AUs' Region of Interest (RoI).

Our model achieves comparable performances with respect to the current SOTA methods.

1.2 Document Structure

In Chapter 2 we will give some background on the Deep Learning techniques that are used as the base for our model.

In Chapter 3 we will give a brief description of the state-of-the-art approaches for AU detection.

Chapter 4 describes the structure of our model, providing insight into each single component.

In Chapter 5 we detail the training and validation procedures, we describe the available data and the pre-processing techniques that we applied.

Finally, in Chapter 6 we report the achieved results, along with ablation studies on different components of the model.

CHAPTER 2

BACKGROUND

2.1 Convolutional Neural Networks

Convolutional Neural Networks (CNN)[11] are a class of deep neural networks, that allowed advancements of Neural Networks in many Computer Vision (CV) related tasks, such as object recognition and detection, image classification and segmentation, redefining the state-of-the-art for those tasks.

A CNN takes in input an image and assign importance, in the form of weights and biases, to various aspects in the image. The main idea behind CNNs is to replace the general matrix multiplication used in Multilayer Perceptron (MLP), with the linear operator that is the convolution.

Furthermore, CNNs handle some of the problems of MLPs when working with images. In particular:

- RGB images are naturally represented as 3-dimensional vectors, or, equivalently, as a 2-dimensional matrix of pixels, in which each value is given by a tuple representing the values for each channel (red, green and blue). When dealing with MLP, images must be flattened, hence represented with a vector with a unique dimension. This representation is not able to capture all the information of an image, which reside in its spatial structure. On the other hand, since CNNs do not require to apply this transformation, they reduce

the images to an easier form to process, without losing information about the spatial dependencies.

- The operation of flattening the image leads also to an increase of parameters that must be learned by an MLP. For example, a 430x430 image with 3 channels, contains 184.900 values. If we consider a single hidden layer with 1000 neurons, the number of parameters of the model would be 184.900.000. Furthermore, the increase in the number of parameters, causes the MLP model to be more prone to overfitting and with a poor convergence.
- The objective of convolution is to extract the relevant features such as edges, from the input image. Adding multiple convolutional layers, the architecture is capable to identify and extract both high-level features, giving a general understanding of the image, and low-level features, such as edges.

Before the arrival of CNNs, the operation of transforming the image in vectors of relevant features was done with classical feature extractor techniques such as SIFT [12] and HOG [13].

A convolutional layer shifts a 2-dimensional filter, and performs a matrix multiplication between the kernel and a portion of the image, as shown in Figure 2.

In mathematical terms, the convolution operator $*$ can be expressed by the formula

$$(F * k)(p) = \sum_{s+t=p} F(s)k(t)$$

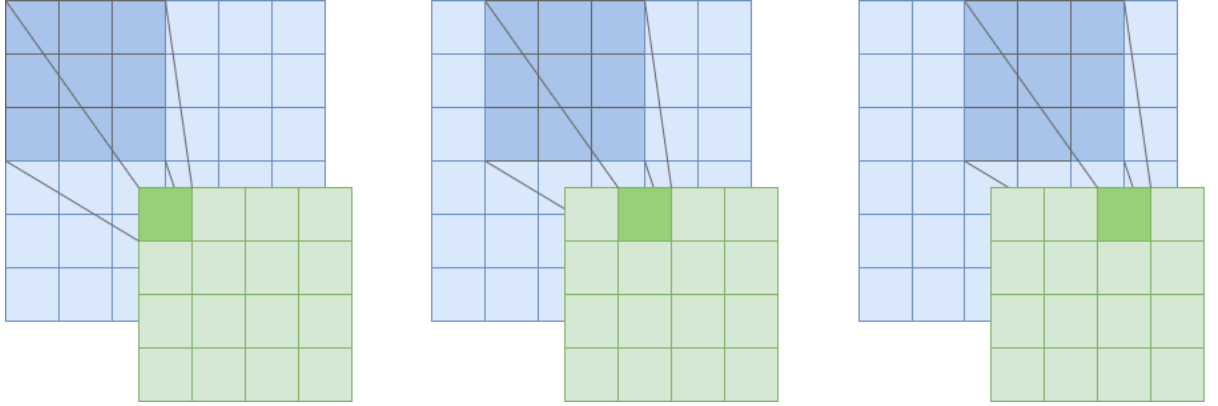


Figure 2: An example of convolution with a filter of size 3x3, stride 1 and no padding

where $F : \mathbb{Z}^2 \rightarrow \mathbb{R}$ is a discrete function and $k : [-r, r]^2 \cap \mathbb{Z}^2 \rightarrow \mathbb{R}$ is a discrete filter of size $(2r + 1)^2$.

A filter is a matrix of parameters which are learned by the architecture. It is relevant to notice that the dimension of the filter, hence the number of parameters, is independent of the dimension of the input.

To reduce the spatial size of the convolved features, hence reducing the computational power required to process the data, a deep convolutional architecture also includes pooling layers. These layers are also useful for extracting dominant features which are rotational and positional invariant.

2.2 ResNet

As stated in the previous paragraph, the introduction of Convolutional Neural Network revolutionized the way of solving many tasks related to Computer Vision.

The first deep convolutional architecture that produced state-of-the-art results in the task of image classification for the ImageNet competition [14] was AlexNet [15]. This architecture was built using 5 convolutional and 3 fully-connected layers.

Despite the good results obtained by this architecture, a study showed that adding more layers to an architecture would create a more complex function, increasing the probability of overfitting, even if some regularization techniques (such as dropout [16] or l2-norms) are applied.

The failure of adding many deep layers was mainly blamed on the vanishing/exploding gradient problem [17], which is the difficulty of updating the values of the weights in a deep architecture.

In ResNet [18] the authors defined a new module for tackling the problem of the vanishing gradient, the Residual Block. A Residual Block consists of multiple layers, in which the input of the block is fed into the next layer and also added to the output of the last layer of the block (skip connection), as depicted in Figure 3. The introduction of the skip connection, allows the gradient to flow through the network.

With the introduction of residual blocks, ResNet defined the new state-of-the-art for the ImageNet classification challenge. Furthermore, the pre-trained model is nowadays widely used as a feature extractor for related tasks, or as a base for transfer learning, as described in 2.7.

2.3 Feature Pyramid Networks

Feature Pyramid Network (FPN) [19] is a feature extractor model introduced for the task of object detection and instance segmentation, in order to detect or segment object at different scales in a fast and accurate way.

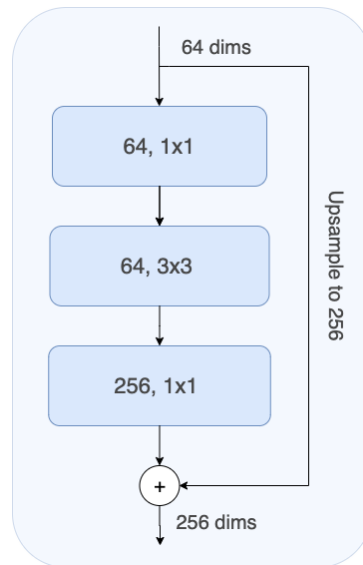


Figure 3: Residual Block introduced in ResNet

The main idea behind this architecture is to create a pyramid of feature maps (multi-scale feature maps) which showed significant improvements with respect to the regular feature maps in various tasks, such as object detection and segmentation.

FPN is composed of two paths:

- **Bottom-Up:** Is a convolutional network for feature extraction, such as ResNet. Going deeper into the model decreases the spatial resolution but increases the semantic information for each layer.
- **Top-Down:** Is used to construct higher resolution layers from a semantic rich layer.

The bottom-up path is responsible for producing dense feature vectors at different resolutions. A dense feature vector is built using the output of a convolution module, that we call C_i , which is used in the top-down path.

In the top-down path, the i -th layer takes as input the sum of C_i (convolved with a 1×1 convolution filter) and the output of the previous layer in the path, upsampled by 2. The output of the different layers in the top-down path is also the output of the architecture.

A similar architecture is the Panoptic Feature Pyramid Network [20]. The only difference is in the top-down path. In particular, the outputs of the convolution module are convolved with a 2×2 convolution filter and upsampled by 2 until it reaches $1/4$ scale. At this point, the top-down path consists of the element-wise summation of the new feature maps.

Figure 4 and Figure 5 show an example of the two architectures.

2.4 Region of Interest Pooling

The Region of Interest (RoI) Pooling layer was first introduced in [21] to solve the burden of dealing with bounding boxes of different sizes in the task of object detection. This layer produces a fixed-sized feature map from inputs with different sizes, applying max-pooling on the inputs.

