# Estimation of BMI from Facial Images using Semantic Segmentation based Region-Aware Pooling

Nadeem Yousaf[a], Sarfaraz Hussein[b] and Waqas Sultani [a,*]

[a]*Intelligent Machine Lab, Information Technology University, Pakistan*
[b]*Machine Learning and Data Science @ The Home Depot, USA*

## ABSTRACT

Body-Mass-Index (BMI) conveys important information about one's life such as health and socio-economic conditions. Large-scale automatic estimation of BMIs can help predict several societal behaviors such as health, job opportunities, friendships, and popularity. The recent works have either employed hand-crafted geometrical face features or face-level deep convolutional neural network features for face to BMI prediction. The hand-crafted geometrical face feature lack generalizability and face-level deep features don't have detailed local information. Although useful, these methods missed the detailed local information which is essential for exact BMI prediction. In this paper, we propose to use deep features that are pooled from different face regions (eye, nose, eyebrow, lips, etc.,) and demonstrate that this explicit pooling from face regions can significantly boost the performance of BMI prediction. To address the problem of accurate and pixel-level face regions localization, we propose to use face semantic segmentation in our framework. Extensive experiments are performed using different Convolutional Neural Network (CNN) backbones including FaceNet and VGG-face on three publicly available datasets: VisualBMI, Bollywood and VIP attributes. Experimental results demonstrate that, as compared to the recent works, the proposed Reg-GAP gives a percentage improvement of 22.4% on VIP-attribute, 3.3% on VisualBMI, and 63.09% on the Bollywood dataset.

## 1. Introduction

Faces depict important information about one's personality e.g., age, gender, race, psychological conditions, poverty level as well as health conditions. To measure the overall health condition of a person, usually, body weight and height are used which are encoded in Body-Mass-Index (BMI). Person's BMI has a significant impact on several aspects of life, including health [28, 17, 3, 29], job opportunities [5], friendships and popularity [19]. Recent studies demonstrate that higher BMI can lead to many diseases such as heart disease [17, 3], diabetes [28, 17, 29], stroke [3], cancer [28, 17, 3, 29], sleep apnea [28, 3, 29], hypertension [28, 3], fatty liver disease [29], kidney disease [29], depression [29], and pregnancy problems [36]. Other than health, BMI values have also been used to estimate and predict the social behaviors of societies. For example, on social media, people with similar BMI values are more likely to make connections as compared to people with dissimilar BMI values [19]. Similarly, people with higher BMI values have fewer followers as compared to people with lower BMI values [19]. Similarly, Caliendo et al [5] pointed out that obese women (higher BMI) observe weight-based discrimination during job interviews. Consequently, obese women find it more difficult to find a job and they get low wages as compared to their less obese counterparts.

In the past, researchers have studied the association between facial measures and body weight. It has been observed that BMI is strongly correlated with eye-detailed information (e.g., intraocular pressure (IOP) and anterior corneal curvature (ACD) [25], neck circumference [31, 1] and face

physical measures such as width-to-height ratio, perimeter-to-area ratio, and cheek-to-jaw-width ratio [7]. However, most of these works have used hand-crafted features and took face measurements manually. Due to the large impact of BMI on a person's life and society behaviors, it is of significant importance to measure the BMI on a large scale. However, this would need a lot of resources to go door-to-door to compute each individual's BMI. Fortunately, computer vision and deep learning provide a non-intrusive, efficient, and cheaper way to estimate people's BMI through their face images on social media. [18, 19]

Although it has been well-established [25, 6] that different facial regions are strongly correlated with BMI values and their measurements can help better prediction of BMI, all the previous deep learning-based methods have used the deep features from the full face images. For example, Pascalia et al., [27] proposed a method for automatic extraction of geometric features using 3D facial data acquired with low-cost depth scanners. They have experimentally shown that these features are highly correlated with weight and BMI. Similarly, authors in [4] employed diffusion tensor imaging on white matter alterations to estimate obesity and BMI. We, on the other hand, explicitly use the information from different facial regions to obtain a robust feature embedding. To exploit the relationship of different facial regions with BMI, we obtain the improved feature vector by pooling the convolution feature maps based on different semantic regions of the face. We obtain different face regions through face semantic segmentation. The face semantic segmentation provides accurate pixel-wise locations of different face regions. Figure 1 demonstrates the pipeline of the proposed approach. Assuming FaceNet [32] as a face feature extractor module, given the input image, we crop the face region employing
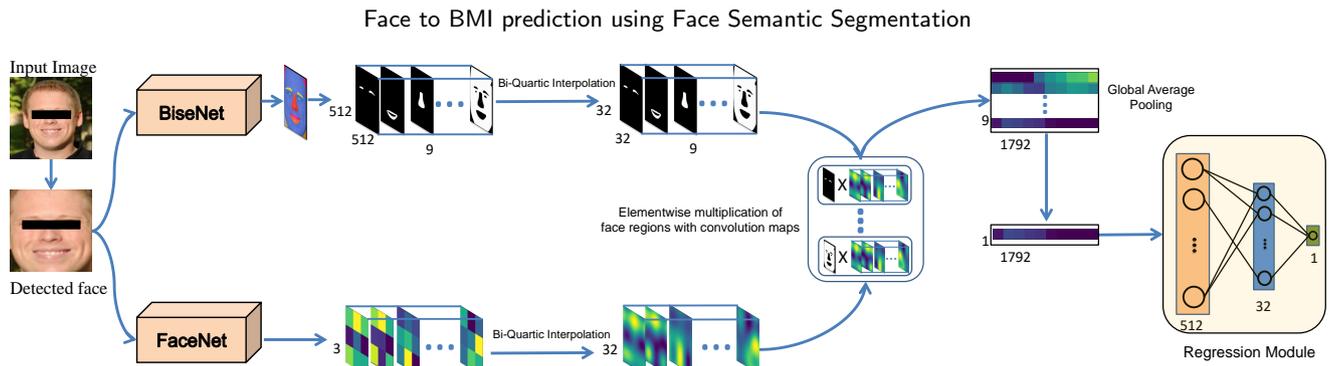
*Corresponding author
ORCID(s):

**Figure 1:** The pipeline of the proposed approach. Given the input image, we crop the face region employing face detection [41]. Each face region mask is obtained through face semantic segmentation. After that, face region masks are element-wise multiplied with the convolution feature maps to give high weights to different face regions. Global average pooling is then applied to each masked convolution map separately. Finally, we employ the regression module to obtain the BMI prediction.

face detection [41]. Face region masks of different face regions such as ear, eyes, eyebrow, hair, lips, neck, nose, and skin are obtained through face semantic segmentation. After resizing face-region masks using bi-quartic interpolation, each mask is multiplied with convolution feature maps to obtain feature map values, specific to different face regions. After that, we perform **Reg**ion aware **G**lobal **A**verage **P**ooling (Reg-GAP) to get the final embedding for training BMI prediction regression module. Finally, the regression module is employed to predict the BMI value of the face. Although simple and straightforward, our experimental results demonstrate that face semantic segmentation-based feature pooling helps improve BMI prediction on two popular publicly available datasets.

