# How are attributes expressed in face DCNNs?

Prithviraj Dhar[1], Ankan Bansal[1], Carlos D. Castillo[1], Joshua Gleason[1], P. Jonathon Phillips[2]
and Rama Chellappa[1]

[1] University of Maryland, College Park, MD, USA

[2] National Institute of Standards and Technology, Gaithersburg, MD, USA

*Abstract*— As deep networks become increasingly accurate at recognizing faces, it is vital to understand how these networks process faces. While these networks are solely trained to recognize identities, they also contain face related information such as sex, age, and pose of the face. The networks are not trained to learn these attributes. We introduce expressivity as a measure of how much a feature vector informs us about an attribute, where a feature vector can be from internal or final layers of a network. Expressivity is computed by a second neural network whose inputs are features and attributes. The output of the second neural network approximates the mutual information between feature vectors and an attribute. We investigate the expressivity for two different deep convolutional neural network (DCNN) architectures: a Resnet-101 and an Inception Resnet v2. In the final fully connected layer of the networks, we found the order of expressivity for facial attributes to be Age > Sex > Yaw. Additionally, we studied the changes in the encoding of facial attributes over training iterations. We found that as training progresses, expressivities of yaw, sex, and age decrease. Our technique can be a tool for investigating the sources of bias in a network and a step towards explaining the network's identity decisions.

## I. INTRODUCTION

Deep convolutional neural networks (DCNN)-based face algorithms are trained to learn the identity of a face; they are not trained to learn attributes of the face. These DCNNs generate representations that encode identity. However, Hill emphet al. [1] found that DCNNs generated identity representations self-organized by sex. Also, identity representations can contain information on pose, age, and illumination direction [1], [2], [3].

Facial attributes, including those mentioned above, affect algorithm accuracy [4], [5]. Assessing bias in algorithms implies measuring the effect of these attributes on accuracy. Explaining how a network comes to an identity decision includes discerning how face representations encode attributes. To gain further understanding on the effects of attributes on bias and assist in developing methods to explain network decisions we address the following two questions. How much information about facial attributes are captured in the internal layers of the network? How does the encoding of facial attributes evolve as training progresses?

In this paper, we explore how attributes are encoded in the internal layers of two different and successful architectures: a Resnet-101 and an Inception Resnet v2 architecture based DCNN [6], [7]. Both networks are trained solely to identify faces. In addition, we examine how the encoding of attributes evolves over training iterations for both these two

networks. To gain a greater understanding of the relationship between attributes, we introduce the concept of *expressivity* of an attribute. Expressivity is a measure of how much a given representation informs us about an attribute, where the representation of a face can be from internal or final layers in a DCNN. The following are the conceptual and experimental contributions of our paper.

1) We are the first to investigate the encoding of facial attributes in the internal layers of DCNNs.
2) For DCNNs, we are the first to monitor the evolution of the encoding of facial attributes during training.
3) In the final fully connected layer of both networks, we observed that the order of expressivity for three facial attributes to be Age > Sex > Yaw.
4) We found that as training progresses, the expressivity of yaw, sex, and age decreases. The observed rate of decrease was Age < Sex < Yaw.
5) The expressivity of identity dramatically increases from the last pooling layer to the final fully connected layer.

Knowing how face attributes are expressed in internal layers and how their representations evolve during training has both scientific and societal importance. Since DCNN-based face recognition systems are fielded in the real-world, the need for these networks to be explainable is pressing. Understanding how internal layers encode attributes introduces new a tool to explain identity decisions and to examine the sources of bias in algorithms. From a scientific perspective, knowing the importance of attributes during training will provide insight into the training process.

## II. RELATED WORK

Significant research has been done in training state-of-the-art face recognition networks in the past few years [6], [7], [8], [9], [10]. Although, the interpretability of such networks has not been widely explored, several existing works explore explainability of deep networks for general visual recognition. These works can be divided into following two categories.

**Methods enforcing interpretability constraints during training**: Yin *et al.* [11] propose a spatial activation diversity loss as a constraint to preserve interpretability while training face recognition networks. Similarly, Kim *et al.* [12] propose a generative technique using the most representative exemplars (prototypes), thus highlighting interpretability of the model. As mentioned in [13], these

methods cannot be used to interpret pre-trained models and hence cannot be applied to networks or models which are already in use.

**Methods interpreting trained models** : TCAV [13] is one of the most influential work in this area. TCAV interprets a network on the basis of its sensitivity to user defined concepts (such as 'stripes'). This is done by learning Concept Activity Vectors (CAVs) by training a linear classifier to distinguish between the activations produced by a concept's examples. While this method works efficiently for discrete physical concepts, such as presence of a specific color or pattern, it cannot be directly modified for checking the sensitivity of a model to a more general continuous concept (such as pose angle, age etc.). This is because for training CAVs, we also need negative example images where the concept being studied is missing. It is not trivial to find such images when the concepts are omnipotent and continuous (facial yaw, age etc.). Also, the method requires the testing images to belong to one of the training classes since the sensitivity computations requires measuring the change in logits of the class being investigated. This cannot be easily modified for our requirement where we use unseen subjects/faces to estimate models' sensitivity to facial attributes. Another important work in this category is [14] where the authors use linear classifiers on different layers of a network to understand the role of intermediate layers. Koh and Liang [15] propose an influence function to measure the model's sensitivity to an infinitesimally-small local perturbation in the training images. However, such a local perturbation-based method cannot be used to estimate models' sensitivity to physical attributes like pose or orientation.