A RoI pooling layer takes two inputs:

- A feature map obtained from a Convolutional Model
- N bounding box (BB) coordinates representing the 'regions of interest' in the image, usually generated by a Region Proposal network.

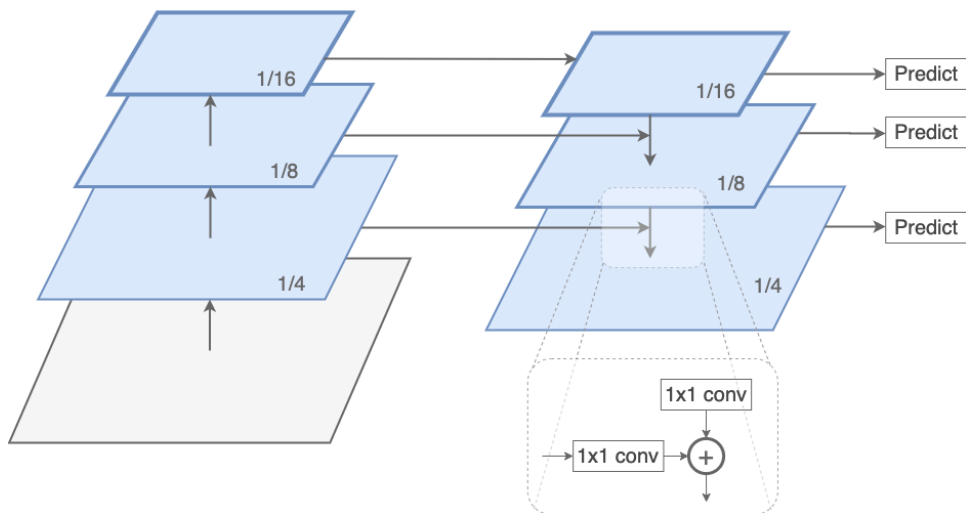


Figure 4: FPN

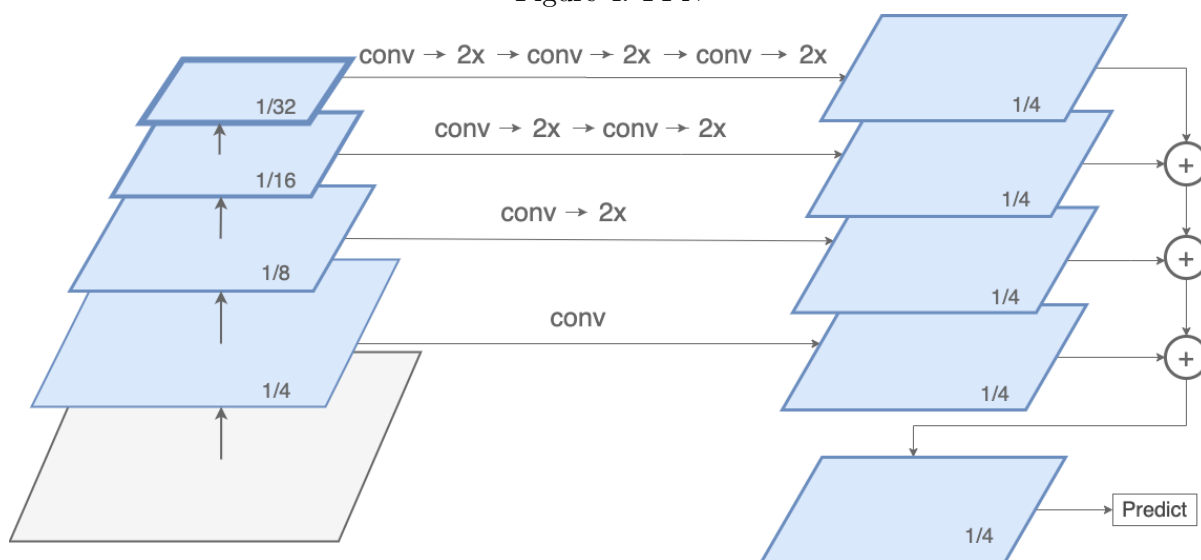


Figure 5: Panoptic FPN

For each bounding box, the RoI pooling layer crop the values of the input feature-map related to that region and convert the cropped section to a fixed-sized feature map. The most important feature of this layer is its capability of producing a fixed-size representation. The size of the new representation is a parameter of the layer itself, hence it is independent both from the dimension of the input feature map and from the RoI size.

Another parameter required by the RoI layer is the spatial scale. In a region proposal network, the RoIs are generated based on the input image size. Hence, there is a need to rescale them in order to extract the relevant region from the input feature map. For example, with an input image of size 160×160 and a feature map of size 20×20 , the spatial scale parameter value is 0.125.

The operations carried out by the RoI pooling layer can be summarized as:

- **Rescaling:** Divide each BB coordinate by k (the spatial scale parameter) and take the integer part, obtaining the new coordinates relative to the input feature map.
- **Quantization:** The cropped part of the feature map is divided into bins, resulting in a $n \times n$ grid, and from each bin, the maximum or average value is extracted. This results in a fixed-sized feature map relative to the RoI.

Figure 6 shows the process of applying RoI pooling for a single RoI.

An alternative to the RoI Pooling layer is the RoI Align layer [22]. This last method rescales the bounding boxes more accurately with the following process

- **Rescaling:** Divide each BB coordinate by k , without rounding up

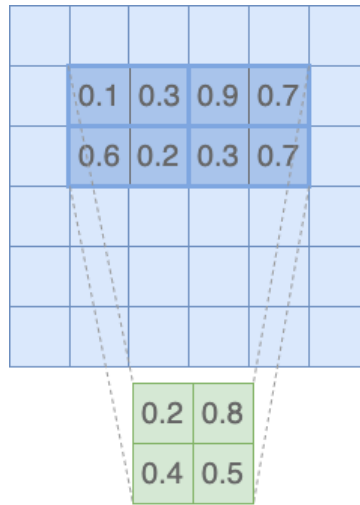


Figure 6: RoI Pooling description

- **Quantization:** The cropped feature map is divided into bins. For each bin are selected 4 points using bilinear interpolation, then the maximum or average value is extracted from these 4 points.

RoI Align is usually preferred to RoI Pooling when dealing with image segmentation, because it does not introduce the misalignment between the RoI and the extracted features.

2.5 Face Alignment

Face alignment, also known as face normalization, is the process of identifying landmarks in images containing human faces, and to use those landmarks to obtain a new face image such that the eyes are aligned and the face is centered. This result is accomplished through scaling, rotating and translating the input image.

Face alignment has been proved to be beneficial for the tasks of face recognition and expression recognition [23] because reduces the variations in face scale and in plane rotations [24].

Two main types of face alignment exists:

- **3D-alignment:** Defines a 3D reference landmark model, and transform the input image in such a way that the landmarks of the input image match the landmarks in the 3D model
- **2D-alignment:** A simpler alignment method, that relies on the scale, rotation and translation of the input image based on the detected facial landmarks

The process of the 2D face alignment can be summarized as:

- **Face detection and Landmark Localization:** This first step consists of determining the location of the face (or faces) in the input image and to extract localized landmark coordinates (usually 5, 44 or 64). Classical face detection algorithms are Haar Cascades Classifier [25] or Histogram of Oriented Gradients (HOG) [13].

A widely used DNN based face and landmark detection model is Multi-Task CNN (MTCNN) [26] which employs multi-task learning and a cascade of learners to detect facial landmarks.

- **Image warping and transformation:** The last phase consists of warping and transforming the input image to an output coordinate space, given a set of facial landmarks.

There are different types of transformations that can be applied in this phase. In our work,



Figure 7: Example of the alignment process

we will use the similarity transform [27] which is a composition of rotation, translations, and magnifications. The points in the new image are defined by the formulas

$$X = s * x * \cos(rotation) - s * y * \sin(rotation) + a_1$$

$$Y = s * x * \sin(rotation) + s * y * \cos(rotation) + b_1$$

in which s is the scale factor, x and y are the translation parameters, rotation is the angle in counter-clockwise as radians and the homogeneous transformation matrix is defined as

$$\begin{bmatrix} a_0 & b_0 & a_1 \\ b_0 & a_0 & b_1 \\ 0 & 0 & 1 \end{bmatrix}$$

Figure 7 shows an example of the alignment process.

2.6 Optimization

Optimization is one of the core components in Machine Learning. In fact, ML algorithms aim at building an optimization model that is capable of learning the parameters in the objective function from the data. The classic optimization method for DNNs is Stochastic Gradient Descent (SGD) [28] and its variants, with the backpropagation algorithm. SGD is a first-order optimization method that iteratively travels down the slopes of the objective function until a minimum is reached. At the end of an iteration of SGD, the backprop algorithm is used to propagate the gradient updates in the network.