The overall organization of the paper is as follows: Section 3 summarizes some of the recent research works for BMI prediction, section 4 provides details of the proposed methodology, section 5 shows experimental results and analysis and finally, section 6 concludes the paper.

## 2. Related work

Due to the significant impact of BMI on health, economic conditions, societal behaviors, several researchers have extensively worked to accurately estimate BMI from people's eyes and neck information, facial dimensions, and face and human body images.

Recent studies have shown the facial features encode useful information about a person's health. Panon et al. [25] investigated the relationship between body mass index and ocular parameters. Employing enhanced depth-imaging optical coherence tomography, several measurements of anterior and posterior segment parameters of the eye measurements were made. The segments include anterior chamber angle, central corneal thickness, macular thickness (MT), anterior chamber depth (ACD), ganglion cell thickness (GCT), retinal nerve fiber layer thickness among many others. Moreover, anterior corneal curvature and intraocular pressure (IOP) were measured by non-contact tonometry. Using data from fifty-three left eyes of normal weight subjects and 67 age-sex matched overweight subjects, they concluded that intraocu-

lar pressure (IOP) and anterior chamber depth (ACD) are positively correlated with BMI.

Coetzee et al. [7] demonstrated that there are two important prerequisites for any health cue. One of them is the perception of weight in the face which can significantly predict perceived health and attractiveness. Authors in [6] tried to spot the facial cues that are associated with BMI. They recruited two groups of African and two groups of Caucasian participants, determined their BMI, and measured their 2-D facial images for perimeter-to-area ratio, cheek-to-jaw-width ratio, and width-to-height ratio, The width-to-height, and cheek-to-jaw-width ratios were found to be significantly associated with BMI in males and females. Mayer et al. [23] studied to assess the association of BMI and waist-to-hip ratio (WHR) with facial shape and texture in females. The females included in the study were middle-aged European women with a BMI between 17-35. They showed that BMI is better predictable than WHR from facial attributes. Saka et al. [31] performed a pilot study on Turkish adults and found that neck circumference can be utilized as an indicator for abdominal obesity and similarly, Atwa et al. [1] have utilized neck circumference to detect children with high BMI.

Segmentation has been used in several medical image analyses to improve detection and classification accuracy. Qayyum et al., [30] proposed a hybrid 3D residual network (RN) with a squeeze and excitation (SE) block for volumetric segmentation of kidney, liver, and their associated tumors. The authors in [21] proposed the NucleiSegNet - a robust deep learning network architecture for the nuclei segmentation of hematoxylin and eosin-stained liver cancer histopathology images. Their proposed deep-learning architecture yielded superior results compared to state-of-the-art nuclei segmentation methods. Similarly, Mussi et al., [24] proposed an algorithm that performs ear depth map segmentation.

The first work to show that the BMI of a person can be automatically predicted from a 2D face image using the geometrical features was done by Wen and Guo [38]. After detecting the face and eyes, they normalized the face based on the eyes' coordinates. Normalization was done to align the face images into common eye coordinates. Next, the active shape model (ASM) was used to detect several key points in

each face image. Seven geometrical features were detected which include width to upper facial height ratio, cheekbone to jaw width, eye size, perimeter to area ratio, lower face to face height ratio, mean of eyebrow height, and face width to lower face height ratio. They normalize these features before applying support vector regression (SVR). They evaluated the method on $14,500$ images from the MORPH-II dataset which is not freely publicly available. Following [38], Jiang et al. [11] extracted geometric features from the whole body to predict BMI from whole-body images. In contrast to whole-body images, face images are more easily available (on National ID cards, driving licenses, etc) and usually have little or no occlusion. Furthermore, to the best of our knowledge, there does not exist any publicly available body to BMI dataset

With the resurgence of deep convolution neural networks, several interesting and new problems, including BMI prediction, have been efficiently addressed with improved accuracy. The first work related to deep learning-based BMI estimation from face photos was done by Enes Kocabey et al. [18]. Instead of using traditional machine learning to extract hand-crafted features, they employed pre-trained deep neural network models to extract features. They have used two models to extract the features. The first model named 'VGG-Face' [26] was trained for the face recognition task and the second model named 'VGG-Net' [35] was trained for general image classification. The features were extracted from the fully connected ($f_{c_6}$) layer. To perform the prediction, the epsilon support vector regression model [9] was employed. Furthermore, they have collected their dataset using Reddit-subreddit called progress-pics. Dantcheva et al. [8] proposed a CNN-based method to estimate the height, weight, and BMI using 50-layers ResNet-architecture. They had also presented a new 'VIP-Attribute' dataset consisting of 1026 subjects. This dataset contains 513 males and 513 females.

Similarly, Jiang et al. [12] introduced a label distribution-based method for BMI estimation from face images. Their proposed approach contains two stages. In the first stage, BMI-related features are computed, and in the second stage, a label distribution-based BMI estimator is learned. Specifically, in the first stage, they utilized a face model [39] which was originally trained for the face recognition task. They used the FIW-BMI dataset [13] to fine-tuned the face-recognition model [39] to a BMI-related face model by replacing the last fully connected layer of 512 dimensions into 1 and use the Euclidean loss. They defined a single BMI value as a discrete probability distribution over the range of BMIs using Gaussian distribution and triangle distribution. After extracting the facial features, five estimators were learned. The estimator includes: Principal Component Analysis (PCA), Support Vector Regression (SVR), Gaussian Process Regression (GPR), Partial Least Square analysis (PLS), Canonical Correlation Analysis (CCA), and two label distribution (LD) based estimator include (LD-CCA, LD-PLS). Authors in [3,9,33] are using deep learning-based models as blackbox i.e., they entirely rely on the network to extract mean-

ingful features to map the face image to the BMI score. On the other hand, our proposed approach explicitly extracts sub-regions of the face and uses these local cues to learn an attention-based feature-space which is a more meaningful representation.

Recently, human face and body semantic segmentation have been used in several computer vision applications. Khalil et al. [16] have used semantic face segmentation for gender and expression analysis. They segmented the facial images into six semantic classes: hair, skin, nose, eyes, mouth, and back-ground using a random decision forest. In their final step, they trained a Support Vector Machine (SVM) classifier for gender using the corresponding probability maps of facial regions. Improved facial attribute prediction based on face semantic segmentation was presented by [15]. The core idea of their research was that facial attributes describe local properties and the probability of an attribute to appear in a face image is far from being uniform in the spatial domain. They obtained an improved facial attributes prediction while using localization cues from facial semantic segmentation. Similarly, Kalayeh et al. [14] presented an approach for improved person re-identification using human semantic parsing.