Another class of works [16], [17] interprets the output of a network by generating saliency/attention maps. While such techniques help to highlight the spatial regions which affected the network's prediction, they do not allow to test the models' sensitivity to user defined concepts. Moreover, this method cannot be applied for concepts which cannot be physically localized (such as facial yaw, age etc.).

Hill *et al.* [18] is one of the few works which interpret trained face recognition networks, where the authors show the following hierarchy: face identity nested under sex, illumination nested under identity, and viewpoint nested under illumination.

We interpret a trained face recognition network by investigating its sensitivity to facial attributes. Previous methods like [15], [13] rely on the change in prediction with respect to a concept/attribute to interpret a network's sensitivity to an attribute. However, we introduce a new measure called expressivity, which quantifies the predictability of an attribute in a given set of features extracted using the model. Moreover, expressivity can be computed for both categorical and continuous attributes, which enables us to compare the predictability of various attributes.

## III. EXPRESSIVITY

Predictability of facial attributes/identities in a given set of face descriptors indicates the attribute-relevant information content encoded in the descriptors. To estimate this information content, we intend to use Mutual Information (MI). MI between two random variables is a measure of the amount of information that can be obtained for one random variable by observing the other variable. Therefore, if we estimate MI between face descriptors and their corresponding identities/attributes, we can estimate the information content of these identities/attributes in the given descriptors. Since MI can be computed for both categorical and continuous attributes, an estimate based on MI provides a measure which is consistent across categorical and continuous attributes.

MI between two random variables $(V_1, V_2)$ is given as:

$$I(V_1, V_2) = D_{KL}(\mathbb{P}_{V_1,V_2} \| \mathbb{P}_{V_1} \otimes \mathbb{P}_{V_2}) \qquad (1)$$

where, $D_{KL}$ represents the Kullback-Leibler divergence, $\mathbb{P}_{V_1,V_2}$ denotes the joint probability distribution, $\mathbb{P}_{V_1}$ and $\mathbb{P}_{V_2}$ denote the marginal distributions, and $\mathbb{P}_{V_1} \otimes \mathbb{P}_{V_2}$ represents the product of the marginal distributions. Tishby and Zaslavsky [19] show that each layer in a deep network can be quantified by the amount of mutual information (MI) it retains on the input variable, on the (desired) output variable. However, as mentioned in [20], computing MI is not a trivial task. Most of the existing non-parametric approaches for estimating MI do not scale with the dimensionality of variables. Belghazi *et al.* [20] propose MINE to estimate MI between high dimensional continuous variables using gradient descent over neural networks. The neural information measure has been defined in [20] as follows.

$$I_\Theta(F, A) = \sup_{\theta \in \Theta} \mathbb{E}_{\mathbb{P}_{FA}}[T_\theta] - \log(\mathbb{E}_{\mathbb{P}_F \otimes \mathbb{P}_A}[e^{T_\theta}]) \qquad (2)$$

where $F, A$ are the variables whose mutual information is to be estimated, $\theta \in \Theta$ represents parameters in a network computing a function $T_\theta : F, A \longrightarrow \mathbb{R}$. As proved in [20], the MINE estimate provides a lower bound estimate to the actual MI.

In the context of our work, we define 'Expressivity' of $A$ in $F$ as the aforementioned information measure (Equation 2). We use MINE to compute the expressivity of face identity and various facial attributes $A$ (discrete and continuous) in a given set of face descriptors $F$. Although face identity is not strictly a face 'attribute', we treat it in the same way as a face attribute. Therefore, in this work, attribute $A$ collectively refers to identity and facial attributes like pose, sex etc.

Following the protocols detailed in [20], we briefly explain how gradient descent over a neural network can be used to compute expressivity. Let $f_i \in \mathbb{R}^m$ denote $i^{th}$ feature in a batch $B$ of size $b$ (i.e. $|B| = b$), and $a_i \in \mathbb{R}$ denote the corresponding attribute value. The set $\{(f_i, a_i)\}_{i=1}^b$ represents the $b$ elements sampled from the joint distribution $(f_i, a_i \sim \mathbb{P}_{FA})$. Similarly $\{\tilde{a}_i\}_{i=1}^b$ represents $b$ attribute values sampled from a marginal distribution $(\tilde{a}_i \sim \mathbb{P}_A)$. To

estimate the the neural information in 2, we compute the expectation over joint and marginal distribution as follows :

$$\mathbb{E}_{\mathbb{P}_{FA}^{(b)}}[T_\theta] = \frac{1}{b} \sum_{i=1}^{b} T_\theta(f_i, a_i)$$

$$\mathbb{E}_{\mathbb{P}_{F}^{(b)} \otimes \mathbb{P}_{A}^{(b)}}[e^{T_\theta}] = \frac{1}{b} \sum_{i=1}^{b} e^{T_\theta(f_i, \tilde{a}_i)}$$

where $b$ is the number of features in a batch $B$ whose mutual information to be computed with their corresponding attributes. We use a network with parameter set $\theta$ (see Fig. 2) to compute the aforementioned arbitrary function $T_\theta(f_i, a_i)$ and $T_\theta(f_i, \tilde{a}_i)$. By substituting $\mathbb{E}_{\mathbb{P}_{FA}}[T_\theta]$ and $\mathbb{E}_{\mathbb{P}_F \otimes \mathbb{P}_A}[e^{T_\theta}]$ in Eq. 2 , the function $\mathcal{V}(\theta)$ is computed as:

$$\mathcal{V}(\theta) = \frac{1}{b} \sum_{i=1}^{b} T_\theta(f_i, a_i) - \log\left(\frac{1}{b} \sum_{i=1}^{b} e^{T_\theta(f_i, \tilde{a}_i)}\right) \quad (3)$$