With Vanilla SGD, however, it is easy to fall in local minima or saddle point, often present in the objective function. For this reason, momentum is usually added to SGD. Momentum [29] is used to accumulate the gradients of the past steps, which will be used to determine the weights' update for the current step. This has been proved to increase the convergency speed and to partially solve the problem of the optimization algorithm being stuck.

The quantity of gradient to inject in the network is regulated by the learning rate (LR), which determines the step size of each iteration. A small value of the LR could lead to slow convergence while a high LR may not lead to a convergence at all.

Some recent methods introduced the capability of adapting the learning rate based on some statistics. For example, RMSProp [30] stores the moving average of the previous squared gradients and divides the learning rate by the root of the mean square. This allows to keep learning even if the value of the gradients is very small.

Another widely used optimization algorithm is Adam [31]. This method combines momentum and RMSProp and has been proved to provide good performances for a wider range of learning rates.

2.7 Transfer Learning

Transfer Learning [32] is a widely used machine learning technique where a model which has been trained on a task can be used as a starting point for a similar task. There are several reasons why this technique has become popular and widely used

- **Data requirement:** DL models usually require a huge amount of high-quality annotated data, which is not often available. With transfer learning, the model may be fine-tuned in a smaller dataset producing high-quality results.
- **Computation:** DL techniques are known to be computationally intensive. With transfer learning, we can save hours of computations without incurring in a decay of the results.
- **Overfitting:** The high number of parameters in a DL model can lead to an increase in variance which leads to a low generalization capability. This problem can be leveraged relying upon information from a related domain. z
- **Model features:** Some of the features learned by a DL model can be transferred to another domain without relevant modifications.

CHAPTER 3

PREVIOUS WORK

In this chapter, we present some of the current works in AU detection which we have taken as inspiration to develop our model.

In recent years, Deep Learning approaches have dominated the scene on multiple CV tasks.

In Facial Action Unit detection, the first Deep Learning based work that achieved considerable results, was developed in 2014. Since then, most of the works are based on deep neural networks.

3.1 Deep Region and Multi-label Learning for Facial Action Unit Detection

Deep Region and Multi-Label Learning (DRML) [33] is one of the first Convolutional Neural Network based model for AU recognition.

The authors address two main problems in AU recognition:

- **Region Learning (RL):** Since AUs are active on sparse facial regions, RL is capable to identify those regions to improve detection performance. A particular example of RL is Patch Learning, which consists of dividing the face image into uniform patches and train a model to assign a value of importance to each patch.
- **Multi-Label Learning (ML):** It is well known there is a strong correlation between some AUs. Previous works on AU detection extracts the correlation information from either FACS heuristics or ground truth labels and uses them as input to the model. On

the other hand, the authors enforce the model to learn those correlations without prior knowledge of the correlation between AUs.

The main component introduced in the model is the region layer, whose main function is to identify the patches considered relevant for detecting the AUs.

The region layer takes as input the feature map produced by the first convolutional layer and divides into 8×8 patches. Each patch is fed into a second convolutional layer and then summed to the original patch. The new re-weighted patches are then concatenated, resulting in a weighted feature map of the same size as the input.

The use of the region layer shows an improvement in identifying AU discriminative regions with respect to a standard CNN, such as AlexNet or ConvNet.

3.2 EAC-Net: Deep Nets with Enhancing and Cropping for Facial Action Unit Detection

The Enhancing And Cropping (EAC) network [34], is based on a VGG [35] model, in which the following layers are introduced:

- **Enhancing layers:** Are attention layers whose aim is to give more importance to the regions of the image that are associated with the AUs.

This attention layer is implemented as a special type of skip connection between the input and the output of two groups of convolution in the backbone model. In detail, the output feature map of a group of convolutions is multiplied element-wise by a handcrafted attention map and in parallel processed by the convolutional layers of the next group. The two output feature maps are then fused by element-wise summation. This process

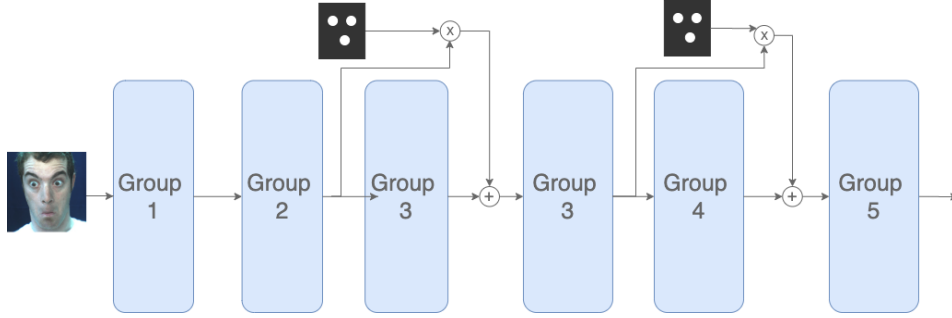


Figure 8: High-level overview of the backbone model with enhancing layers

is repeated for two consecutive groups of convolutional layers. Figure 8 gives a high-level overview of this module.

- **Cropping layers:** Are used to learn features for specific AUs. This module is composed of 20 branches in which each branch takes as input the cropped part of the output feature map of the enhanced VGG. The cropped region is the rescaled center of the 20 AU centers used in the attention map of the enhancing layers.

Each branch upsample the cropped feature map which is then fed into a convolutional layer.

The outputs of the convolutional layers are concatenated and processed by a FC network to produce the final prediction.

In a subsequent work, the author included temporal information to the network, adding two LSTM layers at the end of the previous model, before the dense layer used for the predic-

tion. EAC and EAC with LSTM achieved good results with respect to previous works on AU detection.

3.3 Joint Facial Action Unit Detection and Face Alignment

In Deep Adaptive Attention for Joint Facial Action Unit Detection and Face Alignment [36], the authors defined a deep neural network (JAA-Net) that jointly estimates AUs and facial landmarks (face alignment).

JAA-Net consists of the following modules:

- **Hierarchical and multi-scale region learning:** Is used to extract features at different scales. Similarly to DRML [33], this layer aims at sharing the filter weights within patches with different weights assigned to each patch.

This module is composed of an input layer and three hierarchical convolutional layers. The multi-scale version stands for the way in which the input of a convolution layer is modeled. In particular, the intermediate 8×8 , 4×4 and 2×2 patches of the second, third and fourth convolutional layers are the result of convolution on corresponding patches in the previous layer, input, first and second, respectively. The outputs of the three intermediate layers are concatenated, and the final result is summed element-wise with the output of the input layer.

In Figure 9 is shown the structure of the hierarchical and multi-scale module.

- **Face alignment:** Estimates the locations of facial landmarks and generates the initial attention maps that will be used in the next module.

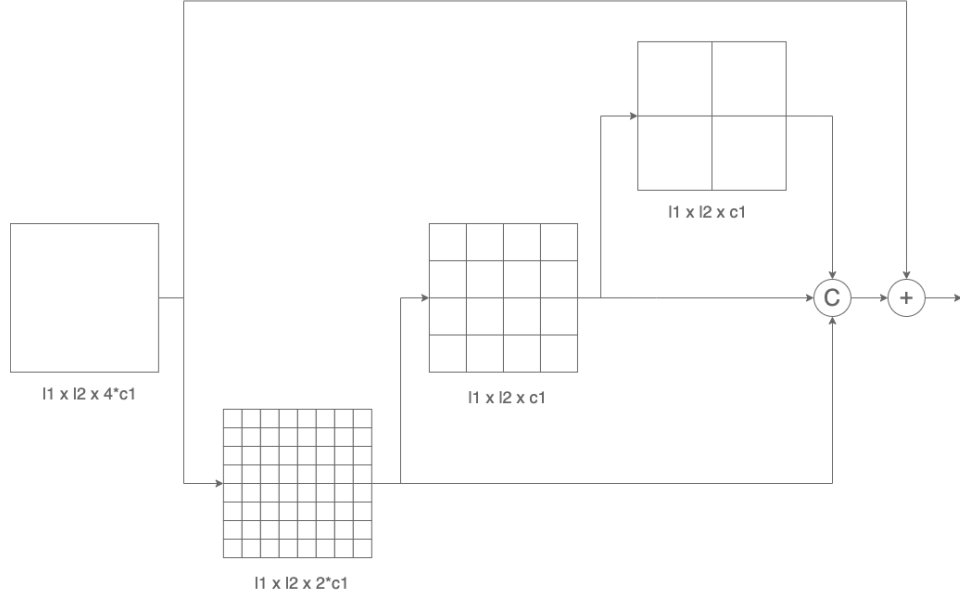


Figure 9: Hierarchical and Multi-Scale module in [36]. In the figure, $l_1 \times l_2 \times c_1$ represent the height, width, and channel of a layer, respectively

This module includes three convolutional layers and a fully connected layer with 2 hidden layers.