Attention-based networks have shown promising results in several vision tasks. Wang et al. [37] put forwarded an attention-based multi-branch network for makeup-invariant face verification. Authors in [33] used the region-wise modelling to predict the human facial skin age and [10] used the multi-step region growing for segmentation of skin cancer images. Similarly, researchers in [2] employed facial regions geometric features for the analysis of in and out-group differences in Western and East Asian facial expression recognition. Similar to us, researchers in [20] have also used Region-based Average Pooling for context-aware object detection. Our proposed approach has several differences from the approach presented in [20]. The main purpose of RAP in [20] is to combine the features of different object regions to achieve improved *object* detection. However, on the other hand, we employ region-aware pooling to explicitly pool features from different face regions to make face to BMI prediction better. The approach in [20] aims at learning the relationship between different regions and use these relationships to better classify and regress each region. Each region in [20] corresponds to a single complete object. That is why they un-pooled the averaged vectors and concatenated them with original representations of each region (object) separately. In our case, each region is a sub-part of one complete object (face) and our purpose of averaging the regions and then combining them is to learn the individual importance of each region for the face to BMI estimation. Therefore, we directly pass the averaged representation to the regression module instead of unpooling it and concatenating it to the region's representations. Furthermore, in [20], since each region is a different object, thus combining features of unrelated classes, such as a car and a cat, might not always result in optimal performance. While in our case, since the regions are a part of the same object, combining

them will always complement the learning process.

In contrast to the above-mentioned research works, in this work, we propose to use face semantic segmentation to extract local facial regions. Our proposed facial regions-based pooling provides robust feature embedding for face to BMI prediction. To the best of our knowledge, we are the first ones to propose facial regions-based pooling for BMI prediction.

## 3. Methodology

The proposed approach for BMI prediction from face images contains three main components. Given the face image, we extract deep features and employ semantic segmentation to obtain pixel-level localization of different face regions. After that, we integrate the semantic segmentation to obtain face-regions based pooling from convolution layers of the neural network which was trained on face images. Finally, the pooled features are used to predict BMI values using a fully-connected regression module. The complete algorithm of our approach is given in Algorithm 1. Below we describe each component of the proposed approach in detail.

---

**Algorithm 1:** Algorithm to Estimate BMI from Facial Images

**Input:** People Images Set I
**Output:** Linear BMI Prediction $B_1...B_k$
**Procedure** Predict_BMI($I$):
    $D_1...D_k \leftarrow Face\_Detection(I)$
    $F_1...F_k \leftarrow VGGFace(D_1...D_K)$
    $S_1...S_k \leftarrow$
    $Face\_Semantic\_Segmentation(D_1...D_k)$
        $S_i \in \mathbb{R}^{h \times w \times j}$
    $M_1...M_k \leftarrow PreProcess\_Masks(S_1...S_k)$
    **for** $k=1$ to $K$ **do**
        **for** $j=1$ to $J$ **do**
            $R_{kj} = F_k \star M_{kj}$
            End
        End
    $v_{f_1}...v_{f_K} \leftarrow Reg\_GAP(R_{kj})$
    $B_1...B_K \leftarrow Regression\_Module(v_{f_1}...v_{f_K})$
    **return**    $B_1...B_k$
    End

---

### 3.1. Face Feature Extraction

We employ two face feature extraction models: FaceNet [32] and VGG-face [26]. Below, we briefly describe both of the feature extraction methods.

**Feature extraction with FaceNet:** FaceNet directly learns a mapping from face images to a compact Euclidean space for tasks such as face recognition and verification. Similarly, face image clustering can be easily implemented using standard techniques by utilizing FaceNet embeddings as feature vectors. FaceNet model is trained on faces detected using multitask cascaded convolutional neural networks (MTCNN)

[41]. The input image shape required for FaceNet is of size 160×160×3. We first resize the images to this shape and then apply MTCNN. The MTCNN approach predicts face and landmark locations in a coarse-to-fine manner by adopting a cascaded structure with three stages of carefully designed deep convolutional networks. Figure 2 represents the detection of facing using MTCNN on VisualBMI dataset [18]. In this work, we extract the features maps of size 3×3× 1792 from last convolution layer of FaceNet, which was resized to $32 \times 32 \times 1792$ using bi-quartic interpolation. Facial-regions feature extraction (see Section 4.3) is applied to these feature maps to extract the region aware features from the FaceNet model.
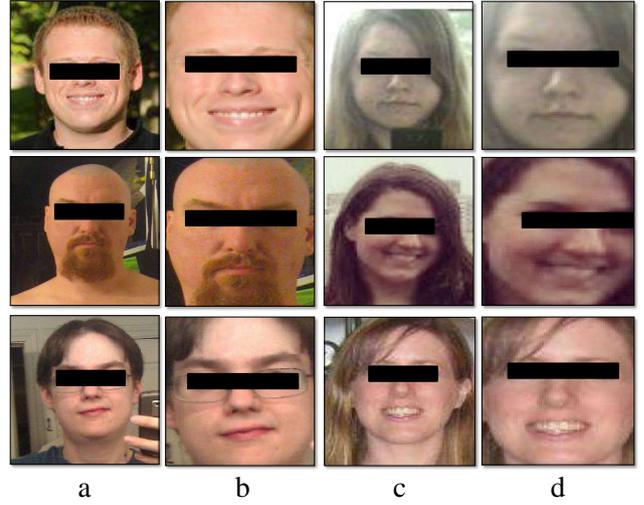


**Figure 2:** Typical examples of face detection to extract FaceNet features are shown in this figure. The face detection is done using Multi-task Cascaded Convolutional Neural Networks [41]. (a) and (c) shows the original face images and (b) and (d) shows the cropped images after face detection.

**Feature extraction with VGGFace:** The second face extraction deep convolution neural network that we use in our experiments is VGGFace [26]. VGGFace is trained as a face classifier with 2.6 million facial images of more than 2600 people. The name 'VGGFace' was given later to the model and it was described by Omkar Parkhi in the 2015 paper titled "Deep Face Recognition [26]. The input shape required for VGGFace is 224×224×3. We resize the images to this shape without applying face detection since this model was originally trained in these settings. We use the features from the $Conv5\_3$ layer of VGGFace which has a shape of 14×14×512. In the experiments, we resize the feature maps to 32×32×512 with the help of bi-quartic interpolation. Finally, facial-regions feature extraction (see section 4.3) is applied to these feature maps to extract the region aware features of the VGGFace model.

### 3.2. Face semantic segmentation

To accurately localize different facial regions, we employ a bilateral segmentation network for face parsing [40].