As mentioned in [20], the supremum of $\mathcal{V}(\theta)$ with respect to parameter set $\theta$ is a lower bound approximation of the mutual information between features $F$ and attributes $A$. Hence we use the following function $L$ as our objective function to train the network $N$.

$$L(\theta) = -\mathcal{V}(\theta) \quad (4)$$

$$\nabla_\theta L(\theta) = -\left( \mathbb{E}_{\mathbb{P}_{FA}^{(b)}}[\nabla T_\theta] - \frac{\mathbb{E}_{\mathbb{P}_{F}^{(b)} \otimes \mathbb{P}_{A}^{(b)}}[\nabla T_\theta e^{T_\theta}]}{\mathbb{E}_{\mathbb{P}_{F}^{(b)} \otimes \mathbb{P}_{A}^{(b)}}[e^{T_\theta}]} \right) \quad (5)$$

At every training iteration, we use a different parameter set $\theta$ to compute the function $T_\theta$. As the training proceeds, the network minimizes (4), thus maximizing $\mathcal{V}(\theta)$ with respect to $\theta$. This is equivalent to computing the supremum of $\mathcal{V}(\theta)$. Hence, the final value of $\mathcal{V}(\theta)$ at convergence is an approximation of Expressivity. (Eq. 2)

In equations (3) and (4), the objective function (and the lower bound of mutual information) is computed in a batch and not on the given set of features and their attribute values, thus making the gradients biased towards the minibatch, rather than the full batch. This issue has been identified in [20], and can be mitigated by replacing the expectation term in the denominator of gradient update (5) by an exponential moving average. More theoretical details related to this approximation of mutual information is provided in [20].

## IV. PROPOSED APPROACH

### A. Networks and datasets used

To compute the expressivity of identity or attribute $A$ in a given set of features $F$, we extract the features $F$ using the following networks :
(1) **Network A** (Resnet-101 architecture) : The architecture is described in [6]. For investigating hierarchical course of the feedforward pass, we use a version of this network trained on a combined dataset of all the MS-Celeb-1M and UMD Faces images. For this trained network, we compute expressivity of attributes using features from these layers : Res4a2b, Res5a2c, Pool5, FC-L2S.

| Layer | Description |
| --- | --- |
| Res4a2b | 2nd convolutional layer of 7th Res-block |
| Res5a2c | 3rd convolutional layer of 9th Res-block |
| Pool5 | Final pooling layer which takes in the output of the 11th Res-block Res_5c |
| FC-L2S | Final fully connected which takes in the output of Pool5 and computes $L_2$ softmax activation |

(a) Network A

| Layer | Description |
| --- | --- |
| a9_concat | Concatenation layer of 10th inception block |
| b5_concat | Concatenation layer of 17th inception block |
| c3_concat | Concatenation layer of 36th inception block |
| Pool8x8 | Final pooling layer which takes in the output of the convolutional layer after 42nd inception block c10 |
| FC-L2S | Final fully connected which takes in the output of Pool8x8 and computes $L_2$ softmax activation |

(b) Network B

TABLE I: Brief descriptions of the layers used to compute expressivity in both the networks. More architectural details are provided in [6].
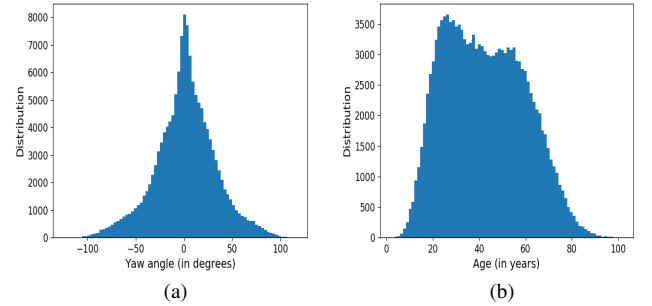


(a)      (b)

Fig. 1: IJB-C dataset [21] shows enough variation with respect to (a.) Yaw, (b.) Age, which is required to compute expressivity of age and yaw, in a given set of IJB-C features.
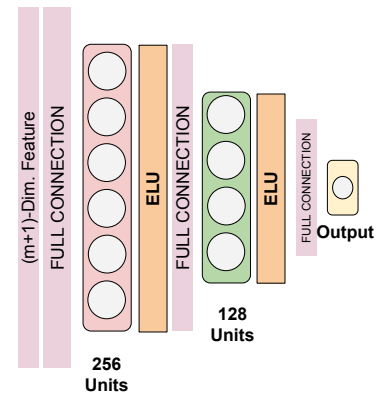


Fig. 2: We use a consistent network architecture to compute expressivity of $m-$dimensional features $F$, augmented with an attribute vector $A$. More details in Sec. IV-C

(2) **Network B** (Inception Resnet v2 architecture). The architecture is described in [6]. The training dataset of this network is same as that of Network A. We explore the following layers in the network: `a9_concat, b5_concat, c3_concat, Pool8x8, FC-L2S`.

