

# Deep Imbalanced Attribute Classification using Visual Attention Aggregation

Nikolaos Sarafianos, Xiang Xu, and Ioannis A. Kakadiaris

Computational Biomedicine Lab  
University of Houston  
{nsarafi a, xxu18, i kakadi a}@central . uh. edu

**Abstract.** For many computer vision applications, such as image description and human identification, recognizing the visual attributes of humans is an essential yet challenging problem. Its challenges originate from its multi-label nature, the large underlying class imbalance and the lack of spatial annotations. Existing methods follow either a computer vision approach while failing to account for class imbalance, or explore machine learning solutions, which disregard the spatial and semantic relations that exist in the images. With that in mind, we propose an effective method that extracts and aggregates visual attention masks at different scales. We introduce a loss function to handle class imbalance both at class and at an instance level and further demonstrate that penalizing attention masks with high prediction variance accounts for the weak supervision of the attention mechanism. By identifying and addressing these challenges, we achieve state-of-the-art results with a simple attention mechanism in both PETA and WIDER-Attribute datasets without additional context or side information.

**Keywords:** Visual Attributes, Deep Imbalanced Learning, Visual Attention

## 1 Introduction

We set out to develop a method that, given an image of a human, predicts its visual attributes. We posed the following questions: (i) what are the challenges of this problem? (ii) what have other people done? and (iii) how should a simple yet effective solution to this problem look like? Human attributes are imbalanced in nature. Bald individuals with a mustache wearing glasses are 14 to 43 times less likely to appear in the CelebA dataset [1] compared to people without these characteristics. Large-scale imbalanced datasets can lead to biased models, optimized to favor the majority classes while failing to identify the subtle discriminant features that are required to recognize the under-represented classes. Setting the class imbalance aside, an additional challenge is identifying which areas in the image provide class-discriminant information. Giving emphasis to the upper part of an image, where the face is located, for attributes such as “glasses” and to the bottom part for attributes such as “long pants” can increase the recognition





