To obtain real-time segmentation with sufficient accuracy, authors in [40] use small stride to preserve the spatial information and employs a fast down-sampling technique to obtain a sufficient receptive field. Specifically, to avoid losing spatial information due to small image size, a 3-layer convolution network is employed that output feature maps that are 1/8 of the original image. Similarly to preserve contextual information through a large receptive field, authors in [40] propose to employ the Xception network along with the global average pooling layer followed by U-structure to fuse the features. Finally, the network is trained end to end using softmax loss function which is given by

$$loss = \frac{1}{N} \sum_k -log\left(\frac{e^{o_k}}{\sum_j e^{o_j}}\right), \qquad (1)$$

where o is the output of the network. We encourage readers to this reference [40] for the model details of segmentation network architecture. The authors in [40] demonstrated the results on Cityscapes, CamVid, and COCO-Stuff datasets. Since we are interested in face image segmentation, we use the model pre-trained on CelebAMask-HQ dataset [22]. We have used modified BiseNet which produces precise segmentation results as shown in Figure 3 and Figure 4. There are several differences between modified Bisnet and the originally proposed BiseNet such as 1) Original BiseNet take the image of the whole scene as input while in our modified version, we first apply face detector (MTCNN) to localize the face and then input it to BiseNet which improves the accuracy for the face parsing, 2) As shown in [40], original BiseNet was trained on two models: Xception and ResNet. Xception has fewer parameters which make it faster but had lower segmentation accuracy of 71.4% as compared to ResNet which had more parameters and has better accuracy of 78.9%. Therefore, our modified version is using ResNet as its backbone to ensure the highest accuracy as compared to its faster version with Xception as a backbone. The typical examples of face semantic segmentation on VisualBMI [18] and VIP attributes [8] datasets are shown in Figure 3. We modified the BiseNet model to generate a separate binary mask for each region instead of a combined mask of all regions as shown in Figure 4. The separate binary masks are later preprocessed according to input image size and are shown in the bottom two rows of Figure 4.

### 3.3. Region-aware Global Average Pooling

We employ face semantic segmentation to pool the deep features from face regions. In Figure 4, in the bottom two rows, we show the masks obtained from different face regions. The face region includes the ear, eyes, eyebrow, hair, lips, neck, nose, and skin, and background. To obtain the face-region aware features, we perform element-wise multiplication of feature maps with that of mask obtained from semantic segmentation. Figure 4, in the top two rows, showed the region-aware features obtained after element-wise multiplication of masked regions with feature maps where feature maps are obtained using FaceNet or VGGFace.
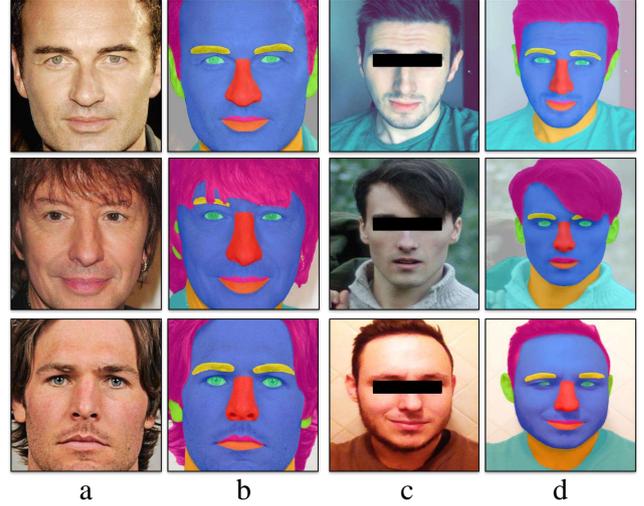


a      b      c      d

**Figure 3:** Examples of the face semantic segmentation. The first and the third column show the face images from the VIP attribute dataset and VisualBMI dataset while second and the fourth column shows their resultant face semantic segmentation.

Formally, assume $\mathcal{F}$ represents convolution feature maps of FaceNet or VGG-Face and $\mathcal{M}_k$ represents stacked binary mask for region $i$ where $i$ can be eye, nose, necks, etc.

$$\mathbf{r}_i = \frac{1}{N} \sum \sum (\mathcal{F} \star \mathcal{M}_i), \qquad (2)$$

where $\mathbf{r}_i$ is the global average pooled vector of face region $i$ and $\star$ represents Hadamard product. We repeat the steps mentioned in Eq 2 for each region separately. Finally, region aware global average pooled feature vector (Reg-GAP) is given by

$$\mathbf{r}_{Reg-Gap} = \frac{1}{K} \sum (\mathbf{r}_i), \qquad (3)$$

where K is the number of regions.

In the experiments, we have also compared Reg-GAP with well known global average pooling (GAP) on the original convolutional feature maps. The GAP is defined as:

$$\mathbf{r}_{Gap} = \frac{1}{N} \sum \sum (\mathcal{F}), \qquad (4)$$

where $\mathcal{F}$ is the original output of the last convolution layer of the model used.

In Figure 5, we show the comparison between the original FaceNet features maps and region-aware FaceNet features maps. The first column shows the cropped image which is input to the FaceNet model. The middle column shows the original feature maps extracted from FaceNet and the last column show the region-aware feature maps. For illustration purposes, the shown feature maps (second column) are generated by taking mean across channels of original feature maps and we overlay it on the face image. The feature maps (third column) are generated by taking max across channels. Finally, to show the region-aware feature maps (last column), we took the max across the channel after taking element-wise multiplication with each face region. It can be seen

**Figure 4:** Given the face semantic segmentation, we extract different face regions and explicitly pool features from those regions. The regions are (Left to Right): ear, eyes, eyebrow, hair, lips, neck, nose, and skin. The bottom two rows show the binary mask obtained from segmentation and the top two rows show region corresponding feature maps. The first and third row samples are from the VisualBMI dataset while the second and fourth-row samples are from the VIP attribute dataset.
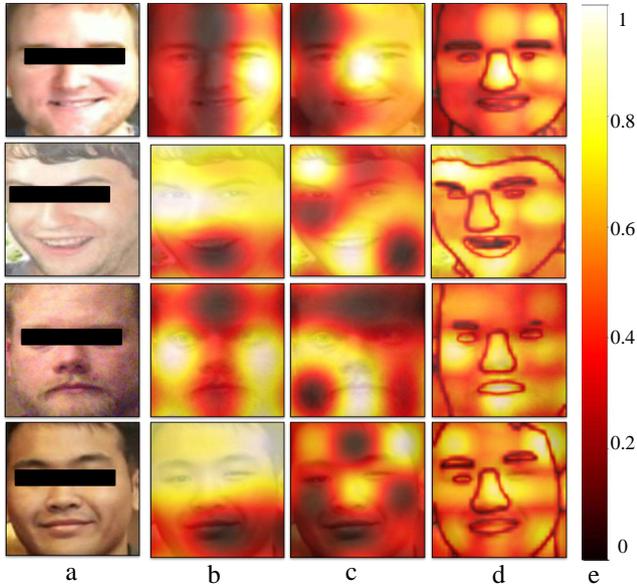


**Figure 5:** This figure shows the comparison of convolution feature maps obtained by the global average pooling (GAP) and by proposed region-aware global averaging pooling (Reg-GAP). (a) shows the input to the FaceNet model, (b) is the mean across channels of the last convolution layer, (c) is the max across channels of last convolution layer, (d) is the max across the channels after Reg-GAP.

that Reg-GAP feature maps capture more details of the face and are invariant to face variations.