Layer-wise details for both the networks are provided in Table I. For computing expressivity, we use IJB-C [21] as our test dataset. IJB-C dataset consists of 3531 identities with a total of 31,334 still images and 117,542 video frames collected in unconstrained settings. We extract the IJB-C features $F$ from different layers of aforementioned trained networks. For computing the corresponding attributes $A$, we use the All-in-one network [22].

### B. Attributes used

We compute the expressivity of identity and three attributes : yaw, sex, age in the extracted IJB-C features. To compute the yaw, sex and age of the corresponding IJB-C images, we use the All-in-one deep CNN proposed in [22], which simultaneously performs face detection, landmarks localization, pose estimation, sex recognition, smile detection, age estimation and face identification and verification. In Figure 1, we verify that the IJB-C datasets show enough variation with respect to yaw and age, so that we can insure that expressivity (which is a lower bound estimate of mutual information) is an accurate model for the corresponding attributes. The IJB-C dataset consists of 2203 male and 1328 female identities, which ensures that sufficient sex variation exists.

We now briefly explain the attribute vector $A$ introduced in Section III. When considering sex, the vector $A$ consists of probability values of the sex being male (`PR_MALE`), which are the outputs of All-in-One network in [22]. When considering identity, $A$ is a discrete vector, where each element is a numerical identity label. We generate the attribute vector $A$ for yaw (in degrees) and age (in years) in a similar manner. The exact methodology to compute the expressivity using the feature descriptors and their respective attribute vector is provided in the next subsection.

---
**Protocol 1** Computing expressivity using flattened features

---
1: **Input**: layer $L$,
2: **Input**: Set of $n$ images $I$
3: **Input**: attribute vector $A \in \mathbb{R}^{n \times 1}$
4: **Initialize** $E = []$
5: For a given image $i \in I$, extract feature $f_i$
6: Augmentation step $X = [F|A]$, where $F = [f_1, f_2 \ldots f_n]^T$
7: **for** $iter$ : 1 to M **do**
8:    Initialize MINE network according to dimensions of $X$
9:    $E \leftarrow$MINE$(X)$
10: **end for**
11: **return** Expressivity = Average$(E)$

---

---
**Protocol 2** Computing expressivity using unflattened feature maps

---
1: **Input**: layer $L$, with $k$ channels, each of dimension $d \times d$
2: **Input**: Set of $n$ images $I$
3: **Input**: attribute vector $A \in \mathbb{R}^{n \times 1}$
4: **Initialize** $E = []$
5: $S =$ Subset of randomly selected $z$ (out of $k$) channels
6: For a given image $i$, vectorize and concatenate all $z$ maps in $S$, to generate vector $f_i$ of dimension $m \times 1$, where $m = d * d * z$
7: Augmentation step $X = [F|A]$ where $F = [f_1, f_2 \ldots f_n]^T$
8: **for** $iter$ : 1 to M **do**
9:    Initialize MINE with input dimensions of $X \in \mathbb{R}^{n \times m+1}$
10:    $E \leftarrow$MINE$(X)$
11: **end for**
12: **return** Expressivity = Average$(E)$

---

### C. Protocols to compute expressivity

In this work, we define expressivity as a lower bound approximation of MINE, as explained in section III. We consider two protocols (described above): Protocol-1, to compute the expressivity of attributes in flattened features (i.e. features from fully connected or pooling layers which are extracted as vectors). Protocol-2, to compute the expressivity of attributes in unflattened feature maps (i.e. feature maps from a convolutional layer which are extracted as 2D channels). For both of these protocols, we need a set of $n$ images, and their corresponding attributes $A \in \mathbb{R}^{n \times 1}$. In step 8 of Protocol-1 and step 9 of Protocol-2, we initialize a MINE approximation network according to the input dimension of the augmented matrix. As explained in Section III, we train the network to compute the lower bound approximation of mutual information between features $F$ and attribute $A$. We use a simple multi layer perceptron (MLP) network, described in Figure 2, for computing $T_\theta$ in (4). The network consists of two hidden layers with 256 and 128 units. These layers are followed by ELU activations. We use this architecture consistent throughout our experiments. When using different sets of features for $F$, the only architectural changes in the network are made in the input layer dimension, according to the feature dimension. The network is trained until the loss in (4) converges and expressivity is computed using the protocols. In step 8 in Protocol-1 and step 9 in Protocol-2 we initialize and train the MINE network multiple times ($M$) to increase number parameter sets $\theta$ in (2), among which the supremum is to be found. In all our experiments, we use $M = 16$. Also, in step 5 of Protocol-2 we select only a subset of features maps (channels), as mentioned in [14]. Note that, for any attribute vector $A$ (yaw/identity/age/sex), we use the same feature subset $S$.