- **Adaptive attention learning:** Used to refine the initial AUs' attention maps and consists of two steps:
 - **AU attention refinement:** Takes as input the handcrafted RoIs for each AU similarly to EAC-Net [34], and applies a series of convolution to the RoI masks using different branches. The output of each branch is multiplied element-wise with the feature maps of the hierarchical region learning.

- **Local AU Feature Learning:** It’s composed by a series of convolutions and pooling layers. This module takes in input the refined AUs from each branch, and outputs a new AU attention map, applying element-wise summation over each convolutional branch.

To estimate the AUs, the output feature maps of the face alignment, and of the hierarchical and multi-scale region learning are concatenated together and fed into a FC network with two hidden layers.

3.4 Facial Action Unit Detection Using Attention and Relation Learning

In Facial Action Unit Detection Using Attention and Relation Learning [37], the authors defined a deep learning framework (ARL) based on attention and relation learning for AU detection.

ARL is composed of 4 modules:

- **Hierarchical and Multi-Scale region learning:** As defined in JAA-Net [36].
- **Channel-Wise Learning:** Consists of a multi-branch architecture whose aim is to select AU related features.

For each action unit, the output of the region learning is fed into a convolutional layer and pooled. This process produces a vector of channel-wise attention weights. This vector is multiplied with the output of the convolution, obtaining a channel-wise weighted feature.

- **Spatial Attention learning:** The output of the previous module is processed by a convolutional layer, obtaining a spatial pixel-wise feature map. The feature map is processed

by a convolutional layer with a single channel and a sigmoid function, obtaining an initial spatial attention weight.

- **Pixel-Level Relation Learning:** This module captures the dependencies at pixel-level used to refine spatial pixel-wise attentions.

This module consists of a Conditional Random Field [38] which takes as input the spatial attention weight and the convolved output of the Channel-Wise learning module, obtaining a refined spatial attention weight. The result is then multiplied with a new convolved version of the output of the Channel-Wise module.

Finally, the output of the Pixel-Level Relation learning module is processed by a FC module to predict the AUs' occurrence.

CHAPTER 4

MODEL STRUCTURE

In this chapter, we give an in-depth description of the proposed architecture.

In particular, we describe the FPN module and how this has been integrated with the Region of Interest pooling module in order to solve the problem of AU detection.

4.1 Model Overview

Our model is based on ResNet-50, which is a ResNet model with 50 layers. In particular, this backbone can be viewed as a collection of 4 groups of convolutions, preceded by a stem module composed of a 7×7 convolution and a 3 max-pooling layer, as described in the Table I.

The architecture is based on two main modules:

- **FPN module:** Is the feature extractor of the model, as detailed in 4.2. We choose to use this structure because of its capability of extracting global features at different scales, which has been proved to be beneficial for detecting AUs [36][37].
- **RoI module:** Is a multi-branch module that processes a feature map in a different way based on the location of the occurrence of the AU in the face. This module extracts the local features related to certain RoIs, and produces an estimation of the occurrence of a set of AU. More details are given in 4.3.

Finally, the outputs provided by the RoI modules generate a vector of AUs estimation.

Figure 10 shows a high-level overview of our model.

Layer Name	Output Size	Operations
conv1	$64 \times 1/4 \times 1/4$	$7 \times 7, 64$, stride 2 3×3 max pool, stride 2
conv2	$256 \times 1/4 \times 1/4$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$
conv3	$512 \times 1/8 \times 1/8$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$
conv4	$1024 \times 1/16 \times 1/16$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$
conv5	$2048 \times 1/32 \times 1/32$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$

TABLE I: ResNet-50. The fractions 1/4, 1/8, 1/16 and 1/32 refer to the height and width dimensions ratio with respect to the input image size

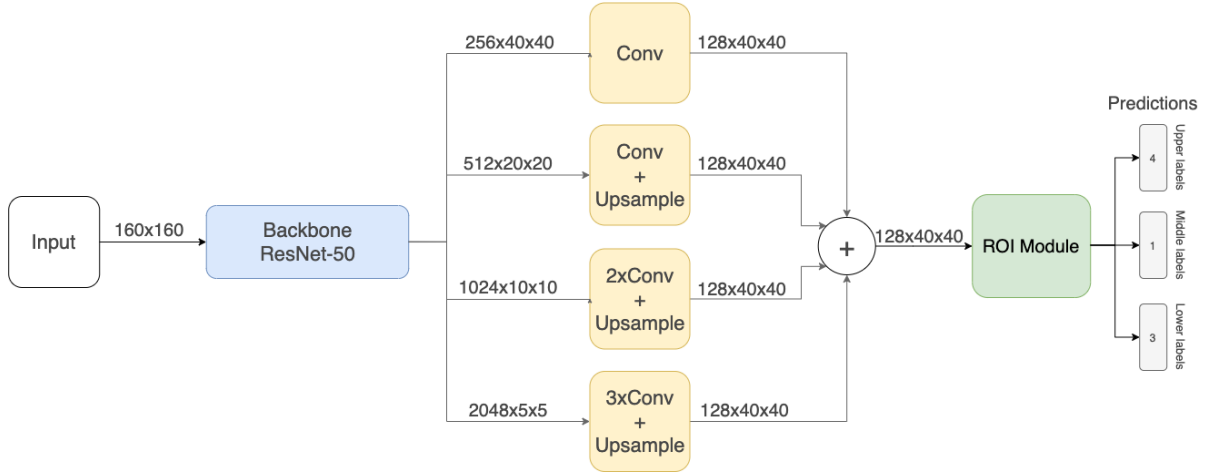


Figure 10: Architecture overview

4.2 Feature Pyramid Network Module

The first module of the proposed architecture is inspired by the Panoptic Feature Pyramid Network described in Section 2.3, and it is used to extract features at different resolutions.

The input of the model is an image with shape $160 \times 160 \times 3$. In the bottom-up path, the model produces 4 feature maps at different resolutions that will be used in the top-down path.

In the top-down path, the output features of each group are fed into an upsampling module, composed by a 3×3 convolution, batch normalization, ReLu, and $2 \times$ bilinear upsampling, as it is illustrated in Figure 11. From the deepest level of the FPN (at $1/32$ scale), we perform three upsampling stages to yield a feature map at $1/4$ scale. This process is repeated for all the remaining feature maps, with progressively fewer upsampling stages.

The upsampled features are then summed element-wise generating the output feature map with 128 channels at $1/4$ scale, which will be used as input for the RoI module, as explained in 4.3. Figure 11, shows the details of the FPN model.

4.3 Regions of Interest Module

The RoI Module is used to extract local features related to the AUs located in different regions of the face, and to estimate the presence of the AU in that part of the face.

In particular, we split the AUs based on their location of occurrence in the face, i.e., upper, middle and lower part of the face. We define three different branches:

- **Upper Branch:** Extracts the local features of the AUs that occurs in the upper region of the face. Particularly, the eyes, the forehead and between the eyebrows. For this branch, we extracted 5 RoIs.

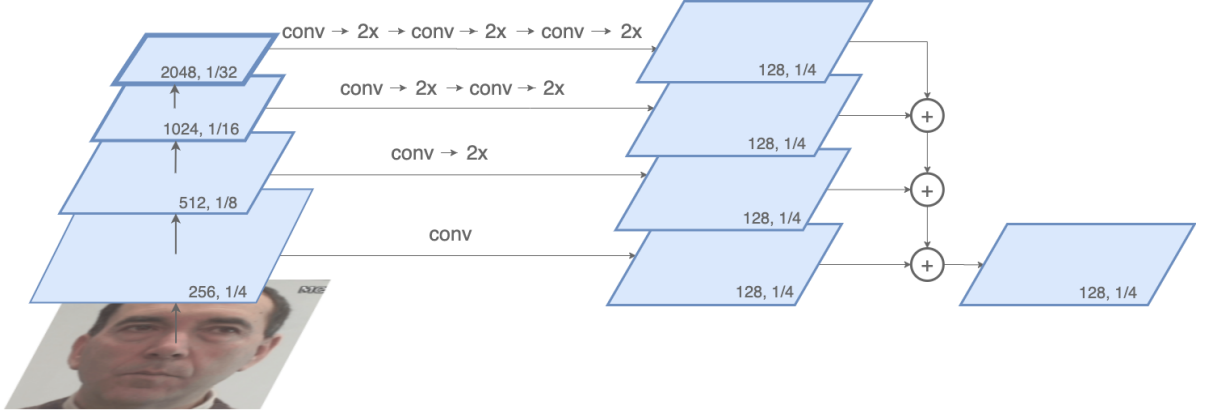


Figure 11: FPN module used in our architecture

- **Middle Branch:** Takes into consideration the middle part of the face. We are interested in detecting the AU related to the cheek's movement. Hence, we defined 2 RoIs, one for each cheek.
- **Lower Branch:** Is tasked with the extraction of the features for the AUs in the lower part of the face. More in detail, we defined the RoIs to extract features related to the mouth, the chin and for the corners of the lips, for a total of 4 RoIs.