.

## 3.4. Regression module

Once the region-aware feature vector is obtained, we employ the regression module to obtain the final Face to BMI prediction. The first layer of the architecture has 512 neurons and a kernel constraint with the max norm of 5. Then we use the dropout of 0.4 to handle the model over-fitting issues. The next layer consists of 256 neurons and a kernel constraint with the max norm of 5. The first two layers have RELU activations. Lastly, we have a single neuron layer with linear activation as we have BMI in the linear range. We employ Adam optimizer with default configurations and the loss function used is Mean Square Error (MSE) which is defined as:

$$MSE = \frac{1}{n} \sum_{i=1}^{n} \left( y_i - x_i \right)^2, \qquad (5)$$

where $n$ is the number of samples, $y_i$ is the ground truth BMI, and $x_i$ is the predicted BMI value.

## 4. Experiments

### 4.1. Dataset

We used three publicly available datasets for evaluation of the proposed methodology: VisualBMI [18], VIP attributes [8] and Bollywood Dataset [34]. We have not used FIW-BMI and Morph II [1] because both datasets are not available free of cost.

**VisualBMI**: The VisualBMI dataset is collected by [18] from Reddit-subreddit called progress-pics where people share th-

---

[1] https://ebill.uncw.edu/C20231_ustores/web/store_main.jsp? STOREID=4

**Figure 6:** Examples of face images in the VisualBMI dataset.

eir images of before and after body transformation. There are a total of 4206 images, out of which 2438 are males and 1768 are females. We followed the split provided by original authors [18]. The first 3368 images are used for training and the rest of the images are used for testing. The dataset is quite challenging due to several low quality and variable size blurry images. Furthermore, some of the images are the images of the picture. The typical example of images are shown in Figure 6.

**VIP Attribute dataset**: The VIP attribute dataset is collected by Bilinski et al. [8]. This dataset consists of 1026 images of males (513) and females (513). For a fair comparison, we have performed experiments using a 78/22 split of data which are provided by Jiang et al. [12] The images are of singers and athletes. Unlike the VisualBMI dataset, these images are mainly frontal and are of high quality. This dataset is challenging due to the presence of makeup, plastic surgery, beard, and mustache. The authors of [8] collected the height and weight from different celebrity websites and calculated the BMI. The typical example of images is shown in Figure 7.

**Bollywood dataset**: The 'Bollywood' dataset is publicly available at GitHub [2]. This dataset contains 237 labeled images. All the images are of Bollywood celebrities. There are a total of 22 identities in this dataset which means that there are multiple images per celebrity. We used a 78/22 split for training and testing. The experimental results in Table 1 show the superiority of our approach on this new dataset as well.

## 4.2. Evaluation metrics

The evaluation metrics used in this papers are mean square error (MSE), root mean square error (RMSE) and Pearson

---

| Model | SVR [34] | RR [34] | Our (GAP) | Our (Reg-GAP) |
|---|---|---|---|---|
| VGG19 | 1.99 | 1.49 | 0.98 | **0.55** |
| VGGFace | 0.96 | 0.97 | 0.40 | **0.32** |

**Table 1**
Results on the Bollywood dataset. Results show that Reg-GAP results are better than that of GAP and [34].

| Model | p | Model | p | Model | p |
|---|---|---|---|---|---|
| VGG16 | 0.45 | ResNet50 | 0.33 | ResNet50v2 | 0.37 |
| VGG19 | 0.42 | ResNet101 | 0.31 | ResNet101v2 | 0.39 |
| **VGGFace** | **0.65** | ResNet152 | 0.27 | ResNet152v2 | 0.4 |
| **FaceNet** | **0.61** | MobileNet | 0.48 | MobileNetV2 | 0.38 |
| DenseNet121 | 0.47 | DenseNet169 | 0.44 | DenseNet201 | 0.49 |

**Table 2**
Model Selection: To select the best pre-trained model for feature extraction, we extract features from different pre-trained models and apply SVR on them to select the best model which gives the highest Pearson correlation for our problem.

correlation. To make our paper self contained, we define each of evaluation metric below.

$$MAE = (\frac{1}{n}) \sum_{i=1}^{n} |y_i - \tilde{y}_i| \tag{6}$$

where $n$ is the number of samples, $y_i$ is the ground truth and $\tilde{y}_i$ is the predicted BMI.

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^{n} \left(y_i - \tilde{y}_i\right)^2} \tag{7}$$

where $n$ is the number of samples, $y_i$ is the ground truth and $\tilde{y}_i$ is the predicted BMI. Finally the Pearson correlation is defined as:

$$Pearson = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{n}(x_i - \bar{x})^2}\sqrt{\sum_{i=1}^{n}(y_i - \bar{y})^2}} \tag{8}$$
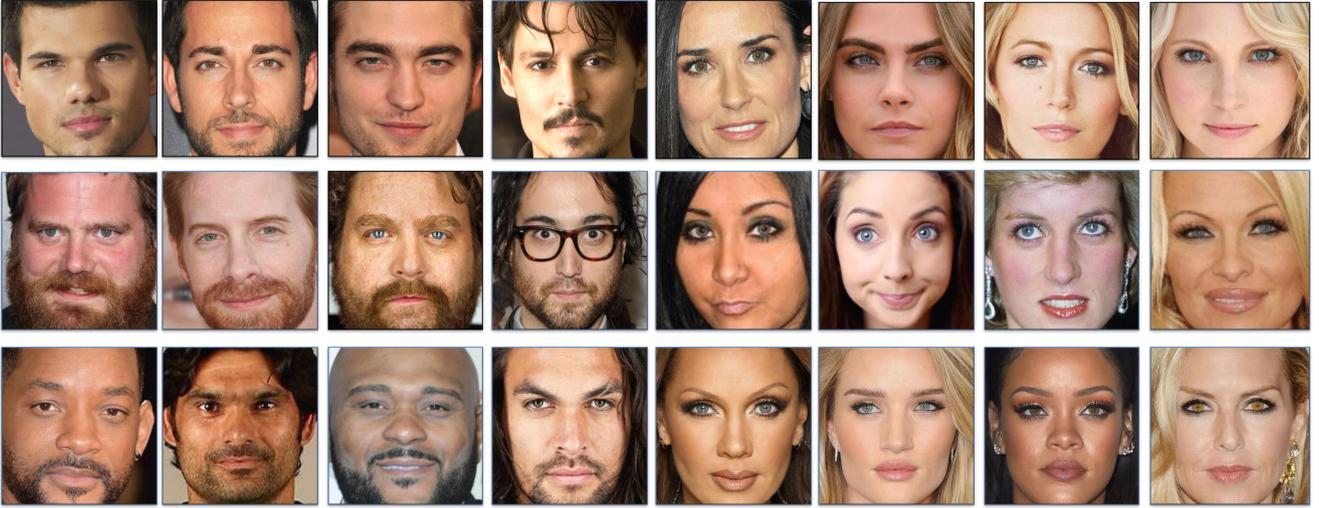
---