Apart from analyzing the network layer features, we also investigate the presence of various attributes in the raw RGB

| Layer | Dim. ($c \times d \times d$) | Protocol | $z$ | Feat. dim |
|---|---|---|---|---|
| Res4a2b | $1024 \times 14 \times 14$ | 2 | 11 | 2156 |
| Res5a2c | $2048 \times 7 \times 7$ | 2 | 42 | 2057 |
| Pool5 | 2048 | 1 | - | 2048 |
| FC-L2S | 512 | 1 | - | 512 |

(a) Network A

| Layer | Dim. ($c \times d \times d$) | Protocol | $z$ | Feat. dim |
|---|---|---|---|---|
| a9_concat | $128 \times 35 \times 35$ | 2 | 3 | 3675 |
| b5_concat | $384 \times 17 \times 17$ | 2 | 13 | 3757 |
| c3_concat | $128 \times 8 \times 8$ | 2 | 59 | 3776 |
| Pool8x8 | 1536 | 1 | - | 1536 |
| FC-L2S | 512 | 1 | - | 512 |

(b) Network B

TABLE II: Network A and B layers used to compute expressivity. $c$ : number of channels in a given layer, $d$: channel dimension, $z$ : Number of channels selected out of $c$ channels, to generate a subset of features

face image in Section V. For this, we average the R,G and B channels to generate a grayscale image. Following this, we vectorize the image and use it as feature $f_i$ in Protocol-1.

### D. Training linear classifiers

To verify that expressivity correctly models the information of attributes in features, we show its correlation with error-rates of linear classifiers trained on the corresponding features. We randomly select a subset of 5000 IJB-C images and extract their features. To train the linear classifier we use 3000 features and test it on 2000 features. This is a trivial task for flattened features. However, to compute the error rate in feature maps from higher layers of the network, we use the same subset $S$ of feature maps as selected in step 5 of Protocol-2. Following this, we vectorize and concatenate them as in Step 7. We provide more specific details in the next section.

### V. EXPERIMENTS

Using the protocols described in Section IV and the expressivity measure defined in III, we extract features from different layers of Networks A and B and use them to train a network (Figure 2) to compute expressivity of various attributes.

### A. Hierarchical course of feedforward pass

Table II shows the network layers explored, along with the final dimension of the features used for computing expressivity for both the networks. The layerwise expressivity values for Networks A and B are shown in Figure 3. It should be noted that both the networks were trained using identity-supervision and no supervision based on pose, sex and age. Our inference is listed as follows:

- In both networks A and B, we find that the expressivity of yaw, sex and age is high and that of identity is the lowest in the shallower layers (Res4a2b, Res5a2c in Network A; a9_concat, b5_concat, c3_concat in Network B) and input image. This shows that yaw and

sex are high level face features as compared to identity, which cannot be extracted using shallow layers.
- As we examine the deepest layer (FC-L2S), the expressivity of yaw and sex attain their lowest values, whereas identity and age have very high expressivity. This shows that identity and age are more fine grained features compared to other attributes.
- There is a rapid increase in the identity expressivity from the pooling to fully connected layers, in both the networks.
- Comparing the expressivity values of all attributes except identity in the final layer, we can infer that for identity recognition, yaw is the least important and age is the most important attribute.

There are three reasons for the relatively high expressivity of the age when compared to the other attributes.

First, in the IJB-C dataset most of the images were acquired over a short period of time, therefore, it can happen that a given age correlates with identity. Second, we used an automated algorithm [22] to estimate the age, and therefore it is computed from the appearance of the face. All attributes were computed automatically, but we can expect that in relative terms the age is the least accurate of all attributes automatically computed. Third, the entropy of the age (Fig. 1(b)) is higher than the entropy of the other attributes and this could increase the mutual information component of the expressivity.

**Discussion of the data processing inequality**: The data processing inequality (DPI) [23] states that for three random variables $P, Q, R$ forming a Markov chain $P \longrightarrow Q \longrightarrow R$,

$$\text{MI}(P, Q) \geq \text{MI}(P, R).$$

The data processing inequality formalizes the concept that no processing of data can increase mutual information. To make this more concrete, let $P$ be any random variable (e.g. sex, yaw, identity), and let $Q, R$ be features for different layers in a network where $R$ is a deterministic function of $Q$, i.e. $R$ is deeper than $Q$. Since $R$ is a function of $Q$ then $P, Q, R$ forms a Markov chain [23]. It follows from the DPI that the information about $P$ contained in the features cannot increase as we go deeper.

The expressivity results in Figure 3 are not monotonically decreasing, which might seem like a contradiction to DPI. However, as pointed out in [14], the features in our context denote *representation*, rather than *information content* described in Information Theory. *Representations* are more closely related to predictability of a specific attribute, as compared to information-theoretic content. Hence, in this work, expressivity refers to the accord between and attribute and attribute, rather than its theoretical information content.

**Relation with linear separability**: Alain and Bengio [14] show that the linear separability of features with respect to output classes, which provided supervision, monotonically increases as we go deeper into the network. From Figure 3, we find this to be true for expressivity identity as well, which provided supervision during training .

| Layer | Yaw regress$^n$ error | Yaw expressivity |
|-------|-----------------------|------------------|
| Res4a2b | 11.42 | 1.36 |
| Res5a2c | 8.64 | 1.48 |
| Pool5 | 11.57 | 1.23 |
| FC-L2S | 11.65 | 0.59 |

TABLE III: Comparison of yaw regression errors with their corresponding expressivity values, in different layers of Network A. The highest accuracy (or lowest error) corresponds to highest expressivity and lowest accuracy (or highest error) corresponds to lowest expressivity.