Each branch takes as input the multi-scale feature map produced by the FPN module and outputs a vector of estimations for the AUs related to that branch.

The branch consists of a RoI Align layer, described in 2.4, that extracts the local features for each of the RoIs, and produces a fixed-size feature map. The output size of this layer is $\#RoIs \times 128 \times \text{pooling size} \times \text{pooling size}$.

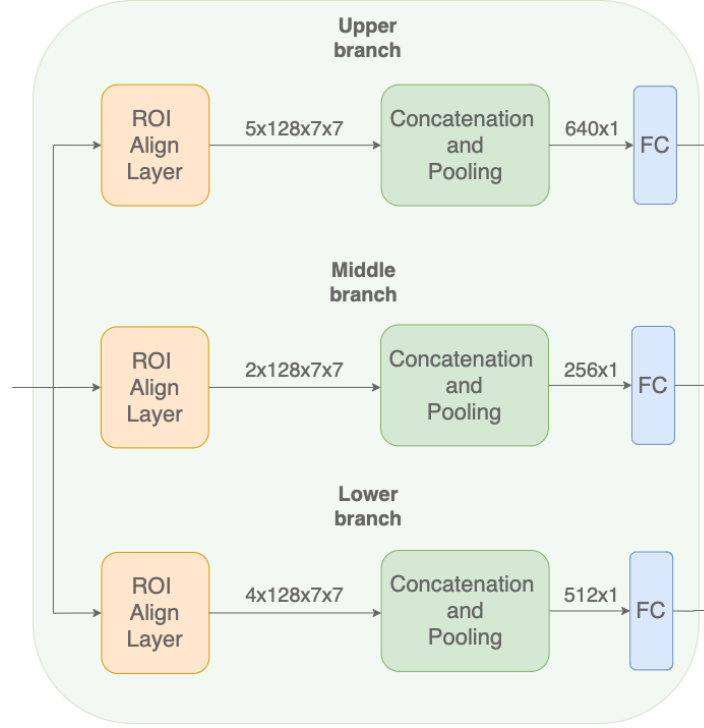


Figure 12: ROI module

These maps are then concatenated on the second dimension and average pooling is applied, before feeding the vector feature map into the bottleneck of each branch. Each bottleneck module is composed of a FC network, with 3 hidden layers. The last layer has a Softmax activation function which produces the estimation of the AUs in that branch.

Figure 12 shows the details for this module.

CHAPTER 5

PRE-PROCESSING, TRAINING AND VALIDATION

This chapter gives more details about the datasets that have been used to train the model, the training procedures, the pre-processing phase and the validation method used to produce the final results.

5.1 Datasets

Six main dataset has been used to train and validate the model:

- **VGGFace2:** VGGFace2 [39] dataset is a face-recognition dataset with more than 3.3 million face images over more than 9,000 entities.

It contains subjects with an equal gender distribution and with a variety of ethnicities and ages. The images are captured "in the wild" such that the images are exposed to pose variations, different lighting, and occlusion conditions.

- **CK+:** The Extended Cohn-Kanade [40] [41] dataset is one of the most widespread laboratory-controlled database for Facial Expression Recognition and AU recognition used.

The dataset contains 593 video sequences from 123 subjects. The duration of each sequence spans from 10 to 60 frames, and it starts from a neutral face to the peak expression. Since the AU is provided for each sequence, usually the last one to three frames are extracted along with the first frame (neutral expression).

- **UNBC-McMaster:** The UNBC-McMaster shoulder pain expression archive [42] database contains videos of subjects who were suffering from shoulder pain, while they were performing some actions.

This dataset contains 200 sequences from 25 subjects. Each frame has been coded by certified FACS coders, resulting in a total of around 47,000 coded frames.

- **CFEE:** The Compound of Facial Expression Emotions [10], is a database of compound emotions images, i.e., emotions constructed by combining basic emotions. The authors defined a total of 21 compounds and basic emotions.

The dataset is a collection of 230 subjects with around 21 images for each one, for a total of 4,816 FACS coded frames.

- **DISFA:** The DISFA [43][44] database, is one of the main benchmark dataset in AU recognition.

It is a collection of 27 sequences, with around 4,800 FACS coded frames for each sequence, from 27 subjects. Each frame is labeled with AUs intensities ranging from 0 to 5. The 8 coded AUs are: 1, 2, 4, 6, 9, 12, 25 and 26.

5.2 Pre-processing

5.2.1 Face Alignment

During this process, we extract 64 facial landmarks in order to normalize the image in terms of scale and rotation. This process aims at aligning the image such that the position of the eyes does not change in all the images of the datasets.

Face portion	AUs	Number of RoIs
Upper	1, 2, 4, 6	5
Middle	9	2
Lower	12, 25, 26	4

TABLE II: RoIs and face locations

The facial landmarks are defined using the framework described in [45] and available in PyTorch [46], which is capable of providing reliable landmarks also in case of partial occlusion of the face.

After retrieving the facial landmarks, we align the image using similarity transformation, previously defined in 2.5. We first identified a set of 64 reference landmarks in an image of size 160×160 , then we applied the transformation using the reference points as destination coordinates. The result of this phase is an RGB image of a size of 160×160 .

5.2.2 Extraction of RoIs

In order to define the RoI regions in which the Action Units occurs, we used the face landmarks previously defined, and we apply the same transformation applied to the image.

In particular, we are using DISFA as a validation dataset so that we defined the RoIs for 8 action units: 1, 2, 4, 6, 9, 12, 25, 26. We mapped the 8 AUs to the portion of the face in which they occur, and we defined a variable number of RoIs for each part of the face, as described in Table II.

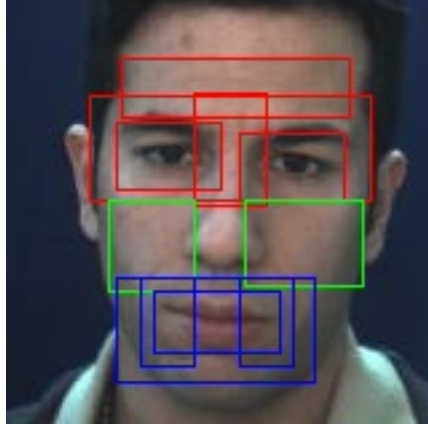


Figure 13: Example of the extracted RoIs. In red, green and blue the RoIs of the upper, middle and lower face are showed, respectively.

Since we are extracting the Regions of Interest based on the facial landmarks, the RoIs are dynamically defined and adapted to the face. Figure 13 shows the RoIs extracted for a single image.

5.3 Training

As previously stated in Section 2.7, transfer learning is a well known method that helps the model to improve generalization and speed up convergence. In our case, we initialize the weights of the ResNet backbone with the weights obtained from training it on the VGGFace2 [39] dataset for face recognition. The weights used are publicly available in [47].

The second training phase aims at fine-tuning the full model for recognizing the facial action units. Each aligned image of size 160×160 is normalized in the same way as the pre-trained backbone with the values shown in Table III.

Channel	μ	σ
R	0.485	0.229
G	0.456	0.224
B	0.406	0.225

TABLE III: Normalization values (mean and standard deviation) for each channel

At training time, we also apply some data augmentation techniques in order to increase and diversify the available data and to improve the robustness of the model. Since we want to capture small movements of the face, we do not apply augmentations that would distort the original features of the face.

In particular, we applied the following augmentations to the training data:

- **Horizontal flip:** Since the AUs are symmetrical, flipping the image allows us to increase the amount of training data without losing information on the location of the action units
- **Color Jitter:** We perform random brightness, contrast, hue, and saturation variations so that the model would be more robust to these variations at inference time
- **Cutout** [48]: Consists of masking out a random part of the input image. This technique simulates occluded examples, encouraging the model to take into consideration a more diverse set of features for classifying images.