**Figure 7:** Examples of face images in the VIP-attribute dataset.

| Model | VisualBMI [18] | | | GAP | | | Reg-GAP | | |
|---|---|---|---|---|---|---|---|---|---|
| Metric | MAE | RMSE | Pearson | MAE | RMSE | Pearson | MAE | RMSE | Pearson |
| FaceNet | 5.38 | 7.51 | 0.61 | 5.23 | 7.04 | 0.645 | **5.03** | **6.92** | **0.663** |
| VGGFace | 5.16 | 7.16 | 0.65 | 5.22 | 7.03 | 0.644 | **4.99** | **6.94** | **0.659** |

**Table 3**
Results on the VisualBMI dataset using VGGFace and FaceNet models. Lower MAE and RMSE is better while the higher Pearson correlation is more useful. Results show that Reg-GAP results are better than that of GAP and [18].

where $n$ is the number of samples, $x_i$ and $y_i$ are the individual's ground truth and predicted BMI for person $i$, $\bar{x}$ and $\bar{y}$ are the mean of ground truth and predicted BMIs. Lower MAE and RMSE and higher Pearson correlation represent the improved results.

### 4.3. Face model selection

Recently several deep convolutional neural networks-based recognition models have been proposed for high accuracy face recognition. We have experimented with several of those models and selected the one which performs better for our problem of face to BMI prediction. Table 2 shows the experimental results. For model selection, we extracted the features from the second last layer of each model and applied $\epsilon$-SVR on the features to predict BMI. As can be seen that among several models, VGGFace [26] and FaceNet[32] have the highest Pearson correlation. Therefore, for our experiments, we have chosen VGG-Face and FaceNet models.

### 4.4. Experimental results on VisualBMI dataset

Table 3 shows the results of region-aware global pooling features on the VisualBMI dataset [18] using VGGFace and FaceNet models. Since the results using MAE and RMSE were not mentioned by the authors of [18], therefore, we have computed the results on these metrics using authors [18] code. The improved face to BMI predictions on all three evaluation metrics (MAE, RMSE, and Pearson correlation) demonstrate the usefulness of the proposed approach. The

superiority of Reg-GAP as compared to GAP shows that explicit feature pooling from different face regions is important and helps in better prediction of BMI from the face image.

#### 4.4.1. Class-level BMI prediction

In this section, we present the BMI class division of the VisualBMI data set. We followed the division of [18] and the details are shown in Table 4. All samples under 18.5 are labeled underweight while all samples above 40 are labeled very severely obese.

Table 5 shows the experimental results of BMI prediction for the different classes of datasets. Improved experimental results of Reg-GAP for all BMI-classes enforce our conjecture that, as compared to extracting only face-level features, face to BMI prediction methods should focus on various facial regions to get a better BMI prediction. Note that the 'very severely obese' class has the highest MAE and RMSE due to the large class variation.

#### 4.4.2. Gender prediction

To demonstrate the discriminative ability of the proposed Reg-GAP, we employ t-SNE to draw the features with GAP and Reg-GAP for the male and females. In Figure 8, red dots show feature vectors for males, and pink dots show feature vectors for females, where the feature vectors are extracted from VGG-Face models for the VisualBMI dataset. The better separation of Reg-GAP features as compared to that of GAP features demonstrates the usefulness of Reg-GAP fea-

| Class | Train Images | Test Images | BMI > | BMI < |
|---|---|---|---|---|
| Under Weight | 7 | 0 | 16 | 18.5 |
| Normal | 555 | 127 | 18.5 | 25 |
| Over Weight | 936 | 215 | 25 | 30 |
| Obese | 772 | 169 | 30 | 35 |
| Severely Obese | 541 | 140 | 35 | 40 |
| Very Severely Obese | 557 | 189 | 40 | - |

**Table 4**
BMI classes with train and test samples for VisualBMI dataset.

| | Class | Normal | Over Weight | Obese | Severely Obese | Very Severely Obese |
|---|---|---|---|---|---|---|
| MAE | GAP | 1.56 | 1.44 | 1.80 | 1.72 | 5.54 |
| | Reg-GAP | **1.25** | **1.28** | **1.22** | **1.23** | **4.91** |
| RMSE | GAP | 1.97 | 1.78 | 2.19 | 2.18 | 8.08 |
| | Reg-GAP | **1.56** | **1.52** | **1.48** | **1.47** | **7.69** |

**Table 5**
Comparisons of GAP with Reg-GAP for different BMI classes. Reg-GAP outpeforms GAP for all the classes for VisualBMI dataset.

| All | ResNet Based[8] | LD-PLS[12] | LD-CCA[12] | GAP | Reg-GAP |
|---|---|---|---|---|---|
| MAE | 2.36 | 2.26 | 2.23 | 1.85 | **1.73** |
| RMSE | - | - | - | 2.74 | **2.61** |
| Pearson | 0.55 | - | - | 0.71 | **0.75** |

**Table 6**
Results on the entire VIP attribute dataset. Lower MAE/RMSE and higher Pearson correlation are better. Results show that Reg-GAP results are better than that of GAP and [12].
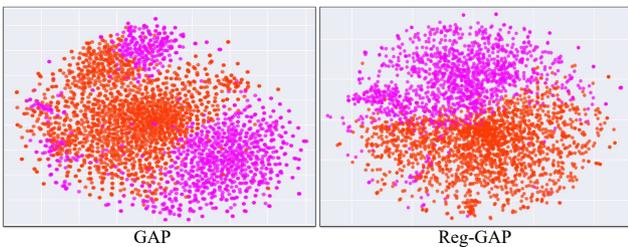
tures.