| Iteration | Yaw error | Yaw expr. | Age error | Age expr. |
|-----------|-----------|-----------|-----------|-----------|
| $T_1$ | 9.40 | 1.13 | 8.06 | 1.22 |
| $T_2$ | 11.27 | 0.68 | 7.70 | 1.39 |
| $T_3$ | 11.65 | 0.59 | 8.11 | 1.19 |

TABLE IV: Comparison of age and yaw regression errors with their corresponding expressivity values in 3 iterations $T_1, T_2, T_3$. For yaw, $T_1, T_2, T_3 = 25k, 100k, 200k$ iterations. For age $50k, 100k, 200k$

In order to ensure that the expressivity values correlate with feature vectors, we compute the accuracy/error rate obtained by training a linear classifier and testing it directly using features from the aforementioned layers in Network A, as explained Section IV. To analyze the yaw expressivity values, we first train a simple linear regression model on 3000 randomly selected features (extracted from IJB-C images) and evaluate its regression error on 2000 IJB-C features. The corresponding results are presented in Table III, from which we can infer that expressivity values do correlate with regression errors for yaw.

### B. Temporal course of training

We also analyze the training process of Network A and B and investigate the changes in the expressivity of yaw, sex, identity and age in the final layer (FC-L2S) of these networks with respect to its training iterations. The features at all iterations ($> 0$) are 512 dimensional and are flattened and Protocol-1 is used for computing attribute-wise expressivity, along with specifications for FC-L2S mentioned in Table II. The features at iteration 0, represent the final layer features of the networks trained on ImageNet [24], without the final fully connected layer for identity recognition. These features are therefore 2048 and 1536 dimensional for networks A and B respectively. The results are presented in Figure 4. Our observations are listed below:

- We find that the expressivity of yaw, age and sex reach their peak values in the first 25000 iterations for Network A and 40000 iterations for Network B, to learn the general concept of facial pose, age and sex.
- Following that, we find that the yaw expressivity decreases rapidly as the training proceeds, showing that making features almost agnostic to pose variance is an essential part of the training process.
- For both networks A and B, the expressivity of age and sex decreases slightly after their corresponding expressivity peaks are attained, during the course of training. However, compared to yaw, the rate of decay in the expressivity of age and sex is low. This shows that age and sex are more important for identity recognition than face yaw.
- Observing the expressivity values in the final iteration, we can infer that for identity recognition, the following is the order of relevance of attributes for which the network does not receive any supervision : Age > Sex > Yaw. The

opposite of this order is observed in the rate by which the expressivity values of yaw, age and sex decreases, i.e. the rate of decrease is : Age < Sex < Yaw. This is true for both networks A and B.
- The expressivity of identity generally increases during training for both networks A and B.
- Features extracted from the final layer of Iteration 0 model (Imagenet features), express identity better than other attributes. This is because the Imagenet features express 'objectness', which is closely related to identity as compared to other attributes (yaw, age, sex).

Similar to what we did in Section V-A, we compare the expressivity values to the corresponding error rates by training and testing linear classifiers directly on the final fully connected layer features. The results are presented in Table IV, where we again find that there exists correlation between expressivity values and age/yaw regression errors.
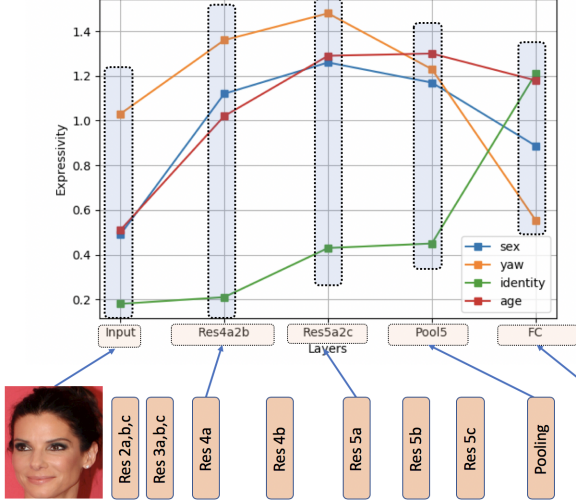
### C. Advantage of expressivity over other techniques

The following are the advantages of Expressivity over existing interpretability techniques:
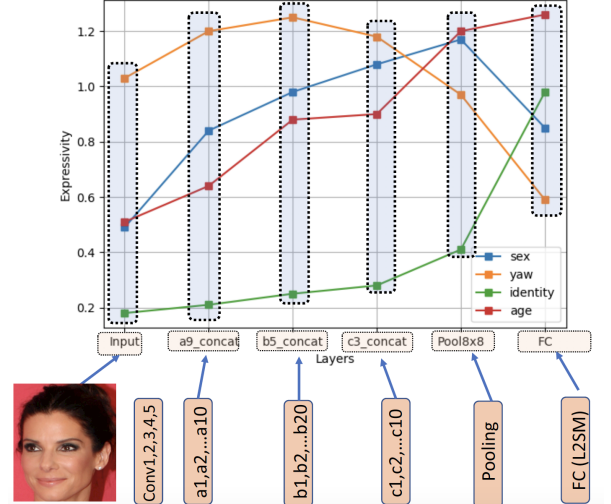
**Comparison of several attributes**: Comparing expressivity of different attributes in a given set of features is important for understanding the information organization in a trained network. As shown in Figures 4 and 3, expressivity helps put the error rates of all attributes on the sample scale, thus enabling their comparisons. This cannot be achieved by directly using the accuracy/error rates of linear classifiers, as different attributes have different scales and evaluation metrics.