For example, to classify a car, the model would be enforced to look at other details of the image, instead of just focusing on the wheels.

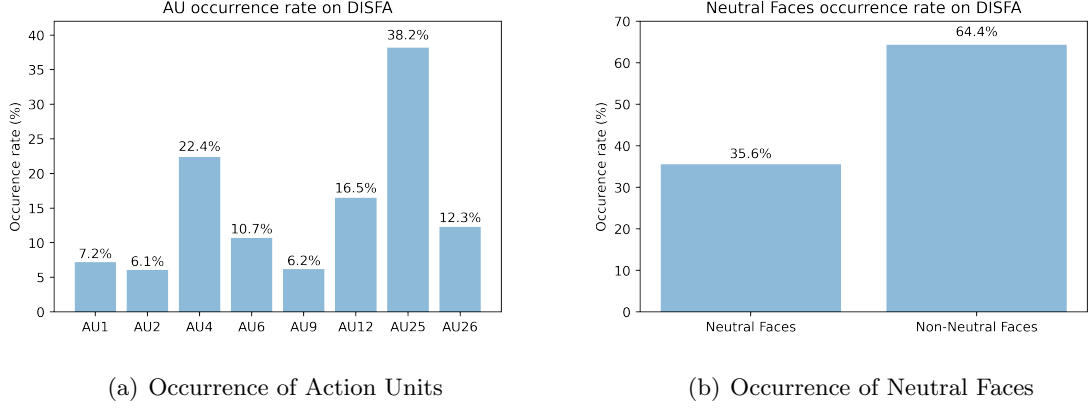


Figure 14: Occurrences of AUs and Neutral Faces in DISFA

Moreover, in order to reduce the imbalance of the training data towards neutral samples, i.e. faces without visible action units, and to increase the diversification of the data, which is reduced by the similarity of subsequent frames, we sub-sampled the images with neutral faces, including in the training data 1 frame over 5.

5.4 Validation

In order to validate the model, we performed 3-fold cross-validation on the DISFA dataset, splitting the data on the subjects. In particular, we used the splits provided by [36]. Figure 14 shows the distribution of AUs and neutral faces on the DISFA dataset.

The training data is then composed of the images in CK+, UNBC and CFEE, and the training folds of the DISFA dataset. Figure 15 shows the AUs' and neutral faces distribution on the training data.

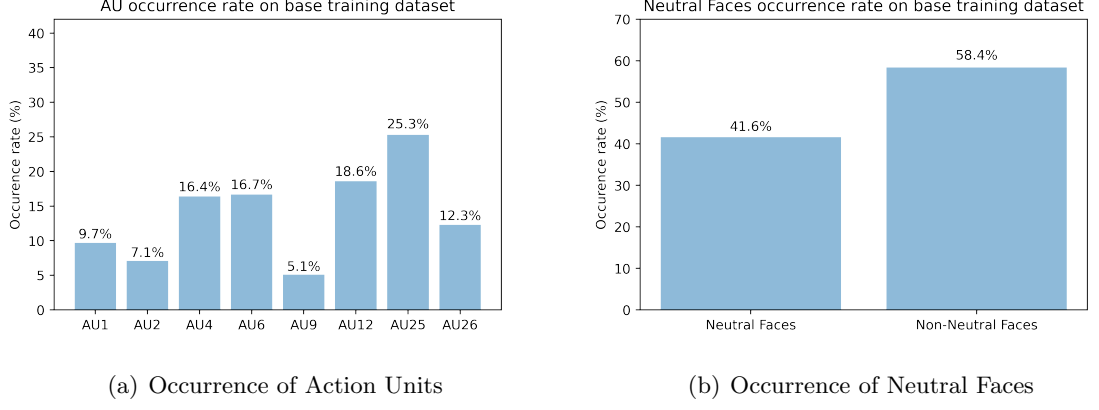


Figure 15: Occurrences of AUs and Neutral Faces in training data

Facial AU recognition can be viewed as a multi-label binary classification problem, which can be solved with the following weighted multi-label softmax loss:

$$E_{softmax} = -\frac{1}{n_{au}} \sum_{i=1}^{n_{au}} w_i [p_i \log(\hat{p}_i) + (1 - p_i) \log(1 - \hat{p}_i)]$$

where

- p_i is the ground truth probability of occurrence of the i -th label (1 or 0)
- \hat{p}_i represents the predicted probability of occurrence
- n_{au} is the total number of action units
- w_i is the weight introduced to alleviate the data imbalance. Particularly, we set

$$w_i = \frac{(1/r_i)n_{au}}{\sum_{i=1}^{n_{au}} (1/r_i)}$$

where r_i is the number of samples with the i -th AU in the training data

Moreover, since some AUs are rarely present in the training data, the network is biased towards absence for this AUs. To mitigate this problem, we introduce a multi-label Dice loss [49], defined as:

$$E_{dice} = \frac{1}{n_{au}} \sum_{i=1}^{n_{au}} w_i \left(1 - \frac{2p_i \hat{p}_i + \epsilon}{p_i^2 \hat{p}_i^2 + \epsilon} \right)$$

where ϵ is a smoothing term.

The Dice coefficient is also known as the F1-score, which is the most popular metric for AU recognition, as described in 6. Introducing this coefficient, we take into account the consistency between the loss function and the evaluation metric.

The final loss function is

$$E = E_{softmax} + \lambda * E_{dice}$$

where λ is an hyper-parameter used to weight the Dice loss.

The final network is trained for up to 2000 iterations in PyTorch [46] using Adam optimizer [31], an initial learning rate of 10^{-4} , mini-batch size of 16 and the λ factor for the Dice Loss has been set to 1.0.

CHAPTER 6

RESULTS

In this chapter, we introduce the evaluation metrics that have been used to validate the model and we show the obtained results in terms of these metrics.

6.1 Evaluation Metrics

- **F1-score:** Is the weighted average of the precision and recall. This metric is used in case of imbalanced classification problems, since the accuracy may not be entirely representative in such cases.

In this work, since we are dealing with a multi-label problem, we take into consideration the macro-averaged F1-score defined as

$$\begin{aligned}\text{F1-score} &= \frac{1}{n_{au}} \sum_{i=1}^{n_{au}} \text{F1-score}_i \\ &= \frac{1}{n_{au}} \sum_{i=1}^{n_{au}} \frac{2 * p_i r_i}{p_i + r_i}\end{aligned}\tag{6.1}$$

in which p_i and r_i refer to the precision and recall of the i -th class, respectively. These measures are defined as

$$\begin{aligned}\text{Precision}_i = p_i &= \frac{\text{TP}_i}{\text{TP}_i + \text{FP}_i} \\ \text{Recall}_i = r_i &= \frac{\text{TP}_i}{\text{TP}_i + \text{FN}_i}\end{aligned}\tag{6.2}$$

where TP are the positive and correctly predicted samples, FP are the negative and correctly predicted samples and the FN are the incorrectly predicted negative samples.

- **Accuracy:** Is the ratio between the correct predictions and the total number of predictions. In a multi-label setting it's computed as the mean of the accuracy for the single class

$$\begin{aligned} \text{Accuracy} &= \frac{1}{n_{au}} \sum_{i=1}^{n_{au}} \text{Accuracy}_i \\ &= \frac{1}{n_{au}} \sum_{i=1}^{n_{au}} \frac{\text{TP}_i}{\text{TP}_i + \text{TN}_i + \text{FP}_i + \text{FN}_i} \end{aligned} \tag{6.3}$$

where TP, FP and FN are defined above and TN represents the correctly predicted negative samples.

6.2 Results on DISFA

We compare our method against some of the SOTA AU detection works under the same 3-fold cross validation setting. In particular, we compare against both traditional methods such as LSVM [50], and DL based ones DRML [33], EAC-Net [34], JAA [36] and ARL [37].

From table Table IV we can see that our method outperforms competing approaches in terms of averaged F1-score. In particular, we can notice that the AU performances are fluctuating in all the methods. A possible explanation for this behavior is that some AUs are easier to recognize with respect to others, resulting in higher F1-score performance.