**Figure 8:** t-SNE of the features of GAP and Reg-GAP for gender classes in VisualBMI dataset.

## 4.5. Experimental results on VIP-Attribute Dataset

In this section, we show the results of the proposed approach on the VIP-attributes dataset and compare our method with the recent state-of-the-art approach. Table 6 shows our results for the entire VIP attribute dataset using the VGGFace model. Table 7 and Table 8 shows the experimental results for males and females separately for the VIP-attribute dataset.

| Female | ResNet Based[8] | LD-PLS[12] | LD-CCA[12] | GAP | Reg-GAP |
|---|---|---|---|---|---|
| MAE | 2.30 | 2.28 | 2.27 | 1.80 | **1.63** |
| RMSE | - | - | - | 2.84 | **2.53** |
| Pearson | 0.55 | - | - | 0.74 | **0.81** |

**Table 7**
Results on the female class of VIP attribute dataset. Results show that Reg-GAP results are better than that of GAP and [12].

| Male | ResNet Based[8] | LD-PLS[12] | LD-CCA[12] | GAP | Reg-GAP |
|---|---|---|---|---|---|
| MAE | 2.32 | 2.25 | 2.19 | 1.89 | **1.77** |
| RMSE | - | - | - | 2.71 | **2.66** |
| Pearson | **0.55** | - | - | 0.37 | 0.41 |

**Table 8**
Results on the male class of VIP-attribute dataset. Results show that Reg-GAP results are better than that of GAP and [12].

The overall and gender-wise experimental results demonstrate the usefulness of the proposed approach.

## 4.6. Run-Time Analysis

In our experiments, we used Google Colab with RAM of 25GB and ROM of 68GB, with the Tensor-flow version 1.14.0, Keras version 2.2, and Python 3.0. The optimizer used is Adam with learning rate of 0.001, and beta1=0.9, beta2=0.999, epsilon=0.48, decay=0.0. For each image, on average, the MTCNN module took 0.76s to detect the faces, the BiseNet took 0.07s for face semantic segmentation, and later reprocessing took 0.14s. The VGGFace module took 0.27s for generating Reg-GAP features and lastly the regression module on average took 2.63E-04s. The entire process of face to BMI prediction for each image took 1.24s.

## 5. Discussion and Concluding remarks

Determining BMI from facial photos is commonly used in recent studies where all of the previous methods have focused on the overall faces but we, on the other hand, have specifically tried to pool features from different face regions. To achieve accurate and pixel-wise localization, we employed face semantic segmentation. Our experimental results show that face to BMI prediction is improved with Reg-GAP as compared to using GAP. The graph in Figure 8 shows that, in addition to the face to BMI prediction problem, Reg-GAP features are more discriminative for the gender prediction problem. To the best of our knowledge, this is one of the first frameworks to utilize facial regions for the BMI prediction while maintaining the state of art accuracy. We have performed experiments on three publicly available datasets and comparisons over several evaluation metrics to validate the proposed ideas. Table 9 shows the result of the significance test on all three datasets with the p-value of 0.0014. It indicates that our results are significant. We have also tested

| Dataset | VIP_Attribute[8] | | | VisualBMI[18] | | Bollywood[34] | | P_Value |
|---|---|---|---|---|---|---|---|---|
| | Overall | Male | Female | FaceNet | VGGFace | VGG19 | VGGFace | |
| Recent Work | 2.23 | 2.19 | 2.27 | 5.38 | 5.16 | 1.49 | 0.97 | |
| Ours | 1.73 | 1.77 | 1.63 | 5.03 | 4.99 | 0.55 | 0.32 | 0.001392827177 |

**Table 9**
Results of significance tests on MAE of all three datasets.

the normality of our samples with the Shapiro–Wilk test and it was also passed with the p-value of 0.0170625

There are still some limitations in our framework. For example, semantic segmentation and feature extraction are being performed separately. Future work could include training of feature extractors and semantic segmentation in a joint framework. For instance, future work may use some student-teacher approach where training is done through the teacher model on semantically segmented facial images and the student model is enforced to predict the same embedding for the same input sample as the teacher model. This approach can let the student model behave as if it was using the segmentation without actually using the semantic segmentation.

The second limitation for the BMI prediction problem is the biases in the BMI-datasets. Mostly the underweight class and very severely obese class are biased because there are very few samples in the underweight class and the very severely obese class has a large range. This can be addressed using losses that handle class imbalance such as focal loss. Lastly, since the datasets are collected from social media or the world wide web, there may be some errors in the annotations. The VisualBMI dataset has many poor quality images as some of the images are images of the pictures and some are blurry. Although the VIP attribute dataset has images in good quality, however, these images are of actors and there is the presence of makeup, plastic surgery which can make the accurate prediction of BMI difficult.

## 6. Appendix

**Declaration of Competing Interest:** We wish to confirm that there are no known conflicts of interest associated with this publication and there has been no significant financial support for this work that could have influenced its outcome. We confirm that the manuscript has been read and approved by all named authors and that there are no other persons who satisfied the criteria for authorship but are not listed. We further confirm that the order of authors listed in the manuscript has been approved by all of us. We confirm that we have given due consideration to the protection of intellectual property associated with this work and that there are no impediments to publication, including the timing of publication, with respect to intellectual property. In so doing we confirm that we have followed the regulations of our institutions concerning intellectual property. We understand that the Corresponding Author is the sole contact for the Editorial process (including Editorial Manager and direct communications with the office). He/she is responsible for communicating with the other authors about progress, submissions of revisions, and final approval of proofs. We confirm that we have provided a current, correct email address which is accessible by the Corresponding Author.

**Credit authorship contribution statement:**
**Nadeem Yousaf:** conceptualization, software, writing - original draft, writing-review & editing **Sarfaraz Hussein:** conceptualization, writing-review, and editing, idea, formal analysis, supervision, project administration. **Waqas Sultani:** conceptualization, writing-review, and editing, idea, formal analysis, supervision, project administration.

## References

[1] Atwa, H., Fiala, L., Handoka, N.M., 2012. Neck circumference as an additional tool for detecting children with high body mass index. J Am Sci 8, 442–446.

[2] Benitez-Garcia, G., Nakamura, T., Kaneko, M., 2017. Analysis of in-and out-group differences between western and east-asian facial expression recognition, in: 2017 Fifteenth IAPR International Conference on Machine Vision Applications (MVA), IEEE. pp. 402–405.

[3] Bray, G.A., 2003. Risks of obesity. Endocrinology and Metabolism Clinics 32, 787–804.

[4] Byeon, K., Park, B.y., Park, H., 2019. Spatially guided functional correlation tensor: A new method to associate body mass index and white matter neuroimaging. Computers in biology and medicine 107, 137–144.

[5] Caliendo, M., Lee, W.S., 2013. Fat chance! obesity and the transition from unemployment to employment. Economics & Human Biology 11, 121–133.

[6] Coetzee, V., Chen, J., Perrett, D.I., Stephen, I.D., 2010. Deciphering faces: Quantifiable visual cues to weight. Perception 39, 51–61.

[7] Coetzee, V., Perrett, D.I., Stephen, I.D., 2009. Facial adiposity: A cue to health? Perception 38, 1700–1711.