**Useful for any physical concept (discrete/continuous)**: Although we used expressivity to analyze face recognition networks in terms of facial attributes, it can be used for a network with respect to any physical attribute. For instance, we can compute the expressivity of 'stripes' concept in a set of features $F$ if we have a binary attribute vector $A$, denoting the presence or absence of 'stripes'. This is similar to TCAVs [13]. However, TCAVs cannot be directly used to quantify the content of continuous concept (like pose angle), since we need images with which demonstrate absence of that concept to train CAVs, and this is not trivial for concepts like pose angle, age etc.

**No dependence on training classes**: Methods like [13] and [16] require computing change of logit values. However, for images not in training classes, this value is not meaningful. Using expressivity allows to measure information about
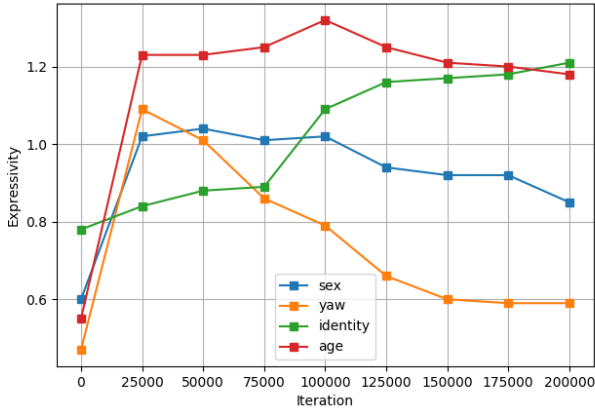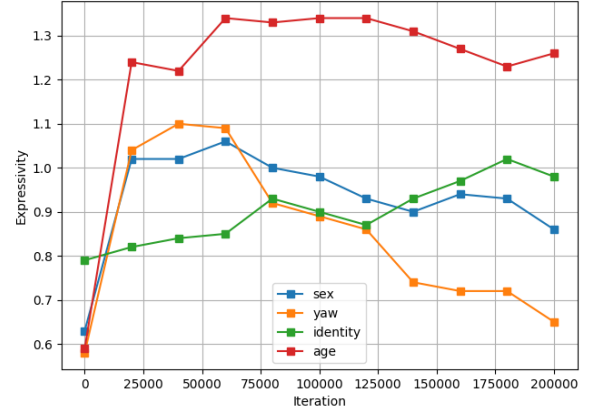
(a) Network A      (b) Network B

Fig. 3: Expressivity of identity, age, sex and yaw in input image and layerwise features from both the networks. The source image for face in this figure is attributed to Eva Rinaldi under the [cc-by-sa-2.0] creative commons licenses respectively. The face was cropped from the source image.



(a)      (b)

Fig. 4: Expressivity of identity, age, sex and yaw in final layer (FC-L2S) features of a.) Network A and b.) Network B. Decreasing expressivity of task irrelevant attributes (yaw, age, sex) is a part of training. Observed rate of decrease : Age < Sex < Yaw

attributes that were not explicitly included in training.

### D. Relation with information bottleneck theory

Saxe et al [25] provide a response to the Information Bottleneck theory [26], and claim that when an input domain consists of a subset of task-relevant and task-irrelevant information, hidden representations do compress the task-irrelevant information. In our context, identity is the task-relevant information, since the networks only receive identity supervision. The task-irrelevant information includes yaw, age and sex information. In our work, we concur with the claims of [25]. It can be seen in Figures 4 the expressivity of task-irrelevant attributes (yaw, age, sex) in the final layer decreases as we train the network, depicting the 'compression'

phase. This occurs while the identity expressivity increases as the training progresses, which corresponds to the 'fitting' phase of the network. Hence, we also verify another result in [25], that this compression happens concurrently with the fitting process rather than during a subsequent compression period.

## VI. DISCUSSION AND FUTURE WORK

We present an approach to quantify the information learned by a face recognition network about several attributes and identities, by computing their expressivity in a given set of features. The scale of this measure is agnostic to the attribute being examined. We use this measure to analyse

layer-wise features, and temporal snapshots of the final fully connected layer in two face recognition networks. We make some important observations in both the networks we investigated: (1) The mid-layer and shallow layer features effectively capture task-irrelevant information (about yaw, sex, age). The deeper layers encode task-relevant identity information. (2.) During the training process, the expressivity of identity increases while that of yaw, sex and age decreases, thus showing that *decreasing expressivity of task-irrelevant attributes is a part of learning*. Expressivity of yaw, especially, decreases very rapidly. (3.) Using the expressivity values in the final layer of trained networks, we find the following order of attribute-wise relevance for identity recognition : Age > Sex > Yaw. This is opposite to the order of the rate by which expressivity of these three attributes decrease during training. We also relate our findings with existing works on interpretability and information bottleneck theory.

There are other face attributes, such as facial expression, presence of eyeglasses, beard, hairstyle etc. which play a crucial role in identity recognition. One future avenue of research is to extend our work for these attributes and obtain a more exhaustive ordering of face attributes in terms of their relevance to face recognition. Also, one could investigate the training processes of layers other than fully connected layers. Finally, one could estimate the expressivity of attributes in face descriptors extracted from other networks, like [27] and [28].