Moreover, is possible to notice in Table V that our accuracy score is almost similar to the one achieved by ARL. Even though we achieve the best accuracy on the AUs 1, 2, 4, 6 and 9



Figure 16: Examples of the AUs in DISFA

AU	Method					
	LSVM	DRML	EAC	JAA	ARL	Ours
AU1	10.8	17.3	41.5	43.7	43.9	40.0
AU2	10.0	17.7	26.4	46.2	42.1	57.0
AU4	21.8	37.4	66.4	56.0	63.6	74.7
AU6	15.7	29.0	50.7	41.4	41.8	49.0
AU9	11.5	10.7	80.5	44.7	40.0	42.7
AU12	70.4	37.7	89.3	69.6	76.2	68.6
AU25	12.0	38.5	88.9	88.3	95.2	88.3
AU26	22.1	20.1	15.6	58.4	66.8	68.3
All	21.8	26.7	48.5	56.0	58.7	61.7

TABLE IV: F1-score (in %) on DISFA dataset

our method performs slightly worse in terms of accuracy in the prediction of the AUs located in the lower part of the face.

6.3 Ablation Studies

In this section, we conduct an analysis of the various part of our model.

In Table VI we show the results obtained removing some components of the model, in which:

- ResNet-50 refers to the pure ResNet model of [1]
- In ResNet+RoI, we removed the FPN bottleneck representation
- FPN+RoI is the final model with FPN and RoI

As we can see, the base model achieves the lowest result in terms of averaged F1-score. In a more in depth analysis, we noticed that the images with an active AU2 (outer brow raiser) were mainly misclassified as AU1 (inner brow raiser). These errors decrease when we include the Region of Interest module.

AU	Method					
	LSVM	DRML	EAC	JAA	ARL	Ours
AU1	21.6	53.3	85.6	93.4	92.1	95.0
AU2	15.8	53.2	84.9	96.1	92.7	96.3
AU4	17.2	60.0	79.1	86.9	88.5	90.8
AU6	8.7	54.9	69.1	91.4	91.6	92.9
AU9	15.0	51.5	88.1	95.8	95.9	96.1
AU12	93.8	54.6	90.0	91.2	93.9	90.1
AU25	3.4	45.6	80.5	93.4	97.3	92.3
AU26	20.1	45.3	64.8	93.2	94.3	92.3
Avg.	27.5	52.3	80.6	92.7	93.3	93.2

TABLE V: Accuracy score (in %) on DISFA dataset

The ResNet with region of interest pooling performs slightly better in almost all the lower face AUs apart for AU26, which shows an improvement of prediction by almost 10 points when using the FPN module as backbone.

Overall, the complete model achieves the best-averaged results outperforming the other models for the AUs 2, 4, 6, and 26.

In Table VII we show the results obtained training the complete model under different settings.

In particular, we can notice that the pre-training phase of the model is fundamental to achieve the obtained result since it brings an F1-score improvement from 49.0 to 58.7.

Furthermore, the weighted loss is beneficial to the model in order to improve the prediction results for the AU1, at the cost of a small decrease in performance for some other AUs.

AU	Method		
	ResNet-50	ResNet+RoI	FPN+RoI
AU1	43.3	44.2	40.0
AU2	27.3	42.3	57.0
AU4	67.7	73.7	74.7
AU6	45.7	43.8	49.0
AU9	48.3	40.3	42.7
AU12	68.0	73.3	68.6
AU25	86.0	90.0	88.3
AU26	56.0	58.0	68.3
Avg.	55.3	58.2	61.7

TABLE VI: F1-score (in %) on DISFA dataset for different model configurations

Method	AU1	AU2	AU4	AU6	AU9	AU12	AU25	AU26	Avg.
No PT	13.2	14.3	57.0	35.1	41.2	74.9	86.7	69.3	49.0
PT	21.2	55.1	71.3	44.2	42.4	81.7	91.3	62.7	58.7
PT + WL	40.2	41.3	67.7	45.3	42.9	78.3	92.7	68.5	59.6
PT + WL + DL	40.0	57.0	74.7	49.0	42.7	68.6	88.3	68.3	61.7

PT: Pre-Training, **WL**: Weighted-Loss, **DL**: Dice-Loss

TABLE VII: F1-score (in %) on the DISFA dataset removing some training techniques

Finally, with the introduction of the Dice Loss, the model achieves the best overall performance, even though the prediction results for the upper faces AUs is lower with respect to the model without it.

CHAPTER 7

CONCLUSIONS AND FUTURE DIRECTIONS

In this thesis, we designed a new deep learning based architecture for Facial Action Unit detection that achieves results superior to other state-of-the-art models, in terms of F1-score.

We showed that the use of a Feature Pyramid network allows the model to learn features at different scales, which has been proven to be beneficial for the task.

Moreover, we defined a new RoI pooling layer that extracts specific features for different regions of the face, and contribute to the detection of AUs in a specific region in the face.

We evaluated the efficacy of our approach on the available benchmark dataset and motivated our choices through ablation studies.

There are possible developments for improving the current model:

- **Increase training data:** In Section 6.3 we show that with a significant amount of pre-training images the model provides better results. Nevertheless, increasing the amount of labeled data might improve the robustness of the model, in particular for the AUs that appear less frequently.

- **Enclose temporal information:** In AU detection, all actions appear in temporal mode.

In particular, an AU can be interpreted as a mutation of the shape in a specific region of the face. This makes us believe that including temporal information in the model, can be beneficial for the detection.

In order to do so, we can introduce an LSTM module before each FC layer on each branch, and training the model with multiple subsequent images, instead of using a simple image. This solution can be beneficial in terms of the obtained result, at the cost of longer training time.

- **Reference frame:** As an alternative to training on sequence data, the use of a Neutral face as reference frame might be beneficial for the model.

With such an approach, the model can be trained to learn the differences in aspects between a neutral face and the one in which the AUs are occurring.

- **RoI extracted by an expert:** The RoI extraction can lack of precision in case of mispredicted facial landmarks, which is mostly due to high head-pose variations. In such a case the model will not be able to extract the correct face region that will be used for the AU detection. For this reason, a manual annotation of the regions of interest for the different AUs could increase the prediction accuracy of the model, but would imply an higher cost related to the data collection.

CITED LITERATURE

1. Menchetti, G., Chen, Z., Wilkie, D. J., Ansari, R., Yardimci, Y., and Çetin, A. E.: Pain detection from facial videos using two-stage deep learning. In 2019 IEEE Global Conference on Signal and Information Processing (GlobalSIP), pages 1–5, Nov 2019.
2. Facial expressions of pain in lung cancer. Analgesia, 1(2), 1995.
3. Rubinow, D. R. and Post, R. M.: Impaired recognition of affect in facial expression in depressed patients. Biological Psychiatry, 31(9):947–953, 2020/03/24 1992.
4. Kapoor, A., Burleson, W., and Picard, R. W.: Automatic prediction of frustration. International Journal of Human-Computer Studies, 65(8):724 – 736, 2007.
5. Lankes, M., Riegler, S., Weiss, A., Mirlacher, T., Pirker, M., and Tscheligi, M.: Facial expressions as game input with different emotional feedback conditions. In Proceedings of the 2008 International Conference on Advances in Computer Entertainment Technology, ACE '08, page 253–256, New York, NY, USA, 2008. Association for Computing Machinery.
6. EKMAN, P.: Facial action coding system. 1977.
7. Lucey, P., Cohn, J. F., Prkachin, K. M., Solomon, P. E., and Matthews, I.: Painful data: The unbc-mcmaster shoulder pain expression archive database. In Face and Gesture 2011, pages 57–64, 2011.
8. Pantic, M., Valstar, M. F., Rademaker, R., and Maat, L.: Web-based database for facial expression analysis. In Proceedings of IEEE Int’l Conf. Multimedia and Expo (ICME’05), pages 317–321, Amsterdam, The Netherlands, July 2005.
9. Lucey, P., Cohn, J. F., Kanade, T., Saragih, J., Ambadar, Z., and Matthews, I.: The extended cohn-kanade dataset (ck+): A complete dataset for action unit and emotion-specified expression. In 2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition - Workshops, pages 94–101, 2010.