[8] Dantcheva, A., Bremond, F., Bilinski, P., 2018. Show me your face and i will tell you your height, weight and body mass index, in: 2018 24th International Conference on Pattern Recognition (ICPR), IEEE. pp. 3555–3560.

[9] Drucker, H., Burges, C.J., Kaufman, L., Smola, A., Vapnik, V., et al., 1997. Support vector regression machines. Advances in neural information processing systems 9, 155–161.

[10] Fondón, I., Serrano, C., Acha, B., 2007. Segmentation of skin cancer images based on multistep region growing., in: MVA, Citeseer. pp. 339–342.

[11] Jiang, M., Guo, G., 2019. Body weight analysis from human body images. IEEE Transactions on Information Forensics and Security 14, 2676–2688.

[12] Jiang, M., Guo, G., Mu, G., 2020. Visual bmi estimation from face images using a label distribution based method. Computer Vision and Image Understanding 197, 102985.

[13] Jiang, M., Shang, Y., Guo, G., 2019. On visual bmi analysis from facial images. Image and Vision Computing 89, 183–196.

[14] Kalayeh, M.M., Basaran, E., Gökmen, M., Kamasak, M.E., Shah, M., 2018. Human semantic parsing for person re-identification, in: Pro-

ceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 1062–1071.

[15] Kalayeh, M.M., Gong, B., Shah, M., 2017. Improving facial attribute prediction using semantic segmentation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 6942–6950.

[16] Khan, K., Mauro, M., Migliorati, P., Leonardi, R., 2017. Gender and expression analysis based on semantic face segmentation, in: International conference on image analysis and processing, Springer. pp. 37–47.

[17] Kissebah, A.H., Freedman, D.S., Peiris, A.N., 1989. Health risks of obesity. Medical Clinics of North America 73, 111–138.

[18] Kocabey, E., Camurcu, M., Ofli, F., Aytar, Y., Marin, J., Torralba, A., Weber, I., 2017. Face-to-bmi: Using computer vision to infer body mass index on social media, in: Proceedings of the International AAAI Conference on Web and Social Media.

[19] Kocabey, E., Ofli, F., Marin, J., Torralba, A., Weber, I., 2018. Using computer vision to study the effects of bmi on online popularity and weight-based homophily, in: International Conference on Social Informatics, Springer. pp. 129–138.

[20] Kuan, K., Manek, G., Lin, J., Fang, Y., Chandrasekhar, V., 2017. Region average pooling for context-aware object detection, in: 2017 IEEE International Conference on Image Processing (ICIP), IEEE. pp. 1347–1351.

[21] Lal, S., Das, D., Alabhya, K., Kanfade, A., Kumar, A., Kini, J., 2021. Nucleisegnet: Robust deep learning architecture for the nuclei segmentation of liver cancer histopathology images. Computers in Biology and Medicine 128, 104075.

[22] Lee, C.H., Liu, Z., Wu, L., Luo, P., 2020. Maskgan: Towards diverse and interactive facial image manipulation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 5549–5558.

[23] Mayer, C., Windhager, S., Schaefer, K., Mitteroecker, P., 2017. Bmi and whr are reflected in female facial shape and texture: a geometric morphometric image analysis. PloS one 12, e0169336.

[24] Mussi, E., Servi, M., Facchini, F., Furferi, R., Governi, L., Volpe, Y., 2021. A novel ear elements segmentation algorithm on depth map images. Computers in Biology and Medicine 129, 104157.

[25] Panon, N., Luangsawang, K., Rugaber, C., Tongchit, T., Thongsepee, N., Cheaha, D., Kongjaidee, P., Changtong, A., Daradas, A., Chotimol, P., 2019. Correlation between body mass index and ocular parameters. Clinical Ophthalmology (Auckland, NZ) 13, 763.

[26] Parkhi, O.M., Vedaldi, A., Zisserman, A., 2015. Deep face recognition .

[27] Pascali, M.A., Giorgi, D., Bastiani, L., Buzzigoli, E., Henríquez, P., Matuszewski, B.J., Morales, M.A., Colantonio, S., 2016. Face morphology: Can it tell us something about body weight and fat? Computers in biology and medicine 76, 238–249.

[28] Pi-Sunyer, F.X., 2002. The medical risks of obesity. Obesity Surgery 12, S6–S11.

[29] Pi-Sunyer, X., 2009. The medical risks of obesity. Postgraduate medicine 121, 21–33.

[30] Qayyum, A., Lalande, A., Meriaudeau, F., 2020. Automatic segmentation of tumors and affected organs in the abdomen using a 3d hybrid model for computed tomography imaging. Computers in Biology and Medicine 127, 104097.

[31] Saka, M., Türker, P., Ercan, A., Kızıltan, G., Baş, M., 2014. Is neck circumference measurement an indicator for abdominal obesity? a pilot study on turkish adults. African health sciences 14, 570–575.

[32] Schroff, F., Kalenichenko, D., Philbin, J., 2015. Facenet: A unified embedding for face recognition and clustering, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 815–823.

[33] Shreve, M., Bala, R., Wu, W., Xu, B., Purwar, A., Matts, P., 2019. Region-wise modeling of facial skin age using deep cnns, in: 2019 16th International Conference on Machine Vision Applications (MVA), IEEE. pp. 1–6.

[34] Siddiqui, H., Rattani, A., Kisku, D.R., Dean, T., 2020. Ai-based bmi inference from facial images: An application to weight monitoring. arXiv preprint arXiv:2010.07442 .

[35] Simonyan, K., Zisserman, A., 2014. Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556 .

[36] Stubert, J., Reister, F., Hartmann, S., Janni, W., 2018. The risks associated with obesity in pregnancy. Deutsches Ärzteblatt International 115, 276.

[37] Wang, W., Fu, Y., Qian, X., Jiang, Y.G., Tian, Q., Xue, X., 2020. Fm2u-net: Face morphological multi-branch network for makeup-invariant face verification, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 5730–5740.

[38] Wen, L., Guo, G., 2013. A computational approach to body mass index prediction from face images. Image and Vision Computing 31, 392–400.

[39] Wen, Y., Zhang, K., Li, Z., Qiao, Y., 2016. A discriminative feature learning approach for deep face recognition, in: European conference on computer vision, Springer. pp. 499–515.

[40] Yu, C., Wang, J., Peng, C., Gao, C., Yu, G., Sang, N., 2018. Bisenet: Bilateral segmentation network for real-time semantic segmentation, in: Proceedings of the European conference on computer vision (ECCV), pp. 325–341.

[41] Zhang, K., Zhang, Z., Li, Z., Qiao, Y., 2016. Joint face detection and alignment using multitask cascaded convolutional networks. IEEE Signal Processing Letters 23, 1499–1503.