## REFERENCES

[1] M Q Hill, C J Parde, C D Castillo, Y I Colon, R Ranjan, J-C Chen, V Blanz, and A J O'Toole. Deep convolutional neural networks in the face of caricature: Identity and image revealed. *Nature Mach Intel*, (in press).

[2] S Nagpal, M Singh, R Singh, M Vatsa, and N K Ratha. Deep learning for face recognition: Pride or prejudiced? *CoRR*, abs/1904.01219, 2019.

[3] C J Parde, C D Castillo, M Q Hill, Y I Colon, S Sankaranarayanan, Jun-Cheng Chen, and Alice J OToole. Face and image representation in deep cnn features. In *2017 12th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2017)*, pages 673–680. IEEE, 2017.

[4] G. H. Givens, J. R. Beveridge, P. J. Phillips, B. A. Draper, Y. M. Lui, and D. S. Bolme. Introduction to face recognition and evaluation of algorithm performance. *Computational Statistics and Data Analysis*, 67:236–247, 2013.

[5] Y. Lee, P. J. Phillips, J. J. Filliben, J. R. Beveridge, and H. Zhang. Generalizing face quality and factor measures to video. In *International Joint Conference on Biometrics (IJCB)*, 2014.

[6] R Ranjan, A Bansal, J Zheng, H Xu, J Gleason, B Lu, A Nanduri, J-C Chen, C D Castillo, and R Chellappa. A fast and accurate system for face detection, identification, and verification. *IEEE Transactions on Biometrics, Behavior, and Identity Science*, 1(2):82–96, 2019.

[7] A Bansal, R Ranjan, C D Castillo, and R Chellappa. Deep features for recognizing disguised faces in the wild. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 10–106. IEEE, 2018.

[8] F Schroff, D Kalenichenko, and J Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 815–823, 2015.

[9] Y Taigman, M Yang, M Ranzato, and L Wolf. Deepface: Closing the gap to human-level performance in face verification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1701–1708, 2014.

[10] J Deng, J Guo, and S Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. *preprint arXiv:1801.07698*, 2018.

[11] B Yin, L Tran, H Li, X Shen, and X Liu. Towards interpretable face recognition. *arXiv preprint arXiv:1805.00611*, 2018.

[12] B Kim, C Rudin, and J A Shah. The bayesian case model: A generative approach for case-based reasoning and prototype classification. In *Advances in Neural Information Processing Systems*, pages 1952–1960, 2014.

[13] B Kim, M Wattenberg, J Gilmer, C Cai, J Wexler, F Viegas, and R Sayres. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). *preprint arXiv:1711.11279*, 2017.

[14] G Alain and Y Bengio. Understanding intermediate layers using linear classifier probes. *preprint arXiv:1610.01644*, 2016.

[15] P W Koh and P Liang. Understanding black-box predictions via influence functions. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1885–1894. JMLR. org, 2017.

[16] R R Selvaraju, M Cogswell, A Das, R Vedantam, D Parikh, and D Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 618–626, 2017.

[17] A Chattopadhay, A Sarkar, P Howlader, and V N Balasubramanian. Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 839–847. IEEE, 2018.

[18] M Q Hill, C J Parde, C D Castillo, Y I Colon, R Ranjan, JC Chen, V Blanz, and A J O'Toole. Deep convolutional neural networks in the face of caricature: Identity and image revealed. *arXiv preprint arXiv:1812.10902*, 2018.

[19] N Tishby and N Zaslavsky. Deep learning and the information bottleneck principle. In *2015 IEEE Information Theory Workshop (ITW)*, pages 1–5. IEEE, 2015.

[20] M I Belghazi, A Baratin, S Rajeswar, S Ozair, Y Bengio, A Courville, and R D Hjelm. Mine: mutual information neural estimation. *arXiv preprint arXiv:1801.04062*, 2018.

[21] Brianna Maze, Jocelyn Adams, James A Duncan, Nathan Kalka, Tim Miller, Charles Otto, Anil K Jain, W Tyler Niggel, Janet Anderson, Jordan Cheney, et al. Iarpa janus benchmark-c: Face dataset and protocol. In *2018 International Conference on Biometrics (ICB)*, pages 158–165. IEEE, 2018.

[22] R Ranjan, S Sankaranarayanan, C D Castillo, and R Chellappa. An all-in-one convolutional neural network for face analysis. In *2017 12th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2017)*, pages 17–24. IEEE, 2017.

[23] T M Cover and J A Thomas. *Elements of information theory*. John Wiley & Sons, 2012.

[24] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and F.-F. Li. Imagenet: A large-scale hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, June 2009.

[25] A M Saxe, Y Bansal, J Dapello, M Advani, A Kolchinsky, B D Tracey, and D D Cox. On the information bottleneck theory of deep learning. 2018.

[26] R Shwartz-Ziv and N Tishby. Opening the black box of deep neural networks via information. *arXiv preprint arXiv:1703.00810*, 2017.

[27] O M Parkhi, A Vedaldi, and A Zisserman. Deep face recognition. In *BMVC*, volume 1, page 6, 2015.

[28] S Sankaranarayanan, A Alavi, C D Castillo, and R Chellappa. Triplet probabilistic embedding for face verification and clustering. In *2016 IEEE 8th International Conference on Biometrics Theory, Applications and Systems (BTAS)*, 2016.