CITED LITERATURE (continued)

10. Du, S., Tao, Y., and Martinez, A. M.: Compound facial expressions of emotion. Proceedings of the National Academy of Sciences, 111(15):E1454–E1462, 2014.
11. LeCun, Y., Haffner, P., Bottou, L., and Bengio, Y.: Object recognition with gradient-based learning. In Shape, Contour and Grouping in Computer Vision, volume 1681 of Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics), pages 319–345. Springer Verlag, 1999. International Workshop on Shape, Contour and Grouping in Computer Vision ; Conference date: 26-05-1998 Through 29-05-1998.
12. Lowe, D. G.: Distinctive image features from scale-invariant keypoints. Int. J. Comput. Vision, 60(2):91–110, November 2004.
13. Dalal, N. and Triggs, B.: Histograms of oriented gradients for human detection. Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on, 1:886–893, 2005.
14. Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L.: ImageNet: A Large-Scale Hierarchical Image Database. In CVPR09, 2009.
15. Krizhevsky, A., Sutskever, I., and Hinton, G. E.: Imagenet classification with deep convolutional neural networks. In Advances in Neural Information Processing Systems 25, eds. F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, pages 1097–1105. Curran Associates, Inc., 2012.
16. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R.: Dropout: A simple way to prevent neural networks from overfitting. J. Mach. Learn. Res., 15(1):1929–1958, January 2014.
17. Hochreiter, S.: The vanishing gradient problem during learning recurrent neural nets and problem solutions. International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems, 6(2):107–116, 1998.
18. He, K., Zhang, X., Ren, S., and Sun, J.: Deep residual learning for image recognition. In 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 770–778, 2016.
19. Lin, T., Dollár, P., Girshick, R., He, K., Hariharan, B., and Belongie, S.: Feature pyramid networks for object detection. In 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 936–944, 2017.

CITED LITERATURE (continued)

20. Kirillov, A., Girshick, R., He, K., and Dollár, P.: Panoptic feature pyramid networks. In 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 6392–6401, 2019.
21. Girshick, R.: Fast r-cnn. In Proceedings of the 2015 IEEE International Conference on Computer Vision (ICCV), ICCV '15, page 1440–1448, USA, 2015. IEEE Computer Society.
22. He, K., Gkioxari, G., Dollár, P., and Girshick, R.: Mask r-cnn. In 2017 IEEE International Conference on Computer Vision (ICCV), pages 2980–2988, 2017.
23. Mollahosseini, A., Chan, D., and Mahoor, M. H.: Going deeper in facial expression recognition using deep neural networks. In 2016 IEEE Winter Conference on Applications of Computer Vision (WACV), pages 1–10, 2016.
24. Li, S. and Deng, W.: Deep facial expression recognition: A survey. IEEE Transactions on Affective Computing, PP, 04 2018.
25. Viola, P. and Jones, M.: Rapid object detection using a boosted cascade of simple features. In Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. CVPR 2001, volume 1, pages I–I, 2001.
26. Zhang, K., Zhang, Z., Li, Z., and Qiao, Y.: Joint face detection and alignment using multitask cascaded convolutional networks. IEEE Signal Processing Letters, 23(10):1499–1503, 2016.
27. Anonymous: Affine transform. <https://scikit-image.org/docs/dev/api/skimage.transform.html#skimage.transform.AffineTransform>, 2018.
28. Rumelhart, D. E., Hinton, G. E., and Williams, R. J.: Learning Representations by Back-propagating Errors. Nature, 323(6088):533–536, 1986.
29. Sutskever, I., Martens, J., Dahl, G., and Hinton, G.: On the importance of initialization and momentum in deep learning. In Proceedings of the 30th International Conference on Machine Learning, eds. S. Dasgupta and D. McAllester, volume 28 of Proceedings of Machine Learning Research, pages 1139–1147, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR.

CITED LITERATURE (continued)

30. Tieleman, T. and Hinton, G.: Lecture 6.5—RmsProp: Divide the gradient by a running average of its recent magnitude. COURSERA: Neural Networks for Machine Learning, 2012.
31. Kingma, D. P. and Ba, J.: Adam: A method for stochastic optimization, 2014. cite arxiv:1412.6980Comment: Published as a conference paper at the 3rd International Conference for Learning Representations, San Diego, 2015.
32. Tan, C., Sun, F., Kong, T., Zhang, W., Yang, C., and Liu, C.: A Survey on Deep Transfer Learning: 27th International Conference on Artificial Neural Networks, Rhodes, Greece, October 4–7, 2018, Proceedings, Part III, pages 270–279. 10 2018.
33. Zhao, K., Chu, W.-S., and Zhang, H.: Deep region and multi-label learning for facial action unit detection. In The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), June 2016.
34. Li, W., Abtahi, F., Zhu, Z., and Yin, L.: Eac-net: Deep nets with enhancing and cropping for facial action unit detection. IEEE Transactions on Pattern Analysis and Machine Intelligence, 40(11):2583–2596, 2018.
35. Simonyan, K. and Zisserman, A.: Very deep convolutional networks for large-scale image recognition. In International Conference on Learning Representations, 2015.
36. Shao, Z., Liu, Z., Cai, J., and Ma, L.: Deep adaptive attention for joint facial action unit detection and face alignment. In European Conference on Computer Vision, pages 725–740. Springer, 2018.
37. Shao, Z., Liu, Z., Cai, J., Wu, Y., and Ma, L.: Facial action unit detection using attention and relation learning. IEEE Transactions on Affective Computing, PP:1–1, 10 2019.
38. Lafferty, J. D., McCallum, A., and Pereira, F. C. N.: Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In Proceedings of the Eighteenth International Conference on Machine Learning, ICML '01, pages 282–289, San Francisco, CA, USA, 2001. Morgan Kaufmann Publishers Inc.
39. Cao, Q., Shen, L., Xie, W., Parkhi, O. M., and Zisserman, A.: Vggface2: A dataset for recognising faces across pose and age. In International Conference on Automatic Face and Gesture Recognition, 2018.

CITED LITERATURE (continued)

40. Kanade, T., Cohn, J. F., and Yingli Tian: Comprehensive database for facial expression analysis. In Proceedings Fourth IEEE International Conference on Automatic Face and Gesture Recognition (Cat. No. PR00580), pages 46–53, 2000.
41. Lucey, P., Cohn, J., Kanade, T., Saragih, J., Ambadar, Z., and Matthews, I.: The extended cohn-kanade dataset (ck+): A complete dataset for action unit and emotion-specified expression. pages 94 – 101, 07 2010.
42. Lucey, P., Cohn, J. F., Prkachin, K. M., Solomon, P. E., Chew, S., and Matthews, I.: Painful monitoring: Automatic pain monitoring using the unbc-mcmaster shoulder pain expression archive database. Image and Vision Computing, 30(3):197 – 205, 2012. Best of Automatic Face and Gesture Recognition 2011.
43. Mavadati, S. M., Mahoor, M. H., Bartlett, K., Trinh, P., and Cohn, J. F.: Disfa: A spontaneous facial action intensity database. IEEE Transactions on Affective Computing, 4(2):151–160, 2013.
44. Mavadati, S. M., Mahoor, M. H., Bartlett, K., and Trinh, P.: Automatic detection of non-posed facial action units. In 2012 19th IEEE International Conference on Image Processing, pages 1817–1820, 2012.
45. Bulat, A. and Tzimiropoulos, G.: How far are we from solving the 2d & 3d face alignment problem? (and a dataset of 230,000 3d facial landmarks). In International Conference on Computer Vision, 2017.
46. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., and Chintala, S.: Pytorch: An imperative style, high-performance deep learning library. In Advances in Neural Information Processing Systems 32, eds. H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, pages 8024–8035. Curran Associates, Inc., 2019.
47. Anonymous: Pytorch face recognizer based on 'vggface2: A dataset for recognising faces across pose and age. <https://github.com/cydonia999/VGGFace2-pytorch>, 2018.
48. DeVries, T. and Taylor, G.: Improved regularization of convolutional neural networks with cutout. 08 2017.

CITED LITERATURE (continued)

49. Milletari, F., Navab, N., and Ahmadi, S.-A.: V-net: Fully convolutional neural networks for volumetric medical image segmentation. pages 565–571, 10 2016.
50. Fan, R.-E., Chang, K.-W., Hsieh, C.-J., Wang, X.-R., and Lin, C.-J.: Liblinear: A library for large linear classification. J. Mach. Learn. Res., 9:1871–1874, June 2008.

VITA

NAME	Guglielmo Menchetti
EDUCATION	<p>Master of Science in “Computer Science”, University of Illinois at Chicago, May 2020, USA</p> <p>Master of Science in “Computer Science and Engineering”, Politecnico di Milano, December 2020, Italy</p> <p>Bachelor’s Degree in ”Computer Science and Engineering””, Università degli Studi di Firenze, April 2017, Italy</p>
LANGUAGE SKILLS	
Italian	Native speaker
English	<p>Full working proficiency</p> <p>2017 - IELTS examination (7/9)</p> <p>A.Y. 2018/20 - One Year of study abroad in Chicago, Illinois</p> <p>A.Y. 2017/18 - Lessons and exams attended exclusively in English</p>
SCHOLARSHIPS	
Spring-Fall 2019	Research Assistantship (RA) position (10 hours/week) with full tuition waiver plus monthly stipend
Spring 2018	Italian scholarship for TOP-UIC students